

Pearson
Education

We work with leading authors to develop the strongest educational materials in chemistry, bringing cutting-edge thinking and best learning practice to a global market.

Under a range of well-known imprints, including Prentice Hall, we craft high-quality print and electronic publications which help readers to understand and apply their content, whether studying or at work.

To find out more about the complete range of our publishing please visit us on the World Wide Web at: www.pearsoneduc.com

Molecular Modelling

PRINCIPLES AND APPLICATIONS

Second edition

Andrew R. Leach

Glaxo Wellcome Research and Development

Prentice
Hall

An imprint of **Pearson Education**

Harlow, England · London · New York · Reading, Massachusetts · San Francisco · Toronto · Don Mills, Ontario · Sydney
Tokyo · Singapore · Hong Kong · Seoul · Taipei · Cape Town · Madrid · Mexico City · Amsterdam · Munich · Paris · Milan

Pearson Education Limited

Edinburgh Gate
Harlow
Essex CM20 2JE
England

and Associated Companies around the world

Visit us on the World Wide Web at:
www.pearsoneduc.com

First published under the Longman imprint 1996
Second edition 2001

© Pearson Education Limited 1996, 2001

The right of Andrew R. Leach to be identified as the author of
this Work has been asserted by him in accordance with
the Copyright, Designs and Patents Act 1988.

All rights reserved. No part of this publication may be reproduced, stored
in a retrieval system, or transmitted in any form or by any means, electronic,
mechanical, photocopying, recording or otherwise without either the prior
written permission of the publisher or a licence permitting restricted copying
in the United Kingdom issued by the Copyright Licensing Agency Ltd,
90 Tottenham Court Road, London W1P 0LP.

ISBN 0-582-38210-6

British Library Cataloguing-in-Publication Data

A catalogue record for this book can be obtained from the British Library

Library of Congress Cataloging-in-Publication Data

Leach, Andrew R.

Molecular modelling: principles and applications / Andrew R. Leach. - 2nd ed.
p. cm.

Includes bibliographical references and index.

ISBN 0-582-38210-6

1. Molecular structure-Computer simulation. 2. Molecules-Models-Computer
simulation. I. Title.

QD480.L43 2001

541.2'2'0113-dc21

00-046480

10 9 8 7 6 5 4 3 2 1
05 04 03 02 01

Top right-hand cover image © American Institute of Physics

Typeset by 60

Printed in Great Britain by Henry Ling Ltd,
at the Dorset Press, Dorchester, Dorset

To Christine

Contents

Preface to the Second Edition	xiii
Preface to the First Edition	xv
Symbols and Physical Constants	xvii
Acknowledgements	xxi
1 Useful Concepts in Molecular Modelling	1
1.1 Introduction	1
1.2 Coordinate Systems	2
1.3 Potential Energy Surfaces	4
1.4 Molecular Graphics	5
1.5 Surfaces	6
1.6 Computer Hardware and Software	8
1.7 Units of Length and Energy	9
1.8 The Molecular Modelling Literature	9
1.9 The Internet	9
1.10 Mathematical Concepts	10
Further Reading	24
References	24
2 An Introduction to Computational Quantum Mechanics	26
2.1 Introduction	26
2.2 One-electron Atoms	30
2.3 Polyelectronic Atoms and Molecules	34
2.4 Molecular Orbital Calculations	41
2.5 The Hartree-Fock Equations	51
2.6 Basis Sets	65
2.7 Calculating Molecular Properties Using <i>ab initio</i> Quantum Mechanics	74
2.8 Approximate Molecular Orbital Theories	86
2.9 Semi-empirical Methods	86
2.10 Hückel Theory	99
2.11 Performance of Semi-empirical Methods	102
Appendix 2.1 Some Common Acronyms Used in Computational Quantum Chemistry	104
Further Reading	105
References	105

3	Advanced <i>ab initio</i> Methods, Density Functional Theory and Solid-state Quantum Mechanics	108
3.1	Introduction	108
3.2	Open-shell Systems	108
3.3	Electron Correlation	110
3.4	Practical Considerations When Performing <i>ab initio</i> Calculations	117
3.5	Energy Component Analysis	122
3.6	Valence Bond Theories	124
3.7	Density Functional Theory	126
3.8	Quantum Mechanical Methods for Studying the Solid State	138
3.9	The Future Role of Quantum Mechanics: Theory and Experiment Working Together	160
	Appendix 3.1 Alternative Expression for a Wavefunction Satisfying Bloch's Function	161
	Further Reading	161
	References	162
4	Empirical Force Field Models: Molecular Mechanics	165
4.1	Introduction	165
4.2	Some General Features of Molecular Mechanics Force Fields	168
4.3	Bond Stretching	170
4.4	Angle Bending	173
4.5	Torsional Terms	173
4.6	Improper Torsions and Out-of-plane Bending Motions	176
4.7	Cross Terms: Class 1, 2 and 3 Force Fields	178
4.8	Introduction to Non-bonded Interactions	181
4.9	Electrostatic Interactions	181
4.10	Van der Waals Interactions	204
4.11	Many-body Effects in Empirical Potentials	212
4.12	Effective Pair Potentials	214
4.13	Hydrogen Bonding in Molecular Mechanics	215
4.14	Force Field Models for the Simulation of Liquid Water	216
4.15	United Atom Force Fields and Reduced Representations	221
4.16	Derivatives of the Molecular Mechanics Energy Function	225
4.17	Calculating Thermodynamic Properties Using a Force Field	226
4.18	Force Field Parametrisation	228
4.19	Transferability of Force Field Parameters	231
4.20	The Treatment of Delocalised π Systems	233
4.21	Force Fields for Inorganic Molecules	234
4.22	Force Fields for Solid-state Systems	236
4.23	Empirical Potentials for Metals and Semiconductors	240
	Appendix 4.1 The Interaction Between Two Drude Molecules	246
	Further Reading	247
	References	247

5	Energy Minimisation and Related Methods for Exploring the Energy Surface	253
5.1	Introduction	253
5.2	Non-derivative Minimisation Methods	258
5.3	Introduction to Derivative Minimisation Methods	261
5.4	First-order Minimisation Methods	262
5.5	Second Derivative Methods: The Newton-Raphson Method	267
5.6	Quasi-Newton Methods	268
5.7	Which Minimisation Method Should I Use?	270
5.8	Applications of Energy Minimisation	273
5.9	Determination of Transition Structures and Reaction Pathways	279
5.10	Solid-state Systems: Lattice Statics and Lattice Dynamics	295
	Further Reading	300
	References	301
6	Computer Simulation Methods	303
6.1	Introduction	303
6.2	Calculation of Simple Thermodynamic Properties	307
6.3	Phase Space	312
6.4	Practical Aspects of Computer Simulation	315
6.5	Boundaries	317
6.6	Monitoring the Equilibration	321
6.7	Truncating the Potential and the Minimum Image Convention	324
6.8	Long-range Forces	334
6.9	Analysing the Results of a Simulation and Estimating Errors	343
	Appendix 6.1 Basic Statistical Mechanics	347
	Appendix 6.2 Heat Capacity and Energy Fluctuations	348
	Appendix 6.3 The Real Gas Contribution to the Virial	349
	Appendix 6.4 Translating Particle Back into Central Box for Three Box Shapes	350
	Further Reading	351
	References	351
7	Molecular Dynamics Simulation Methods	353
7.1	Introduction	353
7.2	Molecular Dynamics Using Simple Models	353
7.3	Molecular Dynamics with Continuous Potentials	355
7.4	Setting up and Running a Molecular Dynamics Simulation	364
7.5	Constraint Dynamics	368
7.6	Time-dependent Properties	374
7.7	Molecular Dynamics at Constant Temperature and Pressure	382
7.8	Incorporating Solvent Effects into Molecular Dynamics: Potentials of Mean Force and Stochastic Dynamics	387
7.9	Conformational Changes from Molecular Dynamics Simulations	392
7.10	Molecular Dynamics Simulations of Chain Amphiphiles	394

Appendix 7.1	Energy Conservation in Molecular Dynamics	405
	Further Reading	406
	References	406
8	Monte Carlo Simulation Methods	410
8.1	Introduction	410
8.2	Calculating Properties by Integration	412
8.3	Some Theoretical Background to the Metropolis Method	414
8.4	Implementation of the Metropolis Monte Carlo Method	417
8.5	Monte Carlo Simulation of Molecules	420
8.6	Models Used in Monte Carlo Simulations of Polymers	423
8.7	'Biased' Monte Carlo Methods	432
8.8	Tackling the Problem of Quasi-ergodicity: J-walking and Multicanonical Monte Carlo	433
8.9	Monte Carlo Sampling from Different Ensembles	438
8.10	Calculating the Chemical Potential	442
8.11	The Configurational Bias Monte Carlo Method	443
8.12	Simulating Phase Equilibria by the Gibbs Ensemble Monte Carlo Method	450
8.13	Monte Carlo or Molecular Dynamics?	452
Appendix 8.1	The Marsaglia Random Number Generator	453
	Further Reading	454
	References	454
9	Conformational Analysis	457
9.1	Introduction	457
9.2	Systematic Methods for Exploring Conformational Space	458
9.3	Model-building Approaches	464
9.4	Random Search Methods	465
9.5	Distance Geometry	467
9.6	Exploring Conformational Space Using Simulation Methods	475
9.7	Which Conformational Search Method Should I Use? A Comparison of Different Approaches	476
9.8	Variations on the Standard Methods	477
9.9	Finding the Global Energy Minimum: Evolutionary Algorithms and Simulated Annealing	479
9.10	Solving Protein Structures Using Restrained Molecular Dynamics and Simulated Annealing	483
9.11	Structural Databases	489
9.12	Molecular Fitting	490
9.13	Clustering Algorithms and Pattern Recognition Techniques	491
9.14	Reducing the Dimensionality of a Data Set	497
9.15	Covering Conformational Space: Poling	499
9.16	A 'Classic' Optimisation Problem: Predicting Crystal Structures	501

	Further Reading	505
	References	506
10	Protein Structure Prediction, Sequence Analysis and Protein Folding	509
10.1	Introduction	509
10.2	Some Basic Principles of Protein Structure	513
10.3	First-principles Methods for Predicting Protein Structure	517
10.4	Introduction to Comparative Modelling	522
10.5	Sequence Alignment	522
10.6	Constructing and Evaluating a Comparative Model	539
10.7	Predicting Protein Structures by 'Threading'	545
10.8	A Comparison of Protein Structure Prediction Methods: CASP	547
10.9	Protein Folding and Unfolding	549
Appendix 10.1	Some Common Abbreviations and Acronyms Used in Bioinformatics	553
Appendix 10.2	Some of the Most Common Sequence and Structural Databases Used in Bioinformatics	555
Appendix 10.3	Mutation Probability Matrix for 1 PAM	556
Appendix 10.4	Mutation Probability Matrix for 250 PAM	557
	Further Reading	557
	References	558
11	Four Challenges in Molecular Modelling: Free Energies, Solvation, Reactions and Solid-state Defects	563
11.1	Free Energy Calculations	563
11.2	The Calculation of Free Energy Differences	564
11.3	Applications of Methods for Calculating Free Energy Differences	569
11.4	The Calculation of Enthalpy and Entropy Differences	574
11.5	Partitioning the Free Energy	574
11.6	Potential Pitfalls with Free Energy Calculations	577
11.7	Potentials of Mean Force	580
11.8	Approximate/'Rapid' Free Energy Methods	585
11.9	Continuum Representations of the Solvent	592
11.10	The Electrostatic Contribution to the Free Energy of Solvation: The Born and Onsager Models	593
11.11	Non-electrostatic Contributions to the Solvation Free Energy	608
11.12	Very Simple Solvation Models	609
11.13	Modelling Chemical Reactions	610
11.14	Modelling Solid-state Defects	622
Appendix 11.1	Calculating Free Energy Differences Using Thermodynamic Integration	630
Appendix 11.2	Using the Slow Growth Method for Calculating Free Energy Differences	631

Appendix 11.3	Expansion of Zwanzig Expression for the Free Energy Difference for the Linear Response Method	631
	Further Reading	632
	References	633
12	The Use of Molecular Modelling and Chemoinformatics to Discover and Design New Molecules	640
12.1	Molecular Modelling in Drug Discovery	640
12.2	Computer Representations of Molecules, Chemical Databases and 2D Substructure Searching	642
12.3	3D Database Searching	647
12.4	Deriving and Using Three-dimensional Pharmacophores	648
12.5	Sources of Data for 3D Databases	659
12.6	Molecular Docking	661
12.7	Applications of 3D Database Searching and Docking	667
12.8	Molecular Similarity and Similarity Searching	668
12.9	Molecular Descriptors	668
12.10	Selecting 'Diverse' Sets of Compounds	680
12.11	Structure-based <i>De Novo</i> Ligand Design	687
12.12	Quantitative Structure-Activity Relationships	695
12.13	Partial Least Squares	706
12.14	Combinatorial Libraries	711
	Further Reading	719
	References	720
	Index	727

Preface to the Second Edition

The impetus for this second edition is a desire to include some of the new techniques that have emerged in recent years and also extend the scope of the book to cover certain areas that were under-represented (even neglected) in the first edition. In this second volume there are three topics that fall into the first category (density functional theory, bioinformatics/protein structure analysis and chemoinformatics) and one main area in the second category (modelling of the solid state). In addition, of course, a new edition provides an opportunity to take a critical view of the text and to re-organise and update the material. Thus whilst much remains from the first edition, and this second book follows much the same path through the subject, readers familiar with the first edition will find some changes which I hope they will agree are for the better.

As with the first edition we initially consider quantum mechanics, but this is now split into two chapters. Thus Chapter 2 provides an introduction to the *ab initio* and semi-empirical approaches together with some examples of the uses of quantum mechanics. Chapter 3 covers more advanced aspects of the *ab initio* approach, density functional theory and the particular problems of the solid state. Molecular mechanics is the subject of Chapter 4 and then in Chapter 5 we consider energy minimisation and other 'static' techniques. Chapters 6, 7 and 8 deal with the two main simulation methods (molecular dynamics and Monte Carlo). Chapter 9 is devoted to the conformational analysis of 'small' molecules but also includes some topics (e.g. cluster analysis, principal components analysis) that are widely used in informatics. In Chapter 10 the problems of protein structure prediction and protein folding are considered; this chapter also contains an introduction to some of the more widely used methods in bioinformatics. In Chapter 11 we draw upon material from the previous chapters in a discussion of free energy calculations, continuum solvent models, and methods for simulating chemical reactions and defects in solids. Finally, Chapter 12 is concerned with modelling and chemoinformatics techniques for discovering and designing new molecules, including database searching, docking, *de novo* design, quantitative structure-activity relationships and combinatorial library design.

As in the first edition, the inexorable pace of change means that what is currently considered 'cutting edge' will soon become routine. The examples are thus chosen primarily because they illuminate the underlying theory rather than because they are the first application of a particular technique or are the most recent available. In a similar vein, it is impossible in a volume such as this to even attempt to cover everything and so there are undoubtedly areas which are under-represented. This is not intended to be a definitive historical account or a review of the current state-of-the-art. Thus, whilst I have tried to include many literature references it is possible that the invention of some technique may appear to be incorrectly attributed or a 'classic' application may be missing. A general guiding principle has been

to focus on those techniques that are in widespread use rather than those which are the province of one particular research group. Despite these caveats I hope that the coverage is sufficient to provide a solid introduction to the main areas and also that those readers who are 'experts' will find something new to interest them.

A Companion Web Site accompanies *Molecular Modelling: Principles and Applications, Second Edition* by Andrew Leach



Visit the *Molecular Modelling* Companion Web Site at www.booksites.net/leach
The website contains general information about the book, up-to-date hyperlinks to related chemistry sources on the web, reference copies of appendices of relevant acronyms, and twenty-six full screen, full-colour graphical representations of molecular structures.

Preface to the First Edition

Molecular modelling used to be restricted to a small number of scientists who had access to the necessary computer hardware and software. Its practitioners wrote their own programs, managed their own computer systems and mended them when they broke down. Today's computer workstations are much more powerful than the mainframe computers of even a few years ago and can be purchased relatively cheaply. It is no longer necessary for the modeller to write computer programs as software can be obtained from commercial software companies and academic laboratories. Molecular modelling can now be performed in any laboratory or classroom.

This book is intended to provide an introduction to some of the techniques used in molecular modelling and computational chemistry, and to illustrate how these techniques can be used to study physical, chemical and biological phenomena. A major objective is to provide, in one volume, some of the theoretical background to the vast array of methods available to the molecular modeller. I also hope that the book will help the reader to select the most appropriate method for a problem and so make the most of his or her modelling hardware and software. Many modelling programs are extremely simple to use and are often supplied with seductive graphical interfaces, which obviously helps to make modelling techniques more accessible, but it can also be very easy to select a wholly inappropriate technique or method.

Most molecular modelling studies involve three stages. In the first stage a model is selected to describe the intra- and inter-molecular interactions in the system. The two most common models that are used in molecular modelling are quantum mechanics and molecular mechanics. These models enable the energy of any arrangement of the atoms and molecules in the system to be calculated, and allow the modeller to determine how the energy of the system varies as the positions of the atoms and molecules change. The second stage of a molecular modelling study is the calculation itself, such as an energy minimisation, a molecular dynamics or Monte Carlo simulation, or a conformational search. Finally, the calculation must be analysed, not only to calculate properties but also to check that it has been performed properly.

The book is organised so that some of the techniques discussed in later chapters refer to material discussed earlier, though I have tried to make each chapter as independent of the others as possible. Some readers may therefore be pleased to know that it is not essential to completely digest the chapters on quantum mechanics and molecular mechanics in order to read about methods for searching conformational space! Readers with experience in one or more areas may, of course, wish to be more selective.

I have tried to provide as much of the underlying theory as seems appropriate to enable the reader to understand the fundamentals of each method. In doing so I have assumed some background knowledge of quantum mechanics, statistical mechanics, conformational analysis and mathematics. A reader with an undergraduate degree in chemistry should

have covered this material, which should also be familiar to many undergraduates in the final year of their degree course. Full discussion can be found in the suggestions for further reading at the end of each chapter. I have also attempted to provide a reasonable selection of original references, though in a book of this scope it is obviously impossible to provide a comprehensive coverage of the literature. In this context, I apologise in advance if any technique is inappropriately attributed.

The range of systems that can be considered in molecular modelling is extremely broad, from isolated molecules through simple atomic and molecular liquids to polymers, biological macromolecules such as proteins and DNA and solids. Many of the techniques are illustrated with examples chosen to reflect the breadth of applications. It is inevitable that, for reasons of space, some techniques must be dealt with in a rudimentary fashion (or not at all), and that many interesting and important applications cannot be described. Molecular modelling is a rapidly developing discipline and has benefited from the dramatic improvements in computer hardware and software of recent years. Calculations that were major undertakings only a few years ago can now be performed using personal computing facilities. Thus, examples used to indicate the 'state of the art' at the time of writing will invariably be routine within a short time.

Symbols and Physical Constants

This list contains the most frequently used symbols and physical constants ordered according to approximate appearance in the text.

λ	Lagrange multiplier
r, θ, ϕ	spherical polar coordinates
$\mathbf{i}, \mathbf{j}, \mathbf{k}$	orthogonal unit vectors along x, y, z axes
ϕ, θ, ψ	Euler angles
$\langle x \rangle$ or \bar{x}	arithmetic mean value of x
\mathbf{I}	unit matrix
i	square root of -1
$\hat{\mathbf{r}}$	unit vector
α	exponent in Gaussian function (normal distribution)
σ	standard deviation
σ^2	variance
h	Planck's constant ($6.626\,18 \times 10^{-34}$ J s)
\hbar	$h/2\pi$ ($1.054\,59 \times 10^{-34}$ J s)
m	particle mass
Ψ	molecular wavefunction
∇^2	$\partial^2/\partial x^2 + \partial^2/\partial y^2 + \partial^2/\partial z^2$ ('del-squared')
\mathcal{H}	Hamiltonian
ψ	spatial orbital
α, β	spin functions ('spin up' and 'spin down')
χ	spin orbital (product of spatial orbital and a spin function)
ϕ	basis function/atomic orbital (usually labelled $\phi_\mu, \phi_\nu, \phi_\lambda, \phi_\sigma$)
$d\nu$ or $d\mathbf{r}$	indicates an integral over all spatial coordinates
$d\sigma$	indicates an integral over all spin coordinates
$d\tau$	indicates an integral over all spatial and spin coordinates
r_{ij}	distances between two particles i and j (usually electrons in quantum mechanics)
R_{AB}	distance between two nuclei A and B
δ_{ij}	Kronecker delta ($\delta_{ij} = 1$ if $i = j$; $\delta_{ij} = 0$ if $i \neq j$)
\mathcal{K}	exchange operator
\mathcal{J}	Coulomb operator
$\mathcal{H}^{\text{core}}$	core Hamiltonian operator
\mathbf{F}	Fock matrix
\mathbf{S}	overlap matrix
S_{ij}	overlap integral between orbitals i and j
f	Fock operator
\mathbf{C}	matrix of basis function coefficients

G	metric matrix (in distance geometry)
p_i	i th principal component
Z	variance-covariance matrix
λ	coupling parameter (used in free energy calculations)
$W(\mathbf{r}^N)$	weighting function used in umbrella sampling
\mathcal{N}	number density ($= N/V$)
S_{AB}	similarity coefficient between two molecules A and B
D_{AB}	'distance' between two molecules A and B
σ	Hammett substitution constant
P	partition coefficient of solute between two solvents
π	$\log(P_X/P_H)$ for a substituent X relative to a hydrogen substituent
r^2	squared correlation coefficient
R^2	squared correlation coefficient in multiple linear regression
Q^2	cross-validated R^2

Acknowledgements

For this second edition I would like to thank Drs Neil Allan, Paul Bamborough, Gianpaolo Bravi, Richard Bryce, Julian Gale, Richard Green, Mike Hann and Alan Lewis, who commented on various parts of the new text. Julian Gale's suggestions were particularly useful for refining the sections concerning materials science and solid-state applications. I would also like to record my thanks once more to those who gave their time to read and comment on draft copies of various chapters of the first edition, upon which this second edition is based (in alphabetical order): Dr D B Adolf, Dr J M Blaney, Professor A V Chadwick, Dr P S Charifson, Dr C-W Chung, Dr A Cleasby, Dr A Emerson, Dr J W Essex, Dr D V S Green, Dr I R Gould, Dr M M Hann, Dr C A Leach, Dr M Pass, Dr D A Pearlman, Dr C A Reynolds, Dr D W Salt, Dr M Saqi, Professor J I Siepmann, Dr W C Swope, Dr N R Taylor, Dr P J Thomas, Professor D J Tildesley and Mr O Warschkow.

Assistance with the illustrations for the second edition was provided by Drs R Groot, S McGrother and V Milman. Many of the figures from the first edition are also included here and so I would like to thank again Dr S E Greasley, Dr M M Hann, Dr H Jhoti, Dr S N Jordan, Professor G R Luckhurst, Dr P M McMeekin, Dr A Nicholls, Dr P Popelier, Dr A Robinson and Dr T E Klein.

Alexandra Seabrook, Pauline Gillet and Julie Knight at Pearson Education provided the foundation of the publishing team, coping with a steady stream of questions and keeping everyone to schedule. Especial thanks are due to Julie, who did a splendid job as editor.

Any errors that remain are of course my own responsibility. If you do find any, I would like to know! I will also be pleased to receive any constructive suggestions, comments or criticisms. We plan to set up a web site that will provide access to various material from the book (such as electronic versions of the colour images) together with email contacts. This can be accessed via www.booksites.net.

Molecular modelling would not be what it is today without the efforts of those who develop computer hardware and software and I would like to acknowledge the authors of the following computer programs which were used to generate figures and/or data described in the text. All calculations were performed using Silicon Graphics computers.

AMBER: D A Pearlman, D A Case, J C Caldwell, G L Seibel, U C Singh, P Weiner and P A Kollman 1991.

Amber 3.0, University of California, San Francisco.

Cambridge Structural Database: F H Allen, S A Bellard, M D Brice, B A Cartwright, A Doubleday, H Higgs, T Hummelink, B G Hummelink-Peters, O Kennard, W D S Motherwell, J R Rodgers and D G Watson 1979. The Cambridge Crystallographic Data Centre: Computer-Based Search, Retrieval, Analysis and Display of Information. *Acta Crystallographica* B35:2331-2339. Cambridge Crystallographic Data Centre, Cambridge, United Kingdom.

CASTEP: Molecular Simulations Inc., 9685 Scranton Road, San Diego, California, USA.

Catalyst: Molecular Simulations Inc., 9685 Scranton Road, San Diego, California, USA.

Cerius2: Molecular Simulations Inc., 9685 Scranton Road, San Diego, California, USA.

- COSMIC: J G Vinter, A Davis, M R Saunders 1987. Strategic approaches to drug design. I. An integrated software framework for molecular modeling. *Journal of Computer-Aided Molecular Design* 1(1):31-51.
- Dials and Windows: G Ravishanker, S Swaminathan, D L Beveridge, R Lavery and H Sklenar 1989. *Journal of Biomolecular Structure and Dynamics* 6:669-699. Wesleyan University, USA.
- Gaussian 92: M J Frisch, G W Trucks, M Head-Gordon, P M W Gill, M W Wong, J B Foresman, B G Johnson, H B Schlegel, M A Robb, E S Replogle, R Gomperts, J L Andres, K Raghavachari, J S Binkley, C Gonzalez, R L Martin, D J Fox, D J DeFrees, J Baker, J J P Stewart and J A Pople. Gaussian Inc., Pittsburgh, Pennsylvania, USA.
- GCG: Genetics Computer Group, Inc., University Research Park, 575 Science Drive, Suite B, Madison, Wisconsin 53711, USA.
- GRASP (Graphical Representation and Analysis of Surface Properties): A Nicholls, Columbia University, New York, USA.
- GRID: P J Goodford 1985. A Computational Procedure for Determining Energetically Favorable Binding Sites on Biologically Important Macromolecules. *Journal of Medicinal Chemistry* 28:849-857. Molecular Discovery Ltd, Oxford, United Kingdom.
- InsightII: Molecular Simulations Inc., 9685 Scranton Road, San Diego, California, USA.
- IsoStar: I J Bruno, J C Cole, J P M Lommerse, R S Rowland, R Taylor and M L Verdonk 1997. IsoStar: a library of information about nonbonded interactions. *Journal of Computer-Aided Molecular Design* 11:525-537. Cambridge Crystallographic Data Centre, Cambridge, United Kingdom.
- Micromol: S M Colwell, A R Marshall, R D Amos and N C Handy 1985. Quantum chemistry on microcomputers, *Chemistry in Britain* 21:655-659.
- Molscrip: P J Kraulis 1991. Molscrip - A program to produce both detailed and schematic plots of protein structures. *Journal of Applied Crystallography* 24:946-950.
- PROCHECK: R Laskowski, M W MacArthur, D S Moss and J M Thornton 1993. Procheck - A program to check the stereochemical quality of protein structures. *Journal of Applied Crystallography* 26:283-291.
- Quanta: Molecular Simulations Inc., 9685 Scranton Road, San Diego, California, USA.
- Spartan: Wavefunction Inc., 18401 Von Karman, Suite 370, Irvine, California, USA.
- SPASMS (San Francisco Package of Applications for the Simulation of Molecular Systems). D A Spellmeyer, W C Swope, E-R Evensen, T Cheatham, D M Ferguson and P A Kollman. University of California, San Francisco, USA.
- Sybyl: Tripos Inc., 1699 South Hanley Road, St. Louis, Missouri, USA.
- The following programs were used to produce draft copies of the manuscript and diagrams: Microsoft Word (Microsoft Corp.), Gnuplot (T Williams and C Kelley), Kaleidagraph (Abelbeck Software), Chem3D (CambridgeSoft Corp.) and Microsoft Excel (Microsoft Corp.).

We are grateful to the following for permission to reproduce copyright material:

- Figure 1.11 from *Mathematical Methods in the Physical Sciences*, 2nd edn, Boas M L, ©1983. Reprinted by permission of John Wiley & Sons, Inc.
- Figure 1.14 from *The FFT Fundamentals and Concepts* by Ramirez ©1985. Reprinted by permission of Prentice-Hall, Inc., Upper Saddle River, NJ.
- Figures 2.7 and 3.3 from *Ab initio Molecular Orbital Theory*, Hehre W J, L Radom, P v R Schleyer, J A Hehre ©1986. Reprinted with permission by John Wiley & Sons, Inc.
- Figure 3.5 from Gerratt J, D L Cooper, P B Karadakov and M Raimondi 1997. Modern valence bond theory. *Chemical Society Reviews* 87-100. Reproduced by permission of The Royal Society of Chemistry.
- Figure 3.22 from Needs R J and Mujica 1995. First-principles pseudopotential study of the structural phases of silicon. *Physical Review B* 51:9652-9660.

- Figure 4.18 from Buckingham A D 1959. Molecular Quadrupole Moments. *Quarterly Reviews of the Chemical Society* 13:183-214. Reproduced by permission of The Royal Society of Chemistry.
- Figure 4.29 from *Computer Simulation in Chemical Physics*, edited by Allen M P and D J Tildesley, 1993. Effective Pair Potentials and Beyond, Sprik M, with kind permission from Kluwer Academic Publishers.
- Figure 4.49 reprinted with permission from Pranata J and W L Jorgensen. Computational Studies on FK506: Conformational Search and Molecular Dynamics Simulations in Water. *The Journal of the American Chemical Society* 113:9483-9493. ©1991 American Chemical Society.
- Figure 4.50 from Molecular Parameters for Organosilicon Compounds Calculated from *Ab Initio* Computations, Grigoras S and T H Lane, *Journal of Computational Chemistry* 9:25-39, ©1988. Reprinted by permission of John Wiley & Sons, Inc.
- Figures 5.4 and 5.8 Press W H, B P Flannery, S A Teukolsky and W T Vetterling, *Numerical Recipes in Fortran*. 1992, Cambridge University Press.
- Figure 5.21 reprinted with permission from Chandrasekhar J, S F Smith and W L Jorgensen. Theoretical Examination of the S_N2 Reaction Involving Chloride Ion and Methyl Chloride in the Gas Phase and Aqueous Solution. *The Journal of the American Chemical Society* 107:154-163. ©1985 American Chemical Society.
- Figure 5.23 reprinted with permission from Doubleday C, J McIver, M Page and T Zielinski. Temperature Dependence of the Transition-State Structure for the Disproportionation of Hydrogen Atom with Ethyl Radical. *The Journal of the American Chemical Society* 107:5800-5801. ©1985 American Chemical Society.
- Figure 5.29 from Gonzalez C and H B Schlegel 1988. An Improved Algorithm for Reaction Path Following. *The Journal of Chemical Physics* 90:2154-2161.
- Figure 5.30 reprinted from *Chemical Physical Letters*, 194, Fischer S and M Karplus. Conjugate Peak Refinement: An Algorithm for Finding Reaction Paths and Accurate Transition States in Systems with Many Degrees of Freedom. 252-261, ©1992, with permission from Elsevier Science.
- Figure 5.35 reprinted with permission from Houk K N, J González and Y Li. Pericyclic Reaction Transition States: Passions and Punctilios 1935-1995. *Accounts of Chemical Research* 28:81-90. ©1995 American Chemical Society.
- Figure 6.25 reprinted from *Chemical Physics Letters*, 196, Ding H-Q, N Karasawa and W A Goddard III, The Reduced Cell Multipole Method for Coulomb Interactions in Periodic Systems with Million-Atom Unit Cells, 6-10, ©1992, with permission of Elsevier Science.
- Figure 7.2 from Alder B J and T E Wainwright 1959. Studies in Molecular Dynamics. I. General Method. *The Journal of Chemical Physics* 31:459-466.
- Figure 7.11 from Alder B J and T E Wainwright 1970. Decay of the Velocity Autocorrelation Function. *Physical Review A* 1:18-21.
- Figure 7.12 from Guillot B 1991. A Molecular Dynamics Study of the Infrared Spectrum of Water. *The Journal of Chemical Physics* 95:1543-1551.
- Figure 7.13 reprinted with permission from Jorgensen W L, R C Binning Jr and B Bigot. Structures and Properties of Organic Liquids: *n*-Butane and 1,2-Dichloroethane and Their Conformational Equilibria. *The Journal of the American Chemical Society* 103:4393-4399. ©1981 American Chemical Society.
- Figure 7.24 (and on cover) from Groot R D and T J Madden 1998. Dynamic simulation of diblock copolymer microphase separation. *The Journal of Chemical Physics* 108:8713-8724. © American Institute of Physics.
- Figure 8.16 from Frantz, D D, D L Freeman and J D Doll 1990. Reducing quasi-ergodic behavior in Monte Carlo simulations by J-walking: applications to atomic clusters. *The Journal of Chemical Physics* 93:2769-2784.

- Figure 8.17 from Cracknell R F, D Nicholson and N Quirke 1994. A Grand Canonical Monte Carlo Study of Lennard-Jones Mixtures in Slit Pores; 2: Mixtures of Two-Centre Ethane with Methane. *Molecular Simulation* 13:161-175. ©1994 OPA (Overseas Publishers Association) N.V. Permission to reproduce granted by Gordon and Breach Publishers.
- Figure 8.22 reprinted with permission from Smit B and J I Siepmann. Simulating the Adsorption of Alkanes in Zeolites. *Science* 264:1118-1120 ©1994 American Association for the Advancement of Science.
- Figure 9.34 from Poling: Promoting Conformational Variation, Smellie A S, S L Teig and P Towbin *Journal of Computational Chemistry* 16:171-187, ©1995. Reprinted by permission of John Wiley & Sons, Inc.
- Figure 10.18 from Pearson W R and D J Lipman 1988. Improved tools for biological sequence comparison. *Proceedings of the National Academy of Sciences USA* 85:2444-2448.
- Figure 10.21 reprinted from *Current Opinion in Structural Biology*, 6, Eddy S R, Hidden Markov Models, 361-365 ©1996, with permission from Elsevier Science.
- Figure 10.26 Reprinted with permission from Lüthy R, J U Bowie and D Eisenberg. Assessment of Protein Models with Three-Dimensional Profiles. *Nature* 356:83-85. ©1992 Macmillan Magazines Limited.
- Figure 11.6 from Lybrand T P, J A McCammon and G Wipff 1986. Theoretical Calculation of Relative Binding Affinity in Host-Guest Systems. *Proceedings of the National Academy of Sciences USA* 83:833-835.
- Figure 11.18 reprinted with permission from Guo Z and C L Brooks III. Rapid Screening of Binding Affinities: Application of the λ -Dynamics Method to a Trypsin-Inhibitor System. *The Journal of the American Chemical Society* 120:1920-1921. ©1998 American Chemical Society.
- Figure 11.24 reprinted with permission from Still W C, A Tempczyk, R C Hawley and T Hendrickson. Semianalytical Treatment of Solvation for Molecular Mechanics and Dynamics. *The Journal of the American Chemical Society* 112:6127-6129. ©1990 American Chemical Society.
- Figure 11.35 reprinted with permission from Chandrasekhar J and W L Jorgensen 1985. Energy Profile for a Nonconcerted S_N2 Reaction in Solution. *The Journal of the American Chemical Society* 107:2974-2975. ©1985 American Chemical Society.
- Figure 11.37 reprinted with permission from Åqvist J, M Fothergill and A Warshel. Computer Simulation of the $\text{CO}_2/\text{HCO}_3^-$ Interconversion Step in Human Carbonic Anhydrase I. *The Journal of the American Chemical Society* 115:631-635. ©1993 American Chemical Society.
- Figure 11.40 reprinted with permission from Saitta A M, P D Sooper, E Wasserman and M L Klein 1999. Influence of a knot on the strength of a polymer strand. *Nature* 399:46-48. ©1999 Macmillan Magazines Limited.
- Figure 11.42 from *NATO ASI Series C 498 (New Trends in Materials Chemistry)*, 1997, 285-318, Defects and Matter Transport in Solid Materials, Chadwick A V and J Corish, with kind permission of Kluwer Academic Publishers.

Whilst every effort has been made to trace the owners of copyright material, we take this opportunity to offer our apologies to any copyright holders whose rights we may have unwittingly infringed.

CHAPTER ONE

Useful Concepts in Molecular Modelling

1.1 Introduction

What is molecular modelling? 'Molecular' clearly implies some connection with molecules. The *Oxford English Dictionary* defines 'model' as 'a simplified or idealised description of a system or process, often in mathematical terms, devised to facilitate calculations and predictions'. Molecular modelling would therefore appear to be concerned with ways to mimic the behaviour of molecules and molecular systems. Today, molecular modelling is invariably associated with computer modelling, but it is quite feasible to perform some simple molecular modelling studies using mechanical models or a pencil, paper and hand calculator. Nevertheless, computational techniques have revolutionised molecular modelling to the extent that most calculations could not be performed without the use of a computer. This is not to imply that a more sophisticated model is necessarily any better than a simple one, but computers have certainly extended the range of models that can be considered and the systems to which they can be applied.

The 'models' that most chemists first encounter are molecular models such as the 'stick' models devised by Dreiding or the 'space filling' models of Corey, Pauling and Koltun (commonly referred to as CPK models). These models enable three-dimensional representations of the structures of molecules to be constructed. An important advantage of these models is that they are interactive, enabling the user to pose 'what if ...' or 'is it possible to ...' questions. These structural models continue to play an important role both in teaching and in research, but molecular modelling is also concerned with more abstract models, many of which have a distinguished history. An obvious example is quantum mechanics, the foundations of which were laid many years before the first computers were constructed.

There is a lot of confusion over the meaning of the terms 'theoretical chemistry', 'computational chemistry' and 'molecular modelling'. Indeed, many practitioners use all three labels to describe aspects of their research, as the occasion demands! 'Theoretical chemistry' is often considered synonymous with quantum mechanics, whereas computational chemistry encompasses not only quantum mechanics but also molecular mechanics, minimisation, simulations, conformational analysis and other computer-based methods for understanding and predicting the behaviour of molecular systems. Molecular modellers use all of these methods and so we shall not concern ourselves with semantics but rather shall consider any theoretical or computational technique that provides insight into the behaviour of molecular systems to be an example of molecular modelling. If a distinction has to be

made, it is in the emphasis that molecular modelling places on the representation and manipulation of the structures of molecules, and properties that are dependent upon those three-dimensional structures. The prominent part that computer graphics has played in molecular modelling has led some scientists to consider molecular modelling as little more than a method for generating 'pretty pictures', but the technique is now firmly established, widely used and accepted as a discipline in its own right.

A closely related subject is molecular informatics. This is a rather new term, making a precise definition difficult, but it is usually considered to encompass two disciplines: chemoinformatics and bioinformatics. Of these two areas, chemoinformatics (also written cheminformatics) is the newer name but the older discipline; chemists have been using computers to store, retrieve and manipulate information about molecules almost since computers were invented. Both chemoinformatics and bioinformatics have risen to prominence primarily as a consequence of the introduction of new experimental techniques. For the chemist these experimental techniques are combinatorial library synthesis and high-throughput screening, which enable very large numbers of molecules to be synthesised and tested; for the biologist they are the automated sequencing machines that are being used to determine the human genome. A characteristic feature of molecular informatics is that it is concerned with information about large numbers of molecules, much larger than is typically the case for a traditional molecular modelling study. For this reason, informatics was initially more concerned with less complex representations of molecules that did not fully represent their three-dimensional properties. However, even this distinction is now being eroded and there is increasing use made of more traditional molecular modelling techniques within informatics.

In the rest of this chapter we shall discuss a number of concepts and techniques that are relevant to many areas of molecular modelling and so do not sit comfortably in any individual chapter. We will also define some of the terms that will be used throughout the book.

1.2 Coordinate Systems

It is obviously important to be able to specify the positions of the atoms and/or molecules in the system to a modelling program*. There are two common ways in which this can be done. The most straightforward approach is to specify the Cartesian (x, y, z) coordinates of all the atoms present. The alternative is to use *internal coordinates*, in which the position of each atom is described relative to other atoms in the system. Internal coordinates are usually written as a Z-matrix. The Z-matrix contains one line for each atom in the system. A sample Z-matrix for the staggered conformation of ethane (see Figure 1.1) is

*For a system containing a large number of independent molecules it is common to use the term 'configuration' to refer to each arrangement; this use of the word 'configuration' is not to be confused with its standard chemical meaning as a different bonding arrangement of the atoms in a molecule.

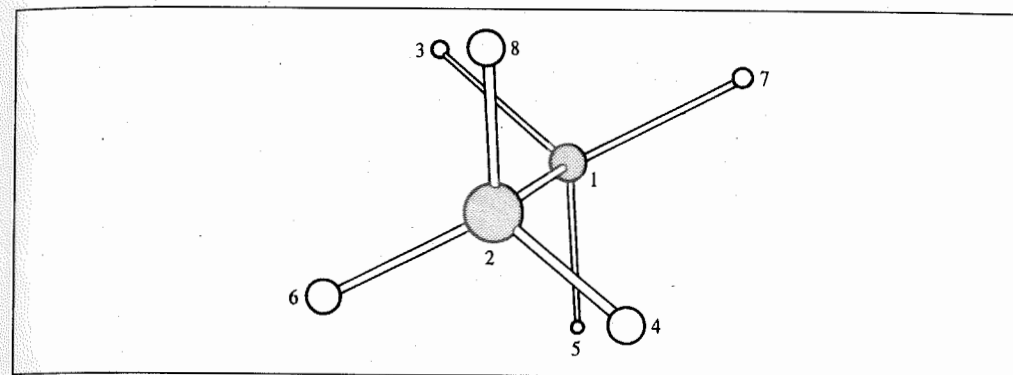


Fig. 1.1: The staggered conformation of ethane.

as follows:

1	C						
2	C	1.54	1				
3	H	1.0	1	109.5	2		
4	H	1.0	2	109.5	1	180.0	3
5	H	1.0	1	109.5	2	60.0	4
6	H	1.0	2	109.5	1	-60.0	5
7	H	1.0	1	109.5	2	180.0	6
8	H	1.0	2	109.5	1	60.0	7

In the first line of the Z-matrix we define atom 1, which is a carbon atom. Atom number 2 is also a carbon atom that is a distance of 1.54 Å from atom 1 (columns 3 and 4). Atom 3 is a hydrogen atom that is bonded to atom 1 with a bond length of 1.0 Å. The angle formed by atoms 2-1-3 is 109.5°, information that is specified in columns 5 and 6. The fourth atom is a hydrogen, a distance of 1.0 Å from atom 2, the angle 4-2-1 is 109.5°, and the torsion angle (defined in Figure 1.2) for atoms 4-2-1-3 is 180°. Thus for all except the first three atoms, each atom has three internal coordinates: the distance of the atom from one of the atoms previously defined, the angle formed by the atom and two of the previous atoms, and the torsion angle defined by the atom and three of the previous atoms. Fewer internal coordinates are required for the first three atoms because the first atom can be placed anywhere in space (and so it has no internal coordinates); for the second atom it is only necessary to specify its distance from the first atom and then for the third atom only a distance and an angle are required.

It is always possible to convert internal to Cartesian coordinates and vice versa. However, one coordinate system is usually preferred for a given application. Internal coordinates can usefully describe the relationship between the atoms in a single molecule, but Cartesian coordinates may be more appropriate when describing a collection of discrete molecules. Internal coordinates are commonly used as input to quantum mechanics programs, whereas calculations using molecular mechanics are usually done in Cartesian coordinates. The total number of coordinates that must be specified in the internal coordinate system is six fewer

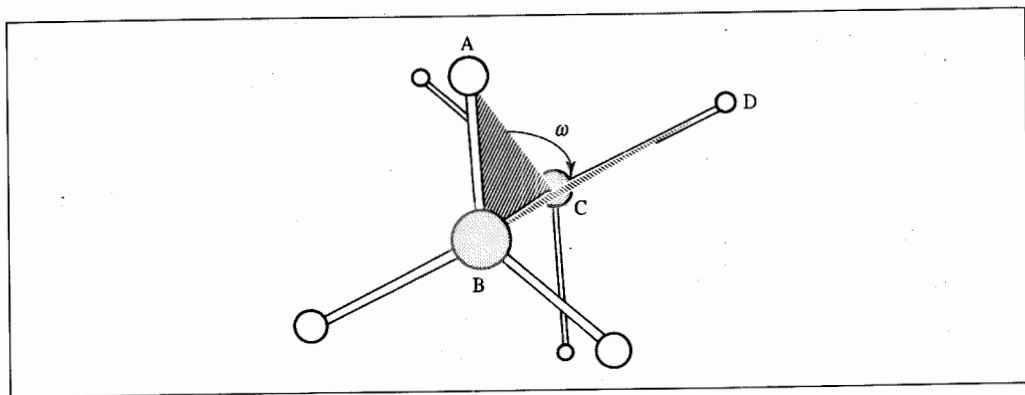


Fig. 1.2: A torsion angle $A-B-C-D$ is defined as the angle between the planes A, B, C and B, C, D . A torsion angle can vary through 360° although the range -180° to $+180^\circ$ is most commonly used. We shall adopt the IUPAC definition of a torsion angle in which an eclipsed conformation corresponds to a torsion angle of 0° and a trans or anti conformation to a torsion angle of 180° . The reader should note that this may not correspond to some of the definitions used in the literature, where the trans arrangement is defined as a torsion angle of 0° . If one looks along the bond $B-C$, then the torsion angle is the angle through which it is necessary to rotate the bond AB in a clockwise sense in order to superimpose the two planes, as shown.

than the number of Cartesian coordinates for a non-linear molecule. This is because we are at liberty to arbitrarily translate and rotate the system within Cartesian space without changing the relative positions of the atoms.

1.3 Potential Energy Surfaces

In molecular modelling the Born-Oppenheimer approximation is invariably assumed to operate. This enables the electronic and nuclear motions to be separated; the much smaller mass of the electrons means that they can rapidly adjust to any change in the nuclear positions. Consequently, the energy of a molecule in its ground electronic state can be considered a function of the nuclear coordinates only. If some or all of the nuclei move then the energy will usually change. The new nuclear positions could be the result of a simple process such as a single bond rotation or it could arise from the concerted movement of a large number of atoms. The magnitude of the accompanying rise or fall in the energy will depend upon the type of change involved. For example, about 3 kcal/mol is required to change the covalent carbon-carbon bond length in ethane by 0.1 \AA away from its equilibrium value, but only about 0.1 kcal/mol is required to increase the non-covalent separation between two argon atoms by 1 \AA from their minimum energy separation. For small isolated molecules, rotation about single bonds usually involves the smallest changes in energy. For example, if we rotate the carbon-carbon bond in ethane, keeping all of the bond lengths and angles fixed in value, then the energy varies in an approximately sinusoidal fashion as shown in Figure 1.3, with minima at the three staggered conformations. The energy in this case can be considered a function of a single coordinate only (i.e. the torsion angle of the carbon-carbon bond), and as such can be displayed graphically, with energy along one axis and the value of the coordinate along the other.

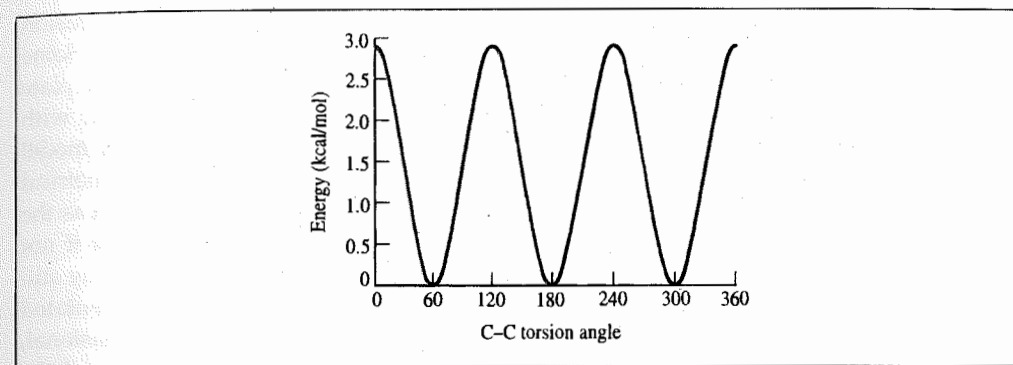


Fig. 1.3: Variation in energy with rotation of the carbon-carbon bond in ethane.

Changes in the energy of a system can be considered as movements on a multidimensional 'surface' called the *energy surface*. We shall be particularly interested in stationary points on the energy surface, where the first derivative of the energy is zero with respect to the internal or Cartesian coordinates. At a stationary point the forces on all the atoms are zero. Minimum points are one type of stationary point; these correspond to stable structures. Methods for locating stationary points will be discussed in more detail in Chapter 5, together with a more detailed consideration of the concept of the energy surface.

1.4 Molecular Graphics

Computer graphics has had a dramatic impact upon molecular modelling. It should always be remembered, however, that there is much more to molecular modelling than computer graphics. It is the interaction between molecular graphics and the underlying theoretical methods that has enhanced the accessibility of molecular modelling methods and assisted the analysis and interpretation of such calculations.

Molecular graphics systems have evolved from delicate and temperamental pieces of equipment that cost hundreds of thousands of pounds and occupied entire rooms, to today's inexpensive workstations that fit on or under a desk and yet are hundreds of times more powerful. Over the years, two different types of molecular graphics display have been used in molecular modelling. First to be developed were vector devices, which construct pictures using an electron gun to draw lines (or dots) on the screen, in a manner similar to an oscilloscope. Vector devices were the mainstay of molecular modelling for almost two decades but have now been largely superseded by raster devices. These divide the screen into a large number of small 'dots', called pixels. Each pixel can be set to any of a large number of colours, and so by setting each pixel to the appropriate colour it is possible to generate the desired image.

Molecules are most commonly represented on a computer graphics screen using 'stick' or 'space-filling' representations, which are analogous to the Dreiding and Corey-Pauling-Koltun (CPK) mechanical models. Sophisticated variations on these two basic types have

been developed, such as the ability to colour molecules by atomic number and the inclusion of shading and lighting effects, which give 'solid' models a more realistic appearance. Some of the commonly used molecular representations are shown in Figure 1.4 (colour plate section). Computer-generated models do have some advantages when compared with their mechanical counterparts. Of particular importance is the fact that a computer model can be very easily interrogated to provide quantitative information, from simple geometrical measures such as the distance between two atoms to more complex quantities such as the energy or surface area. Quantitative information such as this can be very difficult if not impossible to obtain from a mechanical model. Nevertheless, mechanical models may still be preferred in certain types of situation due to the ease with which they can be manipulated and viewed in three dimensions. A computer screen is inherently two-dimensional, whereas molecules are three-dimensional objects. Nevertheless, some impression of the three-dimensional nature of an object can be represented on a computer screen using techniques such as depth cueing (in which those parts of the object that are further away from the viewer are made less bright) and through the use of perspective. Specialised hardware enables more realistic three-dimensional stereo images to be viewed. In the future 'virtual reality' systems may enable a scientist to interact with a computer-generated molecular model in much the same way that a mechanical model can be manipulated.

Even the most basic computer graphics program provides some standard facilities for the manipulation of models, including the ability to translate, rotate and 'zoom' the model towards and away from the viewer. More sophisticated packages can provide the scientist with quantitative feedback on the effect of altering the structure. For example, as a bond is rotated then the energy of each structure could be calculated and displayed interactively.

For large molecular systems it may not always be desirable to include every single atom in the computer image; the sheer number of atoms can result in a very confusing and cluttered picture. A clearer picture may be achieved by omitting certain atoms (e.g. hydrogen atoms) or by representing groups of atoms as single 'pseudo-atoms'. The techniques that have been developed for displaying protein structures nicely illustrate the range of computer graphics representation possible (the use of computational techniques to investigate the structures of proteins is considered in Chapter 10). Proteins are polymers constructed from amino acids, and even a small protein may contain several thousand atoms. One way to produce a clearer picture is to dispense with the explicit representation of any atoms and to represent the protein using a 'ribbon'. Proteins are also commonly represented using the cartoon drawings developed by J Richardson, an example of which is shown in Figure 1.5 (colour plate section). The cylinders in this figure represent an arrangement of amino acids called an α -helix, and the flat arrows an alternative type of regular structure called a β -strand. The regions between the cylinders and the strands have no such regular structure and are represented as 'tubes'.

1.5 Surfaces

Many of the problems that are studied using molecular modelling involve the non-covalent interaction between two or more molecules. The study of such interactions is often facilitated

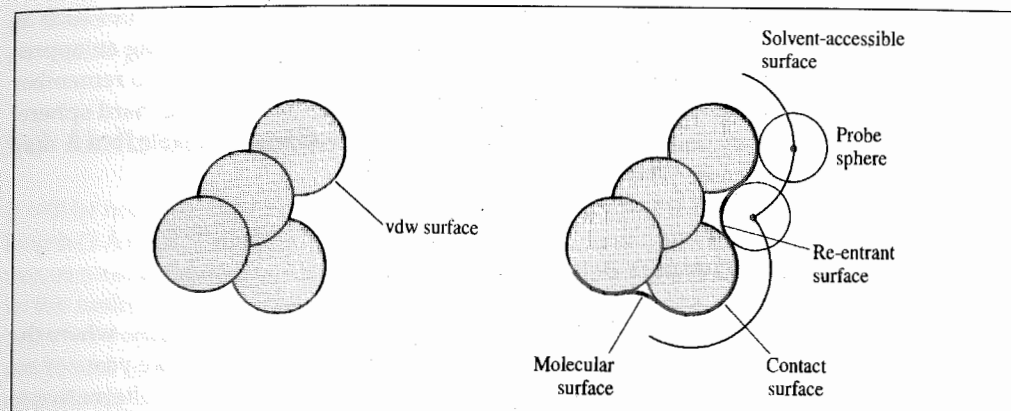


Fig. 1.6: The van der Waals (vdw) surface of a molecule corresponds to the outward-facing surfaces of the van der Waals spheres of the atoms. The molecular surface is generated by rolling a spherical probe (usually of radius 1.4 Å to represent a water molecule) on the van der Waals surface. The molecular surface is constructed from contact and re-entrant surface elements. The centre of the probe traces out the accessible surface.

by examining the van der Waals, molecular or accessible surfaces of the molecule. The *van der Waals surface* is simply constructed from the overlapping van der Waals spheres of the atoms, Figure 1.6. It corresponds to a CPK or space-filling model. Let us now consider the approach of a small 'probe' molecule, represented as a single van der Waals sphere, up to the van der Waals surface of a larger molecule. The finite size of the probe sphere means that there will be regions of 'dead space', crevices that are not accessible to the probe as it rolls about on the larger molecule. This is illustrated in Figure 1.6. The amount of dead space increases with the size of the probe; conversely, a probe of zero size would be able to access all of the crevices. The *molecular surface* [Richards 1977] is traced out by the inward-facing part of the probe sphere as it rolls on the van der Waals surface of the molecule. The molecular surface contains two different types of surface element. The *contact surface* corresponds to those regions where the probe is actually in contact with the van der Waals surface of the 'target'. The *re-entrant surface* regions occur where there are crevices that are too narrow for the probe molecule to penetrate. The molecular surface is usually defined using a water molecule as the probe, represented as a sphere of radius 1.4 Å.

The *accessible surface* is also widely used. As originally defined by Lee and Richards [Lee and Richards 1971] this is the surface that is traced by the centre of the probe molecule as it rolls on the van der Waals surface of the molecule (Figure 1.6). The centre of the probe molecule can thus be placed at any point on the accessible surface and not penetrate the van der Waals spheres of any of the atoms in the molecule.

Widely used algorithms for calculating the molecular and accessible surfaces were developed by Connolly [Connolly 1983a, b], and others [e.g. Richmond 1984] have described formulae for the calculation of exact or approximate values of the surface area. There are many ways to represent surfaces, some of which are illustrated in Figure 1.7 (colour plate section). As shown, it may also be possible to endow a surface with a translucent quality, which enables the molecule inside the surface to be displayed. Clipping can also be used

to cut through the surface to enable the 'inside' to be viewed. In addition, properties such as the electrostatic potential can be calculated on the surface and represented using an appropriate colour scheme. Useful though these representations are, it is important to remember that the electronic distribution in a molecule formally extends to infinity. The 'hard sphere' representation is often very convenient and has certainly proved very valuable, but it may not be appropriate in all cases [Rouvray 1997, 1999, 2000].

1.6 Computer Hardware and Software

One cannot fail to be amazed at the pace of development in the computer industry, where the ratio of performance-to-price has increased by an order of magnitude every five years or so. The workstations that are commonplace in many laboratories now offer a real alternative to centrally maintained 'supercomputers' for molecular modelling calculations, especially as a workstation or even a personal computer can be dedicated to a single task, whereas the super-computer has to be shared with many other users. Nevertheless, in the immediate future there will always be some calculations that require the power that only a supercomputer can offer. The speed of any computer system is ultimately constrained by the speed at which electrical signals can be transmitted. This means that there will come a time when no further enhancements can be made using machines with 'traditional' single-processor serial architectures, and parallel computers will play an ever more important role.

A parallel computer couples processors together in such a way that a calculation is divided into small pieces with the results being combined at the end. Some calculations are more amenable to parallel processing than others, and a significant amount of effort is being spent converting existing algorithms to run efficiently on parallel architectures. In some cases completely new methods have been developed to take maximum advantage of the opportunities of parallel processing. The low cost of personal computer chips means that large 'farms' of processors can be constructed to give significant computing power for relatively small outlay.

To perform molecular modelling calculations one also requires appropriate programs (the software). The software used by molecular modellers ranges from simple programs that perform just a single task to highly complex packages that integrate many different methods. There is also an extremely wide variation in the price of software! Some programs have been so widely used and tested that they can be considered to have reached the status of a 'gold standard' against which similar programs are compared. One hesitates to specify such programs in print, but three items of software have been so widely used and cited that they can safely be afforded the accolade. These are the Gaussian series of programs for performing *ab initio* quantum mechanics, the MOPAC/AMPAC programs for semi-empirical quantum mechanics and the MM2 program for molecular mechanics.

Various pieces of software were used to generate the data for the examples and illustrations throughout this book. Some of these were written specifically for the task; some were freely available programs; others were commercial packages. I have decided not to describe specific programs in any detail, as such descriptions rapidly become outdated. Nevertheless,

all items of software are accredited where appropriate. Please note that the use of any particular piece of software does not imply any recommendation!

1.7 Units of Length and Energy

It will be noted that our Z-matrix for ethane has been defined using the angstrom as the unit of length ($1 \text{ \AA} \equiv 10^{-10} \text{ m} \equiv 100 \text{ pm}$). The ångström is a non-SI unit but is a very convenient one to use, as most bond lengths are of the order of 1–2 Å. One other very common non-SI unit found in the molecular modelling literature is the kilocalorie ($1 \text{ kcal} \equiv 4.1840 \text{ kJ}$). Other systems of units are employed in other types of calculation, such as the atomic units used in quantum mechanics (discussed in Chapter 2). It is important to be aware of, and familiar with, these non-standard units as they are widely used in the literature and throughout this book.

1.8 The Molecular Modelling Literature

The number of scientific papers concerned with molecular modelling methods is rising rapidly, as is the number of journals in which such papers are published. This reflects the tremendous diversity of problems to which molecular modelling can be applied and the ever-increasing availability of molecular modelling methods. It does, however, mean that it can be very difficult to remain up to date with the field. A number of specialist journals are devoted to theoretical chemistry, computational chemistry and molecular modelling, each with their own particular emphasis. Relevant papers are also published in the more 'general' journals, and there are now a number of books covering aspects of molecular modelling, some aimed at the specialist reader, others at the beginner. Many scientists are now fortunate to have access to electronic catalogues of publications which can be searched to find relevant papers. As many journals are now available over the internet it is possible to perform a literature search and obtain copies of the relevant papers without even having to leave the office. Some of the journals which are devoted to short reviews of recent developments often include molecular modelling sections (such as the 'Current Opinion' series); in others, useful review articles appear on an occasional basis. One particularly valuable source of information on molecular modelling methods is the *Reviews in Computational Chemistry*, edited by Lipkowitz and Boyd, beginning in 1990 (see Further Reading). Each of these volumes contains chapters on a variety of subjects, each written by an appropriate expert. A recent addition is the *Encyclopaedia of Computational Chemistry* by Schleyer *et al.* (1998) (see Further Reading), which contains many chapters that cover a wide range of topics.

1.9 The Internet

In the first edition of this book I wrote, 'A major use of the Internet is for electronic mail, but extremely rapid growth is being observed in other areas, particularly the "World-Wide Web" (WWW) ...'. Such a phrase seems an understatement; despite the 'hype', the Internet has certainly made a dramatic impact, not least on the scientific community, where its

origins lie. Anything written about the Internet is almost certain to become obsolete more rapidly than any other topic in this book and so this section will be brief. I will assume that all readers of this book will be familiar with the use of a web browser and the concept of a hyperlink, which enables documents to be linked together. The URL (Uniform Resource Locator) is the currency of the WWW, being the 'electronic address' which enables the particular item to be identified. Most documents are still written using HTML (HyperText Markup Language) but increasingly incorporate more sophisticated features. Given the tremendous growth in the Web it is important to be able to locate relevant information. This is the role of the Internet search engines, which can be used to identify relevant sites of interest via some form of keyword search. Within the molecular modelling context, several trends can be noted. Whilst the Web was initially used to distribute mostly textual information, it is increasingly used for much more sophisticated applications. Interactive molecular graphics are a feature of many sites. Some sites enable calculations or database searches to be performed via the Web, with the results being delivered interactively or via email. This is particularly true for 'intranets' within an organisation. XML (eXtensible Markup Language) is likely to play an increasingly important role in the 'intelligent' exchange of information over the Web, especially in specialist areas such as chemistry [Murray-Rust and Rzepa 1999]. Several 'electronic conferences' have been held with participants from many different countries. Perhaps the only prediction that one can safely make about the Web is that it is here to stay and its use will continue to grow.

1.10 Mathematical Concepts

A full appreciation of all of the techniques of molecular modelling would require a mathematical treatment beyond that appropriate to a book of this size and scope. However, a proper understanding does benefit from some knowledge of mathematical concepts such as vectors, matrices, differential equations, complex numbers, series expansions and Lagrangian multipliers, and some very elementary statistical concepts. There is only space in this book for a cursory introduction to these mathematical concepts and ideas, with very brief descriptions and some key results. The suggestions for further reading provide detailed background information on all of the mathematical topics required.

1.10.1 Series Expansions

There are various series expansions that are useful for approximating functions. Particularly important is the *Taylor series*: if $f(x)$ is a continuous, single-valued function of x with continuous derivatives $f'(x), f''(x), \dots$, then we can expand the function about a point x_0 as follows:

$$f(x_0 + x) = f(x_0) + \frac{x}{1!} f'(x_0) + \frac{x^2}{2!} f''(x_0) + \frac{x^3}{3!} f'''(x_0) + \dots + \frac{x^n}{n!} f^{(n)}(x_0) \quad (1.1)$$

Taylor series are often truncated after the term involving the second derivative, which makes the function vary in a quadratic fashion. This is a common assumption in many of the minimisation algorithms that we will discuss in Chapter 5.

A *Maclaurin series* is a specific form of the Taylor series for which $x_0 = 0$. Some standard expansions in Taylor series form are:

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots \quad (1.2)$$

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots \quad (1.3)$$

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots \quad (1.4)$$

The *binomial expansion* is used for functions of the form $(1+x)^\alpha$:

$$(1+x)^\alpha = 1 + \alpha x + \alpha(\alpha-1)\frac{x^2}{2!} + \alpha(\alpha-1)(\alpha-2)\frac{x^3}{3!} + \dots \quad (1.5)$$

All these series must have $|x| < 1$ to be convergent.

1.10.2 Vectors

A vector is a quantity with both magnitude and direction. For example, the velocity of a moving body is a vector quantity as it defines both the direction in which the body is travelling and the speed at which it is moving. In Cartesian coordinates a vector such as the velocity will have three components, indicating the contribution to the overall motion from the component motions along the x , y and z directions. The addition and subtraction of vectors can be understood using geometrical constructions, as shown in Figure 1.8. Thus, if we want to calculate the force on an atom due to its interactions with all other atoms in the system (as required in molecular dynamics calculations, see Chapter 7), we would perform a vector sum of all the individual forces.

Some of the common manipulations that are performed with vectors include the scalar product, vector product and scalar triple product, which we will illustrate using vectors r_1 , r_2 and r_3 that are defined in a rectangular Cartesian coordinate system:

$$\begin{aligned} r_1 &= x_1\mathbf{i} + y_1\mathbf{j} + z_1\mathbf{k} \\ r_2 &= x_2\mathbf{i} + y_2\mathbf{j} + z_2\mathbf{k} \\ r_3 &= x_3\mathbf{i} + y_3\mathbf{j} + z_3\mathbf{k} \end{aligned} \quad (1.6)$$

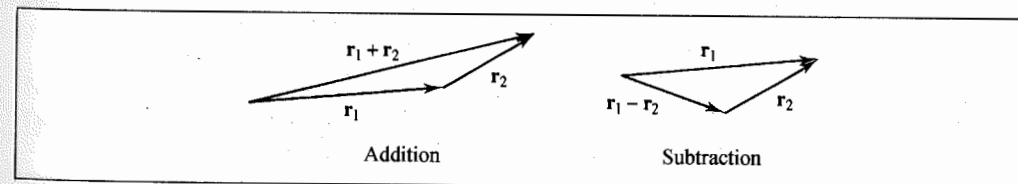


Fig. 1.8: The addition and subtraction of vectors.

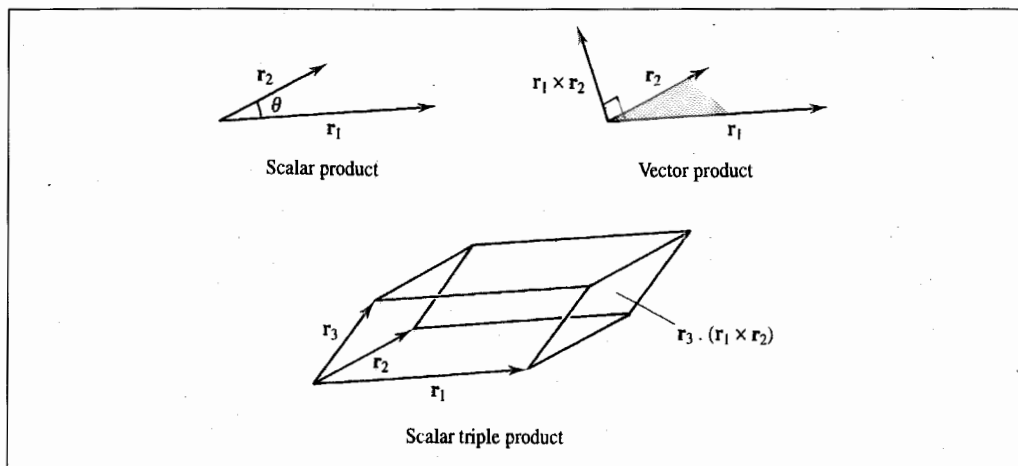


Fig. 1.9: The scalar product, vector product and scalar triple product.

i , j and k are orthogonal unit vectors along the x , y and z axes. The *scalar product* is defined as:

$$\mathbf{r}_1 \cdot \mathbf{r}_2 = |\mathbf{r}_1| |\mathbf{r}_2| \cos \theta \quad (1.7)$$

$|\mathbf{r}_1|$ and $|\mathbf{r}_2|$ are the magnitudes of the two vectors ($|\mathbf{r}_1| = \sqrt{x_1^2 + y_1^2 + z_1^2}$) and θ is the angle between them (Figure 1.9). The angle can be calculated as follows:

$$\cos \theta = \frac{x_1 x_2 + y_1 y_2 + z_1 z_2}{|\mathbf{r}_1| |\mathbf{r}_2|} \quad (1.8)$$

The scalar product of two vectors is thus a scalar.

The *vector product* of two vectors $\mathbf{r}_1 \times \mathbf{r}_2$ (sometimes written $\mathbf{r}_1 \wedge \mathbf{r}_2$) is a new vector (\mathbf{v}), in a direction perpendicular to the plane containing the two original vectors (Figure 1.9). The direction of this new vector is such that \mathbf{r}_1 , \mathbf{r}_2 and the new vector form a right-handed system. If \mathbf{r}_1 and \mathbf{r}_2 are three-component vectors then the components of \mathbf{v} are given by:

$$\mathbf{v} = (y_1 z_2 - z_1 y_2) \mathbf{i} + (z_1 x_2 - x_1 z_2) \mathbf{j} + (x_1 y_2 - y_1 x_2) \mathbf{k} \quad (1.9)$$

Note that the vector product $\mathbf{r}_2 \times \mathbf{r}_1$ is not the same as the vector product $\mathbf{r}_1 \times \mathbf{r}_2$, as it corresponds to a vector in the opposite direction. The vector product is thus not commutative.

The *scalar triple product* $\mathbf{r}_1 \cdot (\mathbf{r}_2 \times \mathbf{r}_3)$ equals the scalar product of \mathbf{r}_1 with the vector product of \mathbf{r}_2 and \mathbf{r}_3 . The result is a scalar. The scalar triple product has a useful geometrical interpretation; it is the volume of the parallelepiped whose sides correspond to the three vectors (Figure 1.9).

1.10.3 Matrices, Eigenvectors and Eigenvalues

A matrix is a set of quantities arranged in a rectangular array. An $m \times n$ matrix has m rows and n columns. A vector can thus be considered to be a one-column matrix. Matrix addition

and subtraction can only be performed with matrices of the same order. For example:

$$\text{If } \mathbf{A} = \begin{pmatrix} 4 & 7 \\ -3 & 5 \\ 8 & -2 \end{pmatrix} \text{ and } \mathbf{B} = \begin{pmatrix} -4 & 3 \\ 5 & 2 \\ -5 & 3 \end{pmatrix}$$

$$\text{Then } \mathbf{A} + \mathbf{B} = \begin{pmatrix} 0 & 10 \\ 2 & 7 \\ 3 & 1 \end{pmatrix}; \quad \mathbf{A} - \mathbf{B} = \begin{pmatrix} 8 & 4 \\ -8 & 3 \\ 12 & -5 \end{pmatrix} \quad (1.10)$$

Multiplication of two matrices (\mathbf{AB}) is only possible if the number of columns in \mathbf{A} is equal to the number of rows in \mathbf{B} . If \mathbf{A} is an $m \times n$ matrix and \mathbf{B} is an $n \times o$ matrix then the product \mathbf{AB} is an $m \times o$ matrix. Each element (i, j) in the matrix \mathbf{AB} is obtained by taking each of the n values in the i th row of \mathbf{A} and multiplying by the corresponding value in the j th column of \mathbf{B} . To illustrate with a simple example:

$$\text{If } \mathbf{A} = \begin{pmatrix} 3 & -2 & 5 \\ -3 & 4 & 1 \end{pmatrix} \text{ and } \mathbf{B} = \begin{pmatrix} 0 & 3 \\ -2 & 4 \\ 1 & 6 \end{pmatrix}$$

Then

$$\mathbf{AB} = \begin{pmatrix} (3 \times 0) + (-2 \times -2) + (5 \times 1) & (3 \times 3) + (-2 \times 4) + (5 \times 6) \\ (-3 \times 0) + (4 \times -2) + (1 \times 1) & (-3 \times 3) + (4 \times 4) + (1 \times 6) \end{pmatrix}$$

$$= \begin{pmatrix} 9 & 31 \\ -7 & 13 \end{pmatrix} \quad (1.11)$$

We shall often encounter square matrices, which have the same number of rows and columns. A diagonal matrix is a square matrix in which all the elements are zero except for those on the diagonal. The *unit* or *identity* matrix \mathbf{I} is a special type of diagonal matrix in which all the non-zero elements are 1; thus the 3×3 unit matrix is:

$$\mathbf{I} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (1.12)$$

A matrix is *symmetric* if it is a square matrix with elements such that the elements above and below the diagonal are mirror images; $A_{ij} = A_{ji}$.

Multiplication of a matrix by its inverse gives the unit matrix:

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{I} \quad (1.13)$$

To compute the inverse of a square matrix it is necessary to first calculate its *determinant*, $|\mathbf{A}|$. The determinants of 2×2 and 3×3 matrices are calculated as follows:

$$\begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc \quad (1.14)$$

$$\begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = a \begin{vmatrix} e & f \\ h & i \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix} \\ = a(ei - hf) - b(di - fg) + c(dh - eg) \quad (1.15)$$

For example:

$$\begin{vmatrix} 3 & 6 \\ -2 & 3 \end{vmatrix} = 21; \quad \begin{vmatrix} 4 & 2 & -2 \\ 2 & 5 & 0 \\ -2 & 0 & 3 \end{vmatrix} = 28 \quad (1.16)$$

As can be seen, the determinant of a 3×3 matrix can be written as a sum of determinants of 2×2 matrices, obtained by first selecting one of the rows or columns in the matrix (the top row was chosen in our example). For each element A_{ij} in this row, the row and column in which that number appears are deleted (i.e. the i th row and the j th column). This leaves a 2×2 matrix whose determinant is calculated and then multiplied by $(-1)^{i+j}$. The result of this calculation is called the *cofactor* of the element A_{ij} . For example, the cofactor of the element A_{12} in the 3×3 matrix

$$\mathbf{A} = \begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & 0 \\ -2 & 0 & 3 \end{pmatrix}$$

is -6 . When calculating the determinant the cofactor is multiplied by the element A_{ij} . The determinants of larger matrices can be obtained by extensions of the scheme illustrated above; thus the determinant of a 4×4 matrix is initially written in terms of 3×3 matrices, which in turn can be expressed in terms of 2×2 matrices.

Determinants have many useful and interesting properties. The determinant of a matrix is zero if any two of its rows or columns are identical. The sign of the determinant is reversed by exchanging any pair of rows or any pair of columns. If all elements of a row (or column) are multiplied by the same number, then the value of the determinant is multiplied by that number. The value of a determinant is unaffected if equal multiples of the values in any row (or column) are added to another row (or column).

The vector product and the scalar triple product can be conveniently written as matrix determinants. Thus:

$$\mathbf{r}_1 \times \mathbf{r}_2 = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \end{vmatrix} \quad (1.17)$$

$$\mathbf{r}_1 \cdot (\mathbf{r}_2 \times \mathbf{r}_3) = \begin{vmatrix} x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \\ x_3 & y_3 & z_3 \end{vmatrix} \quad (1.18)$$

The *transpose* of a matrix, \mathbf{A}^T , is the matrix obtained by exchanging its rows and columns. Thus the transpose of an $m \times n$ matrix is an $n \times m$ matrix:

$$\text{If } \mathbf{A} = \begin{pmatrix} 4 & 7 \\ -3 & 5 \\ 8 & -2 \end{pmatrix} \quad \mathbf{A}^T = \begin{pmatrix} 4 & -3 & 8 \\ 7 & 5 & -2 \end{pmatrix} \quad (1.19)$$

The transpose of a square matrix is, of course, another square matrix. The transpose of a symmetric matrix is itself. One particularly important transpose matrix is the *adjoint* matrix, $\text{adj}\mathbf{A}$, which is the transpose matrix of cofactors. For example, the matrix of cofactors of the 3×3 matrix

$$\mathbf{A} = \begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & 0 \\ -2 & 0 & 3 \end{pmatrix} \quad \text{is} \quad \begin{pmatrix} 15 & -6 & 10 \\ -6 & 8 & -4 \\ 10 & -4 & 16 \end{pmatrix} \quad (1.20)$$

In this case the adjoint matrix is the same as the matrix of cofactors (as \mathbf{A} is a symmetric matrix). The *inverse* of a matrix is obtained by dividing the elements of the adjoint matrix by the determinant:

$$\mathbf{A}^{-1} = \frac{\text{adj}\mathbf{A}}{|\mathbf{A}|} \quad (1.21)$$

Thus the inverse of our 3×3 matrix is

$$\mathbf{A}^{-1} = \begin{pmatrix} 15/28 & -3/14 & 5/14 \\ -3/14 & 2/7 & -1/7 \\ 5/14 & -4 & 4/7 \end{pmatrix} \quad (1.22)$$

One of the most common matrix calculations involves finding its *eigenvalues* and *eigenvectors*. An eigenvector is a column matrix \mathbf{x} such that

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \quad (1.23)$$

λ is the associated eigenvalue. The eigenvector problem can be reformulated as follows:

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}\mathbf{I} \Rightarrow \mathbf{A}\mathbf{x} - \lambda\mathbf{x}\mathbf{I} = 0 \Rightarrow (\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = 0 \quad (1.24)$$

A trivial solution to this equation is $\mathbf{x} = 0$. For a non-trivial solution, we require that the determinant $|\mathbf{A} - \lambda\mathbf{I}|$ equals zero. One way to determine the eigenvalues and their associated eigenvectors is thus to expand the determinant to give a polynomial equation in λ . For our 3×3 symmetric matrix this gives:

$$\begin{vmatrix} 4 - \lambda & 2 & -2 \\ 2 & 5 - \lambda & 0 \\ -2 & 0 & 3 - \lambda \end{vmatrix} \quad (1.25)$$

or:

$$(4 - \lambda)(5 - \lambda)(3 - \lambda) - 2[2(3 - \lambda)] - 2[2(5 - \lambda)] = 0 \quad (1.26)$$

This can be factorised to give:

$$(1 - \lambda)(7 - \lambda)(4 - \lambda) = 0 \quad (1.27)$$

The eigenvalues are thus $\lambda_1 = 1$, $\lambda_2 = 4$, $\lambda_3 = 7$. The corresponding eigenvectors are:

$$\lambda_1 = 1 : \mathbf{x}_1 = \begin{pmatrix} 2/3 \\ -1/3 \\ 2/3 \end{pmatrix} \quad \lambda_2 = 4 : \mathbf{x}_2 = \begin{pmatrix} -1/3 \\ 2/3 \\ 2/3 \end{pmatrix} \quad \lambda_3 = 7 : \mathbf{x}_3 = \begin{pmatrix} 2/3 \\ 2/3 \\ -1/3 \end{pmatrix} \quad (1.28)$$

Here we have expressed the eigenvectors as vectors of unit length; any multiple of each eigenvector would also be a solution. \mathbf{A} is a real, symmetric matrix. The eigenvalues of such matrices are always real and orthogonal (i.e. the scalar products of all pairs of eigenvectors are zero). This can be easily seen in our example.

As can be readily envisaged, expanding the determinant and solving a polynomial in λ is not the most efficient way to determine the eigenvalues and eigenvectors of larger matrices. Matrix diagonalisation methods are much more common. *Diagonalisation* of a matrix \mathbf{A} involves finding a matrix \mathbf{U} such that:

$$\mathbf{U}^{-1}\mathbf{A}\mathbf{U} = \mathbf{D} \quad (1.29)$$

\mathbf{D} is the diagonal matrix of eigenvalues. When \mathbf{A} is a real symmetric matrix, then \mathbf{U} is the matrix of eigenvectors and \mathbf{U}^{-1} is the inverse matrix of eigenvectors. Thus, for our example:

$$\begin{pmatrix} 2/3 & -1/3 & 2/3 \\ -1/3 & 2/3 & 2/3 \\ 2/3 & 2/3 & -1/3 \end{pmatrix} \begin{pmatrix} 4 & 2 & -2 \\ 2 & 5 & 0 \\ -2 & 0 & 3 \end{pmatrix} \begin{pmatrix} 2/3 & -1/3 & 2/3 \\ -1/3 & 2/3 & 2/3 \\ 2/3 & 2/3 & -1/3 \end{pmatrix} \\ = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 7 \end{pmatrix} \quad (1.30)$$

Note that for a real symmetric matrix \mathbf{A} , the inverse \mathbf{U}^{-1} is the same as the transpose, \mathbf{U}^T .

Many methods have been devised for diagonalising matrices; some of these are specific to certain classes of matrices such as the class of real symmetric matrices. Many modelling techniques require us to calculate the eigenvalues and eigenvectors of a matrix, including self-consistent field quantum mechanics (Section 2.5), the distance geometry method for exploring conformational space (Section 9.5) and principal components analysis (Section 9.13.1). The class of *positive definite* matrices is important in energy minimisation and when finding transition structures; the eigenvalues of a positive definite matrix are all positive. A *positive semidefinite* matrix of rank m has m positive eigenvalues.

1.10.4 Complex Numbers

A complex number has two components: a real part (a) and an imaginary part (b), as follows:

$$x = a + bi \quad (1.31)$$

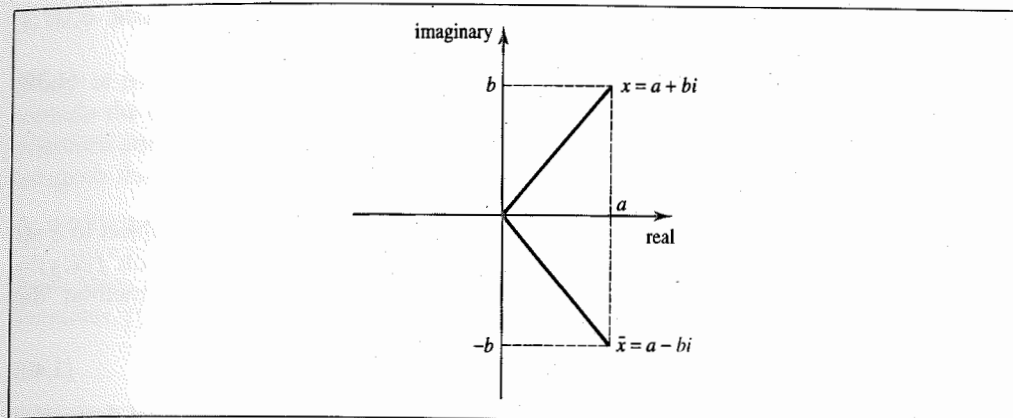


Fig. 1.10: The Argand diagram used to represent complex numbers.

i is the square root of -1 ($i = \sqrt{-1}$). Complex numbers enable certain types of equation that have no real solutions to be solved. For example, the roots of the equation $x^2 - 2x + 3 = 0$ are $x = 1 + \sqrt{2}i$ and $x = 1 - \sqrt{2}i$. A complex number can be considered as a vector in a two-dimensional coordinate system. Complex numbers are commonly represented using an *Argand diagram*, in which the x coordinate corresponds to the real part of the complex number and the y coordinate to the imaginary part (Figure 1.10).

Arithmetical operations on complex numbers are performed much as for vectors. Thus, if $x = a + bi$ and $y = c + di$, then:

$$x + y = (a + c) + (b + d)i \quad (1.32)$$

$$x - y = (a - c) + (b - d)i \quad (1.33)$$

$$xy = (ac - bd) + (ad + bc)i \quad (1.34)$$

The *complex conjugate*, \bar{x} , equals $a - bi$ and is obtained by reflecting x in the real axis in the Argand diagram.

A commonly used relationship involving complex numbers is:

$$e^{i\theta} = \cos \theta + i \sin \theta \quad (1.35)$$

where θ is any real number. This relationship is used in Fourier analysis and can be derived from the expansions of the exponential, cosine and sine functions:

$$e^{i\theta} = 1 + i\theta - \frac{\theta^2}{2!} - \frac{i\theta^3}{3!} + \frac{\theta^4}{4!} + \dots \quad (1.36)$$

$$\sin \theta = \theta - \frac{\theta^3}{3!} + \frac{\theta^5}{5!} - \dots \quad (1.37)$$

$$\cos \theta = 1 - \frac{\theta^2}{2!} + \frac{\theta^4}{4!} - \dots \quad (1.38)$$

Various other relationships can be defined. For example:

$$\cos \theta = \frac{e^{i\theta} + e^{-i\theta}}{2} \quad \sin \theta = \frac{e^{i\theta} - e^{-i\theta}}{2i} \quad (1.39)$$

1.10.5 Lagrange Multipliers

Lagrange multipliers can be used to find the stationary points of functions, subject to a set of constraints. Suppose we wish to find the stationary points of a function $f(x, y) = 4x^2 + 3x + 2y^2 + 6y$ subject to the constraint $y = 4x + 2$. In the Lagrange method the constraint is written in the form $g(x, y) = 0$:

$$g(x, y) = y - 4x - 2 = 0 \quad (1.40)$$

To find stationary points $f(x, y)$ subject to $g(x, y) = 0$ we first determine the total derivative df , which is set equal to zero:

$$df = \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy = (8x + 3) dx + (4y + 6) dy = 0 \quad (1.41)$$

Without the constraint the stationary points would be determined by setting the two partial derivatives $\partial f/\partial x$ and $\partial f/\partial y$ equal to zero, as x and y are independent. With the constraint, x and y are no longer independent but are related via the derivative of the constraint function g :

$$dg = \frac{\partial g}{\partial x} dx + \frac{\partial g}{\partial y} dy = -4dx + dy = 0 \quad (1.42)$$

The derivative of the constraint function, dg , is multiplied by a parameter λ (the *Lagrange multiplier*) and added to the total derivative df :

$$\left(\frac{\partial f}{\partial x} + \lambda \frac{\partial g}{\partial x} \right) dx + \left(\frac{\partial f}{\partial y} + \lambda \frac{\partial g}{\partial y} \right) dy = 0 \quad (1.43)$$

The value of the Lagrange multiplier is obtained by setting each of the terms in parentheses to zero. Thus for our example we have:

$$8x + 3 - 4\lambda = 0 \quad (1.44)$$

$$4y + 6 + \lambda = 0 \quad (1.45)$$

From these two equations we can obtain a further equation linking x and y :

$$\lambda = 2x + 3/4 = -6 - 4y \quad \text{or} \quad x = -27/8 - 2y \quad (1.46)$$

Combining this with the constraint equation enables us to identify the stationary point, which is at $(-59/72, -23/18)$.

This simple example could, of course, have been solved by simply substituting the constraint equation into the original function, to give a function of just one of the variables. However, in many cases this is not possible. The Lagrange multiplier method provides a powerful approach which is widely applicable to problems involving constraints such as in constraint dynamics (Section 7.5) and in quantum mechanics.

1.10.6 Multiple Integrals

Many of the theories used in molecular modelling involve multiple integrals. Examples include the two-electron integrals found in Hartree-Fock theory, and the integral over the positions and momenta used to define the partition function, Q . In fact, most of the multiple integrals that have to be evaluated are double integrals.

A 'traditional' or one-dimensional integral corresponds to the area under the curve between the imposed limit, as illustrated in Figure 1.11. Multiple integrals are simply extensions of these ideas to more dimensions. We shall illustrate the principles using a function of two variables, $f(x, y)$. The double integral

$$\iint_A dx dy f(x, y) \equiv \iint_A f(x, y) dx dy \quad (1.47)$$

is the sum of the volume elements $f(x, y) \delta x \delta y$ (see Figure 1.11) over the area A as δx and δy tend to zero. Note that the ' $dx dy$ ' can be put either immediately after the integral sign or at the end; in this book we often use the first method for multiple integrals.

Some multiple integrals can be written as a product of single integrals. This occurs when $f(x, y)$ is itself a product of functions $g(x)h(y)$, in which case the integral can be separated:

$$\iint_A dx dy g(x)h(y) = \int dx g(x) \int dy h(y) \quad (1.48)$$

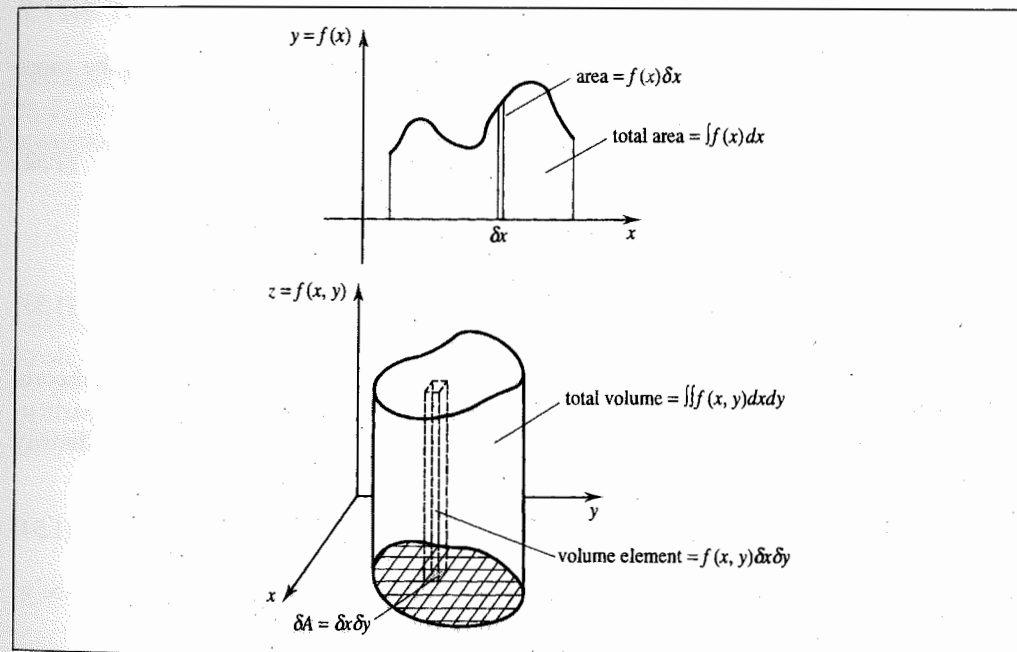


Fig. 1.11: Single and double integrals. (Figure adapted in part from Boas M L, 1983, *Mathematical Methods in the Physical Sciences*. 2nd Edition. New York, Wiley.)

For example:

$$\int_{-1}^1 dx \int_{-\pi/2}^{+\pi/2} dy x^2 \cos y = \int_{-1}^1 x^2 dx [\sin y]_{-\pi/2}^{+\pi/2} = 2 \left(\frac{x^3}{3} \right)_{-1}^{+1} = \frac{4}{3} \quad (1.49)$$

We will use the separation of multiple integrals throughout our discussion of quantum mechanics and computer simulation methods (Chapters 2, 3, 6, 7 and 8).

1.10.7 Some Basic Elements of Statistics

Statistics is concerned with the collection and interpretation of numerical data. The subject is a vast and complex one, and all we shall do here is to state some of the definitions commonly used and to explain some of the terminology.

The *arithmetic mean* of a set of observations is the sum of the observations divided by the number of observations:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad (1.50)$$

N is the number of observations. The mean may also be written $\langle x \rangle$. The *variance*, σ^2 , indicates the extent to which the set of observations cluster around the mean value and equals the average of the squared deviations from the mean:

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (1.51)$$

The variance can also be calculated using the following formula, which may be more convenient:

$$\sigma^2 = \frac{1}{N} \left[\sum_{i=1}^N (x_i^2) - \frac{1}{N} \left(\sum_{i=1}^N x_i \right)^2 \right] \quad (1.52)$$

The *standard deviation*, σ , equals the (positive) square root of the variance:

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} \quad (1.53)$$

It is often desired to compare the distribution of observations in a population with a theoretical distribution. The *normal distribution* (also called the Gaussian distribution) is a particularly important theoretical distribution in molecular modelling. The probability density function for a general normal distribution is:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp[-(x - \bar{x})^2/2\sigma^2] \quad (1.54)$$

The factor before the exponential ensures that the integral of the function $f(x)$ from $-\infty$ to $+\infty$ equals 1. The distribution is often written in terms of a parameter α :

$$f(x) = \sqrt{\frac{\alpha}{\pi}} e^{-\alpha(x - \bar{x})^2} \quad (1.55)$$

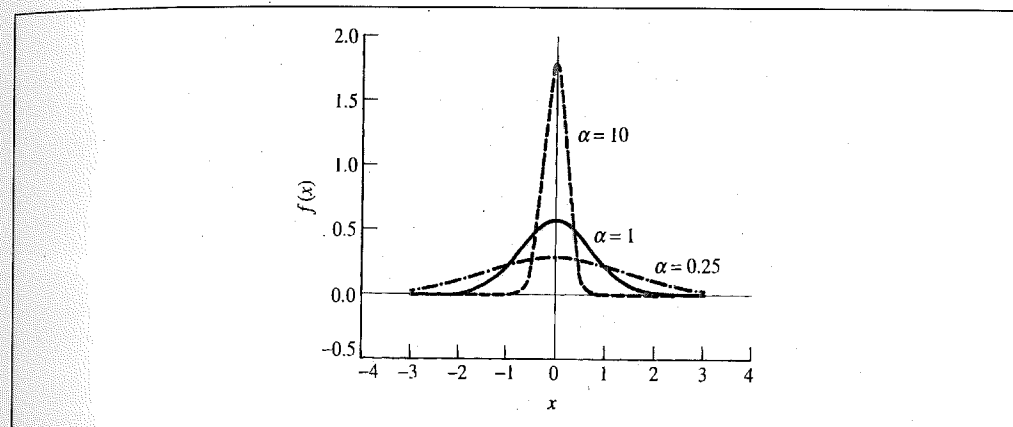


Fig. 1.12: Three normal distributions with different values of α (Equation (1.55)). The functions are normalised, so the area under each curve is the same.

In Figure 1.12 we show three normal distributions that all have zero mean but different values of the variance (σ^2). A variance larger than 1 (small α) gives a flatter function and a variance less than 1 (larger α) gives a sharper function.

1.10.8 The Fourier Series, Fourier Transform and Fast Fourier Transform

Consider a periodic function $x(t)$ that repeats between $t = -\tau/2$ and $t = +\tau/2$ (i.e. has period τ). Even though $x(t)$ may not correspond to an analytical expression it can be written as the superposition of simple sine and cosine functions or *Fourier series*, Figure 1.13.

$$x(t) = a_0 + a_1 \cos \omega_0 t + a_2 \cos 2\omega_0 t + \dots + b_1 \sin \omega_0 t + b_2 \sin 2\omega_0 t + \dots \quad (1.56)$$

$$x(t) = a_0 + \sum_{n=1}^{\infty} (a_n \cos n\omega_0 t + b_n \sin n\omega_0 t) \quad (1.57)$$

ω_0 is related to the period of the function by $\omega_0 = 2\pi/\tau$ and to the frequency of the function by $\omega_0 = 2\pi\nu_0$. The frequencies of the contributing harmonics are thus $n\nu_0$ and are separated by $1/\tau$.

The coefficients a_n and b_n can be obtained as follows:

$$a_0 = \frac{1}{\tau} \int_{-\tau/2}^{\tau/2} x(t) dt \quad (1.58)$$

$$a_n = \frac{2}{\tau} \int_{-\tau/2}^{\tau/2} x(t) \cos(2n\pi x/\tau) dx \quad (1.59)$$

$$b_n = \frac{2}{\tau} \int_{-\tau/2}^{\tau/2} x(t) \sin(2n\pi x/\tau) dx \quad (1.60)$$

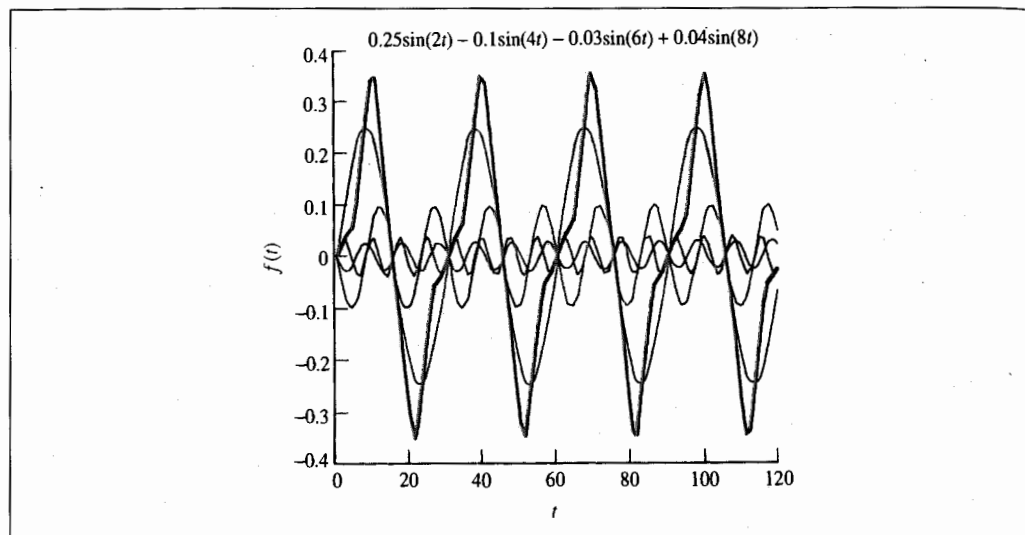


Fig. 1.13: In a Fourier series a periodic function is expressed as a sum of sine and cosine functions.

An alternative way to express a Fourier series makes use of the following relationships:

$$\sin \omega_0 t = [\exp(i\omega_0 t) - \exp(-i\omega_0 t)]/2i \quad (1.61)$$

$$\cos \omega_0 t = [\exp(i\omega_0 t) + \exp(-i\omega_0 t)]/2 \quad (1.62)$$

The Fourier series is then written

$$x(t) = \sum_{-\infty}^{+\infty} c_n \exp(in\omega_0 t) \quad (1.63)$$

with

$$c_n = \frac{1}{\tau} \int_{-\tau/2}^{\tau/2} x(t) \exp(in\omega_0 t) dt \quad (1.64)$$

The Fourier series is used to represent a function that is periodic with period τ in terms of frequencies $n\omega_0 = 2\pi n/\tau$. The *Fourier transform* is used when the function has no periodicity. There is a close relationship between the Fourier series and the Fourier transform. One way to demonstrate the gradual change from a Fourier series to a Fourier transform is to consider how the distribution of contributing frequencies changes as the period increases. This is illustrated in Figure 1.14, where the period of a square wavefunction is gradually increased. Also shown are the frequency contributions. It can be seen that an increasing number of frequency components is needed to describe the function as the period increases, and that when the period is infinite, the frequency spectrum is continuous.

The Fourier transform relationship between a function $x(t)$ and the corresponding frequency function $X(\nu)$ is:

$$x(t) = \int_{-\infty}^{+\infty} X(\nu) \exp(2\pi i\nu t) d\nu \quad (1.65)$$

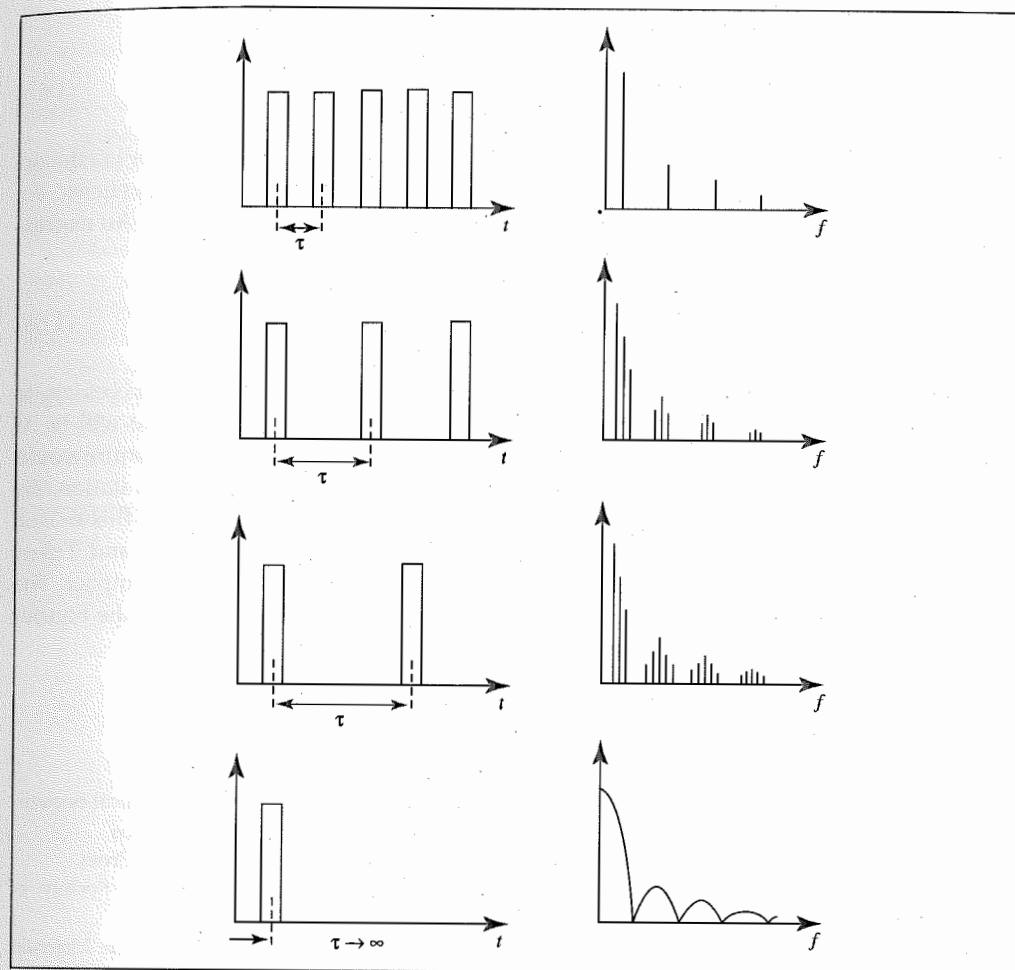


Fig. 1.14: The connection between the Fourier transform and the Fourier series can be established by gradually increasing the period of the function. When the period is infinite a continuous spectrum is obtained. (Figure adapted from Ramirez R W, 1985, *The FFT Fundamentals and Concepts*. Englewood Cliffs, NJ, Prentice Hall.)

The frequency function $X(\nu)$ is given by:

$$X(\nu) = \int_{-\infty}^{+\infty} x(t) \exp(-2\pi i\nu t) dt \quad (1.66)$$

In practical applications, $x(t)$ is not a continuous function, and the data to be transformed are usually discrete values obtained by sampling at intervals. Under such circumstances, the discrete Fourier transform (DFT) is used to obtain the frequency function. Let us suppose that the time-dependent data values are obtained by sampling at regular intervals separated by δt and that a total of M samples are obtained (starting at $t = 0$). From M samples, a total of M frequency coefficients can be obtained using the DFT expression

[Press *et al.* 1992]:

$$X(k\delta\nu) = \delta t \sum_{n=0}^{M-1} x(n\delta t) \exp[-2\pi ink/M] \quad (1.67)$$

Here, $x(n\delta t)$ ($n = 0, 1, \dots, M-1$) are the experimental values obtained and $X(k\delta\nu)$ is the set of Fourier coefficients ($k = 0, 1, \dots, M-1$). The separation between the frequencies, $\delta\nu$, depends on the number of samples and the time between samples: $\delta\nu = 1/M\delta t$. An expression for converting frequency data into the time domain is also possible:

$$x(n\delta t) = \frac{1}{M} \sum_{k=0}^{M-1} X(k\delta\nu) \exp[2\pi ink/M] \quad (1.68)$$

To compute each Fourier coefficient $X(k\delta T)$ (of which there are M) it is therefore necessary to evaluate the summation $\sum_{n=0}^{M-1} x(n\delta t) \exp[-2\pi ink/M]$ for that value of k . There will be M terms in the summation. A simple algorithm to determine the frequency spectrum would scale with the square of the number of measurements, M . This is a severe limitation, for many problems involve an extremely large number of pieces of data. It is for this reason that the fast Fourier transform (FFT) (ascribed to Cooley and Tukey [Cooley and Tukey 1965] but, in fact, using methods developed much earlier) has made such an impact. The FFT algorithm scales as $M \ln M$. With the FFT algorithm it is possible to derive the Fourier transforms, even with a considerable number of data points.

Further Reading

- Bachrach S M 1996. *The Internet: A Guide for Chemists*. Washington, D.C., American Chemical Society.
- Boas M L 1983. *Mathematical Methods in the Physical Sciences*. New York, John Wiley & Sons.
- Grant G H and W G Richards 1995. *Computational Chemistry*. Oxford, Oxford University Press.
- Goodman J M 1998. *Chemical Applications of Molecular Modelling*. Cambridge, Royal Society of Chemistry.
- Leach A R 1999. *Computational Chemistry and the Virtual Laboratory*. In *The Age of the Molecule*. Cambridge, Royal Society of Chemistry.
- Lipkowitz K B and D B Boyd (Editors) 1990-. *Reviews in Computational Chemistry* Vols 1-. New York, VCH.
- Ramirez R W 1985. *The FFT Fundamentals and Concepts*. Englewood Cliffs, NJ, Prentice Hall.
- Schleyer, P v R, N L Allinger, T Clark, J Gasteiger, P A Kollman, H F Schaefer III and P R Schreiner 1998. *The Encyclopedia of Computational Chemistry*. Chichester, John Wiley & Sons.
- Stephenson G 1973. *Mathematical Methods for Science Students*. London, Longman.
- Winter M J, H S Rzepa and B J Whitaker 1995. Surfing the Chemical Net. *Chemistry in Britain* **31**: 685-689 and <http://www.ch.ic.ac.uk/rzepa/cib/>.

References

- Bolin J T, D J Filman, D A Matthews, R C Hamlin and J Kraut 1982. Crystal Structures of *Escherichia coli* and *Lactobacillus casei* Dihydrofolate Reductase Refined at 1.7 Ångstroms Resolution. I. Features and Binding of Methotrexate. *Journal of Biological Chemistry* **257**:13650-13662.

- Connolly M L 1983a. Solvent-accessible Surfaces of Proteins and Nucleic Acids. *Science* **221**:709-713.
- Connolly M L 1983b. Analytical Molecular Surface Calculation. *Journal of Applied Crystallography* **16**:548-558.
- Cooley J W and J W Tukey 1965. An Algorithm for the Machine Calculation of Complex Fourier Series. *Mathematics of Computation* **19**:297-301.
- Lee B and F M Richards 1971. The Interpretation of Protein Structures: Estimation of Static Accessibility. *Journal of Molecular Biology* **55**:379-400.
- Murray-Rust P and H Rzepa 1999. Chemical Markup, XML, and the Worldwide Web. 1. Basic Principles. *Journal of Chemical Information and Computer Science* **39**:923-942.
- Press W H, B P Flannery, S A Teukolsky and W T Vetterling 1992. *Numerical Recipes in Fortran*. Cambridge, Cambridge University Press.
- Richards F M 1977. Areas, Volumes, Packing and Protein Structure. *Annual Review in Biophysics and Bioengineering* **6**:151-176.
- Richmond T J 1984. Solvent Accessible Surface Area and Excluded Volume in Proteins. *Journal of Molecular Biology* **178**:63-88.
- Rouvray D 1997. Do Molecular Models Accurately Reflect Reality? *Chemist in Industry* **15**:587-590.
- Rouvray D 1999. Model Answers. *Chemistry in Britain* **35**:30-32.
- Rouvray D 2000. Atoms as Hard Spheres. *Chemistry in Britain* **36**:25.

CHAPTER TWO

An Introduction to Computational Quantum Mechanics

2.1 Introduction

Our aim in this chapter will be to establish the basic elements of those quantum mechanical methods that are most widely used in molecular modelling. We shall assume some familiarity with the elementary concepts of quantum mechanics as found in most 'general' physical chemistry textbooks, but little else other than some basic mathematics (see Section 1.10). There are also many excellent introductory texts to quantum mechanics. In Chapter 3 we then build upon this chapter and consider more advanced concepts. Quantum mechanics does, of course, predate the first computers by many years, and it is a tribute to the pioneers in the field that so many of the methods in common use today are based upon their efforts. The early applications were restricted to atomic, diatomic or highly symmetrical systems which could be solved by hand. The development of quantum mechanical techniques that are more generally applicable and that can be implemented on a computer (thereby eliminating the need for much laborious hand calculation) means that quantum mechanics can now be used to perform calculations on molecular systems of real, practical interest. Quantum mechanics explicitly represents the electrons in a calculation, and so it is possible to derive properties that depend upon the electronic distribution and, in particular, to investigate chemical reactions in which bonds are broken and formed. These qualities, which differentiate quantum mechanics from the empirical force field methods described in Chapter 4, will be emphasised in our discussion of typical applications.

There are a number of quantum theories for treating molecular systems. The first we shall examine, and the one which has been most widely used, is *molecular orbital theory*. However, alternative approaches have been developed, some of which we shall also describe, albeit briefly. We will be primarily concerned with the *ab initio* and semi-empirical approaches to quantum mechanics but will also mention techniques such as Hückel theory and valence bond theory. An alternative approach to quantum mechanics, density functional theory, is considered in Chapter 3. Density functional theory has always enjoyed significant support from the materials science community but is increasingly used for molecular systems.

Quantum mechanics is often considered to be a difficult subject, and a cursory glance at the following pages in this chapter may simply serve to reinforce that view! However, if followed carefully it is possible to see how models that are developed for very simple

systems can be applied to much more complex systems. As a consequence our treatment does require some consideration of the mathematical background to the simplest and most common types of calculation. Our strategy in developing the underlying theory of molecular orbital quantum mechanical calculations is as follows. First, we revise some key features of quantum mechanics, including the hydrogen atom. We then discuss the functional form of an acceptable wavefunction for a molecular system and show how to calculate the energy of such a system from the wavefunction. This leads to the problem of determining the wavefunction itself and how this can be done using routine mathematical methods. We will then be in a position to understand how quantum mechanical calculations can be performed for 'real' systems and will have the background necessary to consider more advanced topics.

The starting point for any discussion of quantum mechanics is, of course, the Schrödinger equation. The full, time-dependent form of this equation is

$$\left\{ -\frac{\hbar^2}{2m} \left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} \right) + \mathcal{V} \right\} \Psi(\mathbf{r}, t) = i\hbar \frac{\partial \Psi(\mathbf{r}, t)}{\partial t} \quad (2.1)$$

Equation (2.1) refers to a single particle (e.g. an electron) of mass m which is moving through space (given by a position vector $\mathbf{r} = x\mathbf{i} + y\mathbf{j} + z\mathbf{k}$) and time (t) under the influence of an external field \mathcal{V} (which might be the electrostatic potential due to the nuclei of a molecule). \hbar is Planck's constant divided by 2π and i is the square root of -1 . Ψ is the *wavefunction* which characterises the particle's motion; it is from the wavefunction that we can derive various properties of the particle. When the external potential \mathcal{V} is independent of time then the wavefunction can be written as the product of a spatial part and a time part: $\Psi(\mathbf{r}, t) = \psi(\mathbf{r})T(t)$. We shall only consider situations where the potential is independent of time, which enables the time-dependent Schrödinger equation to be written in the more familiar, time-independent form:

$$\left\{ -\frac{\hbar^2}{2m} \nabla^2 + \mathcal{V} \right\} \Psi(\mathbf{r}) = E\Psi(\mathbf{r}) \quad (2.2)$$

Here, E is the energy of the particle and we have used the abbreviation ∇^2 (pronounced 'del-squared'):

$$\nabla^2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} \quad (2.3)$$

It is usual to abbreviate the left-hand side of Equation (2.1) to $\mathcal{H}\Psi$, where \mathcal{H} is the *Hamiltonian operator*:

$$\mathcal{H} = -\frac{\hbar^2}{2m} \nabla^2 + \mathcal{V} \quad (2.4)$$

This reduces the Schrödinger equation to $\mathcal{H}\Psi = E\Psi$. To solve the Schrödinger equation it is necessary to find values of E and functions Ψ such that, when the wavefunction is operated upon by the Hamiltonian, it returns the wavefunction multiplied by the energy. The Schrödinger equation falls into the category of equations known as partial differential eigenvalue equations in which an operator acts on a function (the eigenfunction) and returns the

function multiplied by a scalar (the eigenvalue). A simple example of an eigenvalue equation is:

$$\frac{d}{dx}(y) = ry \quad (2.5)$$

The operator here is d/dx . One eigenfunction of this equation is $y = e^{ax}$ with the eigenvalue r being equal to a . Equation (2.5) is a first-order differential equation. The Schrödinger equation is a second-order differential equation as it involves the second derivative of Ψ . A simple example of an equation of this type is

$$\frac{d^2 y}{dx^2} = ry \quad (2.6)$$

The solutions of Equation (2.6) have the form $y = A \cos kx + B \sin kx$, where A , B and k are constants. In the Schrödinger equation Ψ is the eigenfunction and E the eigenvalue.

2.1.1 Operators

The concept of an operator is an important one in quantum mechanics. The *expectation value* (which we can consider to be the average value) of a quantity such as the energy, position or linear momentum can be determined using an appropriate operator. The most commonly used operator is that for the energy, which is the Hamiltonian operator itself, \mathcal{H} . The energy can be determined by calculating the following integral:

$$E = \frac{\int \Psi^* \mathcal{H} \Psi d\tau}{\int \Psi^* \Psi d\tau} \quad (2.7)$$

The two integrals in Equation (2.7) are performed over all space (i.e. from $-\infty$ to $+\infty$ in the x , y and z directions). Note the use of the complex conjugate notation (Ψ^*), which reminds us that the wavefunction may be a complex number. This equation can be derived by pre-multiplying both sides of the Schrödinger equation, $\mathcal{H}\Psi = E\Psi$, by the complex conjugate of the wavefunction, Ψ^* , and integrating both sides over all space. Thus:

$$\int \Psi^* \mathcal{H} \Psi d\tau = \int \Psi^* E \Psi d\tau \quad (2.8)$$

E is a scalar and so can be taken outside the integral, thus leading to Equation (2.7). If the wavefunction is normalised then the denominator in Equation (2.7) will equal 1.

The Hamiltonian operator is composed of two parts that reflect the contributions of kinetic and potential energies to the total energy. The kinetic energy operator is

$$-\frac{\hbar^2}{2m} \nabla^2 \quad (2.9)$$

and the operator for the potential energy simply involves multiplication by the appropriate expression for the potential energy. For an electron in an isolated atom or molecule the potential energy operator comprises the electrostatic interactions between the electron and the nucleus and the interactions between the electron and the other electrons. For a

single electron and a single nucleus with Z protons the potential energy operator is thus:

$$V = -\frac{Ze^2}{4\pi\epsilon_0 r} \quad (2.10)$$

Another operator is that for linear momentum along the x direction, which is

$$\frac{\hbar}{i} \frac{\partial}{\partial x} \quad (2.11)$$

The expectation value of this quantity can thus be obtained by evaluating the following integral:

$$p_x = \frac{\int \Psi^* \frac{\hbar}{i} \frac{\partial}{\partial x} \Psi d\tau}{\int \Psi^* \Psi d\tau} \quad (2.12)$$

2.1.2 Atomic Units

Quantum mechanics is primarily concerned with atomic particles: electrons, protons and neutrons. When the properties of such particles (e.g. mass, charge, etc.) are expressed in 'macroscopic' units then the value must usually be multiplied or divided by several powers of 10. It is preferable to use a set of units that enables the results of a calculation to be reported as 'easily manageable' values. One way to achieve this would be to multiply each number by an appropriate power of 10. However, further simplification can be achieved by recognising that it is often necessary to carry quantities such as the mass of the electron or electronic charge all the way through a calculation. These quantities are thus also incorporated into the atomic units. The atomic units of length, mass and energy are as follows:

- 1 unit of charge equals the absolute charge on an electron, $|e| = 1.60219 \times 10^{-19} \text{ C}$
- 1 mass unit equals the mass of the electron, $m_e = 9.10593 \times 10^{-31} \text{ kg}$
- 1 unit of length (1 Bohr) is given by $a_0 = \hbar^2/4\pi^2 m_e e^2 = 5.29177 \times 10^{-11} \text{ m}$
- 1 unit of energy (1 Hartree) is given by $E_a = e^2/4\pi\epsilon_0 a_0 = 4.35981 \times 10^{-18} \text{ J}$

The atomic unit of length is the radius of the first orbit in Bohr's treatment of the hydrogen atom. It also turns out to be the most probable distance of a 1s electron from the nucleus in the hydrogen atom. The atomic unit of energy corresponds to the interaction between two electronic charges separated by the Bohr radius. The total energy of the 1s electron in the hydrogen atom equals -0.5 Hartree. In atomic units Planck's constant $\hbar = 2\pi$ and so $\hbar \equiv 1$.

2.1.3 Exact Solutions to the Schrödinger Equation

The Schrödinger equation can be solved exactly for only a few problems, such as the particle in a box, the harmonic oscillator, the particle on a ring, the particle on a sphere and the hydrogen atom, all of which are dealt with in introductory textbooks. A common feature of these problems is that it is necessary to impose certain requirements (often called *boundary*

conditions) on possible solutions to the equation. Thus, for a particle in a box with infinitely high walls, the wavefunction is required to go to zero at the boundaries. For a particle on a ring the wavefunction must have a periodicity of 2π because it must repeat every traversal of the ring. An additional requirement on solutions to the Schrödinger equation is that the wavefunction at a point r , when multiplied by its complex conjugate, is the probability of finding the particle at the point (this is the Born interpretation of the wavefunction). The square of an electronic wavefunction thus gives the electron density at any given point. If we integrate the probability of finding the particle over all space, then the result must be 1 as the particle must be somewhere:

$$\int \Psi^* \Psi d\tau = 1 \quad (2.13)$$

$d\tau$ indicates that the integration is over all space. Wavefunctions which satisfy this condition are said to be *normalised*. It is usual to require the solutions to the Schrödinger equation to be orthogonal:

$$\int \Psi_m^* \Psi_n d\tau = 0 \quad (m \neq n) \quad (2.14)$$

A convenient way to express both the orthogonality of different wavefunctions and the normalisation conditions uses the *Kronecker delta*:

$$\int \Psi_m^* \Psi_n d\tau = \delta_{mn} \quad (2.15)$$

When used in this context, the Kronecker delta can be taken to have a value of 1 if m equals n and zero otherwise. Wavefunctions that are both orthogonal and normalised are said to be *orthonormal*.

2.2 One-electron Atoms

In an atom that contains a single electron, the potential energy depends upon the distance between the electron and the nucleus as given by the Coulomb equation. The Hamiltonian thus takes the following form:

$$\mathcal{H} = -\frac{\hbar^2}{2m} \nabla^2 - \frac{Ze^2}{4\pi\epsilon_0 r} \quad (2.16)$$

In atomic units the Hamiltonian is:

$$\mathcal{H} = -\frac{1}{2} \nabla^2 - \frac{Z}{r} \quad (2.17)$$

For the hydrogen atom, the nuclear charge, Z , equals +1. r is the distance of the electron from the nucleus. The helium cation, He^+ , is also a one-electron atom but has a nuclear charge of +2. As atoms have spherical symmetry it is more convenient to transform the Schrödinger equation to polar coordinates r , θ and ϕ , where r is the distance from the nucleus (located at the origin), θ is the angle to the z axis and ϕ is the angle from the x axis in the xy plane (Figure 2.1). The solutions can be written as the product of a radial function $R(r)$, which depends only on r , and an angular function $Y(\theta, \phi)$ called a *spherical harmonic*, which

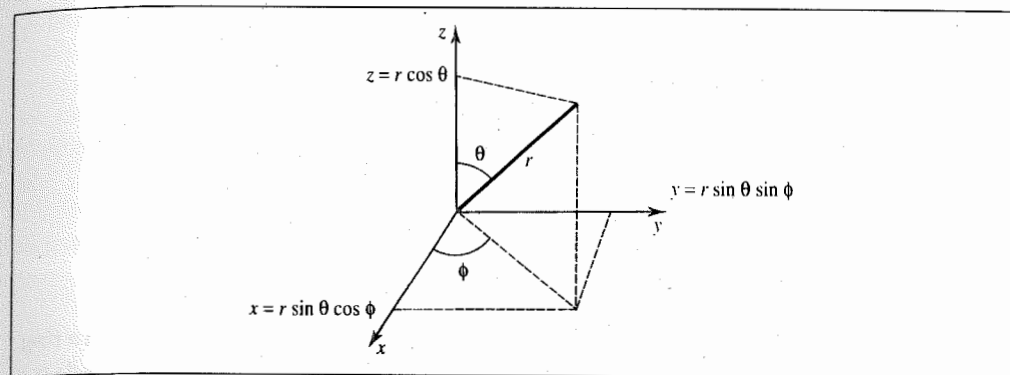


Fig. 2.1: The relationship between spherical polar and Cartesian coordinates.

depends on θ and ϕ :

$$\Psi_{nlm} = R_{nl}(r) Y_{lm}(\theta, \phi) \quad (2.18)$$

The wavefunctions are commonly referred to as *orbitals* and are characterised by three quantum numbers n , m and l . The quantum numbers can adopt values as follows:

n : principal quantum number: 0, 1, 2, ...

l : azimuthal quantum number: 0, 1, ... ($n-1$)

m : magnetic quantum number: $-l, -(l-1), \dots, 0, \dots, (l-1), l$.

The full radial function is:

$$R_{nl}(r) = -\left[\left(\frac{2Z}{na_0} \right)^3 \frac{(n-l-1)!}{2n[(n+l)!]^3} \right]^{1/2} \exp\left(-\frac{\rho}{2}\right) \rho^l L_{n+l}^{2l+1}(\rho) \quad (2.19)$$

$\rho = 2Zr/na_0$, where a_0 is the Bohr radius.* The term in square brackets is a normalising factor. $L_{n+l}^{2l+1}(\rho)$ is a special type of function called a Laguerre polynomial. We shall rarely be interested in any other than the first few members of the series; moreover, they simplify considerably if atomic units are used and we write them in terms of the *orbital exponent* $\zeta = Z/n$. The first few members of the series for low values of n are given in Table 2.1 and are illustrated graphically in Figure 2.2. As can be seen, the radial part of the wavefunction is a polynomial multiplied by a decaying exponential.

The angular part of the wavefunction is the product of a function of θ and a function of ϕ :

$$Y_{lm}(\theta, \phi) = \Theta_{lm}(\theta) \Phi_m(\phi) \quad (2.20)$$

These functions are:

$$\Phi_m(\phi) = \frac{1}{\sqrt{2\pi}} \exp(im\phi) \quad (2.21)$$

$$\Theta_{lm}(\theta) = \left[\frac{(2l+1)(l-|m|)!}{2(l+|m|)!} \right]^{1/2} P_l^{|m|}(\cos\theta) \quad (2.22)$$

* Strictly, a_0 in this case is given by $a_0 = \hbar^2/\pi^2\mu e$, where μ is the reduced mass, $\mu = m_e M/(m_e + M)$; M is the mass of the nucleus.

n	l	$R_{nl}(r)$
1	0	$2\zeta^{3/2} \exp(-\zeta r)$
2	0	$2\zeta^{3/2}(1 - \zeta r) \exp(-\zeta r)$
2	1	$(4/3)^{1/2} \zeta^{5/2} r \exp(-\zeta r)$
3	0	$(2/3)^{1/2} \zeta^{3/2} (3 - 6\zeta r + 2\zeta^2 r^2) \exp(-\zeta r)$
3	1	$(8/9)^{1/2} \zeta^{5/2} (2 - \zeta r) r \exp(-\zeta r)$
3	2	$(8/45)^{1/2} \zeta^{7/2} r^2 \exp(-\zeta r)$

Table 2.1 Radial function for one-electron atoms.

The functions $\Phi_m(\phi)$ are just the solutions to the Schrödinger equation for a particle on a ring. The term in square brackets for the function $\Theta_{lm}(\theta)$ is a normalising factor. $P_l^{|m|}(\cos \theta)$ is a member of a series of functions called the associated Legendre polynomials (the 'Legendre polynomials' are functions for which $|m| = 0$). The total orbital angular momentum of an electron in the orbital is given by $l(l+1)\hbar$ and the component of the angular momentum along the $\theta = 0$ axis is given by $l\hbar$. The energy of each solution is a function of the principal quantum number only; thus orbitals with the same value of n but different l and m are degenerate. The orbitals are often represented as shown in Figure 2.3. These graphical representations are not necessarily the same as the solutions given above. For example, the 'correct' solutions for the 2p orbitals comprise one real and two complex functions:

$$2p(+1) = \sqrt{3/4\pi} R(r) \sin \theta e^{i\phi} \quad (2.23)$$

$$2p(0) = \sqrt{3/4\pi} R(r) \cos \theta \quad (2.24)$$

$$2p(-1) = \sqrt{3/4\pi} R(r) \sin \theta e^{-i\phi} \quad (2.25)$$

$R(r)$ is the radial part of the wavefunction and $\sqrt{3/4\pi}$ is a normalisation factor for the angular part. The $2p(0)$ function is real and corresponds to the $2p_z$ orbital that is pictured in Figure 2.3. A linear combination of the two remaining 2p solutions is used to generate two 'real' 2p wavefunctions, making use of the relationship $\exp(i\phi) = \cos \phi + i \sin \phi$ (Section 1.10.4). These linear combinations are the $2p_x$ and $2p_y$ orbitals shown in Figure 2.3.

$$2p_x = 1/2[2p(+1) + 2p(-1)] = \sqrt{3/4\pi} R(r) \sin \theta \cos \phi \quad (2.26)$$

$$2p_y = -1/2[2p(+1) - 2p(-1)] = \sqrt{3/4\pi} R(r) \sin \theta \sin \phi \quad (2.27)$$

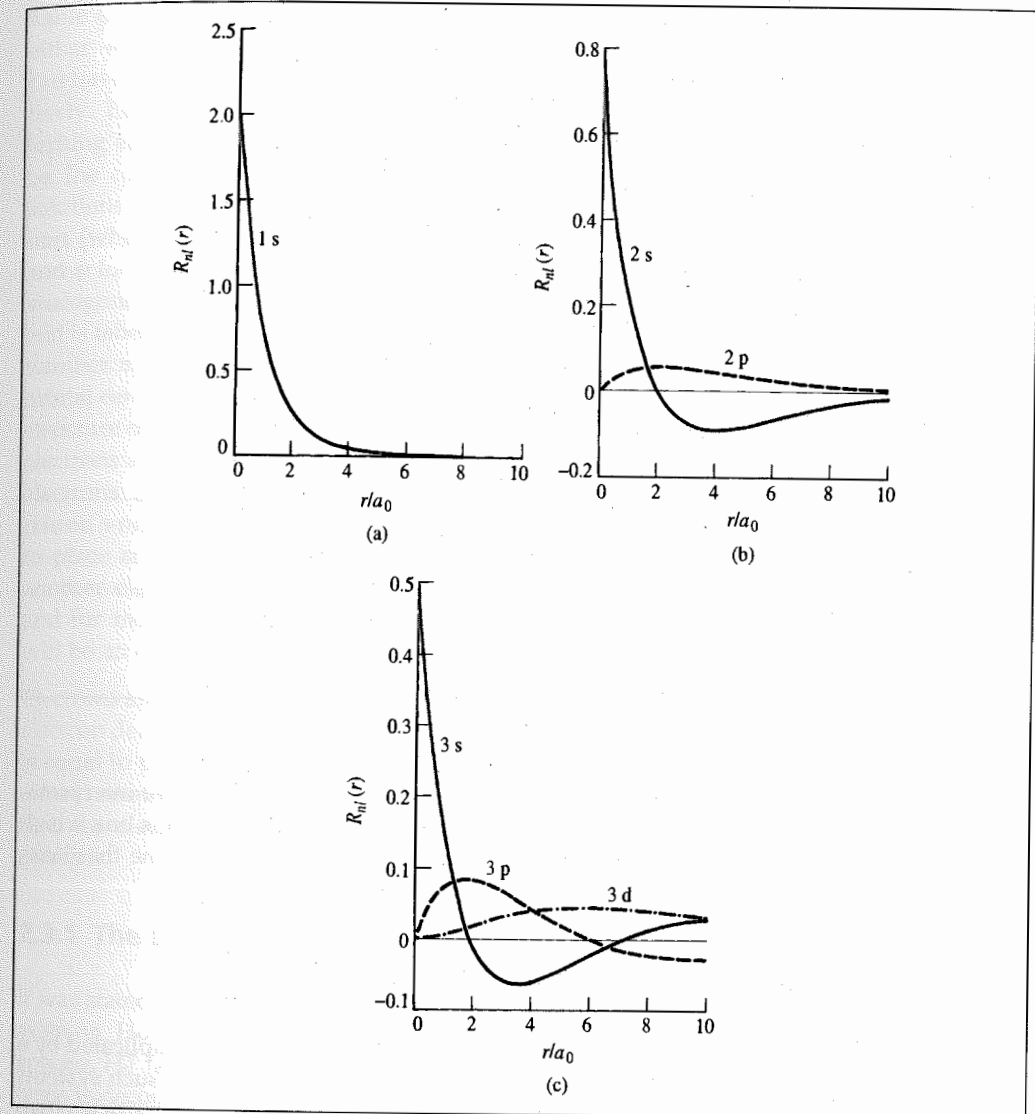
These linear combinations still have the same energy as the original complex wavefunctions. This is a general property of degenerate solutions of the Hamiltonian operator. The reason why they are labelled $2p_x$ and $2p_y$ is that in polar coordinates the Cartesian coordinates x , y and z have the same angular dependence as the orbitals in Figure 2.3:

$$x = r \sin \theta \cos \phi \quad (2.28)$$

$$y = r \sin \theta \sin \phi \quad (2.29)$$

$$z = r \cos \theta \quad (2.30)$$

The solutions of the Schrödinger equation are either real or occur in degenerate pairs. These pairs are complex conjugates that can then be combined to give energetically equivalent real solutions. It is only when dealing with certain types of operator that it is necessary to retain a

Fig. 2.2: The functions $R_{nl}(r)$ for the first three values of the principal quantum number. (a) 1s; (b) 2s and 2p; (c) 3s, 3p and 3d.

complex wavefunction (for the 2p functions, the operator that corresponds to angular momentum about the z axis falls into this category). In fact, to simplify matters we will almost always ignore the complex notation from now on and will deal with real orbitals.

Finally, we should note that the solutions are all orthogonal to each other; if the product of any pair of orbitals is integrated over all space, the result is zero unless the two orbitals are the same. Orthonormality is achieved by multiplying by an appropriate normalisation constant.

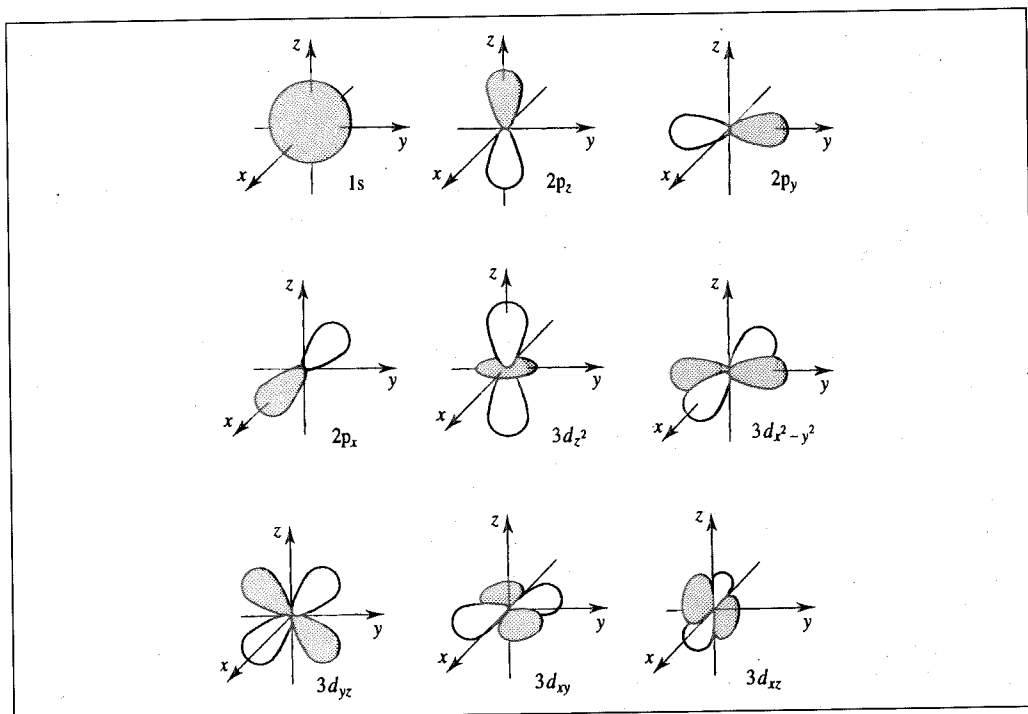


Fig. 2.3: The common graphical representations of s, p and d orbitals.

The orbital picture has proved invaluable for providing insight and qualitative interpretations into the nature of the bonding in and reactivity of chemical systems. It is one which we would like to retain for polyelectronic systems to provide a unifying theme that links the simplest systems with much more complicated ones.

2.3 Polyelectronic Atoms and Molecules

Solving the Schrödinger equation for atoms with more than one electron is complicated by a number of factors. The first complication is that the Schrödinger equation for such systems cannot be solved exactly, even for the helium atom. The helium atom has three particles (two electrons and one nucleus) and is an example of a *three-body problem*. No exact solutions can be found for systems that involve three (or more) interacting particles. Thus, any solutions we might find for polyelectronic atoms or molecules can only be approximations to the real, true solutions of the Schrödinger equation. One consequence of there being no exact solution is that the wavefunction may adopt more than one functional form; no form is necessarily more 'correct' than another. In fact, the most general form of the wavefunction will be an infinite series of functions.

A second complication with multi-electron species is that we must account for electron spin. Spin is characterised by the quantum number s , which for an electron can only take the

value $\frac{1}{2}$. The spin angular momentum is quantised such that its projection on the z axis is either $+\hbar$ or $-\hbar$. These two states are characterised by the quantum number m_s , which can have values of $+\frac{1}{2}$ or $-\frac{1}{2}$, and are often referred to as 'up spin' and 'down spin' respectively. Electron spin is incorporated into the solutions to the Schrödinger equation by writing each one-electron wavefunction as the product of a spatial function that depends on the coordinates of the electron and a spin function that depends on its spin. Such solutions are called *spin orbitals*, which we will represent using the symbol χ . The spatial part (which will be referred to as an orbital and represented using ϕ for atomic orbitals and ψ for molecular orbitals) describes the distribution of electron density in space and is analogous to the orbital diagrams in Figure 2.3. The spin part defines the electron spin and is labelled α or β . These spin functions have the value 0 or 1 depending on the quantum number m_s of the electron. Thus $\alpha(\frac{1}{2}) = 1$, $\alpha(-\frac{1}{2}) = 0$, $\beta(+\frac{1}{2}) = 0$, $\beta(-\frac{1}{2}) = 1$. Each spatial orbital can accommodate two electrons, with paired spins. In order to predict the electronic structure of a polyelectronic atom or a molecule, the *Aufbau principle* is employed, in which electrons are assigned to the orbitals, two electrons per orbital. We need to remember that electrons occupy degenerate states with a maximum number of unpaired electrons (Hund's rules), and that there are certain situations where it is energetically more favourable to place an unpaired electron in a higher-energy spatial orbital rather than pair it with another electron. However, such situations are rare, particularly for molecular systems, and for most of the situations that we shall be interested in the number of electrons, N , will be an even number that occupy the $N/2$ lowest-energy orbitals.

Electrons are indistinguishable. If we exchange any pair of electrons, then the distribution of electron density remains the same. According to the Born interpretation, the electron density is equal to the square of the wavefunction. It therefore follows that the wavefunction must either remain unchanged when two electrons are exchanged, or else it must change sign. In fact, for electrons the wavefunction is required to change sign: this is the *antisymmetry principle*.

2.3.1 The Born–Oppenheimer Approximation

It was stated above that the Schrödinger equation cannot be solved exactly for any molecular systems. However, it is possible to solve the equation exactly for the simplest molecular species, H_2^+ (and isotopically equivalent species such as HD^+), when the motion of the electrons is decoupled from the motion of the nuclei in accordance with the Born–Oppenheimer approximation. The masses of the nuclei are much greater than the masses of the electrons (the resting mass of the lightest nucleus, the proton, is 1836 times heavier than the resting mass of the electron). This means that the electrons can adjust almost instantaneously to any changes in the positions of the nuclei. The electronic wavefunction thus depends only on the positions of the nuclei and not on their momenta. Under the Born–Oppenheimer approximation the total wavefunction for the molecule can be written in the following form:

$$\Psi_{\text{tot}}(\text{nuclei, electrons}) = \Psi(\text{electrons})\Psi(\text{nuclei}) \quad (2.31)$$

The total energy equals the sum of the nuclear energy (the electrostatic repulsion between the positively charged nuclei) and the electronic energy. The electronic energy comprises

the kinetic and potential energy of the electrons moving in the electrostatic field of the nuclei, together with electron–electron repulsion: $E_{\text{tot}} = E(\text{electrons}) + E(\text{nuclei})$.

When the Born–Oppenheimer approximation is used we concentrate on the electronic motions; the nuclei are considered to be fixed. For each arrangement of the nuclei the Schrödinger equation is solved for the electrons alone in the field of the nuclei. If it is desired to change the nuclear positions then it is necessary to add the nuclear repulsion to the electronic energy in order to calculate the total energy of the configuration.

2.3.2 The Helium Atom

We now return to the helium atom, our objective being to find a wavefunction that describes the behaviour of the electrons. The Born–Oppenheimer approximation is not, of course, relevant to systems with just one nucleus, and the wavefunction will be a function of the two electrons (which we shall label 1 and 2 with positions in space \mathbf{r}_1 and \mathbf{r}_2). As noted above, for polyelectronic systems any solution we find can only ever be an approximation to the true solution. There are a number of ways in which approximate solutions to the Schrödinger equation can be found. One approach is to find a simpler but related problem that can be more easily solved and then consider how the differences between the two problems change the Hamiltonian and thereby affect the solutions. This is called *perturbation theory* and is most appropriate when the differences between the real and simple problems are small. For example, a perturbation approach to tackling the helium atom might choose as the related system a ‘pseudo atom’, containing two electrons that interact with the nucleus but not with each other. Although this is a ‘three-body’ problem, the lack of any interaction between the electrons means that it can be solved exactly using the method of the separation of variables. The separation of variables technique can be applied whenever the Hamiltonian can be divided into parts that are themselves dependent solely upon subsets of the coordinates. The equation to be solved in this case is:

$$\left\{ -\frac{\hbar^2}{2m} \nabla_1^2 - \frac{Ze^2}{4\pi\epsilon_0 r_1} - \frac{\hbar^2}{2m} \nabla_2^2 - \frac{Ze^2}{4\pi\epsilon_0 r_2} \right\} \Psi(\mathbf{r}_1, \mathbf{r}_2) = E\Psi(\mathbf{r}_1, \mathbf{r}_2) \quad (2.32)$$

Or, in atomic units,

$$\left\{ -\frac{1}{2} \nabla_1^2 - \frac{Z}{r_1} - \frac{1}{2} \nabla_2^2 - \frac{Z}{r_2} \right\} \Psi(\mathbf{r}_1, \mathbf{r}_2) = E\Psi(\mathbf{r}_1, \mathbf{r}_2) \quad (2.33)$$

We can abbreviate this equation to

$$\{\mathcal{H}_1 + \mathcal{H}_2\} \Psi(\mathbf{r}_1, \mathbf{r}_2) = E\Psi(\mathbf{r}_1, \mathbf{r}_2) \quad (2.34)$$

\mathcal{H}_1 and \mathcal{H}_2 are the individual Hamiltonians for electrons 1 and 2. Let us assume that the wavefunction can be written as a product of individual one-electron wavefunctions, $\phi_1(\mathbf{r}_1)$ and $\phi_2(\mathbf{r}_2)$: $\Psi(\mathbf{r}_1, \mathbf{r}_2) = \phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2)$. Then we can write:

$$\{\mathcal{H}_1 + \mathcal{H}_2\} \phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2) = E\phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2) \quad (2.35)$$

Premultiplying by $\phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2)$ and integrating over all space gives:

$$\iint d\tau_1 d\tau_2 \phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2) [\mathcal{H}_1 + \mathcal{H}_2] \phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2) = \iint d\tau_1 d\tau_2 \phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2) \phi_1(\mathbf{r}_1)\phi_2(\mathbf{r}_2) \quad (2.36)$$

or

$$\begin{aligned} & \int d\tau_1 \phi_1(\mathbf{r}_1) \mathcal{H}_1 \phi_1(\mathbf{r}_1) \int d\tau_2 \phi_2(\mathbf{r}_2) \phi_2(\mathbf{r}_2) + \int d\tau_1 \phi_1(\mathbf{r}_1) \phi_1(\mathbf{r}_1) \int d\tau_2 \phi_2(\mathbf{r}_2) \mathcal{H}_2 \phi_2(\mathbf{r}_2) \\ & = E \int d\tau_1 \phi_1(\mathbf{r}_1) \phi_1(\mathbf{r}_1) \int d\tau_2 \phi_2(\mathbf{r}_2) \phi_2(\mathbf{r}_2) \end{aligned} \quad (2.37)$$

If we assume that the wavefunctions are normalised then it can easily be seen that the total energy E is the sum of the individual orbital energies E_1 and E_2 ($E_1 = \int d\tau_1 \phi_1(\mathbf{r}_1) \mathcal{H}_1 \phi_1(\mathbf{r}_1)$ and $E_2 = \int d\tau_2 \phi_2(\mathbf{r}_2) \mathcal{H}_2 \phi_2(\mathbf{r}_2)$). When the separation of variables method is used the solutions for each electron are just those of the hydrogen atom (1s, 2s, etc.) in Equation (2.19) with $Z = 2$.

We now wish to establish the general functional form of possible wavefunctions for the two electrons in this pseudo helium atom. We will do so by considering first the spatial part of the wavefunction. We will show how to derive functional forms for the wavefunction in which the exchange of electrons is independent of the electron labels and does not affect the electron density. The simplest approach is to assume that each wavefunction for the helium atom is the product of the individual one-electron solutions. As we have just seen, this implies that the total energy is equal to the sum of the one-electron orbital energies, which is not correct as it ignores electron–electron repulsion. Nevertheless, it is a useful illustrative model. The wavefunction of the lowest energy state then has each of the two electrons in a 1s orbital:

$$1s(1)1s(2) \quad (2.38)$$

‘1s(1)’ indicates a 1s function that depends on the coordinates of electron 1 (\mathbf{r}_1) and ‘1s(2)’ indicates a 1s function that depends upon the coordinates of electron 2 (\mathbf{r}_2). This wavefunction satisfies the indistinguishability criterion, for we obtain the same function when we exchange the electrons – 1s(1)1s(2) is the same as 1s(2)1s(1). Its energy is twice that of a single electron in a 1s orbital. What of the first excited state, in which one electron is promoted to the 2s orbital? Two possible wavefunctions for this state are:

$$1s(1)2s(2) \quad (2.39)$$

$$1s(2)2s(1) \quad (2.40)$$

Do these wavefunctions satisfy the indistinguishability criterion? In other words, do we get the same function (or its negative) when we exchange the electrons? We do not, for when the two electrons (1 and 2) are exchanged then a different wavefunction is obtained: ‘1s(1)2s(2)’ and ‘1s(2)2s(1)’ are not the same, nor is one simply minus the other. However, linear combinations of these two wavefunctions do not suffer from the labelling problem and so we might anticipate that functional forms such as the following might constitute acceptable solutions to the Schrödinger equation for the pseudo helium atom:

$$(1/\sqrt{2})[1s(1)2s(2) + 1s(2)2s(1)] \quad (2.41)$$

$$(1/\sqrt{2})[1s(1)2s(2) - 1s(2)2s(1)] \quad (2.42)$$

The factor $(1/\sqrt{2})$ ensures that the wavefunction is normalised. Of the three acceptable spatial forms that we have described so far, two are symmetric (i.e. do not change sign when the electron labels are exchanged) and one is antisymmetric (the sign changes when the electrons are exchanged):

$$1s(1)1s(2) \quad \text{symmetric} \quad (2.43)$$

$$(1/\sqrt{2})[1s(1)2s(2) + 1s(2)2s(1)] \quad \text{symmetric} \quad (2.44)$$

$$(1/\sqrt{2})[1s(1)2s(2) - 1s(2)2s(1)] \quad \text{antisymmetric} \quad (2.45)$$

We now need to consider the effects of electron spin. For two electrons 1 and 2 there are four spin states; $\alpha(1)$, $\beta(1)$, $\alpha(2)$, $\beta(2)$. The indistinguishability criterion holds for the spin components as well, and so the following combinations of spin wavefunctions are possible:

$$\alpha(1)\alpha(2) \quad \text{symmetric} \quad (2.46)$$

$$\beta(1)\beta(2) \quad \text{symmetric} \quad (2.47)$$

$$(1/\sqrt{2})[\alpha(1)\beta(2) + \alpha(2)\beta(1)] \quad \text{symmetric} \quad (2.48)$$

$$(1/\sqrt{2})[\alpha(1)\beta(2) - \alpha(2)\beta(1)] \quad \text{antisymmetric} \quad (2.49)$$

When we combine the spatial and spin wavefunctions, the overall wavefunction must be antisymmetric with respect to exchange of electrons. It is therefore only admissible to combine a symmetric spatial part with an antisymmetric spin part, or an antisymmetric spatial part with a symmetric spin part. The following functional forms are therefore permissible for the wavefunctions of the ground and first few excited states of the helium atom:

$$(1/\sqrt{2})1s(1)1s(2)[\alpha(1)\beta(2) - \alpha(2)\beta(1)] \quad (2.50)$$

$$(1/2)[1s(1)2s(2) + 1s(2)2s(1)][\alpha(1)\beta(2) - \alpha(2)\beta(1)] \quad (2.51)$$

$$(1/\sqrt{2})[1s(1)2s(2) - 1s(2)2s(1)]\alpha(1)\alpha(2) \quad (2.52)$$

$$(1/\sqrt{2})[1s(1)2s(2) - 1s(2)2s(1)]\beta(1)\beta(2) \quad (2.53)$$

$$(1/2)[1s(1)2s(2) - 1s(2)2s(1)][\alpha(1)\beta(2) + \alpha(2)\beta(1)] \quad (2.54)$$

2.3.3 General Polyelectronic Systems and Slater Determinants

We now turn to the general case. What is an appropriate functional form of the wavefunction for a polyelectronic system (not necessarily an atom) with N electrons that satisfies the antisymmetry principle? First, we note that the following functional form of the wavefunction is inappropriate:

$$\Psi(1, 2, \dots, N) = \chi_1(1)\chi_2(2) \dots \chi_N(N) \quad (2.55)$$

This product of spin orbitals is unacceptable because it does not satisfy the antisymmetry principle; exchanging pairs of electrons does not give the negative of the wavefunction. This formulation of the wavefunction is known as a *Hartree product*. The energy of a system described by a Hartree product equals the sum of the one-electron spin orbitals. A key conclusion of the Hartree product description is that the probability of finding an electron at a particular point in space is independent of the probability of finding any

other electron at that point in space. In fact, it turns out that the motions of the electrons are correlated. In addition, the Hartree product assumes that specific electrons have been assigned to specific orbitals, whereas the antisymmetry principle requires that the electrons are indistinguishable. Recall that for the helium atom, an acceptable functional form for the lowest-energy state, is:

$$\begin{aligned} \psi &= 1s(1)1s(2)[\alpha(1)\beta(2) - \alpha(2)\beta(1)] \\ &\equiv 1s(1)1s(2)\alpha(1)\beta(2) - 1s(1)1s(2)\alpha(2)\beta(1) \end{aligned} \quad (2.56)$$

This can be written in the form of a 2×2 determinant:

$$\begin{vmatrix} 1s(1)\alpha(1) & 1s(1)\beta(1) \\ 1s(2)\alpha(2) & 1s(2)\beta(2) \end{vmatrix} \quad (2.57)$$

The two spin orbitals are

$$\chi_1 = 1s(1)\alpha(1) \quad \text{and} \quad \chi_2 = 1s(1)\beta(1) \quad (2.58)$$

A determinant is the most convenient way to write down the permitted functional forms of a polyelectronic wavefunction that satisfies the antisymmetry principle. In general, if we have N electrons in spin orbitals $\chi_1, \chi_2, \dots, \chi_N$ (where each spin orbital is the product of a spatial function and a spin function) then an acceptable form of the wavefunction is:

$$\Psi = \frac{1}{\sqrt{N!}} \begin{vmatrix} \chi_1(1) & \chi_2(1) & \dots & \chi_N(1) \\ \chi_1(2) & \chi_2(2) & \dots & \chi_N(2) \\ \vdots & \vdots & & \vdots \\ \chi_1(N) & \chi_2(N) & \dots & \chi_N(N) \end{vmatrix} \quad (2.59)$$

As before, $\chi_1(1)$ is used to indicate a function that depends on the space and spin coordinates of the electron labelled '1'. The factor $1/\sqrt{N!}$ ensures that the wavefunction is normalised; we shall see later why the normalisation factor has this particular value. This functional form of the wavefunction is called a *Slater determinant* and is the simplest form of an orbital wavefunction that satisfies the antisymmetry principle. The Slater determinant is a particularly convenient and concise way to represent the wavefunction due to the special properties of determinants. Exchanging any two rows of a determinant, a process which corresponds to exchanging two electrons, changes the sign of the determinant and therefore directly leads to the antisymmetry property. If any two rows of a determinant are identical, which would correspond to two electrons being assigned to the same spin orbital, then the determinant vanishes. This can be considered a manifestation of the Pauli principle, which states that no two electrons can have the same set of quantum numbers. The Pauli principle also leads to the notion that each spatial orbital can accommodate two electrons of opposite spins.

When the Slater determinant is expanded, a total of $N!$ terms results. This is because there are $N!$ different permutations of N electrons. For example, for a three-electron system with spin orbitals χ_1, χ_2 and χ_3 the determinant is

$$\Psi = \frac{1}{\sqrt{12}} \begin{vmatrix} \chi_1(1) & \chi_2(1) & \chi_3(1) \\ \chi_1(2) & \chi_2(2) & \chi_3(2) \\ \chi_1(3) & \chi_2(3) & \chi_3(3) \end{vmatrix} \quad (2.60)$$

Expansion of the determinant gives the following expression (ignoring the normalisation constant):

$$\begin{aligned} & \chi_1(1)\chi_2(2)\chi_3(3) - \chi_1(1)\chi_3(2)\chi_2(3) + \chi_2(1)\chi_3(2)\chi_1(3) \\ & - \chi_2(1)\chi_1(2)\chi_3(3) + \chi_3(1)\chi_1(2)\chi_2(3) - \chi_3(1)\chi_2(2)\chi_1(3) \end{aligned} \quad (2.61)$$

This expansion contains six terms ($\equiv 3!$). The six possible permutations of three electrons are: 123, 132, 213, 231, 312, 321. Some of these permutations involve single exchanges of electrons; others involve the exchange of two electrons. For example, the permutation 132 can be generated from the initial permutation by exchanging electrons 2 and 3. If we do so then the following wavefunction is obtained:

$$\begin{aligned} & \chi_1(1)\chi_2(3)\chi_3(2) - \chi_1(1)\chi_3(3)\chi_2(2) + \chi_2(1)\chi_3(3)\chi_1(2) \\ & - \chi_2(1)\chi_1(3)\chi_3(2) + \chi_3(1)\chi_1(3)\chi_2(2) - \chi_3(1)\chi_2(3)\chi_1(2) \\ & = -\chi_1(1)\chi_2(2)\chi_3(3) + \chi_1(1)\chi_3(2)\chi_2(3) - \chi_2(1)\chi_3(2)\chi_1(3) \\ & + \chi_2(1)\chi_1(2)\chi_3(3) - \chi_3(1)\chi_1(2)\chi_2(3) + \chi_3(1)\chi_2(2)\chi_1(3) \\ & = -\Psi \end{aligned} \quad (2.62)$$

By contrast, the permutation 312 requires that electrons 1 and 3 are exchanged and then electrons 1 and 2 are exchanged. This gives rise to an unchanged wavefunction. In general, an odd permutation involves an odd number of electron exchanges and leads to a wavefunction with a changed sign; an even permutation involves an even number of electron exchanges and returns the wavefunction unchanged.

For any sizeable system the Slater determinant can be tedious to write out, let alone the equivalent full orbital expansion, and so it is common to use a shorthand notation. Various notation systems have been devised. In one system the terms along the diagonal of the matrix are written as a single-row determinant. For the 3×3 determinant we therefore have:

$$\begin{vmatrix} \chi_1(1) & \chi_2(1) & \chi_3(1) \\ \chi_1(2) & \chi_2(2) & \chi_3(2) \\ \chi_1(3) & \chi_2(3) & \chi_3(3) \end{vmatrix} \equiv |\chi_1 \ \chi_2 \ \chi_3| \quad (2.63)$$

The normalisation factor is assumed. It is often convenient to indicate the spin of each electron in the determinant; this is done by writing a bar when the spin part is β (spin down); a function without a bar indicates an α spin (spin up). Thus, the following are all commonly used ways to write the Slater determinantal wavefunction for the beryllium atom (which has the electronic configuration $1s^2 2s^2$):

$$\begin{aligned} \Psi &= \frac{1}{\sqrt{24}} \begin{vmatrix} \phi_{1s}(1) & \bar{\phi}_{1s}(1) & \phi_{2s}(1) & \bar{\phi}_{2s}(1) \\ \phi_{1s}(2) & \bar{\phi}_{1s}(2) & \phi_{2s}(2) & \bar{\phi}_{2s}(2) \\ \phi_{1s}(3) & \bar{\phi}_{1s}(3) & \phi_{2s}(3) & \bar{\phi}_{2s}(3) \\ \phi_{1s}(4) & \bar{\phi}_{1s}(4) & \phi_{2s}(4) & \bar{\phi}_{2s}(4) \end{vmatrix} \\ &\equiv |\phi_{1s} \ \bar{\phi}_{1s} \ \phi_{2s} \ \bar{\phi}_{2s}| \\ &\equiv |1s \ \bar{1}s \ 2s \ \bar{2}s| \end{aligned} \quad (2.64)$$

An important property of determinants is that a multiple of any column can be added to another column without altering the value of the determinant. This means that the spin orbitals are not unique; other linear combinations give the same energy. To illustrate this, consider the first excited state configuration of the helium atom ($1s^2 2s^2$), which can be written as the following 2×2 determinant:

$$\begin{vmatrix} 1s(1)\alpha(1) & 2s(1)\alpha(1) \\ 1s(2)\alpha(2) & 2s(2)\alpha(2) \end{vmatrix} = 1s(1)\alpha(1)2s(2)\alpha(2) - 1s(2)\alpha(2)2s(1)\alpha(1) \quad (2.65)$$

We now introduce two new 'spin orbitals':

$$\chi'_1 = \frac{1s + 2s}{\sqrt{2}}\alpha; \quad \chi'_2 = \frac{1s - 2s}{\sqrt{2}}\alpha \quad (2.66)$$

With these new orbitals the value of the determinant is as follows:

$$\begin{aligned} \begin{vmatrix} \chi'_1(1) & \chi'_2(1) \\ \chi'_1(2) & \chi'_2(2) \end{vmatrix} &= \frac{[1s(1) + 2s(1)][1s(2) - 2s(2)]\alpha(1)\alpha(2)}{2} \\ &\quad - \frac{[1s(1) - 2s(1)][1s(2) + 2s(2)]\alpha(1)\alpha(2)}{2} \\ &\equiv -\Psi \end{aligned} \quad (2.67)$$

This can be helpful because it may enable more meaningful sets of orbitals to be generated from the original solutions. Molecular orbital calculations may give solutions that are 'smeared out' throughout the entire molecule, whereas we may find orbitals that are localised in specific regions (e.g. in the bonds between atoms) to be more useful.

2.4 Molecular Orbital Calculations

2.4.1 Calculating the Energy from the Wavefunction: the Hydrogen Molecule

In our treatment of molecular systems we first show how to determine the energy for a given wavefunction, and then demonstrate how to calculate the wavefunction for a specific nuclear geometry. In the most popular kind of quantum mechanical calculations performed on molecules each molecular spin orbital is expressed as a linear combination of atomic orbitals (the LCAO approach*). Thus each molecular orbital can be written as a summation of the following form:

$$\psi_i = \sum_{\mu=1}^K c_{\mu i} \phi_{\mu} \quad (2.68)$$

ψ_i is a (spatial) molecular orbital, ϕ_{μ} is one of K atomic orbitals and $c_{\mu i}$ is a coefficient. In a simple LCAO picture of the lowest energy state of molecular hydrogen, H_2 , there are two electrons with opposite spins in the lowest energy spatial orbital (labelled $1\sigma_g$), which is

* Computational quantum chemistry is well endowed with acronyms and abbreviations. A list of some of the more common ones can be found in Appendix 2.1.

formed from a linear combination of two hydrogen-atom 1s orbitals:

$$1\sigma_g = A(1s_A + 1s_B) \quad (2.69)$$

A is the normalisation factor, whose value is not important in our present discussion. To calculate the energy of the ground state of the hydrogen molecule for a fixed internuclear distance we first write the wavefunction as a 2×2 determinant:

$$\Psi = \begin{vmatrix} \chi_1(1) & \chi_2(1) \\ \chi_1(2) & \chi_2(2) \end{vmatrix} = \chi_1(1)\chi_2(2) - \chi_1(2)\chi_2(1) \quad (2.70)$$

where

$$\begin{aligned} \chi_1(1) &= 1\sigma_g(1)\alpha(1) \\ \chi_2(1) &= 1\sigma_g(1)\beta(1) \\ \chi_1(2) &= 1\sigma_g(2)\alpha(2) \\ \chi_2(2) &= 1\sigma_g(2)\beta(2) \end{aligned} \quad (2.71)$$

For the hydrogen molecule, the Hamiltonian comprises the kinetic energy operator for each electron plus the potential energy operator due to the Coulomb attraction between the two electrons and the two nuclei, and the repulsion between the two electrons. In atomic units the Hamiltonian is thus

$$\mathcal{H} = -\frac{1}{2}\nabla_1^2 - \frac{1}{2}\nabla_2^2 - \frac{Z_A}{r_{1A}} - \frac{Z_B}{r_{1B}} - \frac{Z_A}{r_{2A}} - \frac{Z_B}{r_{2B}} + \frac{1}{r_{12}} \quad (2.72)$$

The electrons have been labelled 1 and 2 and the nuclei have been labelled A and B. For H_2 the nuclear charges Z_A and Z_B are both equal to 1. First we need to consider how to calculate the energy of this hydrogen molecule. This is obtained using Equation (2.7):

$$E = \frac{\int \Psi \mathcal{H} \Psi \, d\tau}{\int \Psi \Psi \, d\tau} \quad (2.73)$$

In general, a quantum mechanical calculation provides molecular orbitals that are normalised but the total wavefunction is not. The normalisation constant for the wavefunction of the two-electron hydrogen molecule is $1/\sqrt{2}$ and so the denominator in Equation (2.73) is equal to 2.

We now substitute the hydrogen molecule wavefunction into Equation (2.73) to provide the following:

$$E = \frac{1}{2} \iint d\tau_1 d\tau_2 \{ [\chi_1(1)\chi_2(2) - \chi_2(1)\chi_1(2)] [-\frac{1}{2}\nabla_1^2 - \frac{1}{2}\nabla_2^2 - (1/r_{1A}) - (1/r_{1B}) - (1/r_{2A}) - (1/r_{2B}) + (1/r_{12})] [\chi_1(1)\chi_2(2) - \chi_2(1)\chi_1(2)] \} \quad (2.74)$$

$d\tau_i$ indicates that the integration is over the spatial and spin coordinates of electron i . It is useful to separate the Hamiltonian operator into two H_2^+ Hamiltonians plus the inter-electronic repulsion term:

$$E = \frac{1}{2} \iint d\tau_1 d\tau_2 \{ [\chi_1(1)\chi_2(2) - \chi_2(1)\chi_1(2)] [\mathcal{H}_1 + \mathcal{H}_2 + (1/r_{12})] \times [\chi_1(1)\chi_2(2) - \chi_2(1)\chi_1(2)] \} \quad (2.75)$$

where

$$\mathcal{H}_1 = -\frac{1}{2}\nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \quad \text{and} \quad \mathcal{H}_2 = -\frac{1}{2}\nabla_2^2 - \frac{1}{r_{2A}} - \frac{1}{r_{2B}} \quad (2.76)$$

We can now start to separate the integral in Equation (2.74) into individual terms and identify the various contributions to the electronic energy:

$$\begin{aligned} E &= \iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) (\mathcal{H}_1) \chi_1(1)\chi_2(2) \\ &\quad - \iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) (\mathcal{H}_1) \chi_2(1)\chi_1(2) + \dots \\ &\quad + \iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) (\mathcal{H}_2) \chi_1(1)\chi_2(2) \\ &\quad - \iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) (\mathcal{H}_2) \chi_2(1)\chi_1(2) + \dots \\ &\quad + \iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) \left(\frac{1}{r_{12}} \right) \chi_1(1)\chi_2(2) \\ &\quad - \iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) \left(\frac{1}{r_{12}} \right) \chi_2(1)\chi_1(2) + \dots \end{aligned} \quad (2.77)$$

Each of these individual terms can be simplified if we recognise that terms dependent upon electrons other than those in the operator can be separated out. For example, the first term in the expansion, Equation (2.77), is:

$$\iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) (\mathcal{H}_1) \chi_1(1)\chi_2(2) \quad (2.78)$$

The operator \mathcal{H}_1 is a function of the coordinates of electron 1 only, so terms involving electron 2 can be separated out as follows:

$$\begin{aligned} &\iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) (\mathcal{H}_1) \chi_1(1)\chi_2(2) \\ &= \int d\tau_2 \chi_2(2)\chi_2(2) \int d\tau_1 \chi_1(1) \left(-\frac{1}{2}\nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \right) \chi_1(1) \end{aligned} \quad (2.79)$$

If the molecular orbitals are normalised, the integral $\int d\tau_2 \chi_2(2)\chi_2(2)$ equals 1. Further simplification can be achieved by splitting the integral involving electron 1 into separate integrals over the spatial and spin parts; the integral over spin orbitals is equal to the product of an integral over the spatial coordinates and an integral over the spin coordinates:

$$\begin{aligned} &\int d\tau_1 \chi_1(1) \left(-\nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \right) \chi_1(1) \\ &= \int d\nu_1 1\sigma_g(1) \left(-\frac{1}{2}\nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \right) 1\sigma_g(1) \int d\sigma_1 \alpha(1)\alpha(1) \end{aligned} \quad (2.80)$$

$d\nu$ indicates integration over spatial coordinates and $d\sigma$ indicates integration over the spin coordinates. The integral over the spin coordinates equals 1. This expression corresponds

to the sum of the kinetic and potential energy of an electron in the orbital $1\sigma_g$ in the electrostatic field of the two bare nuclei. This integral can in turn be expanded by substituting the atomic orbital combination for $1\sigma_g$:

$$\begin{aligned} & \int d\nu_1 1\sigma_g(1) \left(-\frac{1}{2} \nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \right) 1\sigma_g(1) \\ &= A^2 \int d\nu_1 \{1s_A(1) + 1s_B(1)\} \left(-\frac{1}{2} \nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \right) \{1s_A(1) + 1s_B(1)\} \quad (2.81) \end{aligned}$$

A is the normalisation constant. The integral in Equation (2.81) can in turn be factorised to give a sum of integrals, each of which involves a pair of atomic orbitals:

$$\begin{aligned} & \int d\nu_1 \{1s_A(1) + 1s_B(1)\} \left(-\frac{1}{2} \nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \right) \{1s_A(1) + 1s_B(1)\} \\ &= \int d\nu_1 1s_A(1) \left(-\frac{1}{2} \nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \right) 1s_A(1) \\ &+ \int d\nu_1 1s_B(1) \left(-\frac{1}{2} \nabla_1^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \right) 1s_B(1) + \dots \quad (2.82) \end{aligned}$$

Let us now apply the same procedure to the second term in Equation (2.77):

$$\iint d\tau_1 d\tau_2 \chi_1(1) \chi_2(2) (\mathcal{H}_1) \chi_2(1) \chi_1(2) = \int d\tau_1 \chi_1(1) (\mathcal{H}_1) \chi_2(1) \int d\tau_2 \chi_2(2) \chi_1(2) \quad (2.83)$$

This particular integral is zero because the molecular orbitals are orthogonal and so the integral over the coordinates of electron 2 equals zero:

$$\int d\tau_2 \chi_2(2) \chi_1(2) = 0 \quad (2.84)$$

A similar procedure can be applied to the other integrals involving electron-nuclear interactions; it turns out that there are four non-zero integrals, each of which is equal to the energy of a single electron in the field of the two hydrogen nuclei.

There remain four integrals arising from electron-electron interactions. These are:

$$\begin{aligned} & \iint d\tau_1 d\tau_2 \chi_1(1) \chi_2(2) \left(\frac{1}{r_{12}} \right) \chi_1(1) \chi_2(2) + \iint d\tau_1 d\tau_2 \chi_2(1) \chi_1(2) \left(\frac{1}{r_{12}} \right) \chi_2(1) \chi_1(2) \\ & - \iint d\tau_1 d\tau_2 \chi_1(1) \chi_2(2) \left(\frac{1}{r_{12}} \right) \chi_2(1) \chi_1(2) - \iint d\tau_1 d\tau_2 \chi_2(1) \chi_1(2) \left(\frac{1}{r_{12}} \right) \chi_1(1) \chi_2(2) \quad (2.85) \end{aligned}$$

The first two of these can be simplified as follows:

$$\begin{aligned} & \iint d\tau_1 d\tau_2 \chi_1(1) \chi_2(2) \left(\frac{1}{r_{12}} \right) \chi_1(1) \chi_2(2) = \iint d\nu_1 d\nu_2 1\sigma_g(1) 1\sigma_g(2) \left(\frac{1}{r_{12}} \right) 1\sigma_g(1) 1\sigma_g(2) \\ & \quad \times \int d\sigma_1 \alpha(1) \alpha(1) \int d\sigma_2 \beta(2) \beta(2) \\ &= \iint d\nu_1 d\nu_2 1\sigma_g(1) 1\sigma_g(1) \left(\frac{1}{r_{12}} \right) 1\sigma_g(2) 1\sigma_g(2) \quad (2.86) \end{aligned}$$

According to the Born interpretation of the wavefunction, $1\sigma_g(\mathbf{r}_1)1\sigma_g(\mathbf{r}_1)$ equals the electron density of electron 1 in orbital $1\sigma_g$ at a position \mathbf{r}_1 . Similarly, $1\sigma_g(\mathbf{r}_2)1\sigma_g(\mathbf{r}_2)$ is the electron density of electron 2. The electrostatic repulsion between these regions of electron density thus equals $1\sigma_g(\mathbf{r}_1)1\sigma_g(\mathbf{r}_1) \times (1/r_{12}) \times 1\sigma_g(\mathbf{r}_2)1\sigma_g(\mathbf{r}_2)$, where r_{12} is the distance between the two electrons. The integral of this function over all space thus corresponds to the electrostatic (Coulomb) repulsion between the two orbitals.

If we substitute the atomic orbital expansion, we obtain a series of two-electron integrals, each of which involves four atomic orbitals:

$$\begin{aligned} & \iint d\nu_1 d\nu_2 1\sigma_g(1) 1\sigma_g(2) \left(\frac{1}{r_{12}} \right) 1\sigma_g(1) 1\sigma_g(2) \\ &= \iint d\nu_1 d\nu_2 1s_A(1) 1s_A(2) \left(\frac{1}{r_{12}} \right) 1s_A(1) 1s_A(2) \\ &+ \iint d\nu_1 d\nu_2 1s_A(1) 1s_A(2) \left(\frac{1}{r_{12}} \right) 1s_A(1) 1s_B(2) + \dots \quad (2.87) \end{aligned}$$

The remaining two integrals from Equation (2.85) are:

$$\begin{aligned} & \iint d\tau_1 d\tau_2 \chi_1(1) \chi_2(2) \left(\frac{1}{r_{12}} \right) \chi_2(1) \chi_1(2) = \iint d\nu_1 d\nu_2 1\sigma_g(1) 1\sigma_g(2) \left(\frac{1}{r_{12}} \right) 1\sigma_g(1) 1\sigma_g(2) \\ & \quad \times \int d\sigma_1 \alpha(1) \beta(1) \int d\sigma_2 \beta(2) \alpha(2) \quad (2.88) \end{aligned}$$

$$\begin{aligned} & \iint d\tau_1 d\tau_2 \chi_2(1) \chi_1(2) \left(\frac{1}{r_{12}} \right) \chi_1(1) \chi_2(2) = \iint d\nu_1 d\nu_2 1\sigma_g(1) 1\sigma_g(2) \left(\frac{1}{r_{12}} \right) 1\sigma_g(1) 1\sigma_g(2) \\ & \quad \times \int d\sigma_1 \beta(1) \alpha(1) \int d\sigma_2 \alpha(2) \beta(2) \quad (2.89) \end{aligned}$$

Both of these integrals are zero due to the orthogonality of the electron spin states α and β .

The triplet excited state of H_2 is obtained by promoting an electron to a higher-energy molecular orbital. This higher-energy (antibonding) orbital is written $1\sigma_u$ and can be considered to arise from two $1s$ orbitals as follows:

$$1\sigma_u = A(1s_A - 1s_B) \quad (2.90)$$

The triplet state has two unpaired electrons with the same spin (α) and so the wavefunction state is:

$$\begin{vmatrix} 1\sigma_g \alpha(1) & 1\sigma_u \alpha(1) \\ 1\sigma_g \alpha(2) & 1\sigma_u \alpha(2) \end{vmatrix} \quad (2.91)$$

If we now expand the expression for the energy as for the ground state, terms analogous to the electron-nucleus and electron-electron interactions can again be obtained. However, the cross-terms are no longer equal to zero as was the case for the ground state, because the

electron spins are now the same (both α). For example, compare with Equation (2.88):

$$\iint d\tau_1 d\tau_2 \chi_1(1)\chi_2(2) \left(\frac{1}{r_{12}}\right) \chi_2(1)\chi_1(2) = \iint d\nu_1 d\nu_2 1\sigma_g(1)1\sigma_u(2) \left(\frac{1}{r_{12}}\right) 1\sigma_g(2)1\sigma_u(1) \\ \times \int d\sigma_1 \alpha(1)\alpha(1) \int d\sigma_2 \alpha(2)\alpha(2) \quad (2.92)$$

This contribution is called the *exchange interaction*. This appears with a minus sign in the expression for the total energy and so acts to stabilise the triplet $1s^2 2s^1$ state over the analogous singlet state. The exchange term is only non-zero for electrons of the same spin. It has the effect of making electrons of the same spin 'avoid' each other. As a result of this each electron can be considered to have a 'hole' associated with it. This hole is known as the *exchange hole* or the *Fermi hole*.

2.4.2 The Energy of a General Polyelectronic System

The hydrogen molecule is such a small problem that all of the integrals can be written out in full. This is rarely the case in molecular orbital calculations. Nevertheless, the same principles are used to determine the energy of a polyelectronic molecular system. For an N -electron system, the Hamiltonian takes the following general form:

$$\mathcal{H} = \left(-\frac{1}{2} \sum_{i=1}^N \nabla_i^2 - \frac{1}{r_{1A}} - \frac{1}{r_{1B}} \cdots + \frac{1}{r_{12}} + \frac{1}{r_{13}} + \cdots \right) \quad (2.93)$$

As with the hydrogen molecule, we have adopted the convention that the nuclei are labelled using capital letters A, B, C, etc., and the electrons are labelled 1, 2, 3,

Recall that the Slater determinant for a system of N electrons in N spin orbitals can be written:

$$\begin{vmatrix} \chi_1(1) & \chi_2(1) & \chi_3(1) & \cdots & \chi_N(1) \\ \chi_1(2) & \chi_2(2) & \chi_3(2) & \cdots & \chi_N(2) \\ \chi_1(3) & \chi_2(3) & \chi_3(3) & \cdots & \chi_N(3) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \chi_1(N) & \chi_2(N) & \chi_3(N) & \cdots & \chi_N(N) \end{vmatrix} \quad (2.94)$$

Each term in the determinant can thus be written $\chi_i(1)\chi_j(2)\chi_k(3)\cdots\chi_u(N-1)\chi_v(N)$ where i, j, k, \dots, u, v is a series of N integers.

As usual, the energy can be calculated from $E = \int \Psi \mathcal{H} \Psi / \int \Psi \Psi$:

$$\int \Psi \mathcal{H} \Psi = \int \cdots \int d\tau_1 d\tau_2 \cdots d\tau_N \left\{ [\chi_i(1)\chi_j(2)\chi_k(3)\cdots] \right. \\ \times \left(-\frac{1}{2} \sum_i \nabla_i^2 - (1/r_{1A}) - (1/r_{1B}) \cdots + (1/r_{12}) + (1/r_{13}) + \cdots \right) \\ \left. \times [\chi_i(1)\chi_j(2)\chi_k(3)\cdots] \right\} \quad (2.95)$$

$$\int \Psi \Psi = \int \cdots \int d\tau_1 d\tau_2 \cdots d\tau_N \{ [\chi_i(1)\chi_j(2)\chi_k(3)\cdots] [\chi_i(1)\chi_j(2)\chi_k(3)\cdots] \} \quad (2.96)$$

We can now see why the normalisation factor of the Slater determinantal wavefunction is $1/\sqrt{N!}$. If each determinant contains $N!$ terms then the product of two Slater determinants, [determinant][determinant], contains $(N!)^2$ terms. However, if the spin orbitals form an orthonormal set then only products of identical terms from the determinant will be non-zero when integrated over all space. We can illustrate this with the three-electron example. Considering just the first two terms in the expansion we obtain the following:

$$\iiint d\tau_1 d\tau_2 d\tau_3 [\chi_1(1)\chi_2(2)\chi_3(3) - \chi_1(1)\chi_3(2)\chi_2(3) + \cdots] \\ \times [\chi_1(1)\chi_2(2)\chi_3(3) - \chi_1(1)\chi_3(2)\chi_2(3) + \cdots] \quad (2.97)$$

When multiplied out this gives:

$$\iiint d\tau_1 d\tau_2 d\tau_3 [\chi_1(1)\chi_2(2)\chi_3(3)][\chi_1(1)\chi_2(2)\chi_3(3)] \\ - \iiint d\tau_1 d\tau_2 d\tau_3 [\chi_1(1)\chi_2(2)\chi_3(3)][\chi_1(1)\chi_3(2)\chi_2(3)] + \cdots \\ + \iiint d\tau_1 d\tau_2 d\tau_3 [\chi_1(1)\chi_3(2)\chi_2(3)][\chi_1(1)\chi_3(2)\chi_2(3)] + \cdots \quad (2.98)$$

The first of the integrals in Equation (2.98) equals 1 (if the spin orbitals are normalised). The second term is zero because the terms involving both electrons 2 and 3 are different (for example, the integral $\int d\tau_2 \chi_2(2)\chi_3(2)$ will be zero due to the orthogonality of the spin orbitals χ_2 and χ_3). The third term in Equation (2.98) will be equal to 1, and so on. It turns out that there are $N!$ such non-zero terms. Thus if each individual term in the determinant is normalised, then:

$$\int \Psi \Psi = N! \quad (2.99)$$

Hence the normalisation factor for the determinantal wavefunction is $1/\sqrt{N!}$.

Turning now to the numerator in the energy expression (Equation (2.95)), this can be broken down into a series of one-electron and two-electron integrals, as for the hydrogen molecule. Each of these individual integrals has the general form:

$$\int \cdots \int d\tau_1 d\tau_2 \cdots [term1]operator[term2] \quad (2.100)$$

[term1] and [term2] each represent one of the $N!$ terms in the Slater determinant. To simplify this integral, we first recognise that all spin orbitals involving an electron that does not appear in the operator can be taken outside the integral. For example, if the operator is $1/r_{1A}$, then all spin orbitals other than those that depend on the coordinates of electron 1 can be separated from the integral. The orthogonality of the spin orbitals means that the integral will be zero unless all indices involving these other electrons are the same in

[term1] and [term2]. Again, to use our three-electron system as an example:

$$\begin{aligned} & \iiint d\tau_1 d\tau_2 d\tau_3 [\chi_1(1)\chi_2(2)\chi_3(3)] \left(-\frac{1}{r_{1A}}\right) [\chi_1(1)\chi_2(2)\chi_3(3)] \\ &= \iint d\tau_2 d\tau_3 [\chi_2(2)\chi_3(3)] [\chi_2(2)\chi_3(3)] \int d\tau_1 \chi_1(1) \left(-\frac{1}{r_{1A}}\right) \chi_1(1) \\ &= \int d\tau_1 \chi_1(1) \left(-\frac{1}{r_{1A}}\right) \chi_1(1) \end{aligned} \quad (2.101)$$

But:

$$\begin{aligned} & \iiint d\tau_1 d\tau_2 d\tau_3 [\chi_1(1)\chi_2(2)\chi_3(3)] \left(-\frac{1}{r_{1A}}\right) [\chi_1(1)\chi_3(2)\chi_2(3)] \\ &= \iint d\tau_2 d\tau_3 [\chi_2(2)\chi_3(3)] [\chi_3(2)\chi_2(3)] \int d\tau_1 \chi_1(1) \left(-\frac{1}{r_{1A}}\right) \chi_1(1) \\ &= 0 \end{aligned} \quad (2.102)$$

For integrals that involve two-electron operators (i.e. $1/r_{ij}$), only those terms that do not involve the coordinates of the two electrons can be taken outside the integral. For example:

$$\begin{aligned} & \iiint d\tau_1 d\tau_2 d\tau_3 [\chi_1(1)\chi_2(2)\chi_3(3)] \left(\frac{1}{r_{12}}\right) [\chi_1(1)\chi_2(2)\chi_3(3)] \\ &= \iint d\tau_1 d\tau_2 [\chi_1(1)\chi_2(2)] \left(\frac{1}{r_{12}}\right) [\chi_1(2)\chi_2(2)] \int d\tau_3 \chi_3(3)\chi_3(3) \\ &= \iint d\tau_1 d\tau_2 [\chi_1(1)\chi_2(2)] \left(\frac{1}{r_{12}}\right) [\chi_1(2)\chi_2(2)] \end{aligned} \quad (2.103)$$

But:

$$\begin{aligned} & \iiint d\tau_1 d\tau_2 d\tau_3 [\chi_1(1)\chi_2(2)\chi_3(3)] \left(\frac{1}{r_{12}}\right) [\chi_1(1)\chi_3(2)\chi_2(3)] \\ &= \iint d\tau_1 d\tau_2 [\chi_1(1)\chi_2(2)] \left(\frac{1}{r_{12}}\right) [\chi_1(2)\chi_3(2)] \int d\tau_3 \chi_3(3)\chi_2(3) \\ &= 0 \end{aligned} \quad (2.104)$$

As a consequence of these results, most of the individual integrals in the expansion will be zero. Nevertheless, it can be readily envisaged that there will still be an extremely large number of integrals to consider for all except the smallest problems. It is thus more convenient to write the energy expression in a concise form that recognises the three types of interaction that contribute to the total electronic energy of the system.

First, there is the kinetic and potential energy of each electron moving in the field of the nuclei. The energy associated with this contribution for the molecular orbital χ_i is often written H_{ii}^{core} and for M nuclei is given by:

$$H_{ii}^{\text{core}} = \int d\tau_1 \chi_i(1) \left(-\frac{1}{2} \nabla_i^2 - \sum_{A=1}^M \frac{Z_A}{r_{iA}}\right) \chi_i(1) \quad (2.105)$$

For N electrons in N molecular orbitals this contribution to the total energy is:

$$E_{\text{total}}^{\text{core}} = \sum_{i=1}^N \int d\tau_1 \chi_i(1) \left(-\frac{1}{2} \nabla_i^2 - \sum_{A=1}^M \frac{Z_A}{r_{iA}}\right) \chi_i(1) = \sum_{i=1}^N H_{ii}^{\text{core}} \quad (2.106)$$

Here we have followed convention and have used the label '1' wherever there is an integral involving the coordinates of a single electron, even though the actual electron may not be 'electron 1'. Similarly, when it is necessary to consider two electrons then the labels 1 and 2 are conventionally employed. H_{ii}^{core} makes a favourable (i.e. negative) contribution to the electronic energy.

The second contribution to the energy arises from the electrostatic repulsion between pairs of electrons. This interaction depends on the electron-electron distance and, as we have seen, is calculated from integrals such as:

$$J_{ij} = \iint d\tau_1 d\tau_2 \chi_i(1)\chi_j(2) \left(\frac{1}{r_{12}}\right) \chi_i(1)\chi_j(2) \quad (2.107)$$

The symbol J_{ij} is often used to represent this Coulomb interaction between electrons in spin orbitals i and j , and is unfavourable (i.e. positive). The total electrostatic interaction between the electron in orbital χ_i and the other $N - 1$ electrons is a sum of all such integrals, where the summation index j runs from 1 to N , excluding i :

$$\begin{aligned} E_i^{\text{Coulomb}} &= \sum_{j \neq i}^N \int d\tau_1 d\tau_2 \chi_i(1)\chi_j(2) \frac{1}{r_{12}} \chi_j(2)\chi_i(1) \\ &\equiv \sum_{j \neq i}^N \int d\tau_1 d\tau_2 \chi_i(1)\chi_i(1) \frac{1}{r_{12}} \chi_j(2)\chi_j(2) \end{aligned} \quad (2.108)$$

The total Coulomb contribution to the electronic energy of the system is obtained as a double summation over all electrons, taking care to count each interaction just once:

$$E_{\text{total}}^{\text{Coulomb}} = \sum_{i=1}^N \sum_{j=i+1}^N \int d\tau_1 d\tau_2 \chi_i(1)\chi_i(1) \frac{1}{r_{12}} \chi_j(2)\chi_j(2) = \sum_{i=1}^N \sum_{j=i+1}^N J_{ij} \quad (2.109)$$

The third contribution to the energy is the exchange 'interaction'. This has no classical counterpart and arises because the motions of electrons with parallel spins are correlated: whereas there is a finite probability of finding two electrons with opposite (i.e. paired) spins at the same point in space, where the spins are the same then the probability is zero. This can be considered a manifestation of the Pauli principle, for if two electrons occupied the same region of space and had parallel spins then they could be considered to have the same set of quantum numbers. Electrons with the same spin thus tend to 'avoid' each other, and they experience a lower Coulombic repulsion, giving a lower (i.e. more favourable) energy. The exchange interaction involves integrals of the form:

$$K_{ij} = \iint d\tau_1 d\tau_2 \chi_i(1)\chi_j(2) \left(\frac{1}{r_{12}}\right) \chi_i(2)\chi_j(1) \quad (2.110)$$

This integral is only non-zero if the spins of the electrons in the spin orbitals χ_i and χ_j are the same. The energy due to exchange is often represented as K_{ij} . The exchange energy between the electron in spin orbital χ_i and the other $N - 1$ electrons is:

$$E_i^{\text{exchange}} = \sum_{j \neq i}^N \iint d\tau_1 d\tau_2 \chi_i(1) \chi_j(2) \left(\frac{1}{r_{12}} \right) \chi_i(2) \chi_j(1) \quad (2.111)$$

The total exchange energy is calculated thus:

$$E_{\text{total}}^{\text{exchange}} = \sum_{i=1}^N \sum_{j=i+1}^N \iint d\tau_1 d\tau_2 \chi_i(1) \chi_j(2) \left(\frac{1}{r_{12}} \right) \chi_i(2) \chi_j(1) = \sum_{j=1}^N \sum_{i'=i+1}^N K_{ij} \quad (2.112)$$

The prime on the counter j' indicates that the summation is only over electrons with the same spin as electron i .

2.4.3 Shorthand Representations of the One- and Two-electron Integrals

Various shorthand ways have been devised to represent the integrals involved in an electronic structure calculation. The two-electron integrals J_{ij} and K_{ij} are particularly long-winded to write out. In one scheme the Coulomb interaction J_{ij} is written as:

$$\left\langle \chi_i^* \chi_j^* \left| \frac{1}{r_{12}} \right| \chi_i \chi_j \right\rangle \quad (2.113)$$

In this notation the complex parts are written on the left-hand side and the real parts on the right. Sometimes the χ symbol is eliminated:

$$\left\langle ij \left| \frac{1}{r_{12}} \right| ij \right\rangle \quad (2.114)$$

The exchange integrals would be written:

$$\left\langle ij \left| \frac{1}{r_{12}} \right| ji \right\rangle \quad (2.115)$$

in this notation.

A notation that is widely used in the chemical literature writes the orbitals that are functions of electron 1 on the left-hand side (with the complex conjugate orbital first, if appropriate) and the orbitals that are functions of electron 2 on the right-hand side (again with the complex conjugate orbital first). In this notation, which is the one that we will adopt, the Coulomb integral is written $(ii|jj)$ and the exchange integral $(ij|ji)$. The one-electron integrals such as Equation (2.105) are written as follows:

$$\left(i \left| -\frac{1}{2} \nabla_i^2 - \sum_{A=1}^M \frac{Z_A}{r_{iA}} \right| j \right) \equiv \int d\tau_1 \chi_i(1) \left(-\frac{1}{2} \nabla_i^2 - \sum_{A=1}^M \frac{Z_A}{r_{iA}} \right) \chi_j(1) \quad (2.116)$$

When calculating the total energy of the system, we should not forget the Coulomb interaction between the nuclei; this is constant within the Born-Oppenheimer approximation for a given spatial arrangement of nuclei. When it is desired to change the nuclear positions,

it is of course necessary to take the internuclear repulsion energy into account, which is calculated using the Coulomb equation:

$$\sum_{A=1}^M \sum_{B=A+1}^M \frac{Z_A Z_B}{R_{AB}} \quad (2.117)$$

2.4.4 The Energy of a Closed-shell System

In molecular modelling we are usually concerned with the ground states of molecules, most of which have closed-shell configurations. In a closed-shell system containing N electrons in $N/2$ orbitals, there are two spin orbitals associated with each spatial orbital ψ_i : $\psi_i\alpha$ and $\psi_i\beta$. The electronic energy of such a system can be calculated in a manner analogous to that for the hydrogen molecule. First, there is the energy of each electron moving in the field of the bare nuclei. For an electron in a molecular orbital χ_i , this contributes an energy H_{ii}^{core} . If there are two electrons in the orbital then the energy is $2H_{ii}^{\text{core}}$ and for $N/2$ orbitals the total contribution to the energy will be:

$$\sum_{i=1}^{N/2} 2H_{ii}^{\text{core}} \quad (2.118)$$

If we consider the electron-electron terms, the interaction between each pair of orbitals ψ_i and ψ_j involves a total of four electrons. There are four ways in which two electrons in one orbital can interact in a Coulomb sense with two electrons in a second orbital, thus giving $4J_{ij}$. However, there are just two ways to obtain paired electrons from this arrangement, giving a total exchange contribution of $-2K_{ij}$. Finally, the Coulomb interaction between each pair of electrons in the same orbital must be included; there is no exchange interaction because the electrons have paired spins. The total energy is thus given as:

$$E = 2 \sum_{i=1}^{N/2} H_{ii}^{\text{core}} + \sum_{i=1}^{N/2} \sum_{j=i+1}^{N/2} (4J_{ij} - 2K_{ij}) + \sum_{i=1}^{N/2} J_{ii} \quad (2.119)$$

A more concise form of this equation can be obtained if we recognise that $J_{ii} = K_{ii}$:

$$E = 2 \sum_{i=1}^{N/2} H_{ii}^{\text{core}} + \sum_{i=1}^{N/2} \sum_{j=1}^{N/2} (2J_{ij} - K_{ij}) \quad (2.120)$$

2.5 The Hartree-Fock Equations

In our hydrogen molecule calculation in Section 2.4.1 the molecular orbitals were provided as input, but in most electronic structure calculations we are usually trying to calculate the molecular orbitals. How do we go about this? We must remember that for many-body problems there is no 'correct' solution; we therefore require some means to decide whether one proposed wavefunction is 'better' than another. Fortunately, the *variation theorem* provides us with a mechanism for answering this question. The theorem states that the

energy calculated from an approximation to the true wavefunction will always be greater than the true energy. Consequently, the better the wavefunction, the lower the energy. The 'best' wavefunction is obtained when the energy is a minimum. At a minimum, the first derivative of the energy, δE , will be zero. The Hartree-Fock equations are obtained by imposing this condition on the expression for the energy, subject to the constraint that the molecular orbitals remain orthonormal. The orthonormality condition is written in terms of the *overlap integral*, S_{ij} , between two orbitals i and j . Thus

$$S_{ij} = \int \chi_i \chi_j d\tau = \delta_{ij} \quad (\delta_{ij} \text{ is the Kronecker delta}) \quad (2.121)$$

This type of constrained minimisation problem can be tackled using the method of Lagrange multipliers. In this approach (see Section 1.10.5 for a brief introduction to Lagrange multipliers) the derivative of the function to be minimised is added to the derivatives of the constraint(s) multiplied by a constant called a Lagrange multiplier. The sum is then set equal to zero. If the Lagrange multiplier for each of the orthonormality conditions is written λ_{ij} , then:

$$\delta E + \delta \sum_i \sum_j \lambda_{ij} S_{ij} = 0 \quad (2.122)$$

In the Hartree-Fock equations the Lagrange multipliers are actually written $-2\varepsilon_{ij}$ to reflect the fact that they are related to the molecular orbital energies. The equation to be solved is thus:

$$\delta E - 2\delta \sum_i \sum_j \varepsilon_{ij} S_{ij} = 0 \quad (2.123)$$

We will not describe in detail how this equation is solved, as it is rather complicated. However, a qualitative picture is possible. The major difference between polyelectronic systems and systems with single electrons is the presence of interactions between the electrons, which, as we have seen, are expressed as Coulomb and exchange integrals. Suppose we are given the task of finding the 'best' (i.e. lowest energy) wavefunction for a polyelectronic system. We wish to retain the orbital picture of the system, in which single electrons are assigned to individual spin orbitals. The problem is to find a solution which simultaneously enables all the electronic motions to be taken into account, as a change in the spin orbital for one electron will influence the behaviour of an electron in another spin orbital due to the coupling of the electronic motions. We concentrate on a single electron in a spin orbital χ_i in the field of the nuclei and the other electrons in their (fixed) spin orbitals χ_j . The Hamiltonian operator for the electron in χ_i contains three terms appropriate to the three different contributions to the energy that were identified above (core, Coulomb, exchange). The result can be written as an integro-differential equation for χ_i that has the following form:

$$\left[-\frac{1}{2}\nabla_i^2 - \sum_{A=1}^M \frac{Z_A}{r_{iA}} \right] \chi_i(1) + \sum_{j \neq i} \left[\int d\tau_2 \chi_j(2) \chi_j(2) \frac{1}{r_{12}} \right] \chi_i(1) - \sum_{j \neq i} \left[\int d\tau_2 \chi_j(2) \chi_i(2) \frac{1}{r_{12}} \right] \chi_j(1) = \sum_j \varepsilon_{ij} \chi_j(1) \quad (2.124)$$

This expression can be tidied up by introducing three operators that represent the contributions to the energy of the spin orbital χ_i in the 'frozen' system:

The core Hamiltonian operator, $\mathcal{H}^{\text{core}}(1)$:

$$\mathcal{H}^{\text{core}}(1) = -\frac{1}{2}\nabla_1^2 - \sum_{A=1}^M \frac{Z_A}{r_{1A}} \quad (2.125)$$

In the absence of any interelectronic interactions this would be the only operator present, corresponding to the motion of a single electron moving in the field of the bare nuclei.

The Coulomb operator, $\mathcal{J}_j(1)$:

$$\mathcal{J}_j(1) = \int d\tau_2 \chi_j(2) \frac{1}{r_{12}} \chi_j(2) \quad (2.126)$$

This operator corresponds to the average potential due to an electron in χ_j .

The exchange operator $\mathcal{K}_j(1)$:

$$\mathcal{K}_j(1) \chi_i(1) = \left[\int d\tau_2 \chi_j(2) \frac{1}{r_{12}} \chi_i(2) \right] \chi_j(1) \quad (2.127)$$

The form of this operator is rather unusual, insofar as it must be defined in terms of its effect when acting on the spin orbital χ_i .

Equation (2.124) can thus be written:

$$\mathcal{H}^{\text{core}}(1) \chi_i(1) + \sum_{j \neq i}^N \mathcal{J}_j(1) \chi_i(1) - \sum_{j \neq i}^N \mathcal{K}_j(1) \chi_i(1) = \sum_j \varepsilon_{ij} \chi_j(1) \quad (2.128)$$

Making use of the fact that $\{\mathcal{J}_i(1) - \mathcal{K}_i(1)\} \chi_i(1) = 0$ leads to the following form:

$$\left[\mathcal{H}^{\text{core}}(1) + \sum_{j=1}^N \{\mathcal{J}_j(1) - \mathcal{K}_j(1)\} \right] \chi_i(1) = \sum_{j=1}^N \varepsilon_{ij} \chi_j(1) \quad (2.129)$$

Or, more simply:

$$f_i \chi_i = \sum_j \varepsilon_{ij} \chi_j \quad (2.130)$$

f_i is called the *Fock operator*:

$$f_i(1) = \mathcal{H}^{\text{core}}(1) + \sum_{j=1}^N \{\mathcal{J}_j(1) - \mathcal{K}_j(1)\} \quad (2.131)$$

For a closed-shell system, the Fock operator has the following form:

$$f_i(1) = \mathcal{H}^{\text{core}}(1) + \sum_{j=1}^{N/2} \{2\mathcal{J}_j(1) - \mathcal{K}_j(1)\} \quad (2.132)$$

The Fock operator is an effective one-electron Hamiltonian for the electron in the polyelectronic system. However, written in this form of Equation (2.130), the Hartree-Fock

equations do not seem to be particularly useful: on the left-hand side we have the Fock operator acting on the molecular orbital χ_i , but this returns, not the molecular orbital multiplied by a constant as in a normal eigenvalue equation, but rather a series of orbitals χ_j multiplied by some unknown constants ε_{ij} . This is because the solutions to the Hartree–Fock equations are not unique. We have already seen that the value of a determinant is unaffected when the multiple of any column is added to another column. If such a transformation is performed on the Slater determinant, then a different set of constants ε'_{ij} would be obtained with the spin orbitals χ'_i being linear combinations of the first set. Certain transformations give rise to localised orbitals, which are particularly useful for understanding the chemical nature of the system. These localised orbitals are no more ‘correct’ than a delocalised set. Fortunately, it is possible to manipulate Equations (2.130) mathematically so that the Lagrangian multipliers are zero unless the indices i and j are the same. The Hartree–Fock equations then take on the standard eigenvalue form:

$$f_i \chi_i = \varepsilon_i \chi_i \quad (2.133)$$

Recall that in setting up these equations, each electron has been assumed to move in a ‘fixed’ field comprising the nuclei and the other electrons. This has important implications for the way in which we attempt to find a solution, for any solution that we might find by solving the equation for one electron will naturally affect the solutions for the other electrons in the system. The general strategy is called a *self-consistent field* (SCF) approach. One way to solve these equations is as follows. First, a set of trial solutions χ_i to the Hartree–Fock eigenvalue equations are obtained. These are used to calculate the Coulomb and exchange operators. The Hartree–Fock equations are solved, giving a second set of solutions χ_i , which are used in the next iteration. The SCF method thus gradually refines the individual electronic solutions that correspond to lower and lower total energies until the point is reached at which the results for all the electrons are unchanged, when they are said to be *self-consistent*.

2.5.1 Hartree–Fock Calculations for Atoms and Slater’s Rules

The Hartree–Fock equations are usually solved in different ways for atoms and for molecules. For atoms, the equations can be solved numerically if it is assumed that the electron distribution is spherically symmetrical. However, these numerical solutions are not particularly useful. Fortunately, analytical approximations to these solutions, which are very similar to those obtained for the hydrogen atom, can be used with considerable success. These approximate analytical functions thus have the form:

$$\psi = R_{nl}(r) Y_{lm}(\theta, \phi) \quad (2.134)$$

Y is a spherical harmonic (as for the hydrogen atom) and R is a radial function. The radial functions obtained for the hydrogen atom cannot be used directly for polyelectronic atoms due to the screening of the nuclear charge by the inner shell electrons, but the hydrogen atom functions are acceptable if the orbital exponent is adjusted to account for the screening effect. Even so, the hydrogen atom functions are not particularly convenient to use in molecular orbital calculations due to their complicated functional form. Slater [Slater 1930] suggested

a simpler analytical form for the radial functions:

$$R_{nl}(r) = (2\zeta)^{n+1/2} [(2n)!]^{-1/2} r^{n-1} e^{-\zeta r} \quad (2.135)$$

These functions are universally known as *Slater type orbitals* (STOs) and are just the leading term in the appropriate Laguerre polynomials. The first three Slater functions are as follows:

$$R_{1s}(r) = 2\zeta^{3/2} e^{-\zeta r} \quad (2.136)$$

$$R_{2s}(r) = R_{2p}(r) = \left(\frac{4\zeta^5}{3}\right)^{1/2} r e^{-\zeta r} \quad (2.137)$$

$$R_{3s}(r) = R_{3p}(r) = R_{3d}(r) = \left(\frac{8\zeta^7}{45}\right)^{1/2} r^2 e^{-\zeta r} \quad (2.138)$$

To obtain the whole orbital we must multiply $R(r)$ by the appropriate angular part. For example, we would use the following expressions for the 1s, 2s and 2p_z orbitals:

$$\phi_{1s}(r) = \sqrt{\zeta^3/\pi} \exp(-\zeta r) \quad (2.139)$$

$$\phi_{2s}(r) = \sqrt{\zeta^5/3\pi} r \exp(-\zeta r) \quad (2.140)$$

$$\phi_{2p_z}(r) = \sqrt{\zeta^5/\pi} \exp(-\zeta r) \cos \theta \quad (2.141)$$

Slater provided a series of empirical rules for choosing the orbital exponents ζ , which are given by:

$$\zeta = \frac{Z - \sigma}{n^*} \quad (2.142)$$

Z is the atomic number and σ is a *shielding constant*, determined as below. n^* is an effective principal quantum number, which takes the same value as the true principal quantum number for $n = 1, 2$ or 3 , but for $n = 4, 5, 6$ has the values 3.7, 4.0, 4.2, respectively. The shielding constant is obtained as follows:

First, divide the orbitals into the following groups:

$$(1s); (2s, 2p); (3s, 3p); (3d); (4s, 4p); (4d); (4f); (5s, 5p); (5d) \quad (2.143)$$

For a given orbital, σ is obtained by adding together the following contributions:

- zero from an orbital further from the nucleus than those in the group;
- 0.35 from each other electron in the same group, but if the other orbital is the 1s then the contribution is 0.3;
- 1.0 for each electron in a group with a principal quantum number 2 or more fewer than the current orbital;
- for each electron with a principal quantum number 1 fewer than the current orbital: 1.0 if the current orbital is d or f; 0.85 if the current orbital is s or p.

The shielding constant for the valence electrons of silicon is obtained using Slater’s rules as follows. The electronic configuration of Si is $(1s^2)(2s^2 2p^6)(3s^2 3p^2)$. We therefore count

3×0.35 under rule (b), 2.0 under rule (c) and 8×0.85 under rule (d), giving a total of 9.85 . When subtracted from the atomic number (14) this gives 4.15 for the value of $Z - \sigma$.

2.5.2 Linear Combination of Atomic Orbitals (LCAO) in Hartree–Fock Theory

Direct solution of the Hartree–Fock equations is not a practical proposition for molecules and so it is necessary to adopt an alternative approach. The most popular strategy is to write each spin orbital as a linear combination of single electron orbitals:

$$\psi_i = \sum_{\nu=1}^K c_{\nu i} \phi_{\nu} \quad (2.144)$$

The one-electron orbitals ϕ_{ν} are commonly called *basis functions* and often correspond to the atomic orbitals. We will label the basis functions with the Greek letters μ, ν, λ and σ . In the case of Equation (2.144) there are K basis functions and we should therefore expect to derive a total of K molecular orbitals (although not all of these will necessarily be occupied by electrons). The smallest number of basis functions for a molecular system will be that which can just accommodate all the electrons in the molecule. More sophisticated calculations use more basis functions than a minimal set. At the *Hartree–Fock limit* the energy of the system can be reduced no further by the addition of any more basis functions; however, it may be possible to lower the energy below the Hartree–Fock limit by using a functional form of the wavefunction that is more extensive than the single Slater determinant.

In accordance with the variation theorem we require the set of coefficients $c_{\nu i}$ that gives the lowest-energy wavefunction, and some scheme for changing the coefficients to derive that wavefunction. For a given basis set and a given functional form of the wavefunction (i.e. a Slater determinant) the best set of coefficients is that for which the energy is a minimum, at which point

$$\frac{\partial E}{\partial c_{\nu i}} = 0 \quad (2.145)$$

for all coefficients $c_{\nu i}$. The objective is thus to determine the set of coefficients that gives the lowest energy for the system.

2.5.3 Closed-shell Systems and the Roothaan–Hall Equations

We shall initially consider a closed-shell system with N electrons in $N/2$ orbitals. The derivation of the Hartree–Fock equations for such a system was first proposed by Roothaan [Roothaan 1951] and (independently) by Hall [Hall 1951]. The resulting equations are known as the Roothaan equations or the Roothaan–Hall equations. Unlike the integro-differential form of the Hartree–Fock equations, Equation (2.124), Roothaan and Hall recast the equations in matrix form, which can be solved using standard techniques and can be applied to systems of any geometry. We shall identify the major steps in the Roothaan approach,

starting with the expression for the Hartree–Fock energy for our closed-shell system, Equation (2.120):

$$E = 2 \sum_{i=1}^{N/2} H_{ii}^{\text{core}} + \sum_{i=1}^{N/2} \sum_{j=1}^{N/2} (2J_{ij} - K_{ij}) \quad (2.146)$$

The corresponding Fock operator is (Equation (2.132)):

$$f_i(1) = \mathcal{H}^{\text{core}}(1) + \sum_{j=1}^{N/2} \{2\mathcal{J}_j(1) - \mathcal{K}_j(1)\} \quad (2.147)$$

We now introduce the atomic orbital expansion for the orbitals ψ_i and substitute for the corresponding spin orbital χ_i into the Hartree–Fock equation, $f_i(1)\chi_i(1) = \varepsilon_i\chi_i(1)$:

$$f_i(1) \sum_{\nu=1}^K c_{\nu i} \phi_{\nu}(1) = \varepsilon_i \sum_{\nu=1}^K c_{\nu i} \phi_{\nu}(1) \quad (2.148)$$

Pre-multiplying each side by $\phi_{\mu}(1)$ (where ϕ_{μ} is also a basis function) and integrating gives the following matrix equation:

$$\sum_{\nu=1}^K c_{\nu i} \int d\nu_1 \phi_{\mu}(1) f_i(1) \phi_{\nu}(1) = \varepsilon_i \sum_{\nu=1}^K c_{\nu i} \int d\nu_1 \phi_{\mu}(1) \phi_{\nu}(1) \quad (2.149)$$

$\int d\nu_1 \phi_{\mu}(1) \phi_{\nu}(1)$ is the overlap integral between the basis functions μ and ν , written $S_{\mu\nu}$. Unlike the molecular orbitals, which will be required to be orthonormal, the overlap between two basis functions is not necessarily zero (for example, they may be located on different atoms).

The elements of the *Fock matrix* are given by

$$F_{\mu\nu} = \int d\nu_1 \phi_{\mu}(1) f_i(1) \phi_{\nu}(1) \quad (2.150)$$

The Fock matrix elements for a closed-shell system can be expanded as follows by substituting the expression for the Fock operator:

$$F_{\mu\nu} = \int d\nu_1 \phi_{\mu}(1) \mathcal{H}^{\text{core}}(1) \phi_{\nu}(1) + \sum_{j=1}^{N/2} \int d\nu_1 \phi_{\mu}(1) [2\mathcal{J}_j(1) - \mathcal{K}_j(1)] \phi_{\nu}(1) \quad (2.151)$$

The elements of the Fock matrix can thus be written as the sum of core, Coulomb and exchange contributions. The core contribution is:

$$\int d\nu_1 \phi_{\mu}(1) \mathcal{H}^{\text{core}}(1) \phi_{\nu}(1) = \int d\nu_1 \phi_{\mu}(1) \left[-\frac{1}{2} \nabla^2 - \sum_{A=1}^M \frac{Z_A}{|r_1 - R_A|} \right] \phi_{\nu}(1) \equiv H_{\mu\nu}^{\text{core}} \quad (2.152)$$

The core contributions thus require the calculation of integrals that involve basis functions on up to two centres (depending upon whether ϕ_{μ} and ϕ_{ν} are centred on the same nucleus or not). Each element $H_{\mu\nu}^{\text{core}}$ can in turn be obtained as the sum of a kinetic energy integral and a potential energy integral corresponding to the two terms in the one-electron Hamiltonian.

The Coulomb and exchange contributions to the Fock matrix element $F_{\mu\nu}$ are together given by:

$$\sum_{j=1}^{N/2} \int d\nu_1 \phi_\mu(1) [2\mathcal{J}_j(1) - \mathcal{K}_j(1)] \phi_\nu(1) \quad (2.153)$$

Recall that the Coulomb operator $\mathcal{J}_j(1)$ due to interaction with a spin orbital χ_j is given by:

$$\mathcal{J}_j(1) = \int d\tau_2 \chi_j(2) \frac{1}{r_{12}} \chi_j(2) \quad (2.154)$$

We need to write each of the two occurrences of the spin orbital χ_j in this integral in terms of the appropriate linear combination of basis functions:

$$\mathcal{J}_j(1) = \int d\tau_2 \sum_{\sigma=1}^K c_{\sigma j} \phi_\sigma(2) \frac{1}{r_{12}} \sum_{\lambda=1}^K c_{\lambda j} \phi_\lambda(2) \quad (2.155)$$

We have used the indices σ and λ for the basis functions here. Similarly, the exchange contribution can be written:

$$\mathcal{K}_j(1) \chi_i(1) = \left[\int d\tau_2 \sum_{\sigma=1}^K c_{\sigma j} \phi_\sigma(2) \frac{1}{r_{12}} \chi_i(2) \right] \sum_{\lambda=1}^K c_{\lambda j} \phi_\lambda(2) \quad (2.156)$$

When the Coulomb and exchange operators are expressed in terms of the basis functions and the orbital expansion is substituted for χ_i , then their contributions to the Fock matrix element $F_{\mu\nu}$ take the following form:

$$\begin{aligned} & \sum_{j=1}^{N/2} \int d\nu_1 \phi_\mu(1) [2\mathcal{J}_j(1) - \mathcal{K}_j(1)] \phi_\nu(1) \\ &= \sum_{j=1}^{N/2} \sum_{\lambda=1}^K \sum_{\sigma=1}^K c_{\lambda j} c_{\sigma j} \left[\begin{aligned} & 2 \int d\nu_1 d\nu_2 \phi_\mu(1) \phi_\nu(1) \frac{1}{r_{12}} \phi_\lambda(2) \phi_\sigma(2) \\ & - \int d\nu_1 d\nu_2 \phi_\mu(1) \phi_\lambda(1) \frac{1}{r_{12}} \phi_\nu(2) \phi_\sigma(2) \end{aligned} \right] \\ &\equiv \sum_{j=1}^{N/2} \sum_{\lambda=1}^K \sum_{\sigma=1}^K c_{\lambda j} c_{\sigma j} [2(\mu\nu|\lambda\sigma) - (\mu\lambda|\nu\sigma)] \end{aligned} \quad (2.157)$$

We have used the shorthand notation for the integrals in the final expression. Note that the two-electron integrals may involve up to four different basis functions ($\mu, \nu, \lambda, \sigma$), which may in turn be located at four different centres. This has important consequences for the way in which we try to solve the equations.

It is helpful to simplify Equation (2.157) by introducing the *charge density matrix*, \mathbf{P} , whose elements are defined as:

$$P_{\mu\nu} = 2 \sum_{i=1}^{N/2} c_{\mu i} c_{\nu i} \quad \text{and} \quad P_{\lambda\sigma} = 2 \sum_{i=1}^{N/2} c_{\lambda i} c_{\sigma i} \quad (2.158)$$

Note that the summations are over the $N/2$ occupied orbitals. Other properties can be calculated from the density matrix; for example, the electronic energy is:

$$E = \frac{1}{2} \sum_{\mu=1}^K \sum_{\nu=1}^K P_{\mu\nu} (H_{\mu\nu}^{\text{core}} + F_{\mu\nu}) \quad (2.159)$$

The electron density at a point \mathbf{r} can also be expressed in terms of the density matrix:

$$\rho(\mathbf{r}) = \sum_{\mu=1}^K \sum_{\nu=1}^K P_{\mu\nu} \phi_\mu(\mathbf{r}) \phi_\nu(\mathbf{r}) \quad (2.160)$$

The expression for each element $F_{\mu\nu}$ of the Fock matrix elements for a closed-shell system of N electrons then becomes:

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} + \sum_{\lambda=1}^K \sum_{\sigma=1}^K P_{\lambda\sigma} [(\mu\nu|\lambda\sigma) - \frac{1}{2}(\mu\lambda|\nu\sigma)] \quad (2.161)$$

This is the standard form for the expression for the Fock matrix in the Roothaan–Hall equations.

2.5.4 Solving the Roothaan–Hall Equations

The Fock matrix is a $K \times K$ square matrix that is symmetric if real basis functions are used. The Roothaan–Hall equations (2.149) can be conveniently written as a matrix equation:

$$\mathbf{FC} = \mathbf{SCE} \quad (2.162)$$

The elements of the $K \times K$ matrix \mathbf{C} are the coefficients $c_{\nu i}$:

$$\mathbf{C} = \begin{pmatrix} c_{1,1} & c_{1,2} & \dots & c_{1,K} \\ c_{2,1} & c_{2,2} & \dots & c_{2,K} \\ \vdots & \vdots & & \vdots \\ c_{K,1} & c_{K,2} & \dots & c_{K,K} \end{pmatrix} \quad (2.163)$$

\mathbf{E} is a diagonal matrix whose elements are the orbital energies:

$$\mathbf{E} = \begin{pmatrix} \varepsilon_1 & 0 & \dots & 0 \\ 0 & \varepsilon_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \varepsilon_K \end{pmatrix} \quad (2.164)$$

Let us consider how we might solve the Roothaan–Hall equations and thereby obtain the molecular orbitals. The first point we must note is that the elements of the Fock matrix, which appear on the left-hand side of Equation (2.162), depend on the molecular orbital coefficients $c_{\nu i}$, which also appear on the right-hand side of the equation. Thus an iterative procedure is required to find a solution.

The one-electron contributions $H_{\mu\nu}^{\text{core}}$ due to the electrons moving in the field of the bare nuclei do not depend on the basis set coefficients and remain unchanged throughout the calculation. However, the Coulomb and exchange contributions do depend on the coefficients and we would expect these to vary throughout the calculation. The individual two-electron integrals $(\mu\nu|\lambda\sigma)$ are, however, constant throughout the calculation. An obvious strategy is thus to calculate and store these integrals for later use.

Having written the Roothaan–Hall equations in matrix form we would obviously like to solve them using standard matrix eigenvalue methods (discussed in Section 1.10.3). However, standard eigenvalue methods would require an equation of the form $\mathbf{FC} = \mathbf{CE}$. The Roothaan–Hall equations only adopt such a form if the overlap matrix, \mathbf{S} , is equal to the unit matrix, \mathbf{I} (in which all diagonal elements are equal to 1 and all off-diagonal elements are zero). The functions ϕ are usually normalised but they are not necessarily orthogonal (for example, because they are located on different atoms) and so there will invariably be non-zero off-diagonal elements of the overlap matrix. To solve the Roothaan–Hall equations using standard methods they must be transformed. This corresponds to transforming the basis functions so that they form an orthonormal set. We seek a matrix \mathbf{X} , such that $\mathbf{X}^T \mathbf{S} \mathbf{X} = \mathbf{I}$. \mathbf{X}^T is the transpose of \mathbf{X} , obtained by interchanging rows and columns. There are various ways in which \mathbf{X} can be calculated; in *symmetric orthogonalisation*, the overlap matrix is diagonalised. Diagonalisation involves finding the matrix \mathbf{U} such that

$$\mathbf{U}^T \mathbf{S} \mathbf{U} = \mathbf{D} = \text{diag}(\lambda_1 \dots \lambda_K) \quad (2.165)$$

\mathbf{D} is the diagonal matrix containing the eigenvalues λ_i of \mathbf{S} , and \mathbf{U} contains the eigenvectors of \mathbf{S} . \mathbf{U}^T is the transpose of the matrix \mathbf{U} . (This expression is often written $\mathbf{U}^{-1} \mathbf{S} \mathbf{U} = \mathbf{D}$ since for real basis functions $\mathbf{U}^{-1} = \mathbf{U}^T$.) Then the matrix \mathbf{X} is given by $\mathbf{X} = \mathbf{U} \mathbf{D}^{-1/2} \mathbf{U}^T$, where $\mathbf{D}^{-1/2}$ is formed from the inverse square roots of \mathbf{D} . We shall write \mathbf{X} as $\mathbf{S}^{-1/2}$, as it can be considered to be the inverse square root of the overlap matrix: $\mathbf{S}^{-1/2} \mathbf{S} \mathbf{S}^{-1/2} = \mathbf{I}$.

The Roothaan–Hall equations can now be manipulated as follows. Both sides of Equation (2.162) are pre-multiplied by the matrix $\mathbf{S}^{-1/2}$:

$$\mathbf{S}^{-1/2} \mathbf{F} \mathbf{C} = \mathbf{S}^{-1/2} \mathbf{S} \mathbf{C} \mathbf{E} = \mathbf{S}^{1/2} \mathbf{C} \mathbf{E} \quad (2.166)$$

Inserting the unit matrix, in the form $\mathbf{S}^{-1/2} \mathbf{S}^{1/2}$, into the left-hand side gives:

$$\mathbf{S}^{-1/2} \mathbf{F} (\mathbf{S}^{-1/2} \mathbf{S}^{1/2}) \mathbf{C} = \mathbf{S}^{1/2} \mathbf{C} \mathbf{E} \quad (2.167)$$

or

$$\mathbf{S}^{-1/2} \mathbf{F} \mathbf{S}^{-1/2} (\mathbf{S}^{1/2} \mathbf{C}) = (\mathbf{S}^{1/2} \mathbf{C}) \mathbf{E} \quad (2.168)$$

Equation (2.168) can be written $\mathbf{F}' \mathbf{C}' = \mathbf{C}' \mathbf{E}$, where $\mathbf{F}' = \mathbf{S}^{-1/2} \mathbf{F} \mathbf{S}^{-1/2}$ and $\mathbf{C}' = \mathbf{S}^{1/2} \mathbf{C}$.

The matrix equation $\mathbf{F}' \mathbf{C}' = \mathbf{C}' \mathbf{E}$ can be solved using standard methods; a solution only exists if the determinant $|\mathbf{F}' - \mathbf{E} \mathbf{I}|$ equals zero. In simple cases this can be done by multiplying out the determinant to give a polynomial (the secular equation) whose roots are the eigenvalues ϵ_i , but for large matrices a much more practical approach involves the diagonalisation of \mathbf{F}' . The matrix of coefficients, \mathbf{C}' , are the eigenvectors of \mathbf{F}' . The basis function coefficients \mathbf{C} can then be obtained from \mathbf{C}' using $\mathbf{C} = \mathbf{S}^{-1/2} \mathbf{C}'$. A common scheme for

solving the Roothaan–Hall equations is thus as follows:

1. Calculate the integrals to form the Fock matrix, \mathbf{F} .
2. Calculate the overlap matrix, \mathbf{S} .
3. Diagonalise \mathbf{S} .
4. Form $\mathbf{S}^{-1/2}$.
5. Guess, or otherwise calculate, an initial density matrix, \mathbf{P} .
6. Form the Fock matrix using the integrals and the density matrix \mathbf{P} .
7. Form $\mathbf{F}' = \mathbf{S}^{-1/2} \mathbf{F} \mathbf{S}^{-1/2}$.
8. Solve the secular equation $|\mathbf{F}' - \mathbf{E} \mathbf{I}| = 0$ to give the eigenvalues \mathbf{E} and the eigenvectors \mathbf{C}' by diagonalising \mathbf{F}' .
9. Calculate the molecular orbital coefficients, \mathbf{C} from $\mathbf{C} = \mathbf{S}^{-1/2} \mathbf{C}'$.
10. Calculate a new density matrix, \mathbf{P} , from the matrix \mathbf{C} .
11. Check for convergence. If the calculation has converged, stop. Otherwise repeat from step 6 using the new density matrix, \mathbf{P} .

This procedure requires an initial guess of the density matrix, \mathbf{P} . The simplest approach is to use the null matrix, which corresponds to ignoring all the electron–electron terms so that the electrons just experience the bare nuclei. This can sometimes lead to convergence problems, which may be prevented if a lower level of theory (such as semi-empirical or extended Hückel) is used to provide the initial guess. Moreover, a better guess may enable the calculation to be performed more quickly. A variety of criteria can be used to establish whether the calculation has converged or not. For example, the density matrix can be compared with that from the previous iteration, and/or the change in energy can be monitored together with the basis set coefficients.

The result of a Hartree–Fock calculation is a set of K molecular orbitals, where K is the number of basis functions in the calculation. The N electrons are then fed into these orbitals in accordance with the Aufbau principle, two electrons per orbital, starting with the lowest-energy orbitals. The remaining orbitals do not contain any electrons; these are known as the *virtual orbitals*. Alternative electronic configurations can be generated by exciting electrons from the occupied orbitals to the virtual orbitals; these excited configurations are used in more advanced calculations that will be discussed in Chapter 3.

A Hartree–Fock calculation provides a set of orbital energies, ϵ_i . What is the significance of these? The energy of an electron in a spin orbital is calculated by adding the core interaction H_{ii}^{core} to the Coulomb and exchange interactions with the other electrons in the system:

$$\epsilon_i = H_{ii}^{\text{core}} + \sum_{j=1}^{N/2} (2J_{ij} - K_{ij}) \quad (2.169)$$

The total electronic energy of the ground state is given by Equation (2.120):

$$E = 2 \sum_{i=1}^{N/2} H_{ii}^{\text{core}} + \sum_{i=1}^{N/2} \sum_{j=1}^{N/2} (2J_{ij} - K_{ij}) \quad (2.170)$$

The total energy is therefore not equal to the sum of the individual orbital energies but is related as follows:

$$E = \sum_{i=1}^N \varepsilon_i - \sum_{i=1}^{N/2} \sum_{j=1}^{N/2} (2J_{ij} - K_{ij}) \quad (2.171)$$

The reason for the discrepancy is that the individual orbital energies include contributions from the interaction between that electron and all the nuclei and all other electrons in the system. The Coulomb and exchange interactions between pairs of electrons are therefore counted twice when summing the individual orbital energies.

2.5.5 A Simple Illustration of the Roothaan–Hall Approach

We will illustrate the stages involved in the Roothaan–Hall approach using the helium hydrogen molecular ion, HeH^+ , as an example. This is a two-electron system. Our objective here is to show how the Roothaan–Hall method can be used to derive the wavefunction, for a fixed internuclear distance of 1 Å. We use HeH^+ rather than H_2 as our system as the lack of symmetry in HeH^+ makes the procedure more informative. There are two basis functions, $1s_A$ (centred on the helium atom) and $1s_B$ (on the hydrogen). The numerical values of the integrals that we shall use in our calculation were obtained using a Gaussian series approximation to the Slater orbitals (the STO-3G basis set, which is described in Section 2.6). This detail need not concern us here. Each wavefunction is expressed as a linear combination of the two $1s$ atomic orbitals centred on the nuclei A and B:

$$\psi_1 = c_{1A} 1s_A + c_{1B} 1s_B \quad (2.172)$$

$$\psi_2 = c_{2A} 1s_A + c_{2B} 1s_B \quad (2.173)$$

First, it is necessary to calculate the various one- and two-electron integrals and to formulate the Fock and overlap matrices, each of which will be a 2×2 symmetric matrix (as there are two orbitals in the basis set). The diagonal elements of the overlap matrix, \mathbf{S} , are equal to 1.0 as each basis function is normalised; the off-diagonal elements have smaller, but non-zero, values that are equal to the overlap between $1s_A$ and $1s_B$ for the internuclear distance chosen. The matrix \mathbf{S} is:

$$\mathbf{S} = \begin{pmatrix} 1.0 & 0.392 \\ 0.392 & 1.0 \end{pmatrix} \quad (2.174)$$

The core contributions $H_{\mu\nu}^{\text{core}}$ can be calculated as the sum of three 2×2 matrices comprising the kinetic energy (\mathbf{T}) and nuclear attraction terms for the two nuclei A and B (\mathbf{V}_A and \mathbf{V}_B). The elements of these three matrices are obtained by evaluating the following integrals:

$$\begin{aligned} \mathbf{T}_{\mu\nu} &= \int d\nu_1 \phi_\mu(1) \left(-\frac{1}{2} \nabla^2\right) \phi_\nu(1) \\ \mathbf{V}_{A,\mu\nu} &= \int d\nu_1 \phi_\mu(1) \left(-\frac{Z_A}{r_{1A}}\right) \phi_\nu(1) \\ \mathbf{V}_{B,\mu\nu} &= \int d\nu_1 \phi_\mu(1) \left(-\frac{Z_B}{r_{1B}}\right) \phi_\nu(1) \end{aligned} \quad (2.175)$$

The matrices are:

$$\mathbf{T} = \begin{pmatrix} 1.412 & 0.081 \\ 0.081 & 0.760 \end{pmatrix} \quad \mathbf{V}_A = \begin{pmatrix} -3.344 & -0.758 \\ -0.758 & -1.026 \end{pmatrix} \quad \mathbf{V}_B = \begin{pmatrix} -0.525 & -0.308 \\ -0.308 & -1.227 \end{pmatrix} \quad (2.176)$$

\mathbf{H}^{core} is the sum of these three:

$$\mathbf{H}^{\text{core}} = \begin{pmatrix} -2.457 & -0.985 \\ -0.985 & -1.493 \end{pmatrix} \quad (2.177)$$

As far as the two-electron integrals are concerned, with two basis functions there are a total of 16 possible two-electron integrals. There are however only six unique two-electron integrals, as the indices can be permuted as follows:

- (i) $(1s_A 1s_A | 1s_A 1s_A) = 1.056$
- (ii) $(1s_A 1s_A | 1s_A 1s_B) = (1s_A 1s_A | 1s_B 1s_A) = (1s_A 1s_B | 1s_A 1s_A) = (1s_B 1s_A | 1s_A 1s_A) = 0.303$
- (iii) $(1s_A 1s_B | 1s_A 1s_B) = (1s_A 1s_B | 1s_B 1s_A) = (1s_B 1s_A | 1s_A 1s_B) = (1s_B 1s_A | 1s_B 1s_A) = 0.112$
- (iv) $(1s_A 1s_A | 1s_B 1s_B) = (1s_B 1s_B | 1s_A 1s_A) = 0.496$
- (v) $(1s_A 1s_B | 1s_B 1s_B) = (1s_B 1s_A | 1s_B 1s_B) = (1s_B 1s_B | 1s_A 1s_B) = (1s_B 1s_B | 1s_B 1s_A) = 0.244$
- (vi) $(1s_B 1s_B | 1s_B 1s_B) = 0.775$

To reiterate, these integrals are calculated as follows:

$$(\mu\nu|\lambda\sigma) = \iint d\nu_1 d\nu_2 \phi_\mu(1) \phi_\nu(1) \frac{1}{r_{12}} \phi_\lambda(2) \phi_\sigma(2) \quad (2.178)$$

Having calculated the integrals, we are now ready to start the SCF calculation. To formulate the Fock matrix it is necessary to have an initial guess of the density matrix, \mathbf{P} . The simplest approach is to use the null matrix in which all elements are zero. In this initial step the Fock matrix \mathbf{F} is therefore equal to \mathbf{H}^{core} .

The Fock matrix must next be transformed to \mathbf{F}' by pre- and post-multiplying by $\mathbf{S}^{-1/2}$:

$$\mathbf{S}^{-1/2} = \begin{pmatrix} -1.065 & -0.217 \\ -0.217 & 1.065 \end{pmatrix} \quad (2.179)$$

\mathbf{F}' for this first iteration is thus:

$$\mathbf{F}' = \begin{pmatrix} -2.401 & -0.249 \\ -0.249 & -1.353 \end{pmatrix} \quad (2.180)$$

Diagonalisation of \mathbf{F}' gives its eigenvalues and eigenvectors, which are:

$$\mathbf{E} = \begin{pmatrix} -2.458 & 0.0 \\ 0.0 & -1.292 \end{pmatrix} \quad \mathbf{C}' = \begin{pmatrix} 0.975 & -0.220 \\ 0.220 & 0.975 \end{pmatrix} \quad (2.181)$$

The coefficients \mathbf{C} are obtained from $\mathbf{C} = \mathbf{S}^{-1/2} \mathbf{C}'$ and are thus:

$$\mathbf{C} = \begin{pmatrix} 0.991 & -0.446 \\ 0.022 & 1.087 \end{pmatrix} \quad (2.182)$$

To formulate the density matrix P we need to identify the occupied orbital(s). With a two-electron system both electrons occupy the orbital with the lowest energy (i.e. the orbital with the lowest eigenvalue). At this stage the lowest-energy orbital is:

$$\psi = 0.991 1s_A + 0.022 1s_B \quad (2.183)$$

The orbital is composed largely of the s orbital on the helium nucleus; in the absence of any electron–electron repulsion the electrons tend to congregate near the nucleus with the larger charge. The density matrix corresponding to this initial wavefunction is:

$$P = \begin{pmatrix} 1.964 & 0.044 \\ 0.044 & 0.001 \end{pmatrix} \quad (2.184)$$

The new Fock matrix is formed using P and the two-electron integrals together with H^{core} . For example, the element F_{11} is given by:

$$\begin{aligned} F_{11} = & H_{11}^{\text{core}} + P_{11}[(1s_A 1s_A | 1s_A 1s_A) - \frac{1}{2}(1s_A 1s_A | 1s_A 1s_A)] \\ & + P_{12}[(1s_A 1s_A | 1s_A 1s_B) - \frac{1}{2}(1s_A 1s_A | 1s_A 1s_B)] \\ & + P_{21}[(1s_A 1s_A | 1s_B 1s_B) - \frac{1}{2}(1s_A 1s_B | 1s_A 1s_B)] \\ & + P_{12}[(1s_A 1s_A | 1s_B 1s_B) - \frac{1}{2}(1s_A 1s_B | 1s_A 1s_B)] \end{aligned} \quad (2.185)$$

The complete Fock matrix is:

$$F = \begin{pmatrix} -1.406 & -0.690 \\ -0.690 & -0.618 \end{pmatrix} \quad (2.186)$$

The energy that corresponds to this Fock matrix (calculated using Equation (2.159)) is -3.870 Hartree. In the next iteration, the various matrices are as follows:

$$\begin{aligned} F' &= \begin{pmatrix} -1.305 & -0.347 \\ -0.347 & -0.448 \end{pmatrix} & E &= \begin{pmatrix} -1.427 & 0.0 \\ 0.0 & -0.325 \end{pmatrix} \\ C' &= \begin{pmatrix} 0.943 & -0.334 \\ 0.334 & 0.943 \end{pmatrix} & C &= \begin{pmatrix} 0.931 & -0.560 \\ 0.150 & 1.076 \end{pmatrix} \\ P &= \begin{pmatrix} 1.735 & 0.280 \\ 0.280 & 0.045 \end{pmatrix} & F &= \begin{pmatrix} -1.436 & -0.738 \\ -0.738 & -0.644 \end{pmatrix} \end{aligned} \quad (2.187)$$

Energy = -3.909 Hartree

The calculation proceeds as illustrated in Table 2.2, which shows the variation in the coefficients of the atomic orbitals in the lowest-energy wavefunction and the energy for the first four SCF iterations. The energy is converged to six decimal places after six iterations and the charge density matrix after nine iterations.

The final wavefunction still contains a large proportion of the $1s$ orbital on the helium atom, but less than was obtained without the two-electron integrals.

Iteration	$c(1s_A)$	$c(1s_B)$	Energy
1	0.991	0.022	-3.870
2	0.931	0.150	-3.909
3	0.915	0.181	-3.911
4	0.912	0.187	-3.911

Table 2.2 Variation in basis set coefficients and electronic energy for the HeH^+ molecule.

2.5.6 Application of the Hartree–Fock Equations to Molecular Systems

We are now in a position to consider how the Hartree–Fock theory we have developed can be used to perform practical quantum mechanical calculations on molecular systems. This is an appropriate place in our discussion to distinguish the two major categories of quantum mechanical molecular orbital calculations: the *ab initio* and the semi-empirical methods. *Ab initio* strictly means ‘from the beginning’, or ‘from first principles’, which would imply that a calculation using such an approach would require as input only physical constants such as the speed of light, Planck’s constant, the masses of elementary particles, and so on. *Ab initio* in fact usually refers to a calculation which uses the full Hartree–Fock/Roothaan–Hall equations, without ignoring or approximating any of the integrals or any of the terms in the Hamiltonian. The *ab initio* methods do rely upon calibration calculations, and this has led some quantum chemists, notably Dewar (who has played a large part in the development of semi-empirical methods), to claim that any real difference between the *ab initio* and semi-empirical methods is entirely pedagogical. By contrast, semi-empirical methods simplify the calculations, using parameters for some of the integrals and/or ignoring some of the terms in the Hamiltonian. First we shall consider *ab initio* methods.

2.6 Basis Sets

The basis sets most commonly used in quantum mechanical calculations are composed of atomic functions. An obvious choice would be the Slater type orbitals. Unfortunately, Slater functions are not particularly amenable to implementation in molecular orbital calculations. This is because some of the integrals are difficult, if not impossible, to evaluate, particularly when the atomic orbitals are centred on different nuclei. It is relatively straightforward to calculate integrals involving one or two centres, such as $(\mu\mu|\nu\nu)$, $(\mu\nu|\nu\nu)$ and $(\mu\nu|\mu\nu)$. Three- and four-centre integrals are also feasible with Slater functions if the atomic orbitals are located on the same atom. However, three- and four-centre integrals are very difficult if the atomic orbitals are based on different atoms. It is common in *ab initio* calculations to replace the Slater orbitals by functions based upon Gaussians. A Gaussian function has the form $\exp(-\alpha r^2)$, and *ab initio* calculations use basis functions comprising integral powers of x , y and z multiplied by $\exp(-\alpha r^2)$:

$$x^a y^b z^c \exp(-\alpha r^2) \quad (2.188)$$

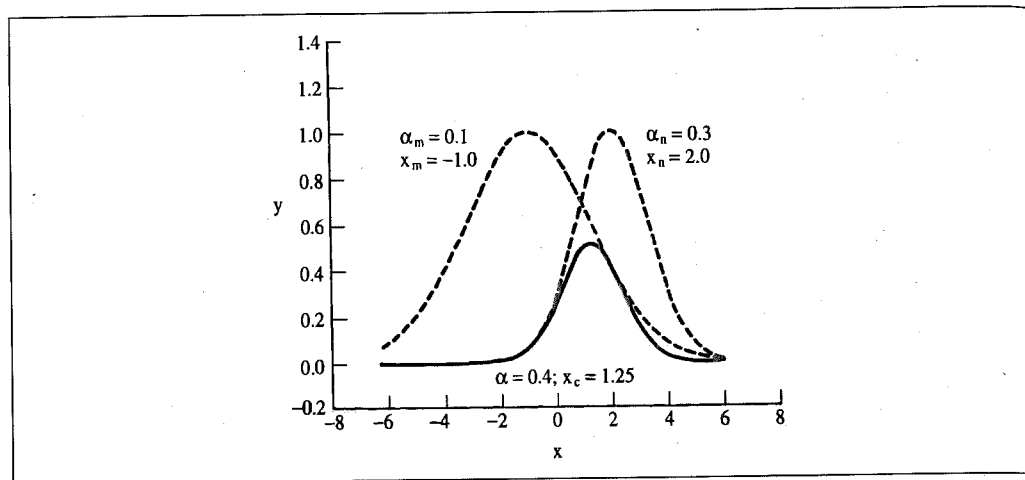


Fig. 2.4: The product of two Gaussian functions is another Gaussian centred along the line joining their centres. In this case the equations of the two functions are $y = \exp[-0.1(x + 1.0)^2]$ and $y = \exp[-0.3(x - 2.0)^2]$ and the equation of the product is $y = \exp(-27/40)[-0.4(x - 1.25)^2]$ (Equation (2.189)).

α determines the radial extent (or 'spread') of a Gaussian function; a function with a large value of α does not spread very far, whereas a small value of α gives a large spread. The order of these Gaussian-type functions is determined by the powers of the Cartesian variables; a zeroth-order function has $a + b + c = 0$; a first-order function has $a + b + c = 1$, and so on. There is thus one zeroth-order function, three first-order functions and six second-order functions. The idea of using Gaussian functions in quantum mechanical calculations is often ascribed to Boys [Boys 1950]. A major advantage of Gaussian functions is that the product of two Gaussians can be expressed as a single Gaussian, located along the line joining the centres of the two Gaussians m and n (Figure 2.4):

$$\exp(-\alpha_m r_m^2) \exp(-\alpha_n r_n^2) = \exp\left(-\frac{\alpha_m \alpha_n}{\alpha_m + \alpha_n} r_{mn}^2\right) \exp(-\alpha r_c^2) \quad (2.189)$$

r_{mn} is the distance between the centres m and n , and the orbital exponent α of the combined function is related to the exponents α_m and α_n by:

$$\alpha = \alpha_m + \alpha_n \quad (2.190)$$

r_c is the distance from point C, which has coordinates:

$$x_c = \frac{\alpha_m x_m + \alpha_n x_n}{\alpha_m + \alpha_n}; \quad y_c = \frac{\alpha_m y_m + \alpha_n y_n}{\alpha_m + \alpha_n}; \quad z_c = \frac{\alpha_m z_m + \alpha_n z_n}{\alpha_m + \alpha_n} \quad (2.191)$$

x_m, y_m, z_m and x_n, y_n, z_n are the centres of the two original Gaussians m and n respectively.

Thus, in a two-electron integral of the form $(\mu\nu|\lambda\sigma)$, the product $\phi_\mu(1)\phi_\nu(1)$ (where ϕ_μ and ϕ_ν may be on different centres) can be replaced by a single Gaussian function that is centred at the appropriate point C. For Cartesian Gaussian functions the calculation is more complicated than for the example we have stated above, due to the presence of the Cartesian functions, but even so, efficient methods for performing the integrals have been devised.

The zeroth-order Gaussian function g_s has s-orbital angular symmetry; the three first-order Gaussian functions have p-orbital symmetry. In normalised form these are:

$$g_s(\alpha, r) = \left(\frac{2\alpha}{\pi}\right)^{3/4} e^{-\alpha r^2} \quad (2.192)$$

$$g_x(\alpha, r) = \left(\frac{128\alpha^5}{\pi^3}\right)^{1/4} x e^{-\alpha r^2} \quad (2.193)$$

$$g_y(\alpha, r) = \left(\frac{128\alpha^5}{\pi^3}\right)^{1/4} y e^{-\alpha r^2} \quad (2.194)$$

$$g_z(\alpha, r) = \left(\frac{128\alpha^5}{\pi^3}\right)^{1/4} z e^{-\alpha r^2} \quad (2.195)$$

The six second-order functions have the following form, exemplified by two of the functions:

$$g_{xx}(\alpha, r) = \left(\frac{2048\alpha^7}{9\pi^3}\right)^{1/4} x^2 e^{-\alpha r^2} \quad (2.196)$$

$$g_{xy}(\alpha, r) = \left(\frac{2048\alpha^7}{9\pi^3}\right)^{1/4} xy e^{-\alpha r^2} \quad (2.197)$$

These second-order functions do not all have the same angular symmetry as the 3d atomic orbitals, but a set comprising g_{xy}, g_{xz} and g_{yz} , together with two linear combinations of the g_{xx}, g_{yy} and g_{zz} , does give the desired result:

$$g_{3zz-rr} = \frac{1}{2}(2g_{zz} - g_{xx} - g_{yy}) \quad (2.198)$$

$$g_{xx-yy} = \sqrt{\frac{3}{4}}(g_{xx} - g_{yy}) \quad (2.199)$$

The remaining sixth linear combination has the symmetry properties of an s function:

$$g_{rr} = \sqrt{5}(g_{xx} + g_{yy} + g_{zz}) \quad (2.200)$$

The advantages of Gaussian functions are countered by some serious shortcomings. This can be readily seen from a graphical comparison of the 1s Slater function and its 'best' Gaussian approximation, Figure 2.5. Unlike the Slater functions the Gaussian functions do not have a cusp at the origin and they also decay towards zero more quickly. It is found that replacing a Slater type orbital by a single Gaussian function leads to unacceptable errors. However, this problem can be overcome if each atomic orbital is represented as a linear combination of Gaussian functions. Each linear combination has the following form:

$$\phi_\mu = \sum_{i=1}^L d_{i\mu} \phi_i(\alpha_{i\mu}) \quad (2.201)$$

$d_{i\mu}$ is the coefficient of the primitive Gaussian function ϕ_i , which has exponent $\alpha_{i\mu}$. L is the number of functions in the expansion. For example, the linear combinations of Gaussian 1s functions that can be used to represent a 1s Slater type orbital with exponent $\xi = 1$ are given in Table 2.3.

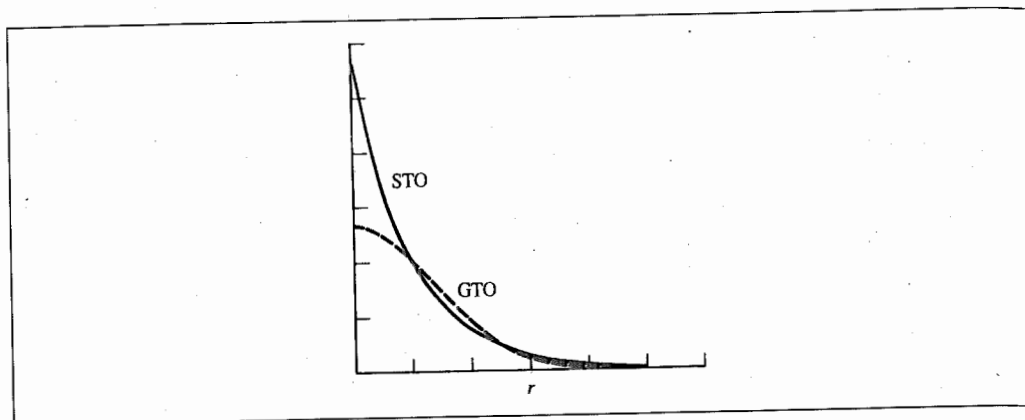


Fig. 2.5: The 1s Slater type orbital and the best Gaussian equivalent.

The coefficients and the exponents are found by least-squares fitting, in which the overlap between the Slater type function and the Gaussian expansion is maximised. Thus, for the 1s Slater type orbital we seek to maximise the following integral:

$$S = \frac{1}{\sqrt{\pi}} \left(\frac{2\alpha}{\pi} \right)^{3/4} \int dx e^{-r} e^{-\alpha r^2} \quad (2.202)$$

A graphical comparison of the 1s Slater type orbital and the four Gaussian expansions in Table 2.3 is shown in Figure 2.6. It is clear that the fit improves as the number of Gaussian functions increases, but even so, the addition of many more Gaussian functions cannot properly describe the exponential tail in the 'true' function and the cusp at the nucleus. This means that Gaussian functions underestimate the long-range overlap between atoms and the charge and spin density at the nucleus.

A Gaussian expansion contains two parameters: the coefficient and the exponent. The most flexible way to use Gaussian functions in *ab initio* molecular orbital calculations permits both of these parameters to vary during the calculation. Such a calculation is said to use

Number of Gaussians	Exponent, α	Expansion coefficient, d
1	0.270 950	1.00
2	0.151 623	0.678 914
	0.851 819	0.430 129
3	0.109 818	0.444 635
	0.405 771	0.535 28
	2.227 66	0.154 329
4	0.088 0187	0.291 626
	0.265 204	0.532 846
	0.954 620	0.260 141
	5.216 86	0.056 7523

Table 2.3 Coefficients and exponents for best-fit Gaussian expansions for the 1s Slater type orbital [Hehre et al. 1969].

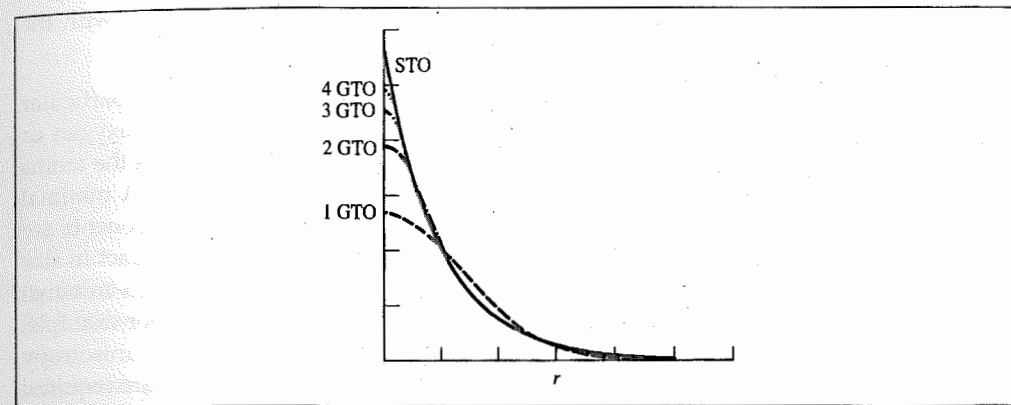


Fig. 2.6: Comparison of 1s Slater type orbital and Gaussian expansions with up to four terms.

uncontracted or *primitive* Gaussians. However, calculations with primitive Gaussians require a significant computational effort and so basis sets that consist of *contracted* Gaussian functions are most commonly employed. In a contracted function the contraction coefficients and exponents are pre-determined and remain constant during the calculation. The series of Gaussian functions in such cases is commonly referred to as a *contraction*, with the *contraction length* being the number of terms in the expansion. A further approximation that is often employed for the sake of computational efficiency is to use the same Gaussian exponents for the s and p orbitals in a given shell. This clearly restricts the flexibility of the basis set, but it does have the advantage of significantly reducing the number of numerically different integrals that need to be calculated.

Quantum chemists have devised efficient short-hand notation schemes to denote the basis set used in an *ab initio* calculation, although this does mean that a proliferation of abbreviations and acronyms are introduced. However, the codes are usually quite simple to understand. We shall concentrate on the notation used by Pople and co-workers in their Gaussian series of programs (see also the appendix to this chapter).

A *minimal basis set* is a representation that, strictly speaking, contains just the number of functions that are required to accommodate all the filled orbitals in each atom. In practice, a minimal basis set normally includes all of the atomic orbitals in the shell. Thus, for hydrogen and helium a single s-type function would be required; for the elements from lithium to neon the 1s, 2s and 2p functions are used, and so on. The basis sets STO-3G, STO-4G, etc. (in general, STO-*n*G) are all minimal basis sets in which *n* Gaussian functions are used to represent each orbital. It is found that at least three Gaussian functions are required to properly represent each Slater type orbital and so the STO-3G basis set is the 'absolute minimum' that should be used in an *ab initio* molecular orbital calculation. In fact, there is often little difference between the results obtained with the STO-3G basis set and the larger minimal basis sets with more Gaussian functions, although for hydrogen-bonded complexes STO-4G can perform significantly better. The STO-3G basis set does perform remarkably well in predicting molecular geometries, though this is due in part to

a fortuitous cancellation of errors. Of course, the computational effort increases with the number of functions in the Gaussian expansion.

The minimal basis sets are well known to have several deficiencies. There are particular problems with compounds containing atoms at the end of a period, such as oxygen or fluorine. Such atoms are described using the same number of basis functions as the atoms at the beginning of the period, despite the fact that they have more electrons. A minimal basis set only contains one contraction per atomic orbital and as the radial exponents are not allowed to vary during the calculation the functions cannot expand or contract in size in accordance with the molecular environment. The third drawback is that a minimal basis set cannot describe non-spherical aspects of the electronic distribution. For example, for a second-row element such as carbon the only functions that incorporate any anisotropy are the $2p_x$, $2p_y$ and $2p_z$ functions. As the radial components of these functions are required to be the same, no one component (x , y or z) can differ from another.

These problems with minimal basis sets can be addressed if more than one function is used for each orbital. A basis set which doubles the number of functions in the minimal basis set is described as a *double zeta* basis. Thus, a linear combination of a 'contracted' function and a 'diffuse' function gives an overall result that is intermediate between the two. The basis set coefficients of the contracted and the diffuse functions are automatically calculated by the SCF procedure, which thus automatically determines whether a more contracted or a more diffuse representation of that particular orbital is required. Such an approach can provide a solution to the anisotropy problem because it is then possible to have different linear combinations for the p_x , p_y and p_z orbitals.

An alternative to the double zeta basis approach is to double the number of functions used to describe the valence electrons but to keep a single function for the inner shells. The rationale for this approach is that the core orbitals, unlike the valence orbitals, do not affect chemical properties very much and vary only slightly from one molecule to another. The notation used for such *split valence* double zeta basis sets is exemplified by 3-21G. In this basis set three Gaussian functions are used to describe the core orbitals. The valence electrons are also represented by three Gaussians: the contracted part by two Gaussians and the diffuse part by one Gaussian. The most commonly used split valence basis sets are 3-21G, 4-31G and 6-31G.

Simply increasing the number of basis functions (triple zeta, quadruple zeta, etc.) does not necessarily improve the model. In fact, it can give rise to wholly erroneous results, particularly for molecules with a strongly anisotropic charge distribution. All of the basis sets we have encountered so far use functions that are centred on atomic nuclei. The use of split valence basis sets can help to surmount the problems with non-isotropic charge distribution but not completely. The charge distribution about an atom in a molecule is usually perturbed in comparison with the isolated atom. For example, the electron cloud in an isolated hydrogen atom is symmetrical, but when the hydrogen atom is present in a molecule the electrons are attracted towards the other nuclei. The distortion can be considered to correspond to mixing p-type character into the 1s orbital of the isolated atom to give a form of sp hybrid. In a similar manner, the unoccupied d orbitals introduce asymmetry into p orbitals (Figure 2.7). The most common solution to this problem is to introduce

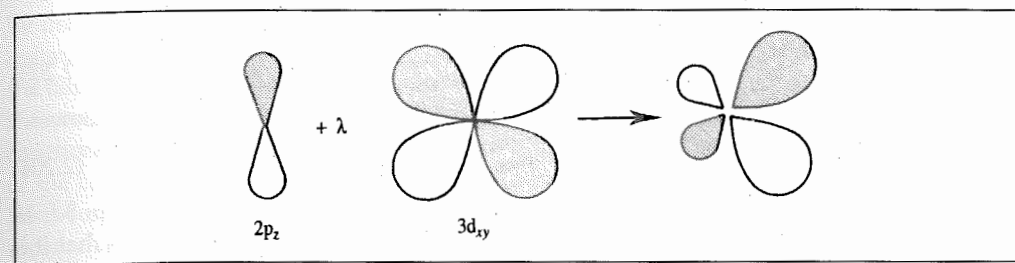


Fig. 2.7: The addition of a $3d_{xy}$ orbital to $2p_z$ gives a distorted orbital. (Figure adapted from Hehre W J, L Radom, P v R Schleyer and J A Hehre 1986. *Ab initio Molecular Orbital Theory*. New York, Wiley.)

polarisation functions into the basis set. The polarisation functions have a higher angular quantum number and so correspond to p orbitals for hydrogen and d orbitals for the first- and second-row elements.

The use of polarisation basis functions is indicated by an asterisk (*). Thus, 6-31G* refers to a 6-31G basis set with polarisation functions on the heavy (i.e. non-hydrogen) atoms. Two asterisks (e.g. 6-31G**) indicate the use of polarisation (i.e. p) functions on hydrogen and helium. The 6-31G** basis set is particularly useful where hydrogen acts as a bridging atom. Partial polarisation basis sets have also been developed. For example, the 3-21G(*) basis set has the same set of Gaussians as the 3-21G basis set (i.e. three functions for the inner shell, two contracted functions and one diffuse function for the valence shell) supplemented by six d-type Gaussians for the second-row elements. This basis set therefore attempts to account for d-orbital effects in molecules containing second-row elements. There are no special polarisation functions on first-row elements, which are described by the 3-21G basis set.

A deficiency of the basis sets described so far is their inability to deal with species such as anions and molecules containing lone pairs which have a significant amount of electron density away from the nuclear centres. This failure arises because the amplitudes of the Gaussian basis functions are rather low far from the nuclei. To remedy this deficiency highly diffuse functions can be added to the basis set. These basis sets are denoted using a '+'; thus the 3-21+G basis set contains an additional single set of diffuse s- and p-type Gaussian functions. '++' indicates that the diffuse functions are included for hydrogen as well as for heavy atoms. At these levels the terminology starts to become a little unwieldy. For example, the 6-311++G(3df, 3pd) basis set uses a single zeta core and triple zeta valence representation with additional diffuse functions on all atoms. The '(3df, 3pd)' indicates three sets of d functions and one set of f functions for first-row atoms and three sets of p functions and one set of d functions for hydrogen. This latter convention is probably the most generic; one commonly encountered example is the 6-31G(d) basis set, which is synonymous with 6-31G*.

The basis sets that we have considered thus far are sufficient for most calculations. However, for some high-level calculations a basis set that effectively enables the basis set limit to be achieved is required. The *even-tempered* basis set is designed to achieve this; each function in this basis set is the product of a spherical harmonic and a Gaussian function multiplied

by a power of the distance from the origin:

$$\chi_{klm}(\rho, \theta, \phi) = \exp(-\zeta_k^2 r) Y_{lm}(\theta, \phi) \quad (2.203)$$

The orbital exponent ζ_k is expressed as a function of two parameters α and β as follows:

$$\zeta_k = \alpha\beta^k \quad k = 1, 2, 3, \dots, N \quad (2.204)$$

The even-tempered basis set consists of the following sequence of functions: 1s, 2p, 3d, 4f, ..., which correspond to increasing values of k . The advantage of this basis set is that it is relatively easy to optimise the exponents for a large sequence of basis functions.

2.6.1 Creating a Basis Set

There is no definitive method for generating basis sets, and the construction of a new basis set is very much an art. Nevertheless, there are a number of well-established approaches that have resulted in widely used basis sets. We have already seen how linear combinations of Gaussian functions can be fitted to Slater type orbitals by minimising the overlap (see Figure 2.6 and Table 2.3). The Gaussian exponents and coefficients are derived by least-squares fitting to the desired functions, such as Slater type orbitals. When using basis sets that have been fitted to Slater orbitals it is often advantageous to use Slater exponents that are different to those obtained from Slater's rules. In general, better results for molecular calculations are obtained if larger Slater exponents are used for the valence electrons; this has the effect of giving a 'smaller', less diffuse orbital. For example, a value of 1.24 is widely used for the Slater exponent of hydrogen rather than the 1.0 that would be suggested by Slater's rules. It is straightforward to derive a basis set for a different Slater exponent if the Gaussian expansion has been fitted to a Slater type orbital with $\zeta = 1.0$. If the Slater exponent ζ is replaced by a new value, ζ' , then the respective Gaussian exponents α and α' are related by:

$$\frac{\alpha'}{\alpha} = \frac{\zeta'^2}{\zeta^2} \quad (2.205)$$

A doubling of the Slater exponent thus corresponds to a quadrupling of the Gaussian exponent. The expansion coefficients remain the same. For example, to obtain the exponents of the Gaussian functions for hydrogen in the STO-3G basis set we need to multiply the appropriate values in Table 2.3 by 1.24^2 , giving exponents of 0.168 856, 0.623 913 and 3.425 25. This strategy can be quite powerful; the STO- n G basis sets were originally defined with exponents that reproduce 'best atom' values for the core orbitals, but the exponents for the valence electrons were values that give optimal performance for a selected set of small molecules. For example, the suggested exponent for the valence orbitals in carbon was 1.72 rather than the 1.625 predicted by Slater's rules. The core orbitals have a Slater exponent of 5.67.

Basis sets can be constructed using an optimisation procedure in which the coefficients and the exponents are varied to give the lowest atomic energies. Some complications can arise when this approach is applied to larger basis sets. For example, in an atomic calculation the diffuse functions can move towards the nucleus, especially if the core region is described

by only a few basis functions. This is contrary to the role of diffuse functions, which is to enhance the description in the internuclear region. It may therefore be necessary to construct the basis set in stages, first determining the diffuse functions, using many basis functions for the core, and then optimising the basis functions for the core region, keeping the diffuse functions fixed. In many of the popular Gaussian basis sets the coefficients and exponents of the core orbitals are designed to reproduce calculations on atoms, whereas the valence basis functions are parametrised to reproduce the properties of a carefully selected set of molecular data.

The basis sets of Dunning [Dunning 1970] are obtained in a rather different way to those of Pople and co-workers. The first step is to perform an atomic SCF calculation using a set of primitive Gaussian functions in which the exponents are optimised to give the lowest energy for the atom. This set of primitive Gaussian functions (usually far too many for general use in molecular calculations) is then contracted to a smaller number of Gaussian functions, so drastically reducing the number of integrals that need be calculated. For example, Huzinga optimised the exponents of an uncontracted basis set that contained nine functions of s symmetry and five functions of p symmetry for the first-row elements [Huzinga 1965]. This (9s5p) basis set represents the 1s, 2s and three 2p orbitals and in fact corresponds to 24 basis functions per atom ($9 + 3 \times 5$). The primitive Gaussians in this uncontracted basis set are then apportioned to the basis functions in the new, contracted basis set, which contains three s functions and two p functions and is written [3s2p]. No primitive is assigned to more than one of the contracted basis functions. The 1s orbital is constructed from six primitives, the 2s orbital from one set of two primitives and one set containing just one primitive, and the 2p orbitals are represented by one contracted function containing five primitives and one contracted function that contains the remaining primitive. The final basis set, which is illustrated in Table 2.4 for nitrogen, contains a total of nine basis functions rather than the original 24. Each of the primitive functions appears

Exponent	Coefficient	Exponent	Coefficient	Exponent	Coefficient
1s		2s		2s	
5900	0.001 190	7.193	-0.160 405	0.2133	1.000 000
887.5	0.009 099	1.707	1.058 215		
204.7	0.044 145				
59.84	0.150 464				
20.00	0.356 741				
7.193	0.446 533				
2.686	0.145 603				
2p		2p			
26.79	0.018 254	0.1654	1.000 000		
5.956	0.116 461				
1.707	0.390 178				
0.5314	0.637 102				

Table 2.4 Exponents and contraction coefficients for the three s -type and the two p -type Gaussian functions in the basis set of Dunning for nitrogen [Dunning 1970].

in just one basis function with its original exponent. The ratios of the coefficients of the primitives in the contracted basis set are equal to the ratios of the coefficients determined in the atomic SCF calculation. The major advantage of this approach is that calculations with the smaller basis set give results that are almost as good as calculations using the full basis set but with much less computational effort.

2.7 Calculating Molecular Properties Using *ab initio* Quantum Mechanics

We have now considered the key features of the *ab initio* approach to quantum mechanical calculations and so, as an antidote to the rather theoretical nature of the chapter so far, it is appropriate to consider how the method might be used in practice. Quantum mechanics can be used to calculate a wide range of properties. In addition to thermodynamic and structural values, quantum mechanics can be used to derive properties dependent upon the electronic distribution. Such properties often cannot be determined by any other method. In this section we shall provide a flavour of the ways in which quantum mechanics is used in molecular modelling. Other applications, such as the location of transition structures and the use of quantum mechanics in deriving force field parameters, will be discussed in later chapters. Many different computer programs are now available for performing *ab initio* calculations; probably the best known of these is the Gaussian series of programs which originated in the laboratory of John Pople, who has made numerous contributions to the field, recognised by the award of the Nobel Prize in 1998.

2.7.1 Setting Up the Calculation and the Choice of Coordinates

The traditional way to provide the nuclear coordinates to a quantum mechanical program is via a Z-matrix, in which the positions of the nuclei are defined in terms of a set of internal coordinates (see Section 1.2). Some programs also accept coordinates in Cartesian format, which can be more convenient for large systems. It can sometimes be important to choose an appropriate set of internal coordinates, especially when locating minima or transition points or when following reaction pathways. This is discussed in more detail in Section 5.7.

2.7.2 Energies, Koopman's Theorem and Ionisation Potentials

The energy of an electron in an orbital (Equation (2.169)) is often equated with the energy required to remove the electron to give the corresponding ion. This is *Koopman's theorem*. Two important caveats must be remembered when applying Koopman's theorem and comparing the results with experimentally determined ionisation potentials. The first of these is that the orbitals in the ionised state are assumed to be the same as in the unionised state; they are 'frozen'. This neglects the fact that the orbitals in the ionised state will be different from those in the unionised state. The energy of the ionised state will thus tend to be higher than it 'should' be, giving too large an ionisation potential. The second caveat is that the Hartree-Fock method does not include the effects of electron correlation.

The correction due to electron correlation would be expected to be greater for the unionised state than for the ionised state, as the former has more electrons. Fortunately, therefore, the effect of electron correlation often opposes the effect of the frozen orbitals, resulting in many cases in good agreement between experimentally determined ionisation potentials and calculated values.

A Hartree-Fock SCF calculation with K basis functions provides K molecular orbitals, but many of these will not be occupied by any electrons; they are the 'virtual' spin orbitals. If we were to add an electron to one of these virtual orbitals then this should provide a means of calculating the electron affinity of the system. Electron affinities predicted by Koopman's theorem are always positive when Hartree-Fock calculations are used, because the virtual orbitals always have a positive energy. However, it is observed experimentally that many neutral molecules will accept an electron to form a stable anion and so have negative electron affinities. This can be understood if one realises that electron correlation would be expected to add to the error due to the 'frozen' orbital approximation, rather than to counteract it as for ionisation potentials.

2.7.3 Calculation of Electric Multipoles

Some of the most important properties that a quantum mechanical calculation provides are the electric multipole moments of the molecule. The electric multipoles reflect the distribution of charge in a molecule. The simplest electric moment (apart from the total net charge on the molecule) is the dipole. The dipole moment of a distribution of charges q_i located at positions \mathbf{r}_i is given by $\sum q_i \mathbf{r}_i$. If there are just two charges $+q$ and $-q$ separated by a distance r then the dipole moment is qr . A dipole moment of 4.8 Debye corresponds to two charges equal in magnitude to the electronic charge e separated by 1 Å. The dipole moment is a vector quantity, with components along the three Cartesian axes. The dipole moment of a molecule has contributions from both the nuclei and the electrons. The nuclear contributions can be calculated using the formula for a system of discrete charges:

$$\mu_{\text{nuclear}} = \sum_{A=1}^M Z_A \mathbf{R}_A \quad (2.206)$$

The electronic contribution arises from a continuous function of electron density and must be calculated using the appropriate operator:

$$\mu_{\text{electronic}} = \int d\tau \Psi_0 \left(\sum_{i=1}^N -\mathbf{r}_i \right) \Psi_0 \quad (2.207)$$

The dipole moment operator is a sum of one-electron operators \mathbf{r}_i , and as such the electronic contribution to the dipole moment can be written as a sum of one-electron contributions. The electronic contribution can also be written in terms of the density matrix, \mathbf{P} , as follows:

$$\mu_{\text{electronic}} = \sum_{\mu=1}^K \sum_{\nu=1}^K P_{\mu\nu} \int d\tau \phi_{\mu}(-\mathbf{r}) \phi_{\nu} \quad (2.208)$$

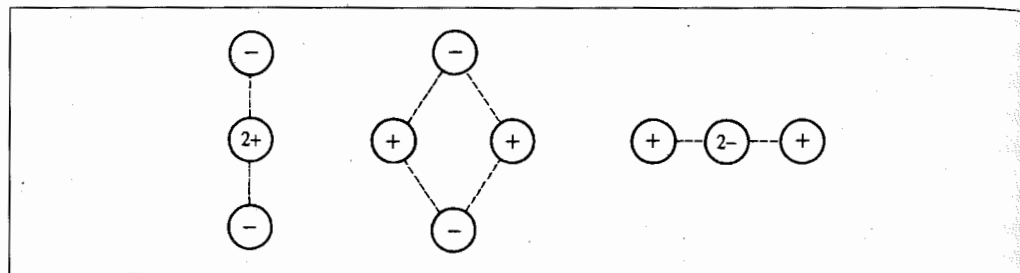


Fig. 2.8: A quadrupole moment can be obtained from various arrangements of two positive and two negative charges.

The electronic contribution to the dipole moment is thus determined from the density matrix and a series of one-electron integrals $\int d\tau \phi_\mu(-\mathbf{r})\phi_\nu$. The dipole moment operator, \mathbf{r} , has components in the x , y and z directions, and so these one-electron integrals are divided into their appropriate components; for example, the x component of the electronic contribution to the dipole moment would be determined using:

$$\mu_x = \sum_{\mu=1}^K \sum_{\nu=1}^K P_{\mu\nu} \int d\tau \phi_\mu(-x)\phi_\nu \quad (2.209)$$

The quadrupole is the next electric moment. A molecule has a non-zero electric quadrupole moment when there is a non-spherically symmetrical distribution of charge. A quadrupole can be considered to arise from four charges that sum to zero which are arranged so that they do not lead to a net dipole. Three such arrangements are shown in Figure 2.8. Whereas the dipole moment has components in the x , y and z directions, the quadrupole has nine components from all pairwise combinations of x and y and is represented by a 3×3 matrix as follows:

$$\Theta = \begin{pmatrix} \sum q_i x_i^2 & \sum q_i x_i y_i & \sum q_i x_i z_i \\ \sum q_i y_i x_i & \sum q_i y_i^2 & \sum q_i y_i z_i \\ \sum q_i z_i x_i & \sum q_i z_i y_i & \sum q_i z_i^2 \end{pmatrix} \quad (2.210)$$

The three moments higher than the quadrupole are the octopole, hexapole and decapole. Methane is an example of a molecule whose lowest non-zero multipole moment is the octopole. The entire set of electric moments is required to completely and exactly describe the distribution of charge in a molecule. However, the series expansion is often truncated after the dipole or quadrupole as these are often the most significant.

Extensive comparisons have been made of experimental and calculated dipole moments (and in some cases the higher moments, though these are difficult to determine accurately by experiment). Factors such as the basis set and electron correlation can have a significant impact on the accuracy of the results, but it is found in many cases that the errors are systematic and that a simple scaling factor can be used to convert the results of a calculation with a small basis set to those obtained from experiment or with a much larger basis set. To illustrate how calculated dipole moments can vary, Table 2.5 provides the dipole moments for formaldehyde calculated at the experimental geometry using a variety of basis sets. It is

STO-3G	1.5258	3-21G	2.2903	4-31G	3.0041
6-31G*	2.7600	6-31G**	2.7576	6-311G**	2.7807
Expt.	2.34				

Table 2.5 Dipole moments calculated for formaldehyde using various basis sets at the experimental geometry.

also important to note that the dipole moment can be very sensitive to the geometry from which it is calculated.

2.7.4 The Total Electron Density Distribution and Molecular Orbitals

The electron density $\rho(\mathbf{r})$ at a point \mathbf{r} can be calculated from the Born interpretation of the wavefunction as a sum of squares of the spin orbitals at the point \mathbf{r} for all occupied molecular orbitals. For a system of N electrons occupying $N/2$ real orbitals, we can write:

$$\rho(\mathbf{r}) = 2 \sum_{i=1}^{N/2} |\psi_i(\mathbf{r})|^2 \quad (2.211)$$

If we express the molecular orbital ψ_i as a linear combination of basis functions, then the electron density at a point \mathbf{r} is given as:

$$\begin{aligned} \rho(\mathbf{r}) &= 2 \sum_{i=1}^{N/2} \left(\sum_{\mu=1}^K c_{\mu i} \phi_\mu(\mathbf{r}) \right) \left(\sum_{\nu=1}^K c_{\nu i} \phi_\nu(\mathbf{r}) \right) \\ &= 2 \sum_{i=1}^{N/2} \sum_{\mu=1}^K c_{\mu i} c_{\mu i} \phi_\mu(\mathbf{r}) \phi_\mu(\mathbf{r}) + 2 \sum_{i=1}^{N/2} \sum_{\mu=1}^K \sum_{\nu=\mu+1}^K 2c_{\mu i} c_{\nu i} \phi_\mu(\mathbf{r}) \phi_\nu(\mathbf{r}) \end{aligned} \quad (2.212)$$

Equation (2.212) can be tidied up considerably if it is written in terms of the elements of the density matrix:

$$\begin{aligned} \left(P_{\mu\nu} = 2 \sum_{i=1}^{N/2} c_{\mu i} c_{\nu i} \right) \\ \rho(\mathbf{r}) &= \sum_{\mu=1}^K \sum_{\nu=1}^K P_{\mu\nu} \phi_\mu(\mathbf{r}) \phi_\nu(\mathbf{r}) \\ &= \sum_{\mu=1}^K P_{\mu\mu} \phi_\mu(\mathbf{r}) \phi_\mu(\mathbf{r}) + 2 \sum_{\mu=1}^K \sum_{\nu=\mu+1}^K P_{\mu\nu} \phi_\mu(\mathbf{r}) \phi_\nu(\mathbf{r}) \end{aligned} \quad (2.213)$$

The integral of $\rho(\mathbf{r})$ over all space equals the number of electrons in the system, N :

$$N = \int d\mathbf{r} \rho(\mathbf{r}) = 2 \sum_{i=1}^{N/2} \int d\mathbf{r} |\psi_i(\mathbf{r})|^2 \quad (2.214)$$

If the overlap between two orbitals ϕ_μ and ϕ_ν is written as $S_{\mu\nu}$, and if the basis functions are

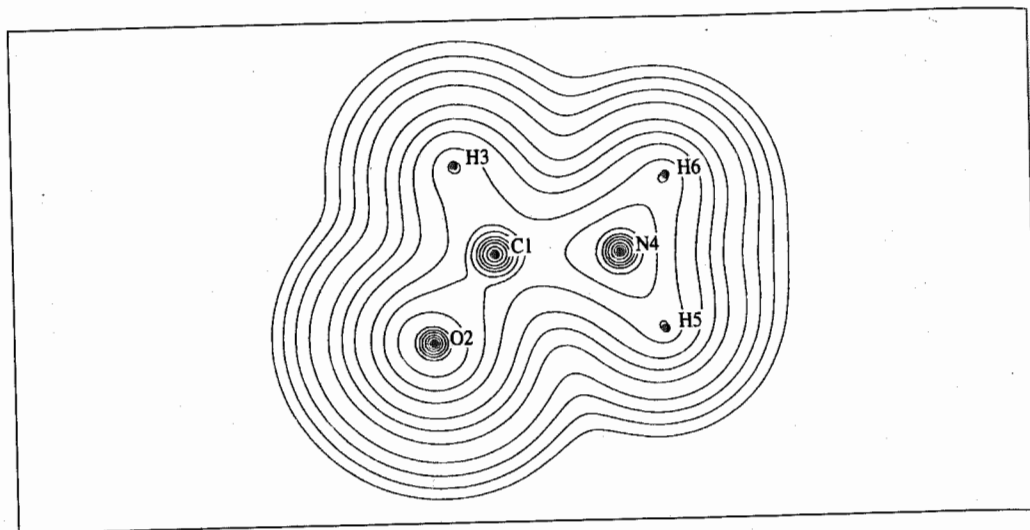


Fig. 2.9: Contour map showing the variation in electron density around formamide.

assumed to be normalised ($S_{\mu\mu} = 1$), then:

$$N = \sum_{\mu=1}^K P_{\mu\mu} + 2 \sum_{\mu=1}^K \sum_{\nu=\mu+1}^K P_{\mu\nu} S_{\mu\nu} \quad (2.215)$$

The electron density can be visualised in several ways. One approach is to construct contours on slices through the molecule, such that each contour connects points of equal density, as shown in Figure 2.9 for formamide. The electron density can also be represented as an isometric projection (or a 'relief map', Figure 2.10), in which the height above the plane

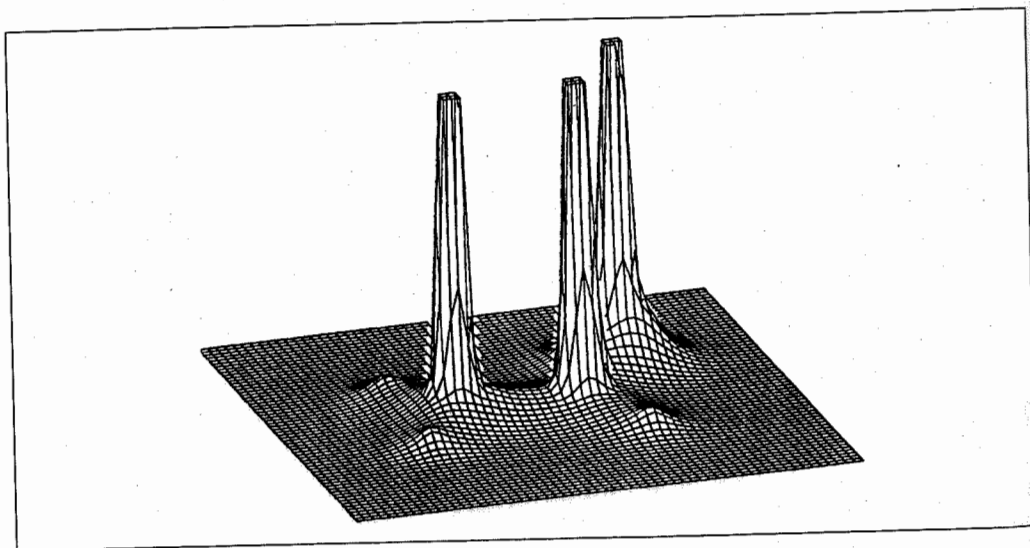


Fig. 2.10: Isometric projection of the electron density around formamide.

represents the magnitude of the electron density. These diagrams show that the electron density tends to be greatest near the nuclei, as would be expected. The electron density can also be represented as a solid object, whose surface connects points of equal density. The surface shown in Figure 2.11 (colour plate section) corresponds to an electron density of 0.0001 a.u. around formamide. Other properties such as the electrostatic potential can be mapped onto this surface, as we shall see in Section 2.7.9.

The electron density distribution of individual molecular orbitals may also be determined and plotted. The highest occupied molecular orbital (HOMO) and the lowest unoccupied molecular orbital (LUMO) are often of particular interest as these are the orbitals most commonly involved in chemical reactions. As an illustration, the HOMO and LUMO for formamide are displayed in Figures 2.12 and 2.13 (colour plate section) as surface pictures.

2.7.5 Population Analysis

Population analysis methods partition the electron density between the nuclei so that each nucleus has a 'number' (not necessarily an integral number) of electrons associated with it. Such a partitioning provides a way to calculate the atomic charge on each nucleus. It should be noted that there is no quantum mechanical operator for the atomic charge and so any partitioning scheme must be arbitrary. Hence many methods have been devised. Here we will consider Mulliken and Löwdin analysis and Bader's theory of atoms in molecules. The alternatives include natural population analysis [Reed *et al.* 1985; Bachrach 1994]. Wiberg and Rablen have compared a number of methods for calculating atomic charges, and we refer to some of their results in the following discussion [Wiberg and Rablen 1993]. To illustrate the variation that can be obtained in the results, for methane they found that the charge on the carbon atom varied from -0.473 to $+0.244$, depending upon the method chosen! We will also consider the problem of calculating atomic charges in more detail in Chapter 4 on molecular mechanics.

2.7.6 Mulliken and Löwdin Population Analysis

RS Mulliken suggested a widely used method for performing population analysis [Mulliken 1955]. The starting point is Equation (2.215), which relates the total number of electrons to the density matrix and to the overlap integrals. In the Mulliken method, all of the electron density ($P_{\mu\mu}$) in an orbital is allocated to the atom on which ϕ_{μ} is located. The remaining electron density is associated with the overlap population, $\phi_{\mu}\phi_{\nu}$. For each element $\phi_{\mu}\phi_{\nu}$ of the density matrix, half of the density is assigned to the atom on which ϕ_{μ} is located and half to the atom on which ϕ_{ν} is located. The net charge on an atom A is then calculated by subtracting the number of electrons from the nuclear charge, Z_A :

$$q_A = Z_A - \sum_{\mu=1; \mu \text{ on } A}^K P_{\mu\mu} - \sum_{\mu=1; \mu \text{ on } A}^K \sum_{\nu=1; \nu \neq \mu}^K P_{\mu\nu} S_{\mu\nu} \quad (2.216)$$

Mulliken population analysis is a trivial calculation to perform once a self-consistent field has been established and the elements of the density matrix have been determined.

However, there are some serious shortcomings to the method, as Mulliken himself pointed out.

A Mulliken analysis depends upon the use of a balanced basis set, in which an equivalent number of basis functions is present on each atom in the molecule. For example, it is possible to calculate a wavefunction for a molecule such as water in which all of the basis functions reside on the oxygen atom; if a large enough basis set is used then a quite reasonable wavefunction for the whole molecule can be obtained. However, the Mulliken analysis would put all of the charge on the oxygen. This is an extreme example of a general problem; p, d and f orbitals are spread quite far from the nucleus with which they are associated and so may be very close to other atoms, yet the charge associated with electron occupation of such orbitals is assigned to the atom on which the orbital is centred. The equal apportioning of electrons between pairs of atoms, even if their electronegativities are very different, can lead in some cases to quite unrealistic values for the net atomic charge. *In extremis*, some orbitals may 'contain' a negative number of electrons and others more than two electrons, in clear contradiction of the Pauli principle. A Mulliken analysis assumes that each basis function can be associated with an atomic centre and so is not applicable if basis functions not centred on the nuclei are used. The atomic charges can be very dependent upon the basis set; for example, Wiberg and Rablen found that the charge on the central carbon in isobutene changed from +0.1 with a 6-31G* basis set to +1.0 for a 6-311++G** basis set.

In the Löwdin approach to population analysis [Löwdin 1970; Cusachs and Politzer 1968] the atomic orbitals are transformed to an orthogonal set, along with the molecular orbital coefficients. The transformed orbitals ϕ'_μ in the orthogonal set are given by:

$$\phi'_\mu = \sum_{\nu=1}^K (\mathbf{S}^{-1/2})_{\nu\mu} \phi_\nu \quad (2.217)$$

The electron population associated with an atom becomes:

$$q_A = Z_A - \sum_{\mu=1; \mu \text{ on } A}^K (\mathbf{S}^{1/2} \mathbf{P} \mathbf{S}^{1/2})_{\mu\mu} \quad (2.218)$$

Löwdin population analysis avoids the problem of negative populations or populations greater than 2. Some quantum chemists prefer the Löwdin approach to that of Mulliken as the charges are often closer to chemically intuitive values and are less sensitive to basis set.

2.7.7 Partitioning Electron Density: The Theory of Atoms in Molecules

R F W Bader's theory of 'atoms in molecules' [Bader 1985] provides an alternative way to partition the electrons between the atoms in a molecule. Bader's theory has been applied to many different problems, but for the purposes of our present discussion we will concentrate on its use in partitioning electron density. The Bader approach is based upon the concept of a *gradient vector path*, which is a curve around the molecule such that it is always perpendicular to the electron density contours. A set of gradient paths is drawn in Figure 2.14 for formamide. As can be seen, some of the gradient paths terminate at the atomic nuclei. Other gradient paths are attracted to points (called critical points) that are

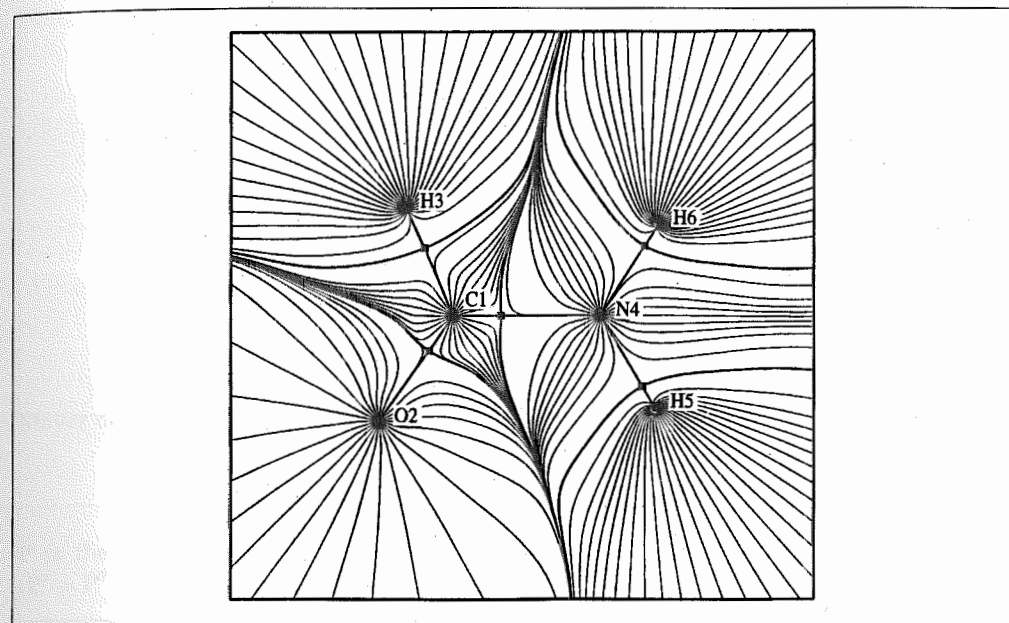


Fig. 2.14: Gradient vector paths around formamide. The paths terminate at atoms or at bond critical points (indicated by squares).

not located at the nuclei; particularly common are the bond critical points, which are located between bonded atoms. Other types of critical point can occur; for example, a *ring critical point* is found in the centre of a benzene ring.

The bond critical points are points of minimum electron charge density between two bonded atoms. If we follow the contour in three-dimensional space from such a point down the gradient path along which the density decreases most rapidly then this gives a means of partitioning the density. This is shown in Figure 2.15 for hydrogen fluoride and in Figure 2.16 for formamide. This procedure can be performed for each bond, resulting in a three-dimensional partitioning of the electron density. The electron population that is assigned to each atom is then calculated by numerically integrating the charge density within the region surrounding that atom.

Wiberg and Rablen found that the charges obtained with the atoms in molecules method were relatively invariant to the basis set. The charges from this method were also consistent with the experimentally determined C-H bond dipoles in methane (in which the carbon is positive) and ethyne (in which the carbon is negative), unlike most of the other methods they examined.

2.7.8 Bond Orders

As with atomic charges, the bond order is not a quantum mechanical observable and so various methods have been proposed for calculating the bond orders in a molecule.

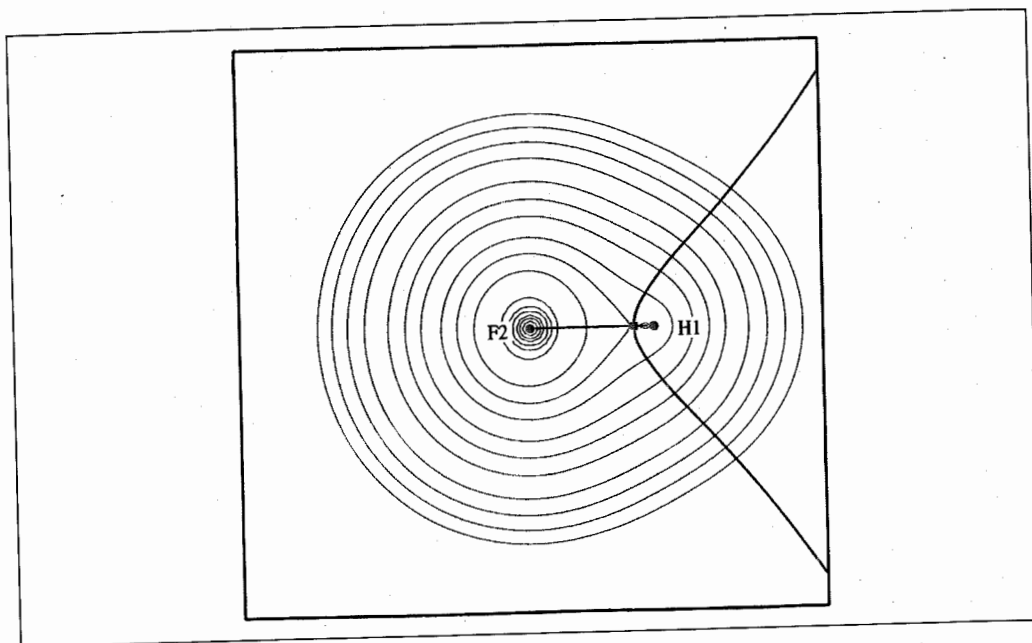


Fig. 2.15: Partitioning the electron density in hydrogen fluoride.

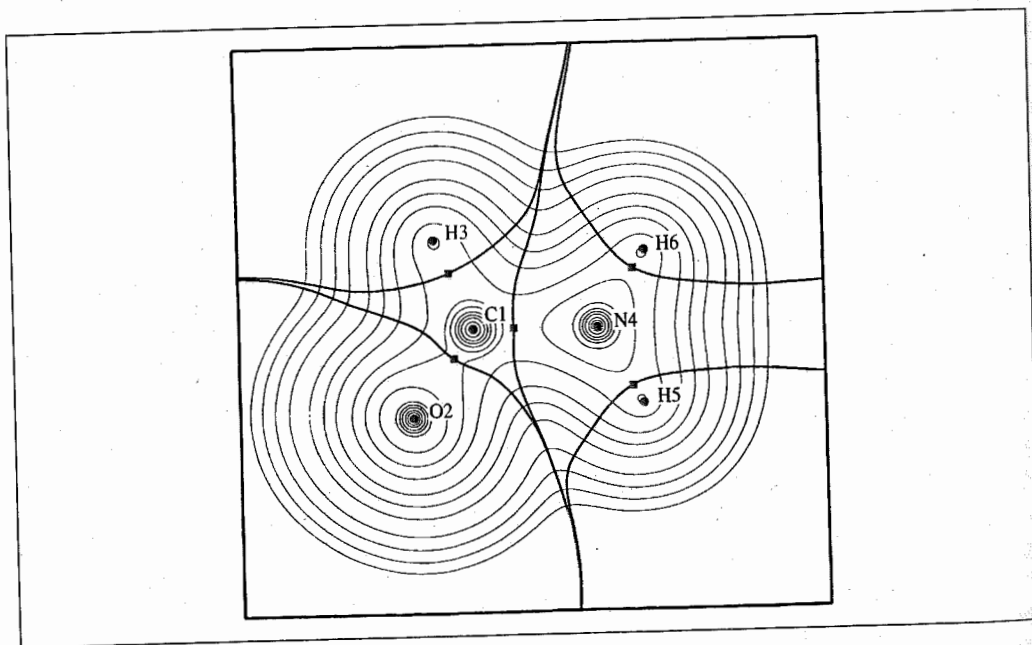


Fig. 2.16: Partitioning the electron density in formamide.

Molecule	Bond	STO-3G	4-31G
H ₂	H-H	1.0	1.0
Methane	C-H	0.99	0.96
Ethene	C=C	2.01	1.96
Ethyne	C-H	0.98	0.96
	C≡C	3.00	3.27
Water	C-H	0.98	0.86
	O-H	0.95	0.80
N ₂	N≡N	3.0	2.67

Table 2.6 Bond order obtained from the Mayer bond order scheme [Mayer 1983].

Mayer defined the bond order between two atoms as follows [Mayer 1983]:

$$B_{AB} = \sum_{\mu \text{ on A}} \sum_{\nu \text{ on B}} [(\mathbf{PS})_{\mu\nu}(\mathbf{PS})_{\nu\mu} + (\mathbf{P}^s\mathbf{S})_{\mu\nu}(\mathbf{P}^s\mathbf{S})_{\nu\mu}] \quad (2.219)$$

\mathbf{P} is the total spinless density matrix ($\mathbf{P} = \mathbf{P}^\alpha + \mathbf{P}^\beta$) and \mathbf{P}^s is the spin density matrix ($\mathbf{P}^s = \mathbf{P}^\alpha - \mathbf{P}^\beta$). For a closed-shell system Mayer's definition of the bond order reduces to:

$$B_{AB} = \sum_{\mu \text{ on A}} \sum_{\nu \text{ on B}} (\mathbf{PS})_{\mu\nu}(\mathbf{PS})_{\nu\mu} \quad (2.220)$$

The bond orders obtained from Mayer's formula often seem intuitively reasonable, as illustrated in Table 2.6 for some simple molecules. The method has also been used to compute the bond orders for intermediate structures in reactions of the form $\text{H} + \text{XH} \rightarrow \text{HX} + \text{H}$ and $\text{X} + \text{H}_2 \rightarrow \text{XH} + \text{H}$ ($\text{X} = \text{F}, \text{Cl}, \text{Br}$). The results suggested that bond orders were a useful way to describe the similarity of the transition structure to the reactants or to the products. Moreover, the bond orders were approximately conserved along the reaction pathway.

As with methods for allocating electron density to atoms, the Mayer method is not necessarily 'correct', though it appears to be a useful measure of the bond order that conforms to accepted pictures of bonding in molecules.

2.7.9 Electrostatic Potentials

The electrostatic potential at a point \mathbf{r} , $\phi(\mathbf{r})$, is defined as the work done to bring unit positive charge from infinity to the point. The electrostatic interaction energy between a point charge q located at \mathbf{r} and the molecule equals $q\phi(\mathbf{r})$. The electrostatic potential has contributions from both the nuclei and from the electrons, unlike the electron density, which only reflects the electronic distribution. The electrostatic potential due to the M nuclei is:

$$\phi_{\text{nucl}}(\mathbf{r}) = \sum_{A=1}^M \frac{Z_A}{|\mathbf{r} - \mathbf{R}_A|} \quad (2.221)$$

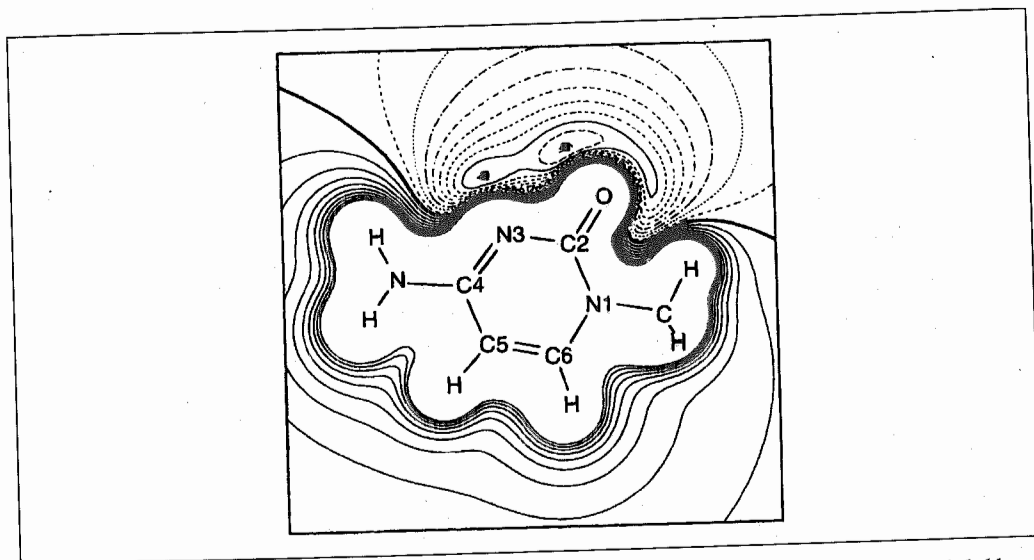


Fig. 2.17: Electrostatic potential contours around cytosine. Negative contours are dashed, the zero contour is bold. The minima near N3 and O are marked.

The potential due to the electrons is obtained from the appropriate integral of the electron density:

$$\phi_{\text{elec}}(\mathbf{r}) = - \int \frac{d\mathbf{r}' \rho(\mathbf{r}')}{|\mathbf{r}' - \mathbf{r}|} \quad (2.222)$$

The total electrostatic potential equals the sum of the nuclear and the electronic contributions:

$$\phi(\mathbf{r}) = \phi_{\text{nucl}}(\mathbf{r}) + \phi_{\text{elec}}(\mathbf{r}) \quad (2.223)$$

The electrostatic potential has proved to be particularly useful for rationalising the interactions between molecules and molecular recognition processes. This is because electrostatic forces are primarily responsible for long-range interactions between molecules. The electrostatic potential varies through space, and so it can be calculated and visualised in the same way as the electron density. Electrostatic potential contours can be used to propose where electrophilic attack might occur; electrophiles are often attracted to regions where the electrostatic potential is most negative. For example, the experimentally determined position of electrophilic attack at the nucleic acid cytosine is at N3 (Figure 2.17). This atom is next to a minimum in the electrostatic potential (also shown in Figure 2.17), as pointed out by Politzer and Murray [Politzer and Murray 1991].

Non-covalent interactions between molecules often occur at separations where the van der Waals radii of the atoms are just touching and so it is often most useful to examine the electrostatic potential in this region. For this reason, the electrostatic potential is often calculated at the molecular surface (defined in Section 1.5) or the equivalent isodensity surface as shown in Figure 2.18 (colour plate section). Such pictorial representations

can be used to qualitatively assess the degree of electrostatic similarity between two molecules.

2.7.10 Thermodynamic and Structural Properties

The total energy of a system is equal to the sum of the electronic energy and the Coulombic nuclear repulsion energy:

$$E_{\text{tot}} = E_{\text{elec}} + \sum_{A=1}^M \sum_{B=A+1}^M \frac{Z_A Z_B}{R_{AB}} \quad (2.224)$$

A more useful quantity for comparison with experiment is the heat of formation, which is defined as the enthalpy change when one mole of a compound is formed from its constituent elements in their standard states. The heat of formation can thus be calculated by subtracting the heats of atomisation of the elements and the atomic ionisation energies from the total energy. Unfortunately, *ab initio* calculations that do not include electron correlation (which we will discuss in Chapter 3) provide uniformly poor estimates of heats of formation with errors in bond dissociation energies of 25–40 kcal/mol, even at the Hartree–Fock limit for diatomic molecules.

When combined with an energy minimisation algorithm, quantum mechanics can be used to calculate equilibrium geometries of molecules. The results of such calculations can be compared with the structures obtained from gas-phase experiments using microwave spectroscopy, electronic spectroscopy and electron diffraction. Extensive tables listing comparisons between calculations and experiment for many molecules have been published in several reviews. Not surprisingly, the agreement between theory and experiment for *ab initio* calculations generally improves as one increases the size of the basis set. Hehre *et al.* suggest that the 3-21G basis set offers a good compromise between performance and applicability [Hehre *et al.* 1986]. It is often found that errors in structural predictions are systematic rather than random. For example, STO-3G bond lengths are generally too long, whilst 6-31G* bond lengths tend to be too short. By analysing the trends in such calculations it can be possible to derive scaling factors which enable more accurate predictions to be made for each level of theory.

Quantum mechanics can be used to calculate the relative energies of conformations and the energy barriers between them. Experimental data is available for both relative stabilities and barrier heights in some cases, though this tends to be limited to relatively simple molecules. Butane is one molecule that has been investigated in great detail, with its *gauche* and *anti* conformations and the barriers that separate them. The energy difference between the *syn* and *anti* conformations of butane (Figure 2.19) was found to fall significantly with increasing basis set size, particularly when correlated levels of theory were employed [Wiberg and Murcko 1988; Allinger *et al.* 1990; Smith and Jaffe 1996]. However, the smaller energy difference between the minimum energy *anti* and *gauche* conformations can be calculated quite accurately even with a relatively small basis set. Quantum mechanics calculations of the change in energy as a bond is rotated are often used to parametrise the torsional terms in molecular mechanics force fields, as will be discussed in Section 4.18.

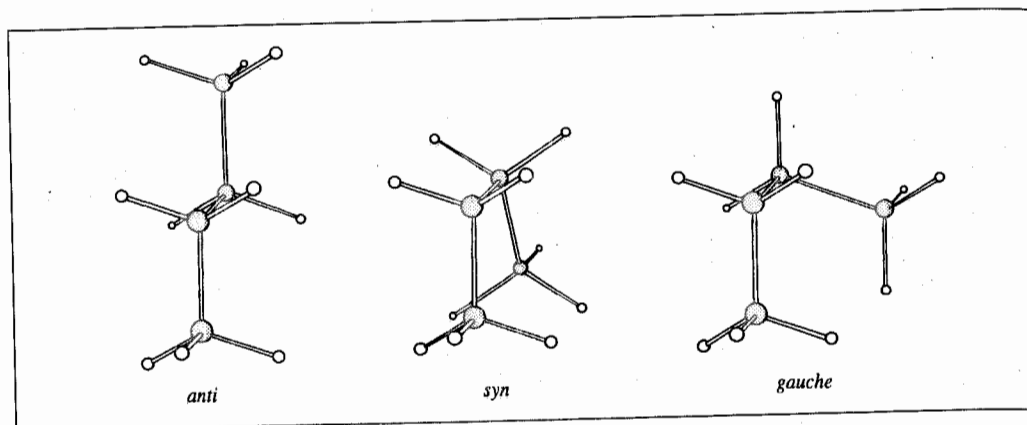


Fig. 2.19: syn, anti and gauche conformations of butane (C-C-C torsion angles 0° , 180° and $\pm 60^\circ$ respectively).

2.8 Approximate Molecular Orbital Theories

Ab initio calculations can be extremely expensive in terms of the computer resources required. Nevertheless, improvements in computer hardware and the availability of easy-to-use programs have helped to make *ab initio* methods a widely used computational tool. The approximate quantum mechanical methods require significantly less computational resources. Indeed, the earliest approximate methods such as Hückel theory predate computers by many years. Moreover, by their incorporation of parameters derived from experimental data some approximate methods can calculate certain properties more accurately than even the highest level of *ab initio* methods.

Many approximate molecular orbital theories have been devised. Most of these methods are not in widespread use today in their original form. Nevertheless, the more widely used methods of today are derived from earlier formalisms, which we will therefore consider where appropriate. We will concentrate on the semi-empirical methods developed in the research groups of Pople and Dewar. The former pioneered the CNDO, INDO and NDDO methods, which are now relatively little used in their original form but provided the basis for subsequent work by the Dewar group, whose research resulted in the popular MINDO/3, MNDO and AM1 methods. Our aim will be to show how the theory can be applied in a practical way, not only to highlight their successes but also to show where problems were encountered and how these problems were overcome. We will also consider the Hückel molecular orbital approach and the extended Hückel method. Our discussion of the underlying theoretical background of the approximate molecular orbital methods will be based on the Roothaan-Hall framework we have already developed. This will help us to establish the similarities and the differences with the *ab initio* approach.

2.9 Semi-empirical Methods

A discussion of semi-empirical methods starts most appropriately with the key components

of the Roothaan-Hall equations, which for a closed-shell system are:

$$\mathbf{FC} = \mathbf{SCE} \quad (2.225)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} + \sum_{\lambda=1}^K \sum_{\sigma=1}^K P_{\lambda\sigma} [(\mu\nu|\lambda\sigma) - \frac{1}{2}(\mu\lambda|\nu\sigma)] \quad (2.226)$$

$$P_{\lambda\sigma} = 2 \sum_{i=1}^{N/2} c_{\lambda i} c_{\sigma i} \quad (2.227)$$

$$H_{\mu\nu}^{\text{core}} = \int d\nu_1 \phi_\mu(1) \left[-\frac{1}{2} \nabla^2 - \sum_{A=1}^M \frac{Z_A}{|r_1 - R_A|} \right] \phi_\nu(1) \quad (2.228)$$

In *ab initio* calculations all elements of the Fock matrix are calculated using Equation (2.226), irrespective of whether the basis functions ϕ_μ , ϕ_ν , ϕ_λ and ϕ_σ are on the same atom, on atoms that are bonded or on atoms that are not formally bonded. To discuss the semi-empirical methods it is useful to consider the Fock matrix elements in three groups: $F_{\mu\mu}$ (the diagonal elements), $F_{\mu\nu}$ (where ϕ_μ and ϕ_ν are on the same atom) and $F_{\mu\nu}$ (where ϕ_μ and ϕ_ν are on different atoms).

We have mentioned several times that the greatest proportion of the time required to perform an *ab initio* Hartree-Fock SCF calculation is invariably spent calculating and manipulating integrals. The most obvious way to reduce the computational effort is therefore to neglect or approximate some of these integrals. Semi-empirical methods achieve this in part by explicitly considering only the valence electrons of the system; the core electrons are subsumed into the nuclear core. The rationale behind this approximation is that the electrons involved in chemical bonding and other phenomena that we might wish to investigate are those in the valence shell. By considering all the valence electrons the semi-empirical methods differ from those theories (e.g. Hückel theory) that explicitly consider only the π electrons of a conjugated system and which are therefore limited to specific classes of molecule. The semi-empirical calculations invariably use basis sets comprising Slater type s, p and sometimes d orbitals. The orthogonality of such orbitals enables further simplifications to be made to the equations.

A feature common to the semi-empirical methods is that the overlap matrix, \mathbf{S} (in Equation (2.225)), is set equal to the identity matrix \mathbf{I} . Thus all diagonal elements of the overlap matrix are equal to 1 and all off-diagonal elements are zero. Some of the off-diagonal elements would naturally be zero due to the use of orthogonal basis sets on each atom, but in addition the elements that correspond to the overlap between two atomic orbitals on different atoms are also set to zero. The main implication of this is that the Roothaan-Hall equations are simplified: $\mathbf{FC} = \mathbf{SCE}$ becomes $\mathbf{FC} = \mathbf{CE}$ and so is immediately in standard matrix form. It is important to note that setting \mathbf{S} equal to the identity matrix does not mean that all overlap integrals are set to zero in the calculation of Fock matrix elements. Indeed, it is important specifically to include some of the overlaps in even the simplest of the semi-empirical models.

2.9.1 Zero-differential Overlap

Many semi-empirical theories are based upon the zero-differential overlap approximation (ZDO). In this approximation, the overlap between pairs of different orbitals is set to zero for all volume elements $d\nu$:

$$\phi_\mu \phi_\nu d\nu = 0 \quad (2.229)$$

This directly leads to the following result for the overlap integrals:

$$S_{\mu\nu} = \delta_{\mu\nu} \quad (2.230)$$

If the two atomic orbitals ϕ_μ and ϕ_ν are located on different atoms then the differential overlap is referred to as diatomic differential overlap; if ϕ_μ and ϕ_ν are on the same atom then we have monatomic differential overlap. If the ZDO approximation is applied to the two-electron repulsion integral $(\mu\nu|\lambda\sigma)$ then the integral will equal zero if $\mu \neq \nu$ and/or if $\lambda \neq \sigma$. This can be written concisely using the Kronecker delta:

$$(\mu\nu|\lambda\sigma) = (\mu\mu|\lambda\lambda)\delta_{\mu\nu}\delta_{\lambda\sigma} \quad (2.231)$$

It can immediately be seen that all three- and four-centre integrals are set to zero under the ZDO approximation. If the ZDO approximation is applied to all orbital pairs then the Roothaan-Hall equations for a closed-shell molecule (Equation (2.226)) simplify considerably to give the following for $\mu \equiv \nu$:

$$F_{\mu\mu} = H_{\mu\mu}^{\text{core}} + \sum_{\lambda=1}^K P_{\lambda\lambda}(\mu\mu|\lambda\lambda) - \frac{1}{2}P_{\mu\mu}(\mu\mu|\mu\mu) \quad (2.232)$$

The summation over λ includes $\lambda = \mu$, and the terms in $(\mu\mu|\mu\mu)$ can be separated to give:

$$F_{\mu\mu} = H_{\mu\mu}^{\text{core}} + \frac{1}{2}P_{\mu\mu}(\mu\mu|\mu\mu) + \sum_{\lambda=1; \lambda \neq \mu}^K P_{\lambda\lambda}(\mu\mu|\lambda\lambda) \quad (2.233)$$

For $\nu \neq \mu$ we have:

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} - \frac{1}{2}P_{\mu\nu}(\mu\mu|\nu\nu) \quad (2.234)$$

Sensible results cannot be obtained by simply applying the ZDO approximation to all pairs of orbitals *carte blanche*. There are two major reasons for this.

The first consideration is that the total wavefunction and the molecular properties calculated from it should be the same when a transformed basis set is used. We have already encountered this requirement in our discussion of the transformation of the Roothaan-Hall equations to an orthogonal set. To reiterate: suppose a molecular orbital is written as a linear combination of atomic orbitals:

$$\psi_i = \sum_{\mu} c_{\mu i} \phi_{\mu} \quad (2.235)$$

If an alternative basis set is used in which the basis functions are just linear combinations of the original basis functions, then the same wavefunction can be written as a linear

combination of these new transformed functions:

$$\psi_i = \sum_{\alpha} c_{\alpha i} \phi'_{\alpha} \quad (2.236)$$

$$\phi'_{\alpha} = \sum_{\mu_{\alpha}} t_{\mu_{\alpha}} \phi_{\mu} \quad (2.237)$$

$t_{\mu_{\alpha}}$ are the coefficients of the original basis functions in the linear expansion of the transformed basis set. Different types of transformation are possible; for example, some transformations mix orbitals with the same principal and azimuthal quantum numbers (e.g. mixing $2p_x$, $2p_y$ and $2p_z$); others mix orbitals with the same principal quantum number but different azimuthal quantum numbers (e.g. mixing $2s$, $2p_x$, $2p_y$ and $2p_z$ orbitals to give sp^3 hybrid orbitals); yet other transformations mix orbitals located on different atoms. Suppose we mix $2p_x$ and $2p_y$ atomic orbitals on the same atom. The differential overlap between these two orbitals is $2p_x 2p_y$. We now introduce the following two new coordinates, which correspond to a rotation in the xy plane:

$$x' = \frac{1}{\sqrt{2}}(x + y) \quad (2.238)$$

$$y' = \frac{1}{\sqrt{2}}(-x + y) \quad (2.239)$$

The overlap between the $2p'_x$ and $2p'_y$ orbitals in this new coordinate system is $\frac{1}{2}(2p_y^2 - 2p_x^2)$. If the zero differential overlap approximation were applied, then different results would be obtained for the two coordinate systems unless the overlap in the new, transformed system was also ignored.

The second reason why the ZDO approximation is not applied to all pairs of orbitals is that the major contributors to bond formation are the electron-core interactions between pairs of orbitals and the nuclear cores (i.e. $H_{\mu\nu}^{\text{core}}$). These interactions are therefore not subjected to the ZDO approximation (and so do not suffer from any transformation problems).

2.9.2 CNDO

The complete neglect of differential overlap (CNDO) approach of Pople, Santry and Segal was the first method to implement the zero-differential overlap approximation in a practical fashion [Pople *et al.* 1965]. To overcome the problems of rotational invariance, the two-electron integrals $(\mu\mu|\lambda\lambda)$, where μ and λ are on different atoms A and B, were set equal to a parameter γ_{AB} which depends only on the nature of the atoms A and B and the internuclear distance, and not on the type of orbital. The parameter γ_{AB} can be considered to be the average electrostatic repulsion between an electron on atom A and an electron on atom B. When both atomic orbitals are on the same atom the parameter is written γ_{AA} and represents the average electron-electron repulsion between two electrons on an atom A.

With this approximation we can divide the elements of the Fock matrix into three groups: $F_{\mu\mu}$ (the diagonal elements), $F_{\mu\nu}$ (where μ and ν are on different atoms) and $F_{\mu\nu}$ (where μ

and ν are on the same atom). To obtain $F_{\mu\mu}$ we substitute γ_{AB} for the two-electron integrals $(\mu\mu|\lambda\lambda)$ where μ and λ are on different atoms and γ_{AA} where μ and λ are on the same atom into the Fock matrix equations, Equations (2.240)–(2.242):

$$F_{\mu\mu} = H_{\mu\mu}^{\text{core}} + \sum_{\lambda=1; \lambda \text{ on A}}^K P_{\lambda\lambda}\gamma_{AA} - \frac{1}{2}P_{\mu\mu}\gamma_{AA} + \sum_{\lambda=1; \lambda \text{ not on A}}^K P_{\lambda\lambda}\gamma_{AB} \quad (2.240)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} - \frac{1}{2}P_{\mu\nu}\gamma_{AA}; \quad \mu \text{ and } \nu \text{ both on atom A} \quad (2.241)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} - \frac{1}{2}P_{\mu\nu}\gamma_{AB}; \quad \mu \text{ and } \nu \text{ on different atoms, A and B} \quad (2.242)$$

Equation (2.240) is rather untidy, involving summations over basis functions on atom A and basis functions not on atom A. It is often simplified by writing P_{AA} as the total electron density on atom A, where:

$$P_{AA} = \sum_{\lambda \text{ on A}}^A P_{\lambda\lambda} \quad (2.243)$$

A similar expression can also be introduced for P_{BB} . With this notation $F_{\mu\mu}$ simplifies to:

$$F_{\mu\mu} = H_{\mu\mu}^{\text{core}} + (P_{AA} - \frac{1}{2}P_{\mu\mu})\gamma_{AA} + \sum_{B \neq A} P_{BB}\gamma_{AB} \quad (2.244)$$

The core Hamiltonian expressions, $H_{\mu\mu}^{\text{core}}$ and $H_{\mu\nu}^{\text{core}}$, correspond to electrons moving in the field of the parent nucleus and the other nuclei. In semi-empirical methods the core electrons are subsumed into the nucleus and so the nuclear charges are altered accordingly (for example, carbon has a nuclear 'charge' of +4).

In CNDO $H_{\mu\mu}^{\text{core}}$ is separated into an integral involving the atom on which ϕ_μ is situated (labelled A), and all the others (labelled B). Thus:

$$H_{\mu\mu}^{\text{core}} = U_{\mu\mu} - \sum_{B \neq A} V_{AB} \quad (2.245)$$

where:

$$U_{\mu\mu} = \left(\mu \left| -\frac{1}{2}\nabla^2 - \frac{Z_A}{|\mathbf{r}_1 - \mathbf{R}_A|} \right| \mu \right) \quad \text{and} \quad V_{AB} = \left(\mu \left| \frac{Z_B}{|\mathbf{r}_1 - \mathbf{R}_B|} \right| \mu \right) \quad (2.246)$$

$U_{\mu\mu}$ is thus the energy of the orbital ϕ_μ in the field of its own nucleus (A) and core electrons; $-V_{AB}$ is the energy of the electron in the field of another nucleus (B). To maintain consistency with the way in which the two-electron integrals are treated, the terms

$$\left(\mu \left| \frac{Z_B}{|\mathbf{r}_1 - \mathbf{R}_B|} \right| \mu \right) \quad (2.247)$$

must be the same for all orbitals ϕ_μ on atom A (i.e. the interaction energy between any electron in an orbital on atom A with the core of atom B is equal to V_{AB}).

We next consider $H_{\mu\nu}^{\text{core}}$, where ϕ_μ and ϕ_ν are both on the same atom, A. In this case the core Hamiltonian has the following form:

$$\begin{aligned} H_{\mu\nu}^{\text{core}} &= \left(\mu \left| -\frac{1}{2}\nabla^2 - \frac{Z_A}{|\mathbf{r}_1 - \mathbf{R}_A|} \right| \nu \right) - \sum_{B \neq A} \left(\mu \left| \frac{Z_B}{|\mathbf{r}_1 - \mathbf{R}_B|} \right| \nu \right) \\ &= U_{\mu\nu} - \sum_{B \neq A} \left(\mu \left| \frac{Z_B}{|\mathbf{r}_1 - \mathbf{R}_B|} \right| \nu \right) \end{aligned} \quad (2.248)$$

As ϕ_μ and ϕ_ν are on the same atom, $U_{\mu\nu}$ is zero due to the orthogonality of atomic orbitals. The term

$$\left(\mu \left| \frac{Z_B}{|\mathbf{r}_1 - \mathbf{R}_B|} \right| \nu \right) \quad (2.249)$$

is zero in accordance with the zero-differential overlap approximation. Thus $H_{\mu\nu}^{\text{core}}$ is zero in CNDO.

Finally, if ϕ_μ and ϕ_ν are on two different atoms A and B, then we can write:

$$H_{\mu\nu}^{\text{core}} = \left(\mu \left| -\frac{1}{2}\nabla^2 - \frac{Z_A}{|\mathbf{r}_1 - \mathbf{R}_A|} - \frac{Z_B}{|\mathbf{r}_1 - \mathbf{R}_B|} \right| \nu \right) - \sum_{C \neq A, B} \left(\mu \left| -\frac{Z_C}{|\mathbf{r}_1 - \mathbf{R}_C|} \right| \nu \right) \quad (2.250)$$

The second term corresponds to the interaction of the distribution $\phi_\mu\phi_\nu$ with the atoms C ($\neq A, B$). These interactions are ignored. The first part (known as the *resonance integral* and commonly written $\beta_{\mu\nu}$) is not subject to the ZDO approximation, because it is the main cause of bonding. In CNDO the resonance integral is made proportional to the overlap integral, $S_{\mu\nu}$:

$$H_{\mu\nu}^{\text{core}} = \beta_{AB}^0 S_{\mu\nu} \quad (2.251)$$

where β_{AB}^0 is a parameter which depends on the nature of atoms A and B.

With these approximations the Fock matrix elements for CNDO become:

$$F_{\mu\mu} = U_{\mu\mu} + \sum_{B \neq A} V_{AB} + (P_{AA} - \frac{1}{2}P_{\mu\mu})\gamma_{AA} + \sum_{B \neq A} P_{BB}\gamma_{AB} \quad (2.252)$$

$$F_{\mu\nu} = -\frac{1}{2}P_{\mu\nu}\gamma_{AA}; \quad \mu \text{ and } \nu \text{ on the same atom, A} \quad (2.253)$$

$$F_{\mu\nu} = \beta_{AB}^0 S_{\mu\nu} - \frac{1}{2}P_{\mu\nu}\gamma_{AB}; \quad \mu \text{ on A and } \nu \text{ on B} \quad (2.254)$$

To perform a CNDO calculation requires the following to be calculated or specified: the overlap integrals, $S_{\mu\nu}$, the core Hamiltonians $U_{\mu\mu}$, the electron-core interactions V_{AB} , the electron repulsion integrals γ_{AB} and γ_{AA} and the bonding parameters β_{AB}^0 . The CNDO basis set comprises Slater type orbitals for the valence shell with the exponents being chosen using Slater's rules (except for hydrogen, where an exponent of 1.2 is used as this value is more appropriate to hydrogen atoms in molecules). Thus the basis set comprises 1s for hydrogen and 2s, 2p_x, 2p_y and 2p_z for the first-row elements. The overlap integrals are calculated explicitly (the overlap between two basis functions on the same atom is, of course, zero with an s, p basis set). The electron repulsion integral parameter γ_{AB} is

calculated using valence *s* functions on the two atoms A and B:

$$\gamma_{AB} = \iint d\nu_1 d\nu_2 \phi_{s,A}(1) \phi_{s,A}(1) \left(\frac{1}{r_{12}} \right) \phi_{s,B}(2) \phi_{s,B}(2) \quad (2.255)$$

The use of spherically symmetric *s* orbitals avoids the problems associated with transformations of the axes. The core Hamiltonians ($U_{\mu\mu}$) are not calculated but are obtained from experimental ionisation energies. This is because it is important to distinguish between *s* and *p* orbitals in the valence shell (i.e. the 2*s* and 2*p* orbitals for the first-row elements), and without explicit core electrons this is difficult to achieve. The resonance integrals, β_{AB}^0 , are written in terms of empirical single-atom values as follows:

$$\beta_{AB}^0 = \frac{1}{2}(\beta_A^0 + \beta_B^0) \quad (2.256)$$

The β^0 values are chosen to fit the results of minimal basis set *ab initio* calculations on diatomic molecules.

The electron–core interaction, V_{AB} , is calculated as the interaction between an electron in a valence *s* orbital on atom A with the nuclear core of atom B:

$$V_{AB} = \int d\nu_1 \phi_{s,A}(1) \frac{Z_B}{|r_1 - R_B|} \phi_{s,A}(1) \quad (2.257)$$

CNDO is rightly recognised as the first in a long line of important semi-empirical models. However, there were some important limitations with the model. One especially serious deficiency of the first version of CNDO (introduced in 1965 [Pople and Segal 1965, Pople *et al.* 1965] and now known as CNDO/1) is that two neutral atoms show a significant (and incorrect) attraction, even when separated by several ångströms. The predicted equilibrium distances for diatomic molecules are also too short and the dissociation energies too large. These effects are due to electrons on one atom penetrating the valence shell of another atom and so experiencing a nuclear attraction. This penetration effect can be quantified more explicitly as follows. The net charge on an atom B equals the difference between its nuclear charge and the total electron density: $Q_B = Z_B - P_{BB}$. If we now substitute for P_{BB} ($= Z_B - Q_B$) in the diagonal elements of the Fock matrix, Equation (2.252), we obtain:

$$F_{\mu\mu} = U_{\mu\mu} + (P_{AA} - \frac{1}{2}P_{\mu\mu})\gamma_{AA} + \sum_{B \neq A} [-Q_B\gamma_{AB} + (Z_B\gamma_{AB} - V_{AB})] \quad (2.258)$$

$-Q_B\gamma_{AB}$ is the contribution from the total charge on atom B; this is zero if the atomic charge is exactly balanced by the electron density. $Z_B\gamma_{AB} - V_{AB}$ is called the *penetration integral*. It was this contribution that caused the anomalous results for two neutral atoms at large separation. In the second version of CNDO (CNDO/2 [Pople and Segal 1966]) the penetration integral effect was eliminated by putting $V_{AB} = Z_B\gamma_{AB}$. The core Hamiltonian $U_{\mu\mu}$ was also defined differently in CNDO/2, using both ionisation energies and electron affinities.

2.9.3 INDO

CNDO makes no allowance for the fact that the interaction between two electrons depends upon their relative spins. This effect can be particularly severe for electrons on the same

atom. Thus, in CNDO all two-electron integrals ($\mu\nu|\lambda\nu$) are set to zero, and integrals ($\mu\mu|\nu\nu$) and ($\mu\mu|\mu\mu$) are forced to be equal (to γ_{AA}). The next development was the intermediate neglect of differential overlap model (INDO [Pople *et al.* 1967]), which includes monatomic differential overlap for one-centre integrals (i.e. for integrals involving basis functions centred on the same atom). This enables the interaction between two electrons on the same atom with parallel spins to have a lower energy than the comparable interaction between electrons with paired spins. For this reason the Fock matrix elements are usually written with the spin (α or β) explicitly specified. The elements $F_{\mu\mu}$ and $F_{\mu\nu}$ (where μ and ν are located on atom A) then change from their CNDO/2 values as follows:

$$F_{\mu\mu}^\alpha = U_{\mu\mu} + \sum_{\lambda \text{ on A}} \sum_{\sigma \text{ on A}} [P_{\lambda\sigma}(\mu\mu|\lambda\sigma) - P_{\lambda\sigma}^\alpha(\mu\lambda|\mu\sigma)] + \sum_{B \neq A} (P_{BB} - Z_B)\gamma_{AB} \quad (2.259)$$

$$F_{\mu\nu}^\alpha = U_{\mu\nu} + \sum_{\lambda \text{ on A}} \sum_{\sigma \text{ on A}} [P_{\lambda\sigma}(\mu\nu|\lambda\sigma) - P_{\lambda\sigma}^\alpha(\mu\lambda|\nu\sigma)]; \quad \mu \text{ and } \nu \text{ both on atom A} \quad (2.260)$$

In Equation (2.259) we have included the CNDO/2 approximation $V_{AB} = Z_B\gamma_{AB}$. The matrix element $F_{\mu\nu}$, where μ and ν are on different atoms, is the same as in CNDO/2:

$$F_{\mu\nu}^\alpha = \frac{1}{2}(\beta_A^0 + \beta_B^0)S_{\mu\nu} - P_{\mu\nu}^\alpha\gamma_{AB} \quad (2.261)$$

In a closed-shell system, $P_{\mu\nu}^\alpha = P_{\mu\nu}^\beta = \frac{1}{2}P_{\mu\nu}$ and the Fock matrix elements can be obtained by making this substitution. If a basis set containing *s*, *p* orbitals is used, then many of the one-centre integrals nominally included in INDO are equal to zero, as are the core elements $U_{\mu\nu}$. Specifically, only the following one-centre, two-electron integrals are non-zero: ($\mu\mu|\mu\mu$), ($\mu\mu|\nu\nu$) and ($\mu\nu|\mu\nu$). The elements of the Fock matrix that are affected can then be written as follows:

$$F_{\mu\mu} = U_{\mu\mu} + \sum_{\nu \text{ on A}} [P_{\nu\nu}(\mu\mu|\nu\nu) - \frac{1}{2}P_{\nu\nu}(\mu\nu|\mu\nu)] + \sum_{B \neq A} (P_{BB} - Z_B)\gamma_{AB} \quad (2.262)$$

$$F_{\mu\nu} = \frac{3}{2}P_{\mu\nu}(\mu\nu|\mu\nu) - \frac{1}{2}P_{\mu\nu}(\mu\mu|\nu\nu); \quad \mu, \nu \text{ on the same atom} \quad (2.263)$$

Some of the one-centre two-electron integrals in INDO are semi-empirical parameters, obtained by fitting to atomic spectroscopic data. The core integrals $U_{\mu\mu}$ are obtained in a slightly different fashion to that of CNDO/2, to take into account the new electronic configurations under the INDO model for atoms and their cations and anions. An INDO calculation requires little additional computational effort compared with the corresponding CNDO calculation and has the key advantage that states of different multiplicities can be distinguished. For example, in CNDO the singlet and triplet configurations $1s^2 2s^2 2p^2$ of carbon have the same energy, whereas these can be distinguished using INDO. Two of the systems considered in the original INDO publication were the methyl and ethyl radicals, the unpaired electron density being compared with experimentally determined hyperfine coupling constants. INDO gave a much more favourable result for these systems than CNDO.

2.9.4 NDDO

The next level of approximation is the neglect of diatomic differential overlap model (NDDO [Pople *et al.* 1965]); this theory only neglects differential overlap between atomic orbitals on

different atoms. Thus all of the two-electron, two-centre integrals of the form $(\mu\nu|\lambda\sigma)$, where μ and ν are on the same atom and λ and σ are also on the same atom, are retained. The Fock matrix elements become:

$$F_{\mu\mu} = H_{\mu\mu}^{\text{core}} + \sum_{\lambda \text{ on A}} \sum_{\sigma \text{ on A}} [P_{\lambda\sigma}(\mu\mu|\lambda\sigma) - \frac{1}{2}P_{\lambda\sigma}(\mu\lambda|\mu\sigma)] + \sum_{B \neq A} \sum_{\lambda \text{ on B}} \sum_{\sigma \text{ on B}} P_{\lambda\sigma}(\mu\mu|\lambda\sigma) \quad (2.264)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} + \sum_{\lambda \text{ on A}} \sum_{\sigma \text{ on A}} [P_{\lambda\sigma}(\mu\nu|\lambda\sigma) - \frac{1}{2}P_{\lambda\sigma}(\mu\lambda|\nu\sigma)] + \sum_{B \neq A} \sum_{\lambda \text{ on B}} \sum_{\sigma \text{ on B}} P_{\lambda\sigma}(\mu\nu|\lambda\sigma); \quad \mu \text{ and } \nu \text{ both on A} \quad (2.265)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} - \frac{1}{2} \sum_{\lambda \text{ on B}} \sum_{\sigma \text{ on A}} P_{\lambda\sigma}(\mu\sigma|\nu\lambda); \quad \mu \text{ on A and } \nu \text{ on B} \quad (2.266)$$

It is again possible to tidy up equations (2.264) and (2.265) when an s, p basis set is used:

$$F_{\mu\mu} = H_{\mu\mu}^{\text{core}} + \sum_{\nu \text{ on A}} [P_{\nu\nu}(\mu\mu|\nu\nu) - \frac{1}{2}P_{\nu\nu}(\mu\nu|\mu\nu)] + \sum_{B \neq A} \sum_{\lambda \text{ on B}} \sum_{\sigma \text{ on B}} P_{\lambda\sigma}(\mu\mu|\lambda\sigma) \quad (2.267)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} + \frac{3}{2}P_{\mu\nu}(\mu\nu|\mu\nu) - \frac{1}{2}P_{\mu\nu}(\mu\mu|\nu\nu) + \sum_{B \neq A} \sum_{\lambda \text{ on B}} \sum_{\sigma \text{ on B}} P_{\lambda\sigma}(\mu\nu|\lambda\sigma) \quad (2.268)$$

Whereas the computation required for an INDO calculation is little more than for the analogous CNDO calculation, in NDDO the number of two-electron, two-centre integrals is increased by a factor of approximately 100 for each pair of heavy atoms in the system.

2.9.5 MINDO/3

The CNDO, INDO and NDDO methods, as originally devised and implemented, are now little used, in comparison with the methods subsequently developed by Dewar and colleagues, but they were of considerable importance in showing how a systematic series of approximations could be used to develop methods of real practical value. Moreover, the calculations could be performed in a fraction of the time required to solve the full Roothaan-Hall equations. However, they did not produce very accurate results, largely because they were parametrised upon the results from relatively low-level *ab initio* calculations, which themselves agreed poorly with experiment. They were also limited to small classes of molecule, and they often required a good experimental geometry to be supplied as input because their geometry optimisation algorithms were not very sophisticated.

It was through the introduction of the MINDO/3 method by Bingham, Dewar and Lo [Bingham *et al.* 1975a-d] that a wider audience was able to apply semi-empirical methods in their own research. MINDO/3 was not so much a significant change in the theory, being based upon INDO (MINDO stands for modified INDO), but it did differ significantly in the way in which the method was parametrised, making much more use of experimental data. It also incorporated a geometry optimisation routine (the Davidon-Fletcher-Powell method; see Chapter 5), which enabled the program to accept crude initial geometries as input and derive the associated minimum energy structures.

MINDO/3 uses an s, p basis set and its Fock matrix elements are:

$$F_{\mu\mu} = U_{\mu\mu} + \sum_{\nu \text{ on A}} (P_{\nu\nu}(\mu\mu|\nu\nu) - \frac{1}{2}P_{\nu\nu}(\mu\nu|\mu\nu)) + \sum_{B \neq A} (P_{BB} - Z_B)\gamma_{AB} \quad (2.269)$$

$$F_{\mu\nu} = -\frac{1}{2}P_{\mu\nu}(\mu\nu|\mu\nu); \quad \mu \text{ and } \nu \text{ both on the same atom A} \quad (2.270)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} - \frac{1}{2}P_{\mu\nu}(\mu\nu|\mu\nu) = H_{\mu\nu}^{\text{core}} - \frac{1}{2}P_{\mu\nu}\gamma_{AB}; \quad \mu \text{ on A and } \nu \text{ on B} \quad (2.271)$$

The two-centre repulsion integrals γ_{AB} in MINDO/3 are calculated using the following function:

$$\gamma_{AB} = \frac{e^2}{\left[R_{AB}^2 + \frac{1}{4} \left(\frac{e^2}{\bar{g}_A} + \frac{e^2}{\bar{g}_B} \right)^2 \right]^{1/2}} \quad (2.272)$$

\bar{g}_A is the average of the one-centre, two-electron integrals $g_{\mu\nu}$ on atom A (i.e. $g_{\mu\nu} \equiv (\mu\mu|\nu\nu)$) and \bar{g}_B is the equivalent average for atom B. This seemingly complex function for γ_{AB} is, in fact, quite simple; at large R_{AB} it tends towards the Coulomb's law expression e^2/R_{AB} and as R_{AB} tends to zero it approaches the average of the one-centre integrals on the two atoms. The two-centre, one-electron integrals $H_{\mu\nu}^{\text{core}}$ are given in MINDO/3 by:

$$H_{\mu\nu}^{\text{core}} = S_{\mu\nu}\beta_{AB}(I_\mu + I_\nu) \quad (2.273)$$

$S_{\mu\nu}$ is the overlap integral, I_μ and I_ν are ionisation potentials for the appropriate orbitals and β_{AB} is a parameter dependent upon both of the two atoms A and B.

The core-core interaction between pairs of nuclei was also changed in MINDO/3 from the form used in CNDO/2. One way to correct the fundamental problems with CNDO/2 such as the repulsion between two hydrogen atoms (or indeed any neutral molecules) at all distances is to change the core-core repulsion term from a simple Coulombic expression ($E_{AB} = Z_A Z_B / R_{AB}$) to:

$$E_{AB} = Z_A Z_B \gamma_{AB} \quad (2.274)$$

In fact, while this correction gives the desired behaviour at relatively long separations, it does not account for the fact that as two nuclei approach each other the screening by the core electrons decreases. As the separation approaches zero the core-core repulsion should be described by Coulomb's law. In MINDO/3 this is achieved by making the core-core interaction a function of the electron-electron repulsion integrals as follows:

$$E_{AB} = Z_A Z_B \{ \gamma_{AB} + [(e^2/R_{AB}) - \gamma_{AB}] \exp(-\alpha_{AB} R_{AB}) \} \quad (2.275)$$

α_{AB} is a parameter dependent upon the nature of the atoms A and B. For OH and NH bonds a slightly different core-core interaction was found to be more appropriate:

$$E_{XH} = Z_X Z_H \{ \gamma_{XH} + [(e^2/R_{XH}) - \gamma_{XH}] \alpha_{XH} \exp(-R_{XH}) \} \quad (2.276)$$

The parameters for MINDO/3 were obtained in an entirely different way from previous semi-empirical methods. Some of the values that were fixed in CNDO, INDO and NNDO were permitted to vary during the MINDO/3 parametrisation procedure. For example, the exponents of the Slater atomic orbitals were allowed to vary from the values given by Slater's rules, and indeed the exponents for s and p orbitals were not required to be the

same. $U_{\mu\mu}$ and β_{AB} were also regarded as variable parameters. Another key difference was that the MINDO/3 parametrisation used experimental data such as molecular geometries and heats of formation, rather than theoretical values from *ab initio* calculations or data from atomic spectra. The parametrisation effort was a considerable undertaking, and it was only at the fourth attempt that an acceptable model was obtained (as is implicit in the appearance of the '3' in the name). For example, just to parametrise two atoms such as carbon and hydrogen using a set of 20 molecules required between 30 000 and 50 000 SCF calculations for each parametrisation scheme that was investigated.

2.9.6 MNDO

MINDO/3 proved to be very successful when it was introduced; it is important to realise that even simple *ab initio* calculations were beyond the computational resources of all but a few research groups in the 1970s. However, there were some significant limitations. For example, heats of formation of unsaturated molecules were consistently too positive, the errors in calculated bond angles were often quite large, and the heats of formation for molecules containing adjacent atoms with lone pairs were too negative. Some of these limitations were due to the use of the INDO approximation, and in particular the inability of INDO to deal with systems containing lone pairs. Dewar and Thiel therefore introduced the modified neglect of diatomic overlap (MNDO) method, which was based on NDDO [Dewar and Thiel 1977a, b]. The Fock matrix elements in MNDO were as follows:

$$F_{\mu\mu} = H_{\mu\mu}^{\text{core}} + \sum_{\nu \text{ on A}} [P_{\nu\nu}(\mu\mu|\nu\nu) - \frac{1}{2}P_{\nu\nu}(\mu\nu|\mu\nu)] + \sum_{B \neq A} \sum_{\lambda \text{ on B}} \sum_{\sigma \text{ on B}} P_{\lambda\sigma}(\mu\mu|\lambda\sigma) \quad (2.277)$$

$$\text{where } H_{\mu\mu}^{\text{core}} = U_{\mu\mu} - \sum_{B \neq A} V_{\mu\mu B} \quad (2.278)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} + \frac{3}{2}P_{\mu\nu}(\mu\nu|\mu\nu) - \frac{1}{2}P_{\mu\nu}(\mu\mu|\nu\nu) + \sum_{B \neq A} \sum_{\lambda \text{ on B}} \sum_{\sigma \text{ on B}} P_{\lambda\sigma}(\mu\nu|\lambda\sigma); \quad \mu \text{ and } \nu \text{ both on A} \quad (2.279)$$

$$\text{where } H_{\mu\nu}^{\text{core}} = - \sum_{B \neq A} V_{\mu\nu B} \quad (2.280)$$

$$F_{\mu\nu} = H_{\mu\nu}^{\text{core}} - \frac{1}{2} \sum_{\lambda \text{ on B}} \sum_{\sigma \text{ on A}} P_{\lambda\sigma}(\mu\sigma|\nu\lambda); \quad \mu \text{ on A and } \nu \text{ on B} \quad (2.281)$$

$$\text{where } H_{\mu\nu}^{\text{core}} = \frac{1}{2}S_{\mu\nu}(\beta_{\mu} + \beta_{\nu}) \quad (2.282)$$

The similarity with the NDDO expressions, Equations (2.264)–(2.266), can clearly be seen; the major new features are the appearance of terms $V_{\mu\mu B}$ and $V_{\mu\nu B}$ and a new form for the two-centre, one-electron core resonance integrals, which depend upon the overlap $S_{\mu\nu}$ and parameters β_{μ} and β_{ν} , as shown in Equation (2.282). $V_{\mu\mu B}$ and $V_{\mu\nu B}$ are two-centre, one-electron attractions between an electron distribution $\phi_{\mu}\phi_{\mu}$ or $\phi_{\mu}\phi_{\nu}$, respectively, on atom A and the core of atom B. These are expressed as follows:

$$V_{\mu\mu B} = -Z_B(\mu_A\mu_A|s_Bs_B) \quad (2.283)$$

$$V_{\mu\nu B} = -Z_B(\mu_A\nu_A|s_Bs_B) \quad (2.284)$$

The core–core repulsion terms are also different in MNDO from those in MINDO/3, with OH and NH bonds again being treated separately:

$$E_{AB} = Z_A Z_B (s_A s_A | s_B s_B) \{1 + \exp(-\alpha_A R_{AB}) + \exp(-\alpha_B R_{AB})\} \quad (2.285)$$

$$E_{XH} = Z_X Z_H (s_X s_X | s_H s_H) \{1 + R_{XH} \exp(-\alpha_X R_{XH}) / R_{AB} + \exp(-\alpha_H R_{XH})\} \quad (2.286)$$

Perhaps the most significant advantage of MNDO over MINDO/3 is the use throughout of monatomic parameters; MINDO/3 requires diatomic parameters in the resonance integral (β_{AB}) and the core–core repulsion (α_{AB}). It has been possible to expand MNDO to cover a much wider variety of elements such as aluminium, silicon, germanium, tin, bromine and lead. However, the use of an (s, p) basis set in the original MNDO method did mean that the method could not be applied to most transition metals, which require a basis set containing d orbitals. In addition, hypervalent compounds of sulphur and phosphorus are not modelled well. In more recent versions of the MNDO method d orbitals have been explicitly included for the heavier elements [Thiel and Voityuk 1994]. Another serious limitation of MNDO is its inability to accurately model intermolecular systems involving hydrogen bonds (for example, the heat of formation of the water dimer is far too low in MNDO). This is because of a tendency to overestimate the repulsion between atoms when they are separated by a distance approximately equal to the sum of their van der Waals radii. Conjugated systems can also present difficulties for MNDO. An extreme example of this occurs with compounds such as nitrobenzene in which the nitro group is predicted to be orthogonal to the aromatic ring rather than conjugated with it. In addition, MNDO energies are too positive for sterically crowded molecules and too negative for molecules containing four-membered rings.

2.9.7 AM1

The Austin Model 1 (AM1) model was the next semi-empirical theory produced by Dewar's group [Dewar *et al.* 1985]. AM1 was designed to eliminate the problems with MNDO, which were considered to arise from a tendency to overestimate repulsions between atoms separated by distances approximately equal to the sum of their van der Waals radii. The strategy adopted was to modify the core–core term using Gaussian functions. Both attractive and repulsive Gaussian functions were used; the attractive Gaussians were designed to overcome the repulsion directly and were centred in the region where the repulsions were too large. Repulsive Gaussian functions were centred at smaller internuclear separations. With this modification the expression for the core–core term was related to the MNDO expression by:

$$E_{AB} = E_{\text{MINDO}} + \frac{Z_A Z_B}{R_{AB}} \times \left\{ \sum_i K_{A_i} \exp[-L_{A_i}(R_{AB} - M_{A_i})^2] + \sum_j K_{B_j} \exp[-L_{B_j}(R_{AB} - M_{B_j})^2] \right\} \quad (2.287)$$

The additional terms are spherical Gaussian functions with a width determined by the parameter L . It was found that the values of these parameters were not critical and many

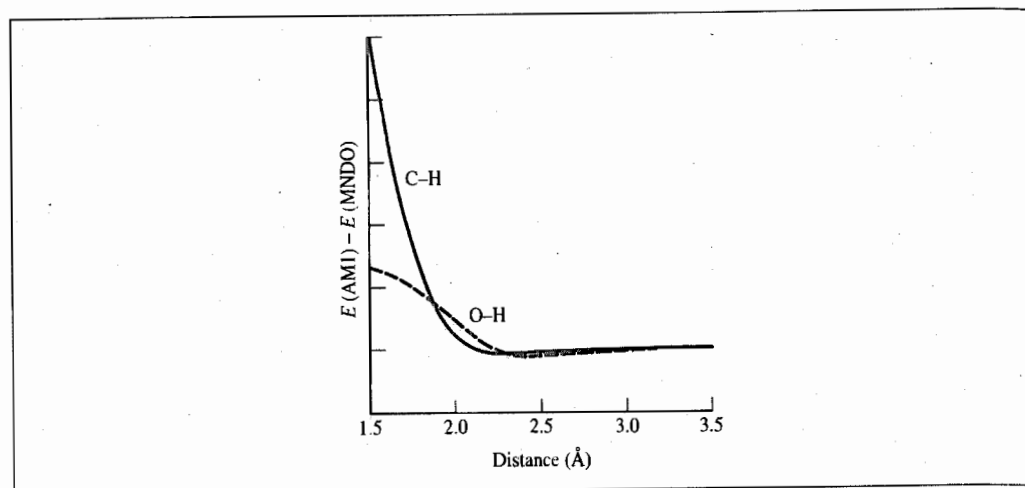


Fig. 2.20: The difference in the core-core energy for AM1 and MNDO for carbon-hydrogen and oxygen-hydrogen interactions.

were set to the same value. The M and K parameters were optimised for each atom, together with the α parameters in the exponential terms in Equations (2.285) and (2.286). In the original parametrisation of AM1 there are four terms in the Gaussian expansion for carbon, three for hydrogen and nitrogen and two for oxygen (both attractive and repulsive Gaussians were used for carbon, hydrogen and nitrogen but only repulsive Gaussians for oxygen). The effect of including these Gaussian functions can be seen in Figure 2.20, which plots the difference in the MNDO and AM1 core-core terms for the carbon-hydrogen and oxygen-hydrogen interactions. The inclusion of these Gaussians significantly increased the number of parameters per atom, from seven in the MNDO to between 13 and 16 per atom in AM1. This, of course, made the parametrisation process considerably more difficult. Overall, AM1 was a significant improvement over MNDO and many of the deficiencies associated with the core repulsion were corrected.

2.9.8 PM3

PM3 is also based on MNDO (the name derives from the fact that it is the third parametrisation of MNDO, AM1 being considered the second) [Stewart 1989a, b]. The PM3 Hamiltonian contains essentially the same elements as that for AM1, but the parameters for the PM3 model were derived using an automated parametrisation procedure devised by JJP Stewart. By contrast, many of the parameters in AM1 were obtained by applying chemical knowledge and 'intuition'. As a consequence, some of the parameters have significantly different values in AM1 and PM3, even though both methods use the same functional form and they both predict various thermodynamic and structural properties to approximately the same level of accuracy. Some problems do remain with PM3. One of the most important of these is the rotational barrier of the amide bond, which is much too low and in some cases almost non-existent. This problem can be corrected through the use of an

empirical torsional potential (see Section 4.5). There has been considerable debate over the relative merits of the AM1 and PM3 approaches to parametrisation.

2.9.9 SAM1

The final offering from the Dewar group* was SAM1, which stands for 'Semi-Ab-initio Model 1' [Dewar *et al.* 1993]. The name was chosen to reflect Dewar's belief that methods like AM1 offer such a significant enhancement over the earlier semi-empirical methods like CNDO/2 that they should be given a different generic name. In SAM1 a standard STO-3G Gaussian basis set is used to evaluate the electron repulsion integrals; close inspection of the results from AM1 and MNDO suggested that steric effects were overestimated because of the way in which the electron repulsion integrals were calculated. The resulting integrals were then scaled, partly to enable some of the effects of electron correlation to be included and partly to compensate for the use of a minimal basis set. The Gaussian terms in the core-core repulsion were retained to fine-tune the model. The number of parameters in SAM1 is no greater than in AM1 and fewer than in PM3. It does take longer to run (by up to two orders of magnitude) though it was felt that with the improvements in computer hardware such an increase was acceptable.

2.9.10 Programs for Semi-empirical Quantum Mechanical Calculations

The popularity of the MNDO, AM1 and PM3 methods is due in large part to their implementation in the MOPAC and AMPAC programs. The programs are able to perform many kinds of calculation and to calculate many different properties.

The contributions of the Dewar group are rightly recognised as particularly significant in the development of semi-empirical methods, but other research groups have also made important contributions. The SINDO1 and ZINDO programs have been developed in the groups of Jug and Zerner, respectively, and both contain novel features. The ZINDO program of Zerner and co-workers can perform a wide variety of semi-empirical calculations and has been particularly useful for calculations on transition metal and lanthanide compounds and for predicting molecular electronic spectra.

2.10 Hückel Theory

Hückel theory can be considered the 'grandfather' of approximate molecular orbital methods, having been formulated in the early 1930s [Hückel 1931]. Hückel theory is limited to conjugated π systems and was originally devised to explain the non-additive nature of certain properties of aromatic compounds. For example, the properties of benzene are much different from those of the hypothetical 'cyclohexatriene' molecule. Although Hückel theory, as originally formulated, is relatively little used in research today, extensions

*Michael Dewar died in 1997.

to it such as extended Hückel theory are still employed and can provide qualitative insights into the electronic structure of important classes of molecule. Hückel theory is also widely used for teaching purposes to introduce a 'real' theory that can be applied to relatively complex systems with little more than pencil and paper or a simple computer program.

Hückel theory separates the π system from the underlying σ framework and constructs molecular orbitals into which the π electrons are then fed in the usual way according to the Aufbau principle. The π electrons are thus considered to be moving in a field created by the nuclei and the 'core' of σ electrons. The molecular orbitals are constructed from linear combinations of atomic orbitals and so the theory is an LCAO method. For our purposes it is most appropriate to consider Hückel theory in terms of the CNDO approximation (in fact, Hückel theory was the first ZDO molecular orbital theory to be developed). Let us examine the three types of Fock matrix element in Equations (2.252)–(2.254). First, $F_{\mu\mu}$. In a neutral species, the net charge on each atom will be approximately zero, and so if we take Equation (2.258), from which penetration effects have been eliminated, then we are left with $U_{\mu\mu} + (P_{AA} - 0.5P_{\mu\mu})\gamma_{AA}$. Now if each nucleus (A) in the π system is the same (i.e. carbon) then this expression will be approximately constant for all nuclei being considered. The matrix elements $F_{\mu\mu}$ are often (confusingly) called Coulomb integrals in Hückel theory and are assigned the symbol α . All off-diagonal elements of the Fock matrix are assumed to be zero with the exception of elements $F_{\mu\nu}$, where μ and ν are π orbitals on two bonded atoms. These $F_{\mu\nu}$ are assumed to be constant, are assigned the symbol β and are known as resonance integrals. The Fock matrix in Hückel theory thus has as many rows and columns as the number of atoms in the π system with diagonal elements that are all set to α . All off-diagonal elements F_{ij} are zero unless there is a bond between the atoms i and j , in which case the element is β . For benzene the Fock matrix is of the following form (atom labelling as in Figure 2.21):

$$\begin{pmatrix} \alpha & \beta & 0 & 0 & 0 & \beta \\ \beta & \alpha & \beta & 0 & 0 & 0 \\ 0 & \beta & \alpha & \beta & 0 & 0 \\ 0 & 0 & \beta & \alpha & \beta & 0 \\ 0 & 0 & 0 & \beta & \alpha & \beta \\ \beta & 0 & 0 & 0 & \beta & \alpha \end{pmatrix} \quad (2.288)$$

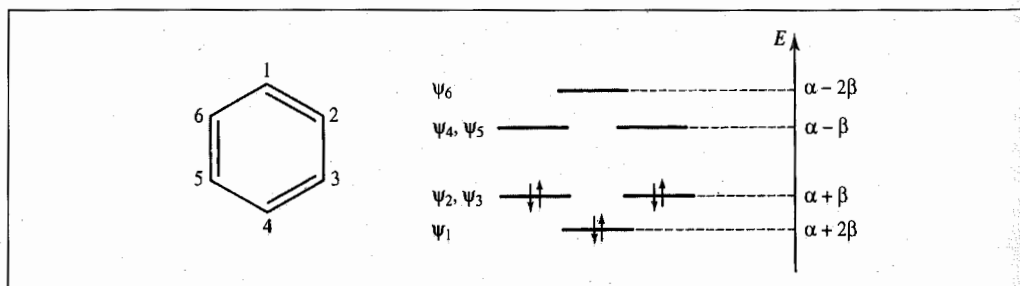


Fig. 2.21: Benzene and its Hückel molecular orbitals.

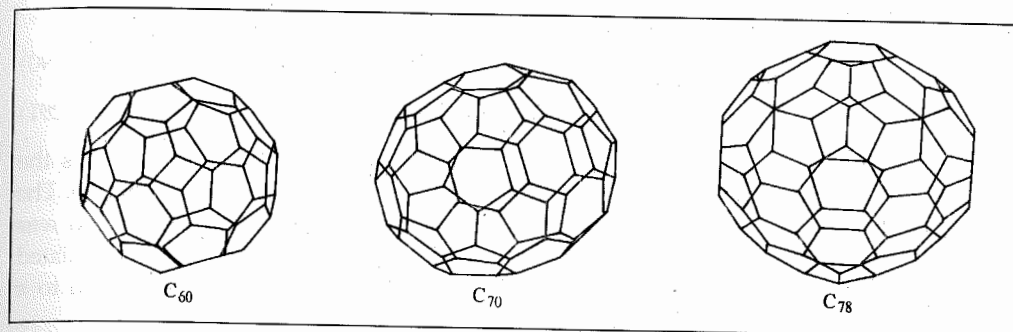


Fig. 2.22: Three fullerenes, C_{60} , C_{70} and C_{78} .

As with the other semi-empirical methods that we have considered so far, the overlap matrix is equal to the identity matrix. The following simple matrix equation must then be solved:

$$FC = CE \quad (2.289)$$

The equation can be solved by standard methods to give the basis set coefficients and the molecular orbital energies E . The orbital energies for benzene are $E_1 = \alpha + 2\beta$; $E_2, E_3 = \alpha + \beta$; $E_4, E_5 = \alpha - \beta$; $E_6 = \alpha - 2\beta$, and so the ground state places two electrons in ψ_1 and two each in the two degenerate orbitals ψ_2 and ψ_3 . The lowest-energy orbital ψ_1 is a linear combination of the six carbon p orbitals.

Hückel theory was extended to cover various other systems, including those with heteroatoms, but it was not particularly successful and has largely been superseded by other semi-empirical methods. Nevertheless, for appropriate problems Hückel theory can be very useful. One example is the calculations of P W Fowler and colleagues, who studied the relationship between geometry and electronic structure for a range of buckminsterfullerenes (the parent molecule of which, C_{60} , was discovered in 1985) [Fowler 1993]. The fullerenes (or 'buckyballs') are excellent candidates for Hückel theory as they are composed of carbon and have extensive π systems; three examples are shown in Figure 2.22.

The results of their calculations were summarised in two rules. The first rule states that at least one isomer C_n with a properly closed p shell (i.e. bonding HOMO, antibonding LUMO) exists for all $n = 60 + 6k$ ($k = 0, 2, 3, \dots$, but not 1). Thus C_{60} , C_{72} , C_{78} , etc., are in this group. The second rule is for carbon cylinders and states that a closed-shell structure is found for $n = 2p(7 + 3k)$ (for all k). C_{70} is the parent of this family. The calculations were extended to cover different types of structure and fullerenes doped with metals.

2.10.1 Extended Hückel Theory

Hückel theory is clearly limited, in part because it is restricted to π systems. The extended Hückel method is a molecular orbital theory that takes account of all the valence electrons in the molecule [Hoffmann 1963]. It is largely associated with R Hoffmann, who received the Nobel Prize for his contributions. The equation to be solved is $FC = SCE$, with the

Fock matrix elements taking the following simple forms:

$$F_{\mu\mu}^{AA} = H_{\mu\mu} = -I_{\mu} \quad (2.290)$$

$$F_{\mu\nu}^{AB} = H_{\mu\nu} = -\frac{1}{2}K(I_{\mu} + I_{\nu})S_{\mu\nu} \quad (2.291)$$

In these equations, μ and ν are two atomic orbitals (e.g. Slater type orbitals), I_{μ} is the ionisation potential of the orbital and K is a constant, which was originally set to 1.75. The formula for the off-diagonal elements $H_{\mu\nu}$ (where μ and ν are on different atoms) was originally suggested by R S Mulliken. These off-diagonal matrix elements are calculated between all pairs of valence orbitals and so extended Hückel theory is not limited to π systems.

The extended Hückel approach has proved to be rather successful for such a simple theory; for example, the famous Woodward-Hoffmann rules (see Section 5.9.4) were based upon calculations using this model. Extended Hückel theory has found particular application in those areas where alternative theories cannot be used. This is largely due to the fact that the basis set requires no more than experimentally determined ionisation potentials. It is particularly useful for studying systems containing metals; these systems are problematic for many other methods due to the lack of suitable basis sets.

2.11 Performance of Semi-empirical Methods

Our discussion of the application of quantum mechanics calculations was not explicitly directed towards any particular quantum mechanical theory but was - implicitly at least - written with *ab initio* methods in mind. All of the properties we considered in Section 2.7 can also be determined using semi-empirical methods. Extensive tables detailing the performance of the popular semi-empirical methods have been published, both in the original papers and in review articles, some of which are listed at the end of this chapter. The parametrisation of the semi-empirical approaches typically includes geometrical variables, dipole moments, ionisation energies and heats of formation. In Table 2.7 we provide a summary of the performance of the MINDO/3, MNDO, AM1, PM3 and SAM1 semi-empirical methods from data supplied in the original publications. The performance of successive semi-empirical methods has gradually improved from one method to another, though one should always remember that anomalous results may be obtained for certain types of system. Some of these limitations were outlined in the discussion of the various semi-empirical methods. It is worth emphasising that some of the major drawbacks with the semi-empirical methods arise simply because one is trying to calculate properties that were not given a major consideration in the parametrisation process. For example, many of the molecules used for the parametrisation of the MNDO, AM1 and PM3 methods had little or no conformational flexibility and it is therefore not so surprising that some rotational barriers are not calculated with the same accuracy as (say) heats of formation. In addition, to achieve optimal performance for specific classes of molecules (e.g. the amino acids) or specific properties (e.g. conformational barriers) then it would be appropriate to include representative systems during the parametrisation procedure.

	MINDO/3	MNDO	AM1	PM3	SAM1	Reference
138 heats of formation (kcal/mol)	11.0	6.3				[Dewar and Thiel 1977b]
228 bond lengths	0.022 Å	0.014 Å				
91 angles	5.6°	2.8°				
57 dipole moments	0.49 D	0.38 D				
58 heats of formation of hydrocarbons (kcal/mol)	9.7	5.87	5.07			[Dewar et al. 1985]
80 heats of formation for species with N and/or O (kcal/mol)	11.69	6.64	5.88			
46 dipole moments	0.54 D	0.32 D	0.26 D			
29 ionisation energies	0.31 eV	0.39 eV	0.29 eV			
406 heats of formation (kcal/mol)			8.82	7.12	5.21	[Dewar et al. 1993]
196 dipole moments			0.35 D	0.40 D	0.32 D	

Table 2.7 Comparison of quantities calculated with various semi-empirical methods.

Appendix 2.1 Some Common Acronyms Used in Computational Quantum Chemistry

AM1	Austin Model 1
AO	Atomic orbital
B3LYP	Scheme for hybrid Hartree-Fock/density functional theory introduced by Becke
BLYP	Becke-Lee-Yang-Parr gradient-corrected functional for use with density functional theory
BSSE	Basis set superposition error
CASSCF	Complete active space self-consistent field
CI	Configuration interaction
CIS	Configuration interaction singles
CISD	Configuration interaction singles and doubles
CNDO	Complete neglect of differential overlap
DFT	Density functional theory
DIIS	Direct inversion of iterative subspace
DVP	Double zeta with polarisation
DZ	Double zeta
EHT	Extended Hückel theory
GVB	Generalised valence bond model
HF	Hartree-Fock
HOMO	Highest occupied molecular orbital
INDO	Intermediate neglect of differential overlap
LCAO	Linear combination of atomic orbitals
LDA	Local density approximation
LSDFT	Local spin density functional theory
LUMO	Lowest unoccupied molecular orbital
MBPT	Many-body perturbation theory
MINDO/3	Modified INDO version 3
MNDO	Modified neglect of diatomic overlap
MO	Molecular orbital
MP	Møller-Plesset
MP2, MP3, etc.	Møller-Plesset theory at second order, third order, etc.
NDDO	Neglect of diatomic differential overlap
PM3	Parametrisation 3 of MNDO
QCISD	Quadratic configuration interaction singles and doubles
QCISD(T)	Configuration interaction method involving single, double and quadratic excitations with an estimated triple excitation
RHF	Restricted Hartree-Fock
SAM1	Semi- <i>Ab initio</i> Model 1
SCF	Self-consistent field
STO	Slater type orbital
STO-3G, STO-4G, etc.	Minimal basis sets in which 3, 4 etc, Gaussian functions are used to represent the atomic orbitals on an atom

UHF	Unrestricted Hartree-Fock
WVN	Correlation functional due to Wilk, Vosko and Nusair
ZDO	Zero differential overlap

Further Reading

- Atkins P W 1991. *Quanta: A Handbook of Concepts*. Oxford, Oxford University Press.
- Atkins P W 1998. *Physical Chemistry*. 6th Edition. Oxford, Oxford University Press.
- Atkins P W and R S Friedman 1996. *Molecular Quantum Mechanics*. Oxford, Oxford University Press.
- Clark T 1985. *A Handbook of Computational Chemistry: A Practical Guide to Chemical Structure and Energy Calculations*. New York, Wiley-Interscience.
- Dewar M J S 1969. *The Molecular Orbital Theory of Organic Chemistry*. New York, McGraw-Hill.
- Hinchliffe A 1988. *Computational Quantum Chemistry*. Chichester, John Wiley & Sons.
- Hinchliffe A 1995. *Modelling Molecular Structures*. Chichester, John Wiley & Sons.
- Hirst D M 1990. *A Computational Approach to Chemistry*. Oxford, Blackwell Scientific.
- Pople J A and D L Beveridge 1970. *Approximate Molecular Orbital Theory*. New York, McGraw-Hill.
- Richards W G and D L Cooper 1983. *Ab initio Molecular Orbital Calculations for Chemists*. 2nd Edition. Oxford, Clarendon Press.
- Schaeffer H F III (Editor) 1977. *Applications of Electronic Structure Theory*. New York, Plenum Press.
- Schaeffer H F III (Editor) 1977. *Methods of Electronic Structure Theory*. New York, Plenum Press.
- Stewart J J P 1990. MOPAC: A Semi-Empirical Molecular Orbital Program. *Journal of Computer-Aided Molecular Design* 4:1-45.
- Stewart J J P 1990. Semi-empirical Molecular Orbital Methods. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry Volume 1*. New York, VCH Publishers, pp 45-82.
- Szabo A and N S Ostlund 1982. *Modern Quantum Chemistry. Introduction to Advanced Electronic Structure Theory*. New York, McGraw-Hill.
- Zerner M C 1991. Semi-empirical Molecular Orbital Methods. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry Volume 2*. New York, VCH Publishers, pp 313-366.

References

- Allinger N L, R S Grev, B F Yates and H F Schaeffer III 1990. The Syn Rotational Barrier in Butane. *Journal of the American Chemical Society* 112:114-118.
- Bachrach S M 1994. Population Analysis and Electron Densities from Quantum Mechanics. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry Volume 5*. New York, VCH Publishers, pp 171-227.
- Bader R F W 1985. Atoms in Molecules. *Accounts of Chemistry Research* 18:9-15.
- Bingham R C, M J S Dewar and D H Lo 1975a. Ground States of Molecules. XXV. MINDO/3. An Improved Version of the MINDO Semi-empirical SCFMO Method. *Journal of the American Chemical Society* 97:1285-1293.
- Bingham R C, M J S Dewar and D H Lo 1975b. Ground States of Molecules. XXVI. MINDO/3. Calculations for Hydrocarbons. *Journal of the American Chemical Society* 97:1294-1301.
- Bingham R C, M J S Dewar and D H Lo 1975c. Ground States of Molecules. XXVII. MINDO/3. Calculations for CHON Species. *Journal of the American Chemical Society* 97:1302-1306.
- Bingham R C, M J S Dewar and D H Lo 1975d. Ground States of Molecules. XXVIII. MINDO/3. Calculations for Compounds Containing Carbon, Hydrogen, Fluorine and Chlorine. *Journal of the American Chemical Society* 97:1307-1310.

- Boys S F 1950. Electronic Wave Functions. I. A General Method of Calculation for the Stationary States of Any Molecular System. *Proceedings of the Royal Society (London)* **A200**:542-554.
- Cusachs L C and Politzer 1968. On the Problem of Defining the Charge on an Atom in a Molecule. *Chemical Physics Letters* **1**:529-531.
- Dewar M J S, C Jie and J Yu 1993. SAM1; The First of a New Series of General Purpose Quantum Mechanical Molecular Models. *Tetrahedron* **49**:5003-5038.
- Dewar M J S and Thiel W 1977a. Ground States of Molecules. 38. The MNDO Method. Approximations and Parameters. *Journal of the American Chemical Society* **99**:4899-4907.
- Dewar M J S and Thiel W 1977b. Ground States of Molecules. 39. MNDO Results for Molecules Containing Hydrogen, Carbon, Nitrogen and Oxygen. *Journal of the American Chemical Society* **99**:4907-4917.
- Dewar M J S, E G Zoebisch, E F Healy and J J P Stewart 1985. AM1: A New General Purpose Quantum Mechanical Model. *Journal of the American Chemical Society* **107**:3902-3909.
- Dunning T H Jr 1970. Gaussian Basis Functions for Use in Molecular Calculations. I. Contraction of (9s5p) Atomic Basis Sets for First-Row Atoms. *Journal of Chemical Physics* **53**:2823-2883.
- Fowler P W 1993. Systematics of Fullerenes and Related Clusters. *Philosophical Transactions of the Royal Society (London)* **A343**:39-52.
- Hall G G 1951. The Molecular Orbital Theory of Chemical Valency VIII. A Method for Calculating Ionisation Potentials. *Proceedings of the Royal Society (London)* **A205**:541-552.
- Hehre W J, R F Stewart and J A Pople 1969. Self-Consistent Molecular-Orbital Methods. I. Use of Gaussian Expansions of Slater-Type Atomic Orbitals. *Journal of Chemical Physics* **51**:2657-2664.
- Hehre W J, L Radom, P v R Schleyer and J A Pople 1986. *Ab initio Molecular Orbital Theory*. New York, John Wiley & Sons.
- Hoffmann R 1963. An Extended Hückel Theory. I. Hydrocarbons. *Journal of Chemical Physics* **39**:1397-1412.
- Hückel Z 1931. Quanten theoretische Beiträge zum Benzolproblem. I. Die Electron enkonfiguration des Benzols. *Zeitschrift für Physik*. **70**:203-286.
- Huzinga S 1965. Gaussian-type Functions for Polyatomic Systems. I. *Journal of Chemical Physics* **42**:1293-1302.
- Löwdin P-Q 1970. On the Orthogonality Problem. *Advances in Quantum Chemistry* **5**:185-199.
- Mayer I 1983. Charge, Bond Order and Valence in the *Ab initio* SCF Theory. *Chemical Physics Letters* **97**:270-274.
- Mulliken R S 1955. Electronic Population Analysis on LCAO-MO Molecular Wave Functions. I. *Journal of Chemical Physics* **23**:1833-1846.
- Politzer P and J S Murray 1991. Molecular Electrostatic Potentials and Chemical Reactivity. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 2. New York, VCH Publishers, pp 273-312.
- Pople J A, D L Beveridge and P A Dobosh 1967. Approximate Self-Consistent Molecular Orbital Theory. V. Intermediate Neglect of Differential Overlap. *Journal of Chemical Physics* **47**:2026-2033.
- Pople J A, D P Santry and G A Segal 1965. Approximate Self-Consistent Molecular Orbital Theory. I. Invariant Procedures. *Journal of Chemical Physics* **43**:S129-S135.
- Pople J A and G A Segal 1965. Approximate Self-Consistent Molecular Orbital Theory. II. Calculations with Complete Neglect of Differential Overlap. *The Journal of Chemical Physics* **43**:S136-S149.
- Pople J A and G A Segal 1966. Approximate Self-Consistent Molecular Orbital Theory. III. CNDO Results for AB₂ and AB₃ systems. *Journal of Chemical Physics* **44**:3289-3296.
- Reed A E, R B Weinstock and F Weinhold 1985. Natural Population Analysis. *Journal of Chemical Physics* **83**: 735-746.
- Roothaan C C J 1951. New Developments in Molecular Orbital Theory. *Reviews of Modern Physics* **23**:69-89.
- Slater J C 1930. Atomic Shielding Constants. *Physical Review* **36**:57-64.

- Smith G D and R L Jaffe 1996. Quantum Chemistry Study of Conformational Energies and Rotational Energy Barriers in *n*-Alkanes. *Journal of Physical Chemistry* **100**:18718-18724.
- Stewart J J P 1989a. Optimisation of Parameters for Semi-empirical Methods I. Method. *Journal of Computational Chemistry* **10**:209-220.
- Stewart J J P 1989b. Optimisation of Parameters for Semi-empirical Methods II. Applications. *Journal of Computational Chemistry* **10**:221-264.
- Thiel W and A A Voityuk 1994. Extension of MNDO to d Orbitals: Parameters and Results for Silicon. *Journal of Molecular Structure (Theochem)* **313**:141-154.
- Wiberg K B and M A Murcko 1988. Rotational Barriers. 2. Energies of Alkane Rotamers. An Examination of Gauche Interactions. *Journal of the American Chemical Society* **110**:8029-8038.
- Wiberg K B and P R Rablen 1993. Comparison of Atomic Charges Derived via Different Procedures. *Journal of Computational Chemistry* **14**:1504-1518.

CHAPTER THREE

Advanced *ab initio* Methods, Density Functional Theory and Solid-state Quantum Mechanics

3.1 Introduction

In Chapter 2 we worked through the two most commonly used quantum mechanical models for performing calculations on ground-state 'organic'-like molecules, the *ab initio* and semi-empirical approaches. We also considered some of the properties that can be calculated using these techniques. In this chapter we will consider various advanced features of the *ab initio* approach and also examine the use of density functional methods. Finally, we will examine the important topic of how quantum mechanics can be used to study the solid state.

3.2 Open-shell Systems

The Roothaan-Hall equations are not applicable to open-shell systems, which contain one or more unpaired electrons. Radicals are, by definition, open-shell systems as are some ground-state molecules such as NO and O₂. Two approaches have been devised to treat open-shell systems. The first of these is *spin-restricted* Hartree-Fock (RHF) theory, which uses combinations of singly and doubly occupied molecular orbitals. The closed-shell approach that we have developed thus far is a special case of RHF theory. The doubly occupied orbitals use the same spatial functions for electrons of both α and β spin. The orbital expansion Equation (2.144) is employed together with the variational method to derive the optimal values of the coefficients. The alternative approach is the *spin-unrestricted Hartree-Fock* (UHF) theory of Pople and Nesbet [Pople and Nesbet 1954], which uses two distinct sets of molecular orbitals: one for electrons of α spin and the other for electrons of β spin. Two Fock matrices are involved, one for each type of spin, with elements as follows:

$$F_{\mu\nu}^{\alpha} = H_{\mu\nu}^{\text{core}} + \sum_{\lambda=1}^K \sum_{\sigma=1}^K [[P_{\lambda\sigma}^{\alpha} + P_{\lambda\sigma}^{\beta}](\mu\nu|\lambda\sigma) - P_{\lambda\alpha}^{\alpha}(\mu\lambda|\nu\sigma)] \quad (3.1)$$

$$F_{\mu\nu}^{\beta} = H_{\mu\nu}^{\text{core}} + \sum_{\lambda=1}^K \sum_{\sigma=1}^K [[P_{\lambda\sigma}^{\alpha} + P_{\lambda\sigma}^{\beta}](\mu\nu|\lambda\sigma) - P_{\lambda\alpha}^{\beta}(\mu\lambda|\nu\sigma)] \quad (3.2)$$

UHF theory also uses two density matrices, the full density matrix being the sum of these two:

$$P_{\mu\nu}^{\alpha} = \sum_{i=1}^{\alpha_{\text{occ}}} c_{\mu i}^{\alpha} c_{\nu i}^{\alpha} \quad P_{\mu\nu}^{\beta} = \sum_{i=1}^{\beta_{\text{occ}}} c_{\mu i}^{\beta} c_{\nu i}^{\beta} \quad (3.3)$$

$$P_{\mu\nu} = P_{\mu\nu}^{\alpha} + P_{\mu\nu}^{\beta} \quad (3.4)$$

The summations in Equations (3.3) and (3.4) are over the occupied orbitals with α and β spin as appropriate. Thus, $\alpha_{\text{occ}} + \beta_{\text{occ}}$ equals the total number of electrons in the system. In a closed-shell Hartree-Fock wavefunction the distribution of electron spin is zero everywhere because the electrons are paired. In an open-shell system, however, there is an excess of electron spin, which can be expressed as the spin density, analogous to the electron density. The spin density $\rho^{\text{spin}}(\mathbf{r})$ at a point \mathbf{r} is given by:

$$\rho^{\text{spin}}(\mathbf{r}) = \rho^{\alpha}(\mathbf{r}) - \rho^{\beta}(\mathbf{r}) = \sum_{\mu=1}^K \sum_{\nu=1}^K [P_{\mu\nu}^{\alpha} - P_{\mu\nu}^{\beta}] \phi_{\mu}(\mathbf{r}) \phi_{\nu}(\mathbf{r}) \quad (3.5)$$

Clearly, the UHF approach is more general and indeed the restricted Hartree-Fock approach is a special case of unrestricted Hartree-Fock. Figure 3.1 illustrates the conceptual difference between the RHF and the UHF models. Unrestricted wavefunctions are also the most appropriate way to deal with other problems such as molecules near the dissociation limit. The simplest example of this type of behaviour is the H₂ molecule, the ground state of which is a singlet with a bond length of approximately 0.75 Å. The restricted wavefunction is the appropriate Hartree-Fock wavefunction, with two paired electrons in a single spatial orbital. As the bond length increases towards the dissociation limit, this description is clearly inappropriate, for hydrogen is experimentally observed to dissociate to two

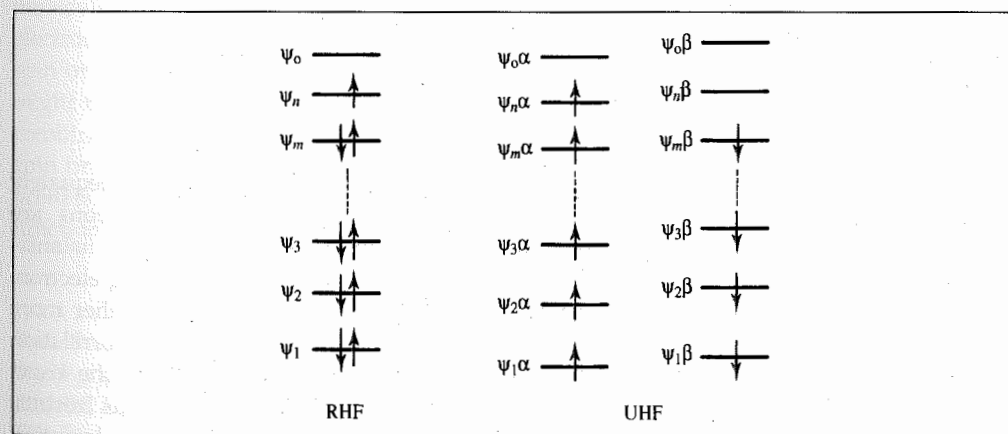


Fig. 3.1: The conceptual difference between the RHF and UHF models.

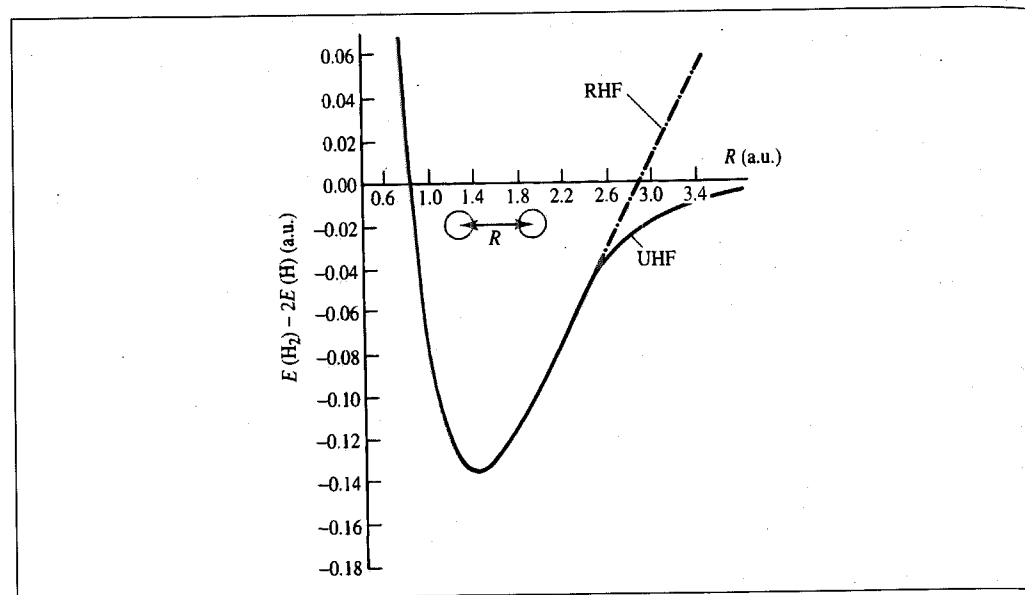


Fig. 3.2: UHF and RHF dissociation curves for H_2 . (Figure adapted from Szabo A, N S Ostlund 1982. Modern Quantum Chemistry. Introduction to Advanced Electronic Structure Theory. New York, McGraw-Hill.)

hydrogen atoms. This behaviour cannot be achieved using a restricted Hartree-Fock wavefunction, which requires the two electrons to occupy the same spatial orbital and leads to H^+ and H^- , but it is appropriately described by a UHF wavefunction. Beyond about 1.2 \AA the 'correct' wavefunction for hydrogen must thus be obtained using UHF theory. The results obtained by calculating the potential energy curves of the hydrogen molecule using the RHF and UHF theories are shown in Figure 3.2. As can be seen, RHF theory gives a dissociation energy that is much too large, whereas the UHF theory shows the correct dissociation behaviour.

3.3 Electron Correlation

The most significant drawback of Hartree-Fock theory is that it fails to adequately represent electron correlation. In the self-consistent field method the electrons are assumed to be moving in an average potential of the other electrons, and so the instantaneous position of an electron is not influenced by the presence of a neighbouring electron. In fact, the motions of electrons are correlated and they tend to 'avoid' each other more than Hartree-Fock theory would suggest, giving rise to a lower energy. The correlation energy is defined as the difference between the Hartree-Fock energy and the exact energy. Neglecting electron correlation can lead to some clearly anomalous results, especially as the dissociation limit is approached. For example, an uncorrelated calculation would predict that the electrons in H_2 spend equal time on both nuclei, even when they are

infinitely separated. Hartree-Fock geometries and relative energies for equilibrium structures are often in good agreement with experiment and as many molecular modelling applications are concerned with species at equilibrium it might be considered that correlation effects are not so important. Nevertheless, there is increasing evidence that the inclusion of correlation effects is warranted, especially when quantitative information is required. Moreover, electron correlation is crucial in the study of dispersive effects (which we shall consider in Section 4.10.1), which play a major role in intermolecular interactions. Electron correlation is most frequently discussed in the context of *ab initio* calculations, but it should be noted that the effects of electron correlation are implicitly included in the semi-empirical methods because of the way in which they are parametrised. However, specific electron correlation methods have also been developed for use with the various levels of semi-empirical calculation; this in turn necessitates the modification of some parameters.

3.3.1 Configuration Interaction

There are a number of ways in which correlation effects can be incorporated into an *ab initio* molecular orbital calculation. A popular approach is configuration interaction (CI), in which excited states are included in the description of an electronic state. To illustrate the principle, let us consider a lithium atom. The ground state of lithium can be written $1s^2 2s^1$ (although we have used the conventional nomenclature here, we should remember that the wavefunction is really a Slater determinant). Excitation of the outer valence electron gives states such as $1s^2 3s^1$. A better description of the overall wavefunction is a linear combination of the ground and excited-state wavefunctions. If a Hartree-Fock calculation is performed with K basis functions then $2K$ spin orbitals are obtained. If these $2K$ spin orbitals are filled with N electrons ($N < 2K$) there will be $2K - N$ unoccupied, virtual orbitals. The wavefunction obtained from the single-determinant approach that we have considered thus far is expressed only in terms of the occupied orbitals. For example, a very simple calculation on H_2 , using as a basis set just the $1s$ orbitals on each hydrogen, results in two molecular orbitals ($1\sigma_g$ and $1\sigma_u$). In the ground state, the $1\sigma_g$ orbital is filled with two electrons. An excited state can be generated by replacing one or more of the occupied spin orbitals with a virtual spin orbital. Possible excited states for the hydrogen molecule might thus include $1\sigma_g^1 \sigma_u^1$ and $1\sigma_u^2$ (in fact, the first of these two configurations cannot be combined with the ground state, as we shall see). In addition to the replacement of single spin orbitals by single virtual orbitals, two spin orbitals can be replaced by two virtual orbitals, three spin orbitals by three virtual orbitals, and so on. In general, the CI wavefunction can be written as:

$$\Psi = c_0 \Psi_0 + c_1 \Psi_1 + c_2 \Psi_2 + \dots \quad (3.6)$$

Ψ_0 is the single-determinant wavefunction obtained by solving the Hartree-Fock equations. Ψ_1, Ψ_2 , etc. are wavefunctions (expressed as determinants) that represent configurations derived by replacing one or more of the occupied spin orbitals by a virtual spin orbital. The energy of the system is then minimised in order to determine the coefficients c_0, c_1 , etc., using a linear variational approach, just as for a single-determinant calculation. A CI

calculation thus involves an additional level of complexity; each configuration is written in terms of molecular orbitals, which in turn are expressed as a linear combination of basis functions. The number of integrals can become extremely large. The total number of ways to permute N electrons and K orbitals is $(2K!)/[N!(2K - N)!]$. This is a very large number for all except small values of K and N , which explains why it is not usual to consider all possibilities (termed *full configuration interaction*) except for very small systems. However, full CI is important because it is the most complete treatment possible within the limitations imposed by the basis set. In the limit of a complete basis set full CI becomes complete CI and virtually exact – but is generally considered impractical as at large K the number of Slater determinants increases exponentially with N as $K^N/N!$. It is common practice to limit the excited states considered. For example, in configuration interaction singles (CIS) only wavefunctions that differ from the Hartree–Fock wavefunction by a single spin orbital are included. The next levels of the theory involve double substitutions (configuration interaction doubles, CID) or both singles and double substitutions (configuration interaction singles and doubles, CISD). Even at the CIS or CID levels, the number of excited states to be included can be very large, and it may be desirable (or necessary) to restrict the spin orbitals that are involved in the substitutions. For example, only excitations involving the highest occupied molecular orbital (HOMO) and the lowest unoccupied molecular orbital (LUMO) may be permitted. Alternatively, the orbitals corresponding to the inner electron core may be neglected (the ‘frozen core’ approximation). Some of these options are illustrated in Figure 3.3.

Not all excitations necessarily help to lower the energy; some determinants do not mix with the ground state. A consequence of *Brillouin’s theorem* is that single excitations do not mix

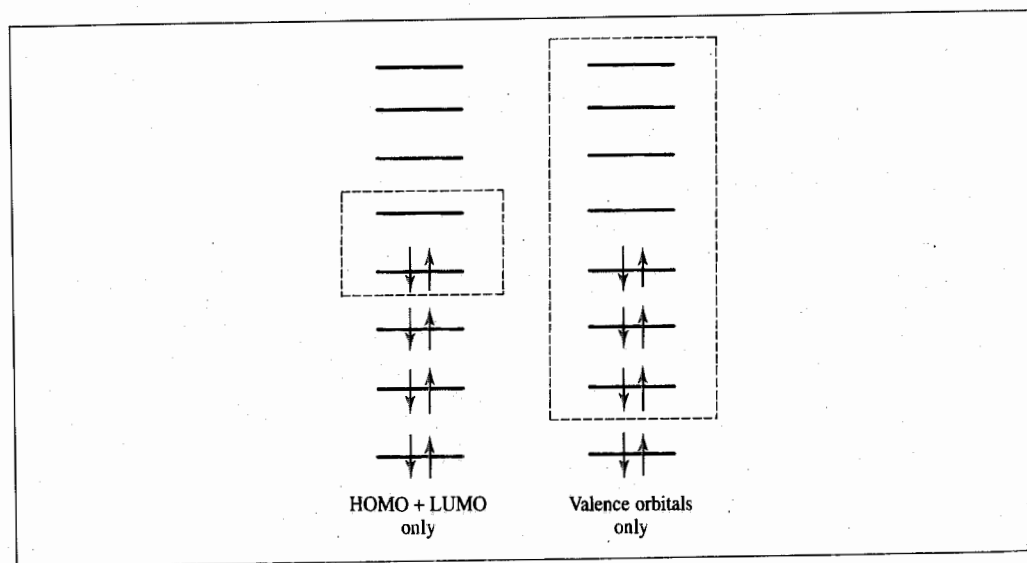


Fig. 3.3: Some of the ways in which excited-state wavefunctions can be included in a configuration interaction calculation. (Figure adapted from Hehre W J, L Radom, P v R Schleyer and J A Hehre 1986. *Ab initio Molecular Orbital Theory*. New York, Wiley.)

directly with the single-determinant, ground-state wavefunction Ψ_0 . It would therefore be anticipated that double excitations would be most important and that single excitations would have no effect on the energy of the ground state. However, the single excitations can interact with the double excitations, which in turn interact with Ψ_0 , and so single excitations do have a small indirect effect on the energy. The determinants of triple and higher excitations also do not interact directly with Ψ_0 (though they may do indirectly via other levels of excitation). This is because the Hamiltonian contains elements involving at most interactions between pairs of electrons, and so if the Slater determinants differ by more than two electron functions, their integral over all space will be zero.

In a ‘traditional’ CI calculation the determinants in the expansion, Equation (3.6), are those obtained from a Hartree–Fock calculation; only the coefficients c_0, c_1 , etc. are permitted to vary. Clearly, a better (i.e. lower-energy) wavefunction should be obtained if the coefficients of the basis functions themselves can vary as well as the coefficients of the determinants. This approach is known as the multiconfiguration self-consistent field method (MCSCF). MCSCF theory is considerably more complicated than the Roothaan–Hall equations and well beyond the scope of our discussion. One MCSCF technique that has attracted considerable attention is the complete active-space SCF method (CASSCF) of Roos [Roos *et al.* 1980]. CASSCF enables very large numbers of configurations to be included in the calculation by dividing the molecular orbitals into three sets: those which are doubly occupied in all configurations, those which are unoccupied in all configurations, and then all the remaining ‘active’ orbitals. The list of configurations is generated by considering all possible arrangements of the active electrons among the active orbitals.

A CI calculation is variational: the energy obtained is guaranteed to be greater than the ‘true’ energy. A drawback of CI calculations other than those performed at the full CI level is that they are not size consistent. Simply put, this means that the energy of a number N of non-interacting atoms or molecules is not equal to N times the energy of a single atom or molecule. Another consequence of size consistency is that, as the bond length in a diatomic molecule increases to infinity, so the energy of the system should become equal to the sum of the energies of the respective atoms. To illustrate why this lack of size consistency arises, consider CID calculations on Be_2 and on two beryllium atoms. The electronic configuration of Be is $1s^2 2s^2$ and so if we label the two atoms A and B, then the wavefunction for each of the two separated atoms will include the configuration $1s_A^2 2p_A^1 1s_B^2 2p_B^2$ ($\equiv 1s_A^2 1s_B^2 2p_A^2 2p_B^2$), in which two electrons have been promoted in each beryllium atom from the $2s$ to the $2p$ orbitals. This configuration represents a *quadruple* excitation for the beryllium dimer, which has the electronic configuration $1s_A^2 1s_B^2 2s_A^2 2s_B^2$. This quadruply excited configuration is not included in the CID wavefunction for the dimer, which is restricted to double excitations. In fact, the energy of a CI calculation including only doubly excited states is expected to scale in proportion to \sqrt{N} , where N is the number of non-interacting species present, rather than N . The Quadratic Configuration Interaction method (QCISD) was introduced to try to deal with this; it can be considered a size-consistent CISD theory [Pople *et al.* 1987]. The procedure involves the addition of higher excitation terms which are quadratic in their expansion coefficients. Higher still in theory is QCISD(T), in which an estimated contribution from the triple excitations can be incorporated, though with extra computational expense.

3.3.2 Many-body Perturbation Theory

Møller and Plesset proposed an alternative way to tackle the problem of electron correlation [Møller and Plesset 1934]. Their method is based upon Rayleigh–Schrödinger perturbation theory, in which the ‘true’ Hamiltonian operator \mathcal{H} is expressed as the sum of a ‘zeroth-order’ Hamiltonian \mathcal{H}_0 (for which a set of molecular orbitals can be obtained) and a perturbation, \mathcal{V} :

$$\mathcal{H} = \mathcal{H}_0 + \mathcal{V} \quad (3.7)$$

The eigenfunctions of the true Hamiltonian operator are Ψ_i with corresponding energies E_i . The eigenfunctions of the zeroth-order Hamiltonian are written $\Psi_i^{(0)}$ with energies $E_i^{(0)}$. The ground-state wavefunction is thus $\Psi_0^{(0)}$ with energy $E_0^{(0)}$. To devise a scheme by which it is possible to gradually improve the eigenfunctions and eigenvalues of \mathcal{H}_0 we can write the true Hamiltonian as follows:

$$\mathcal{H} = \mathcal{H}_0 + \lambda \mathcal{V} \quad (3.8)$$

λ is a parameter that can vary between 0 and 1; when λ is zero then \mathcal{H} is equal to the zeroth-order Hamiltonian, but when λ is 1 then \mathcal{H} equals its true value. The eigenfunctions Ψ_i and eigenvalues E_i of \mathcal{H} are then expressed in powers of λ :

$$\Psi_i = \Psi_i^{(0)} + \lambda \Psi_i^{(1)} + \lambda^2 \Psi_i^{(2)} + \dots = \sum_{n=0} \lambda^n \Psi_i^{(n)} \quad (3.9)$$

$$E_i = E_i^{(0)} + \lambda E_i^{(1)} + \lambda^2 E_i^{(2)} + \dots = \sum_{n=0} \lambda^n E_i^{(n)} \quad (3.10)$$

$E_i^{(1)}$ is the first-order correction to the energy, $E_i^{(2)}$ is the second-order correction, and so on. These energies can be calculated from the eigenfunctions as follows:

$$E_i^{(0)} = \int \Psi_i^{(0)} \mathcal{H}_0 \Psi_i^{(0)} d\tau \quad (3.11)$$

$$E_i^{(1)} = \int \Psi_i^{(0)} \mathcal{V} \Psi_i^{(0)} d\tau \quad (3.12)$$

$$E_i^{(2)} = \int \Psi_i^{(0)} \mathcal{V} \Psi_i^{(1)} d\tau \quad (3.13)$$

$$E_i^{(3)} = \int \Psi_i^{(0)} \mathcal{V} \Psi_i^{(2)} d\tau \quad (3.14)$$

To determine the corrections to the energy it is therefore necessary to determine the wavefunctions to a given order. In Møller–Plesset perturbation theory the unperturbed Hamiltonian \mathcal{H}_0 is the sum of the one-electron Fock operators for the N electrons:

$$\mathcal{H}_0 = \sum_{i=1}^N f_i = \sum_{i=1}^N \left(\mathcal{H}^{\text{core}} + \sum_{j=1}^N (\mathcal{J}_j + \mathcal{K}_j) \right) \quad (3.15)$$

The Hartree–Fock wavefunction, $\Psi_0^{(0)}$, is an eigenfunction of \mathcal{H}_0 , and the corresponding zeroth-order energy $E_0^{(0)}$ is equal to the sum of orbital energies for the occupied molecular

orbitals:

$$E_0^{(0)} = \sum_{i=1}^{\text{occupied}} \varepsilon_i \quad (3.16)$$

In order to calculate higher-order wavefunctions we need to establish the form of the perturbation, \mathcal{V} . This is the difference between the ‘real’ Hamiltonian \mathcal{H} and the zeroth-order Hamiltonian, \mathcal{H}_0 . Remember that the Slater determinant description, based on an orbital picture of the molecule, is only an approximation. The true Hamiltonian is equal to the sum of the nuclear attraction terms and electron repulsion terms:

$$\mathcal{H}_0 = \sum_{i=1}^N (\mathcal{H}^{\text{core}}) + \sum_{i=1}^N \sum_{j=i+1}^N \frac{1}{r_{ij}} \quad (3.17)$$

Hence the perturbation \mathcal{V} is given by:

$$\mathcal{V} = \sum_{i=1}^N \sum_{j=i+1}^N \frac{1}{r_{ij}} - \sum_{j=1}^N (\mathcal{J}_j + \mathcal{K}_j) \quad (3.18)$$

The first-order energy $E_0^{(1)}$ is given by:

$$E_0^{(1)} = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \frac{1}{r_{ij}} [(ii|jj) - (ij|ij)] \quad (3.19)$$

The sum of the zeroth-order and first-order energies thus corresponds to the Hartree–Fock energy (compare with Equation (2.110), which gives the equivalent result for a closed-shell system):

$$E_0^{(0)} + E_0^{(1)} = \sum_{i=1}^N \varepsilon_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N [(ii|jj) - (ij|ij)] \quad (3.20)$$

To obtain an improvement on the Hartree–Fock energy it is therefore necessary to use Møller–Plesset perturbation theory to at least second order. This level of theory is referred to as MP2 and involves the integral $\int \Psi_0^{(0)} \mathcal{V} \Psi_0^{(1)} d\tau$. The higher-order wavefunction $\Psi_0^{(1)}$ is expressed as linear combinations of solutions to the zeroth-order Hamiltonian:

$$\Psi_0^{(1)} = \sum_j c_j^{(1)} \Psi_j^{(0)} \quad (3.21)$$

The $\Psi_j^{(0)}$ in Equation (3.21) will include single, double, etc. excitations obtained by promoting electrons into the virtual orbitals obtained from a Hartree–Fock calculation. The second-order energy is given by:

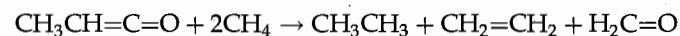
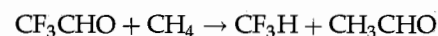
$$E_0^{(2)} = \sum_i^{\text{occupied}} \sum_{j>i}^{\text{virtual}} \sum_a \sum_{b>a} \frac{\int \int d\tau_1 d\tau_2 \chi_i(1) \chi_j(2) \left(\frac{1}{r_{12}} \right) [\chi_a(1) \chi_b(2) - \chi_b(1) \chi_a(2)]}{\varepsilon_a + \varepsilon_b - \varepsilon_i - \varepsilon_j} \quad (3.22)$$

These integrals will be non-zero only for double excitations, according to the Brillouin theorem. Third- and fourth-order Møller–Plesset calculations (MP3 and MP4) are also

available as standard options in many *ab initio* packages. For the fourth-order calculations single, triple and quadruple excitations will also contribute. As the triple substitutions are most difficult to perform computationally a partial theory that involves just single, double and quadruple substitutions (MP4SDQ) is a popular alternative.

The advantage of many-body perturbation theory is that it is size-independent, unlike configuration interaction – even when a truncated expansion is used. However, Møller–Plesset perturbation theory is not variational and can sometimes give energies that are lower than the ‘true’ energy. Møller–Plesset calculations are computationally intensive and so their use is often restricted to ‘single-point’ calculations at a geometry obtained using a lower level of theory. They are at present the most popular way to incorporate electron correlation into molecular quantum mechanical calculations, especially at the MP2 level. A Møller–Plesset calculation is specified using the level of theory used (e.g. MP2, MP3) together with the basis set. Thus MP2/6-31G* indicates a second-order Møller–Plesset calculation with the 6-31G* basis set.

Certain properties benefit more from the use of correlation methods than others do. For example, a single-determinant Hartree–Fock method and a reasonable basis set give geometrical parameters often very close (bond lengths within 0.01–0.02 Å and angles within 1–2°) to the experimental values. This contrasts with the situation for processes which result in the unpairing of electrons. A simple example is the bond dissociation energy of H₂, for which the Hartree–Fock limit is 84 kcal/mol. MP2, MP3 and MP4 calculations using the 6-31G** basis set give results of 101, 105 and 106 kcal/mol, respectively, for this process, much closer to the experimental value of 109 kcal/mol. In these and similar situations, electron correlation is often advised, if the computational resources permit. However, one class of reactions can be well described using single-determinant Hartree–Fock theory. These are known as *isodesmic reactions*, which are transformations in which the number of electron pairs is constant and the chemical bond types are conserved. Such reactions would be expected to benefit from a judicious cancellation of errors as only the environment of the bonds has changed. Examples of isodesmic reactions are:



Even at the STO-3G level quite respectable results can often be obtained.

In an attempt to deal with some of the shortcomings of even the correlated methods a number of correction factors have been developed. The Gaussian-*n* procedures [Pople *et al.* 1989, Curtiss *et al.* 1991, 1998] represent an attempt to develop a protocol for the accurate calculation of various properties such as atomisation energies, ionisation potentials, electron affinities and proton affinities for atoms and molecules containing first-row and second-row elements. Currently, the most recent member of this series is Gaussian-3 (G3) theory [Curtiss *et al.* 1998]. The G3 method involves a defined sequence of calculations involving geometry optimisation first at the Hartree–Fock level with the 6-31G* basis set and then at the MP2/6-31G* level. A single-point calculation is next carried out using this geometry with the full MP4 method (singles, doubles, triples and quadruples). This energy is then

refined through a series of corrections, which deal with the need for higher polarisation functions, for correlation effects beyond fourth-order perturbation theory (i.e. QCISD(T)) and for larger basis set effects. These correction factors are combined, together with a zero-point energy derived from a series of scaled harmonic frequencies determined from the first, HF/6-31G*, geometry optimisation, to give the final G3 energy. When tested on 299 experimental energies the overall average absolute deviation from experiment was 1.02 kcal/mol, with the average deviations for the four different types of data being 0.94 kcal/mol for the enthalpies of formation (148 values), 1.13 kcal/mol for the ionisation energies (85 values), 1.00 kcal/mol for electron affinities (58 values) and 1.34 kcal/mol for proton affinities (eight values). Detailed examination of the results can help to identify systems requiring most attention in subsequent developments of the theory. For example, the enthalpy of formation of both SO₂ and PF₃ have large negative deviations from experiment, perhaps due to the need for a larger basis set to describe the bonding in these molecules. Likewise some of the strained hydrocarbon ring systems (cyclopropene, cyclobutene and bicyclobutane) also show relatively large deviations.

The G3 method is still rather computationally intensive and so some efforts have been made to reduce the computational requirements whilst retaining an acceptable level of error. The G3(MP2) variant [Curtiss *et al.* 1999] replaces the MP4 calculations (which are particularly time-consuming), with comparable calculations at the MP2 level. This leaves the QCISD(T) stage as the most demanding step. The average absolute deviation of the energies calculated using the G3(MP2) method was 1.89 kcal/kmol on the entire 299 test systems, a significantly less accurate result than that of the full G3 method, but still noteworthy.

3.4 Practical Considerations When Performing *ab initio* Calculations

Ab initio calculations can be extremely time-consuming, especially when using the higher levels of theory or when the nuclei are free to move, as in a minimisation calculation (see Chapter 5). Various ‘tricks’ have been developed which can significantly reduce the computational effort involved. Many of these options are routinely available in the major software packages and are invoked by the specification of simple keywords. One common tactic is to combine different levels of theory for the various stages of a calculation. For example, a lower level of theory can be used to provide the initial guess for the density matrix prior to the first SCF iteration. Lower levels of theory can also be used in other ways. Suppose we wish to determine some of the electronic properties of a molecule in a minimum energy structure. Energy minimisation requires that the nuclei move and is typically performed in a series of steps, at each of which the energy (and frequently the gradient of the energy) must be calculated. Minimisation is therefore a computationally expensive procedure, particularly when performed at the high level of theory. To reduce this computational burden a lower level of theory can be employed for the geometry optimisation. A ‘single-point’ calculation using a high level of theory is then performed at the geometry so obtained to give a wavefunction from which the properties are determined. The assumption here of course is that the geometry does not change much between the two levels of

theory. Such calculations are denoted by slashes (/). For example, a calculation that is described as '6-31G*/STO-3G' indicates that the geometry was determined using the STO-3G basis set and the wavefunction was obtained using the 6-31G* basis set. Two slashes are used when each calculation is itself described using a slash, such as when electron correlation methods are used. For example, 'MP2/6-31G**//HF/6-31G*' indicates a geometry optimisation using a Hartree-Fock calculation with a 6-31G* basis set followed by a single-point calculation using the MP2 method for incorporating electron correlation, again using a 6-31G* basis set.

3.4.1 Convergence of Self-consistent Field Calculations

In an SCF calculation the wavefunction is gradually refined until self-consistency is achieved. For closed-shell ground-state molecules this is usually quite straightforward and the energy converges after a few cycles. However, in some cases convergence is a problem, and the energy may oscillate from one iteration to the next or even diverge rapidly. Various methods have been proposed to deal with such situations. A simple strategy is to use an average set of orbital coefficients rather than the set obtained from the immediately preceding iteration. The coefficients in this average set can be weighted according to the energies of each iteration. This tends to weed out those coefficients that give rise to higher energies.

The initial guess of the density matrix may influence the convergence of the SCF calculation; a null matrix is the simplest approach, but better results may be obtained by using a density matrix from a calculation performed at a lower level of theory. For example, the density matrix from a semi-empirical calculation may be used as the starting point for an *ab initio* calculation. Conversely, such an approach may itself lead to problems if there is a significant difference between the density matrices for the lower and higher levels of theory.

A more sophisticated method that is often very successful is Pulay's direct inversion of the iterative subspace (DIIS) [Pulay 1980]. Here, the energy is assumed to vary as a quadratic function of the basis set coefficients. In DIIS the coefficients for the next iteration are calculated from their values in the previous steps. In essence, one is predicting where the minimum in the energy will lie from a knowledge of the points that have been visited and by assuming that the energy surface adopts a parabolic shape.

3.4.2 The Direct SCF Method

An *ab initio* calculation can be logically considered to involve two separate stages. First, the various one- and two-electron integrals are calculated. This is a computationally intensive task and considerable effort has been expended finding ways to make the calculation of the integrals as efficient as possible. In the second stage, the wavefunction is determined using the variation theorem. In a 'traditional' SCF calculation all of the integrals are first calculated and stored on disk, to be retrieved later during the SCF calculation as required. The number of integrals to be stored may run into millions and this inevitably leads to delays in accessing the data, particularly as the retrieval of information from a disk requires

physical movement of the read head and so is slow. Modern computers (both workstations and supercomputers) have much faster (and cheaper) processing units, and many of these machines also have a substantial amount of internal memory, which can be accessed in a fraction of the time it takes to read data from the disk. In a direct SCF calculation, the integrals are not stored on the disk but are kept in memory or recalculated when required [Almlöf *et al.* 1982].

A much-quoted 'fact' is that *ab initio* calculations scale as the fourth power of the number of basis functions for ground-state, closed-shell systems. This scaling factor arises because each two-electron integral ($\mu\nu|\lambda\sigma$) involves four basis functions, so the number of two-electron integrals would be expected to increase in proportion to the fourth power of the number of basis functions. In fact, the number of such integrals is not exactly equal to the fourth power of the number of basis functions because many of the integrals are related by symmetry. We can calculate exactly the number of two-electron integrals that are required in a Hartree-Fock *ab initio* calculation as follows. There are seven different types of two-electron integral:

1. $(ab|cd) \equiv (ab|dc) \equiv (ba|cd) \equiv (ba|dc) \equiv (cd|ab) \equiv (cd|ba) \equiv (dc|ab) \equiv (dc|ba)$
2. $(aa|bc) \equiv (aa|cb) \equiv (bc|aa) \equiv (cb|aa)$
3. $(ab|ac) \equiv (ab|ca) \equiv (ba|ac) \equiv (ba|ca) \equiv (ac|ab) \equiv (ac|ba) \equiv (ca|ab) \equiv (ca|ba)$
4. $(aa|bb) \equiv (bb|aa)$
5. $(ab|ab) \equiv (ab|ba) \equiv (ba|ab) \equiv (ba|ba)$
6. $(aa|ab) \equiv (aa|ba) \equiv (ab|aa) \equiv (ba|aa)$
7. $(aa|aa)$

For a basis set with K basis functions, there are $K(K-1)(K-2)(K-3)$ integrals of type $(ab|cd)$, but due to symmetry only one-eighth of these are unique as shown. Similarly, there are $2K(K-1)(K-2)$ of type (2); $4K(K-1)(K-2)$ of type (3); $K(K-1)$ of type (4); $2K(K-1)$ of type (5); $4K(K-1)$ of type (6) and K of type 7. Thus, a basis set with 200 functions has a total of 202 015 050 unique two-electron integrals. For all but the smallest of basis sets most integrals are of type (1) which is why an *ab initio* problem is often considered to scale as $K^4/8$ ($200^4/8 = 200\,000\,000$). Including electron correlation adds significantly to the computational cost; for example, MP2 calculations scale as the fifth power of the number of basis functions. Electron correlation methods may also require significantly more memory and disk than the comparable SCF calculation; the higher levels scale as the sixth power, and in QCISD(T), one part of the calculation is seventh order.

In practice, *ab initio* calculations often scale as a significantly smaller power than four. It is found that in favourable cases the computational cost of a direct SCF calculation on a large molecule scales as approximately the *square* of the number of basis functions used. This significant reduction (from four to two) is due to several factors. We have already noted some of the ways in which a carefully chosen basis set can reduce the computational effort, for example by making many of the integrals (particularly the two-electron integrals) identical by using the same Gaussian exponents for s and p orbitals in the same shell. Another way in which the calculation time can be significantly reduced is to exploit any symmetry of the system. Many isolated molecules contain symmetry elements such as centres of inversion and mirror planes, information which can be used to reduce the

computational effort required. In the case of an *ab initio* calculation that scales as the fourth power of the number of basis functions then a four-fold reduction in the number of atoms can (in principle at least) result in the computational time being reduced by about 250 times. The most effective way to reduce the computational effort is to identify integrals which are so small that ignoring them (i.e. setting them to zero) will not affect the results. The number of 'important' integrals is believed to scale as $K^2 \ln K$. The negligible integrals are determined by calculating an upper limit for each integral. This can be done rapidly and so those integrals that are guaranteed to be negligible can be identified and so ignored. The cutoff value which determines whether an integral is explicitly calculated or is set to zero can vary from one program to another, so it is always useful to check its value if different programs give different results for a given calculation.

3.4.3 Calculating Derivatives of the Energy

Considerable effort has been spent devising efficient ways of directly calculating the first and second derivatives of the energy with respect to the nuclear coordinates. Derivatives are primarily used during minimisation procedures for finding equilibrium structures (the first derivative of the energy with respect to its coordinates equals the force on an atom) and are also used by methods which locate transition structures and determine reaction pathways.

A self-consistent field wavefunction (and thus its energy) can be considered a complicated function of the nuclear coordinates, basis functions and basis function coefficients (and, for a CI calculation, the coefficients of single determinantal wavefunctions). In order to determine the first, second, etc. derivatives of the energy with respect to the nuclear coordinates [Pulay 1977] it is necessary to consider not only how the energy depends directly on the nuclear coordinates but also whether there is an indirect dependence via other parameters. Indeed, it is only the one-electron part of the Hamiltonian that depends directly upon the nuclear coordinates ($H^{\text{core}}(1)$, Equation (2.125)), to which is added an internuclear Coulomb repulsion term. For the other parameters the derivative with respect to the nuclear coordinates is generally determined via the chain rule (for first derivatives). For example, for a generic nuclear coordinate q_i and a generic parameter x_j we can write:

$$\frac{\partial E}{\partial q_i} = \frac{\partial E}{\partial x_j} \frac{\partial x_j}{\partial q_i} \quad (3.23)$$

In Equation (3.23) q_i would be the x , y or z coordinate of an atom and x_j would be a parameter such as a basis function coefficient or a basis function exponent. An important result is that the terms involving variationally determined parameters (such as basis function coefficients) are equal to zero; the energy is a minimum when $(\partial E/\partial c_j)$ is zero. This greatly reduces the computational effort. Most of the numerical work in calculating the gradient is due to the various basis set parameters (e.g. orbital centres and exponents) which require the derivatives of the various electron integrals. For Gaussian basis sets these derivatives can be obtained analytically and indeed it is relatively straightforward to obtain first derivatives for many levels of theory. The time taken to calculate the derivatives is comparable to that required for the calculation of the total energy. Second (and

higher) derivatives are more difficult and expensive to calculate, even at the lower levels of theory.

A possible alternative approach to the calculation of forces is via the use of the Hellmann-Feynman theorem. If Ψ is an exact wavefunction of a Hamiltonian H with energy E then this theorem states that the derivative of E with respect to some parameter P can be written:

$$\frac{\partial E}{\partial P} = \left\langle \frac{\partial H}{\partial P} \right\rangle \quad (3.24)$$

In the case of the derivative with respect to some nuclear coordinate q_i , we would consider the exact force and the Hellmann-Feynman force to be equal:

$$\frac{\partial}{\partial q_i} \langle \Psi | H | \Psi \rangle = \left\langle \Psi \left| \frac{\partial H}{\partial q_i} \right| \Psi \right\rangle \quad (3.25)$$

Unfortunately, this only holds for the exact wavefunction and certain other types of wavefunction (such as at the Hartree-Fock limit). Moreover, even though the Hellmann-Feynman forces are much easier to calculate they are very unreliable, even for accurate wavefunctions, giving rise to spurious forces (often referred to as 'Pulay forces' [Pulay 1987]).

3.4.4 Basis Set Superposition Error

Suppose we wish to calculate the energy of formation of a bimolecular complex, such as the energy of formation of a hydrogen-bonded water dimer. Such complexes are sometimes referred to as 'supermolecules'. One might expect that this energy value could be obtained by first calculating the energy of a single water molecule, then calculating the energy of the dimer, and finally subtracting the energy of the two isolated water molecules (the 'reactants') from that of the dimer (the 'products'). However, the energy difference obtained by such an approach will invariably be an overestimate of the true value. The discrepancy arises from a phenomenon known as *basis set superposition error* (BSSE). As the two water molecules approach each other, the energy of the system falls not only because of the favourable intermolecular interactions but also because the basis functions on each molecule provide a better description of the electronic structure around the other molecule. It is clear that the BSSE would be expected to be particularly significant when small, inadequate basis sets are used (e.g. the minimal basis STO- n G basis sets) which do not provide for an adequate representation of the electron distribution far from the nuclei, particularly in the region where non-covalent interactions are strongest. One way to estimate the basis set superposition error is via the counterpoise correction method of Boys and Bernardi, in which the entire basis set is included in all calculations [Boys and Bernardi 1970]. Thus, in the general case:



$$\Delta E = E(AB) - [E(A) + E(B)] \quad (3.27)$$

The calculation of the energy of the individual species A is performed in the presence of 'ghost' orbitals of B ; that is, without the nuclei or electrons of B . A similar calculation is

performed for B using ghost orbitals on A. An alternative approach is to use a basis set in which the orbital exponents and contraction coefficients have been optimised for molecular calculations rather than for atoms. The relevance of the basis set superposition error and its dependence upon the basis set and the level of theory employed (i.e. SCF or with electron correlation) remains a subject of much research.

3.5 Energy Component Analysis

The interaction between atoms and molecules can vary from the weak attraction between a pair of closed-shell atoms (e.g. two rare gas atoms in a molecular beam) to the large energy associated with the formation of a chemical bond. Intermediate between these two extremes are interactions due to hydrogen bonding or electron donor-acceptor processes. In these intermediate cases it is often difficult to determine what factors are important in contributing to the interaction. For example, what can a hydrogen bond be ascribed to?

Morokuma analysis is a method for decomposing the energy change on formation of an intermolecular complex into five components: electrostatic, polarisation, exchange repulsion, charge transfer and mixing [Morokuma 1977]. Suppose we have performed *ab initio* SCF calculations on two molecules, X and Y, and on the intermolecular complex (or 'supermolecule') XY. The wavefunctions obtained can be written $A\Psi_X^0$, $A\Psi_Y^0$ and $A\Psi_{XY}^0$. 'A' indicates the use of an antisymmetrised wavefunction (e.g. a Slater determinant). The sum of the energies of the isolated molecules is E_0 and the energy of the supermolecule is E_4 (we follow the original notation of Morokuma). The interaction energy ΔE is thus given by $E_4 - E_0$. The five components are calculated as follows.

The electrostatic contribution equals the interaction between the unperturbed electron distributions of the two isolated species, A and B. It is identical to the classical Coulomb interaction and equals the difference $E_1 - E_0$, where E_1 is the energy associated with the product of the two individual wavefunctions, Ψ_1 :

$$\Psi_1 = A\Psi_A^0 A\Psi_B^0 \quad (3.28)$$

The electronic distributions of both X and Y will be changed by the presence of the other molecule. These polarisation effects cause a dipole to be induced in (say) molecule Y due to the charge distribution in molecule X and vice versa. Polarisation also affects the higher-order multipoles. To calculate the polarisation contribution we first calculate molecular wavefunctions Ψ_A and Ψ_B in the presence of the other molecule. The energy of the wavefunction Ψ_2 is determined as E_2 , where Ψ_2 is:

$$\Psi_2 = A\Psi_X A\Psi_Y \quad (3.29)$$

The polarisation contribution equals $E_2 - E_1$ and is always attractive.

In determining Ψ_1 and Ψ_2 , no electron exchange interactions are considered. The overlap between the electron distributions of X and Y at short range causes a repulsion because to bring together electrons with the same spin into the same region of space ultimately leads to a violation of the Pauli principle.

The exchange repulsion is calculated as $E_3 - E_1$, where E_3 is the energy of the wavefunction Ψ_3 :

$$\Psi_3 = A(\Psi_X^0 \cdot \Psi_Y^0) \quad (3.30)$$

Ψ_3 is derived from the undistorted wavefunctions of X and Y but the exchange of electrons is permitted. The exchange term is always repulsive.

The charge transfer term arises from the transfer of charge (i.e. electrons) from occupied molecular orbitals on one molecule to unoccupied orbitals on the other molecule. This contribution is calculated as the difference between the energy of the supermolecule XY when this charge transfer is specifically allowed to occur, and an analogous calculation in which it is not.

The Morokuma formalism also requires an additional, 'mixing' or 'coupling' term to be included. This equals the difference between the total SCF difference, ΔE , and the sum of the four contributions (electrostatic, polarisation, exchange repulsion and charge transfer). The mixing term has little physical significance and is used because the four components do not completely account for the entire interaction energy (it is a fudge factor!). Fortunately, it is often relatively small.

Morokuma studied a number of hydrogen-bonded complexes using this scheme in order to assess the contribution from each component. The systems studied were typically of intermolecular complexes involving small molecules such as H_2O , HF and NH_3 . In addition, Morokuma and his colleagues also examined a series of electron donor-acceptor complexes such as H_3N-BF_3 , $OC-BH_3$, HF-CIF and benzene- $OC(CN)_2$. He also studied the basis-set dependence of the results and observed that the energy components were more sensitive than the energy differences. For example, a minimal STO-3G basis set overestimates the charge transfer contribution, whereas double zeta basis sets tend to exaggerate the electrostatic interaction.

3.5.1 Morokuma Analysis of the Water Dimer

The water dimer (H_2O)₂ has been subject to perhaps the closest scrutiny of all hydrogen-bonded complexes. A variety of stable geometries are available to the water dimer, in which one or more hydrogen bonds are present. There has been considerable debate over the relative energies of these structures and even some dispute over which structures are actually at minimum points on the energy surface [Smith *et al.* 1990]. As might be expected, the results depend upon the basis set used. A linear geometry is observed experimentally and is also predicted to be the most stable structure by *ab initio* calculations with a wide variety of basis sets (see Figure 3.4). Using a 6-31G** basis set, Umeyama and Morokuma calculated that the -5.6 kcal/mol stabilisation energy was composed of -7.5 kcal/mol electrostatic stabilisation, 4.3 kcal/mol exchange repulsion, -0.5 kcal/mol polarisation and -1.8 kcal/mol charge transfer [Umeyama and Morokuma 1977]. The 'mixing term' contributed -0.1 kcal/mol. Thus the hydrogen bond in the water dimer was considered to arise primarily from electrostatic effects with a smaller charge transfer contribution. Morokuma and Umeyama also extended their analysis of charge transfer to investigate whether this was due to transfer from the proton donor to the acceptor, or from acceptor

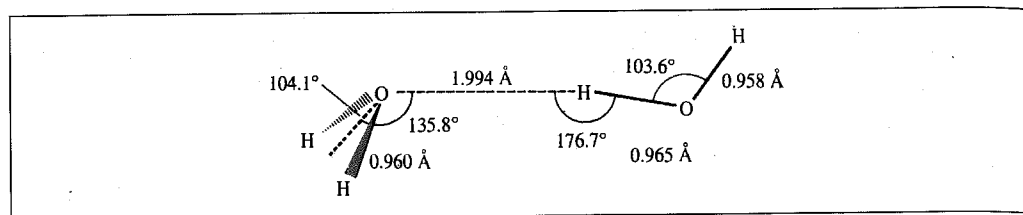


Fig. 3.4: The linear structure of the water dimer [Smith et al. 1990].

to donor. The results showed that approximately 90% of the charge transfer resulted from proton acceptor to proton donor transfer.

Morokuma analysis was widely used in the years after its introduction; it is less popular now as some problems have been encountered when trying to interpret the results with the larger basis sets that are feasible with today's faster computers and improved algorithms. In particular, when diffuse basis sets are used then there is a substantial amount of intermolecular overlap even at relatively large distances, which can make it difficult to factor out the different components. Nevertheless, the approach is certainly a useful way to assess the major causes of a particular type of intermolecular interaction, if only to provide a qualitative picture.

3.6 Valence Bond Theories

An entirely different way to treat the electronic structure of molecules is provided by valence bond theory, which was developed at about the same time as the molecular orbital approach. However, valence bond theory was not so amenable to calculations on large molecules, and molecular orbital theory came to dominate electronic structure theory for such systems. Nevertheless, valence bond theories are often considered to be more appropriate for certain types of problem than molecular orbital theory, especially when dealing with processes that involve bonds being broken and/or formed. Recall from Figure 3.2 that a self-consistent field wavefunction gives a wholly inaccurate picture for the dissociation of H_2 ; by contrast, the correct dissociation behaviour is naturally built into valence bond theories.

Valence bond theory is usually introduced using the famous Heitler-London model of the hydrogen molecule [Heitler and London 1927]. This model considers two non-interacting hydrogen atoms (a and b) in their ground states that are separated by a long distance. The wavefunction for this system is:

$$\Psi = \phi_{1sa}(1)\phi_{1sb}(2) \quad (3.31)$$

As the two hydrogen atoms approach to form a hydrogen molecule, such a wavefunction is inappropriate as it implies that electron 1 remains confined to orbital 1sa and electron 2 to orbital 1sb. This clearly violates the indistinguishability principle, and so a linear combination is used

$$\Psi_{vb} \propto \phi_{1sa}(1)\phi_{1sb}(2) + \phi_{1sa}(2)\phi_{1sb}(1) \quad (3.32)$$

The corresponding molecular orbital function for this system is:

$$\Psi_{mo} \propto \phi_{1sa}(1)\phi_{1sb}(2) + \phi_{1sa}(2)\phi_{1sb}(1) + \phi_{1sa}(1)\phi_{1sa}(2) + \phi_{1sb}(1)\phi_{1sb}(2) \quad (3.33)$$

The additional terms in the molecular orbital wavefunction correspond to states with the two electrons in the same orbital, which endows ionic character to the bond (H^+H^-). The valence bond wavefunction does not include any ionic character and in fact it correctly describes the dissociation into two hydrogen atoms. The simple valence bond and molecular orbital pictures in Equations (3.32) and (3.33) are extremes, with the 'true' wavefunction being somewhere in the middle. The valence bond representation can be improved by including a degree of ionic character as follows:

$$\Psi_{vb} \propto \phi_{1sa}(1)\phi_{1sb}(2) + \phi_{1sa}(2)\phi_{1sb}(1) + \lambda[\phi_{1sa}(1)\phi_{1sa}(2) + \phi_{1sb}(1)\phi_{1sb}(2)] \quad (3.34)$$

λ is a parameter that can be varied to give the 'correct' amount of ionic character. Another way to view the valence bond picture is that the incorporation of ionic character corrects the overemphasis that the valence bond treatment places on electron correlation. The molecular orbital wavefunction underestimates electron correlation and requires methods such as configuration interaction to correct for it. Although the presence of ionic structures in species such as H_2 appears counterintuitive to many chemists, such species are widely used to explain certain other phenomena such as the ortho/para or meta directing properties of substituted benzene compounds under electrophilic attack. Moreover, it has been shown that the ionic structures correspond to the deformation of the atomic orbitals when they are involved in chemical bonds.

One widely used valence bond theory is the generalised valence bond (GVB) method of Goddard and co-workers [Bobrowicz and Goddard 1977]. In the simple Heitler-London treatment of the hydrogen molecule the two orbitals are the non-orthogonal atomic orbitals on the two hydrogen atoms. In the GVB theory the analogous wavefunction is written:

$$\Psi_{GVB} \propto u(1)\nu(2) + u(2)\nu(1) \quad (3.35)$$

u and ν are non-orthogonal orbitals that are each expressed as a basis set expansion with the coefficients being variationally optimised to minimise the energy. The construction of the wavefunction from orbitals that are not necessarily orthogonal is characteristic of many valence bond theories and complicates the computational problem. The GVB approach is particularly successful for describing the electronic nature of systems as they approach dissociation.

Another approach is spin-coupled valence bond theory, which divides the electrons into two sets: 'core' electrons, which are described by doubly occupied orthogonal orbitals, and 'active' electrons, which occupy singly occupied non-orthogonal orbitals. Both types of orbital are expressed in the usual way as a linear combination of basis functions. The overall wavefunction is completed by two spin functions; one that describes the coupling of the spins of the core electrons and one that deals with the active electrons. The choice of spin function for these active electrons is a key component of the theory [Gerratt *et al.* 1997]. One of the distinctive features of this theory is that a considerable amount of chemically significant electronic correlation is incorporated into the wavefunction, giving an accuracy comparable to CASSCF. An additional benefit is that the orbitals tend to be

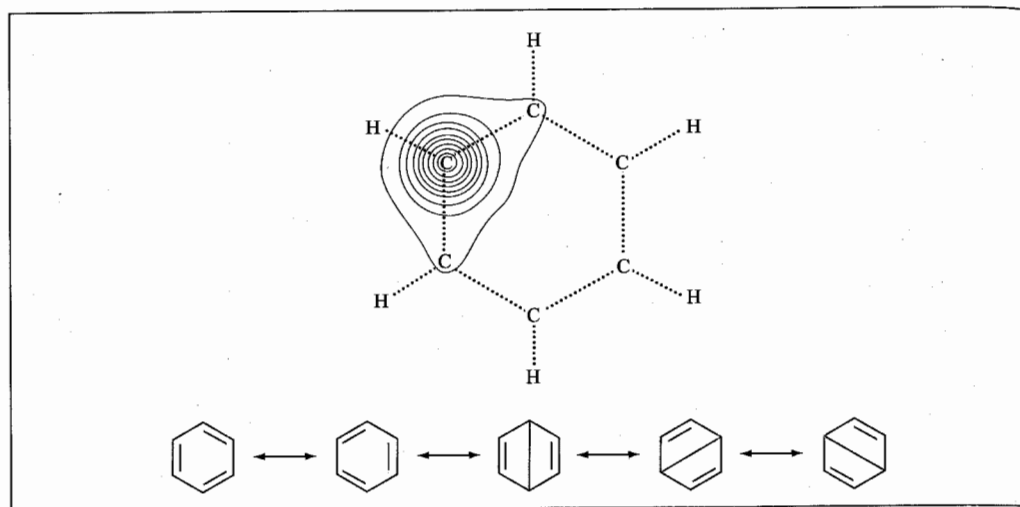


Fig. 3.5: π orbital for benzene obtained from spin-coupled valence bond theory. (Figure redrawn from Gerratt J, D L Cooper, P B Karadakov and M Raimondi 1997. *Modern valence bond theory*. Chemical Society Reviews 87:100.) The figure also shows the two Kekulé and three Dewar benzene forms which contribute to the overall wavefunction; each Kekulé form contributes approximately 40.5% and each Dewar form approximately 6.4%.

localised, closely resembling atomic or hybrid atomic orbitals, and consequently very visual. Various chemical phenomena have been examined using this approach, including dissociation reactions and hypervalence. One particularly interesting study was of the π system of benzene [Cooper *et al.* 1986]. This calculation resulted in six orbitals, each localised on one of the carbon atoms in the ring, though with some deformations towards neighbouring atoms (Figure 3.5). Moreover, the spin-coupling patterns suggested that the bonding was more akin to the Kekulé picture of benzene (with alternating double and single bonds) together with small contributions from Dewar benzene rather than the completely delocalised representation from molecular orbital theory.

3.7 Density Functional Theory

Density functional theory (DFT) is an approach to the electronic structure of atoms and molecules which has enjoyed a dramatic surge of interest since the late 1980s and 1990s [Parr 1983; Wimmer 1997]. Our approach here will be to introduce the key elements of the theory and to identify the similarities and differences between DFT and the Hartree-Fock approach. In Hartree-Fock theory the multi-electron wavefunction is expressed as a Slater determinant which is constructed from a set of N single-electron wavefunctions (N being the number of electrons in the molecule). DFT also considers single-electron functions. However, whereas Hartree-Fock theory does indeed calculate the full N -electron wavefunction, density functional theory only attempts to calculate the total electronic energy and the overall electronic density distribution. The central idea underpinning DFT is that

there is a relationship between the total electronic energy and the overall electronic density. This is not a particularly new idea; indeed an approximate model developed in the late 1920s (the Thomas-Fermi model) contains some of the basic elements. However, the real breakthrough came with a paper by Hohenberg and Kohn in 1964 [Hohenberg and Kohn 1964], who showed that the ground-state energy and other properties of a system were uniquely defined by the electron density. This is sometimes expressed by stating that the energy, E , is a unique *functional* of $\rho(\mathbf{r})$. A functional enables a function to be mapped to a number and is usually written using square brackets. Thus:

$$Q[f(\mathbf{r})] = \int f(\mathbf{r}) d\mathbf{r} \quad (3.36)$$

The function $f(\mathbf{r})$ is usually dependent upon other well-defined functions. A simple example of a functional would be the area under a curve, which takes a function $f(x)$ defining the curve between two points and returns a number (the area, in this case). In the case of DFT the function depends upon the electron density, which would make Q a functional of $\rho(\mathbf{r})$; in the simplest case $f(\mathbf{r})$ would be equivalent to the density (i.e. $f(\mathbf{r}) \equiv \rho(\mathbf{r})$). If the function $f(\mathbf{r})$ were to depend in some way upon the gradients (or higher derivatives) of $\rho(\mathbf{r})$ then the functional is referred to as being 'non-local', or 'gradient-corrected'. By contrast, a 'local' functional would only have a simple dependence upon $\rho(\mathbf{r})$. In DFT the energy functional is written as a sum of two terms:

$$E[\rho(\mathbf{r})] = \int V_{\text{ext}}(\mathbf{r})\rho(\mathbf{r}) d\mathbf{r} + F[\rho(\mathbf{r})] \quad (3.37)$$

The first term arises from the interaction of the electrons with an external potential $V_{\text{ext}}(\mathbf{r})$ (typically due to the Coulomb interaction with the nuclei). $F[\rho(\mathbf{r})]$ is the sum of the kinetic energy of the electrons and the contribution from interelectronic interactions. The minimum value in the energy corresponds to the exact ground-state electron density, so enabling a variational approach to be used (i.e. the 'best' solution corresponds to the minimum of energy and an incorrect density gives an energy above the true energy). There is a constraint on the electron density as the number of electrons (N) is fixed:

$$N = \int \rho(\mathbf{r}) d\mathbf{r} \quad (3.38)$$

In order to minimise the energy we introduce this constraint as a Lagrangian multiplier ($-\mu$), leading to:

$$\frac{\delta}{\delta\rho(\mathbf{r})} \left[E[\rho(\mathbf{r})] - \mu \int \rho(\mathbf{r}) d\mathbf{r} \right] = 0 \quad (3.39)$$

From this we can write:

$$\left(\frac{\delta E[\rho(\mathbf{r})]}{\delta\rho(\mathbf{r})} \right)_{V_{\text{ext}}} = \mu \quad (3.40)$$

Equation (3.40) is the DFT equivalent of the Schrödinger equation. The subscript V_{ext} indicates that this is under conditions of constant external potential (i.e. fixed nuclear positions). It is interesting to note that the Lagrange multiplier, μ , can be identified with the chemical potential of an electron cloud for its nuclei, which in turn is related to the

electronegativity, χ :

$$-\chi = \mu = \left(\frac{\partial E}{\partial N} \right)_{V_{\text{ext}}} \quad (3.41)$$

The second landmark paper in the development of density functional theory was by Kohn* and Sham who suggested a practical way to solve the Hohnberg–Kohn theorem for a set of interacting electrons [Kohn and Sham 1965]. The difficulty with Equation (3.37) is that we do not know what the function $F[\rho(\mathbf{r})]$ is. Kohn and Sham suggested that $F[\rho(\mathbf{r})]$ should be approximated as the sum of three terms:

$$F[\rho(\mathbf{r})] = E_{\text{KE}}[\rho(\mathbf{r})] + E_{\text{H}}[\rho(\mathbf{r})] + E_{\text{XC}}[\rho(\mathbf{r})] \quad (3.42)$$

where $E_{\text{KE}}[\rho(\mathbf{r})]$ is the kinetic energy, $E_{\text{H}}[\rho(\mathbf{r})]$ is the electron–electron Coulombic energy, and $E_{\text{XC}}[\rho(\mathbf{r})]$ contains contributions from exchange and correlation. It is important to note that the first term in Equation (3.42), $E_{\text{KE}}[\rho(\mathbf{r})]$, is defined as the kinetic energy of a system of *non-interacting* electrons with the same density $\rho(\mathbf{r})$ as the real system:

$$E_{\text{KE}}[\rho(\mathbf{r})] = \sum_{i=1}^N \int \psi_i(\mathbf{r}) \left(-\frac{\nabla^2}{2} \right) \psi_i(\mathbf{r}) d\mathbf{r} \quad (3.43)$$

The second term, $E_{\text{H}}(\rho)$, is also known as the Hartree electrostatic energy. The Hartree approach to solving the Schrödinger equation was introduced briefly in Section 2.3.3 and almost immediately dismissed because it fails to recognise that electronic motions are correlated. In the Hartree approach this electrostatic energy arises from the classical interaction between two charge densities, which, when summed over all possible pairwise interactions, gives:

$$E_{\text{H}}[\rho(\mathbf{r})] = \frac{1}{2} \iint \frac{\rho(\mathbf{r}_1)\rho(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \quad (3.44)$$

Combining these two and adding the electron–nuclear interaction leads to the full expression for the energy of an N -electron system within the Kohn–Sham scheme:

$$E[\rho(\mathbf{r})] = \sum_{i=1}^N \int \psi_i(\mathbf{r}) \left(-\frac{\nabla^2}{2} \right) \psi_i(\mathbf{r}) d\mathbf{r} + \frac{1}{2} \iint \frac{\rho(\mathbf{r}_1)\rho(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 + E_{\text{XC}}[\rho(\mathbf{r})] - \sum_{A=1}^M \int \frac{Z_A}{|\mathbf{r} - \mathbf{R}_A|} \rho(\mathbf{r}) d\mathbf{r} \quad (3.45)$$

This equation acts to *define* the exchange–correlation energy functional $E_{\text{XC}}[\rho(\mathbf{r})]$, which thus contains not only contributions due to exchange and correlation but also a contribution due to the difference between the true kinetic energy of the system and $E_{\text{KE}}[\rho(\mathbf{r})]$.

* Walter Kohn, whose name appears on the two key papers which provided the impetus for the development of ‘modern’ density functional theory, was awarded the Nobel Prize for Chemistry in 1998, jointly with John Pople.

Kohn and Sham wrote the density $\rho(\mathbf{r})$ of the system as the sum of the square moduli of a set of one-electron orthonormal orbitals:

$$\rho(\mathbf{r}) = \sum_{i=1}^N |\psi_i(\mathbf{r})|^2 \quad (3.46)$$

By introducing this expression for the electron density and applying the appropriate variational condition the following one-electron Kohn–Sham equations result:

$$\left\{ -\frac{\nabla^2}{2} - \left(\sum_{A=1}^M \frac{Z_A}{r_{1A}} \right) + \int \frac{\rho(\mathbf{r}_2)}{r_{12}} d\mathbf{r}_2 + V_{\text{XC}}[\mathbf{r}_1] \right\} \psi_i(\mathbf{r}_1) = \varepsilon_i \psi_i(\mathbf{r}_1) \quad (3.47)$$

In Equation (3.47) we have written the external potential in the form appropriate to the interaction with M nuclei. ε_i are the orbital energies and V_{XC} is known as the exchange–correlation functional, related to the exchange–correlation energy by:

$$V_{\text{XC}}[\mathbf{r}] = \left(\frac{\delta E_{\text{XC}}[\rho(\mathbf{r})]}{\delta \rho(\mathbf{r})} \right) \quad (3.48)$$

The total electronic energy is then calculated from Equation (3.45).

To solve the Kohn–Sham equations a self-consistent approach is taken. An initial guess of the density is fed into Equation (3.47) from which a set of orbitals can be derived, leading to an improved value for the density, which is then used in the second iteration, and so on until convergence is achieved.

3.7.1 Spin-polarised Density Functional Theory

Local spin density functional theory (LSDFT) is an extension of ‘regular’ DFT in the same way that restricted and unrestricted Hartree–Fock extensions were developed to deal with systems containing unpaired electrons. In this theory both the electron density and the spin density are fundamental quantities with the net spin density being the difference between the density of up-spin and down-spin electrons:

$$\sigma(\mathbf{r}) = \rho_{\uparrow}(\mathbf{r}) - \rho_{\downarrow}(\mathbf{r}) \quad (3.49)$$

The total electron density is just the sum of the densities for the two types of electron. The exchange–correlation functional is typically different for the two cases, leading to a set of spin-polarised Kohn–Sham equations:

$$\left\{ -\frac{\nabla^2}{2} - \left(\sum_{A=1}^M \frac{Z_A}{r_{1A}} \right) + \int \frac{\rho(\mathbf{r}_2)}{r_{12}} d\mathbf{r}_2 + V_{\text{XC}}[\mathbf{r}_1, \sigma] \right\} \psi_i^{\sigma}(\mathbf{r}_1) = \varepsilon_i^{\sigma} \psi_i^{\sigma}(\mathbf{r}_1) \quad \sigma = \alpha, \beta \quad (3.50)$$

This leads to two sets of wavefunctions, one for each spin, similar to UHF theory.

3.7.2 The Exchange–correlation Functional

The exchange–correlation functional is clearly key to the success (or otherwise) of the density functional approach. One reason why DFT is so appealing is that even relatively simple

approximations to the exchange-correlation functional can give favourable results. The simplest way to obtain this contribution uses the so-called *local density approximation* (LDA; the acronym LSDA is also used, for local spin density approximation), which is based upon a model called the uniform electron gas, in which the electron density is constant throughout all space. The total exchange-correlation energy, E_{XC} , for our system can then be obtained by integrating over all space:

$$E_{XC}[\rho(\mathbf{r})] = \int \rho(\mathbf{r})\epsilon_{XC}(\rho(\mathbf{r})) d\mathbf{r} \quad (3.51)$$

$\epsilon_{XC}(\rho(\mathbf{r}))$ is the exchange-correlation energy per electron as a function of the density in the uniform electron gas. The exchange-correlation functional is obtained by differentiation of this expression:

$$V_{XC}[\mathbf{r}] = \rho(\mathbf{r}) \frac{d\epsilon_{XC}(\rho(\mathbf{r}))}{d\rho(\mathbf{r})} + \epsilon_{XC}(\rho(\mathbf{r})) \quad (3.52)$$

In the local density approximation it is assumed that at each point \mathbf{r} in the inhomogeneous electron distribution (i.e. in the system of interest) where the density is $\rho(\mathbf{r})$ then $V_{XC}[\rho(\mathbf{r})]$ and $\epsilon_{XC}(\rho(\mathbf{r}))$ have the same values as in the homogeneous electron gas. In other words, the real electron density surrounding a volume element at position \mathbf{r} is replaced by a constant electron density with the same value as at \mathbf{r} . However, this 'constant' electron density is different for each point in space (Figure 3.6).

The exchange-correlation energy per electron (i.e. the energy density) of the uniform electron gas is known accurately for all densities of practical interest from various approaches such as quantum Monte Carlo methods [Ceperley and Alder 1980]. In order to be of practical use this exchange-correlation energy density is then expressed in an analytical form that makes it amenable to computation. It is usual to express $\epsilon_{XC}[\rho(\mathbf{r})]$ as an analytical function of the electron density and to consider the exchange and correlation contributions separately. However, some analytical expressions for the combined exchange and correlation energy density do exist, such as the following expression of Gunnarsson and

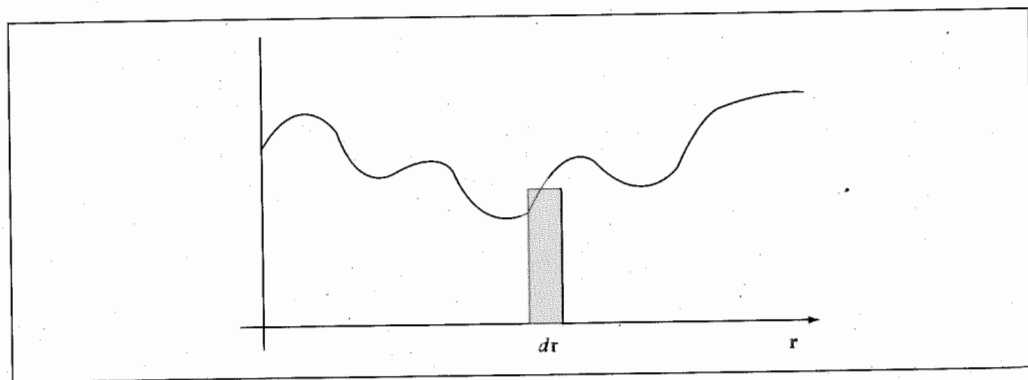


Fig. 3.6: Schematic representation of the way in which the local density approximation assumes that the electron density within a volume element dr surrounding a point \mathbf{r} is assumed to be constant.

Lundqvist [Gunnarsson and Lundqvist 1976]:

$$\epsilon_{XC}(\rho(\mathbf{r})) = -\frac{0.458}{r_s} - 0.0666G\left(\frac{r_s}{11.4}\right);$$

$$G(x) = \frac{1}{2} \left[(1+x) \log(1+x^{-1}) - x^2 + \frac{x}{2} - \frac{1}{3} \right], \quad r_s^3 = \frac{3}{4\pi\rho(\mathbf{r})} \quad (3.53)$$

The following relatively simple expression is commonly used for the exchange-only energy under the local density approximation [Slater 1974]:

$$E_X[\rho_\alpha(\mathbf{r}), \rho_\beta(\mathbf{r})] = -\frac{3}{2} \left(\frac{3}{4\pi} \right)^{1/3} \int (\rho_\alpha^{4/3}(\mathbf{r}) + \rho_\beta^{4/3}(\mathbf{r})) d\mathbf{r} \quad (3.54)$$

where α and β represent up and down spins. In general, more attention has been paid to the correlation contribution, for which there is no such simple functional form. Perdew and Zunger suggested the following parametric relationship for the correlation contribution [Perdew and Zunger 1981]:

$$\epsilon_C(\rho(\mathbf{r})) = \begin{cases} -0.1423/(1 + 1.9529r_s^{1/2} + 0.3334r_s) & r_s \geq 1 \\ -0.0480 + 0.0311 \ln r_s - 0.0116r_s + 0.0020r_s \ln r_s & r_s < 1 \end{cases} \quad (3.55)$$

This result applies when the number of up spins equals the number of down spins and so is not applicable to systems with an odd number of electrons. The correlation energy functional was also considered by Vosko, Wilk and Nusair [Vosko *et al.* 1980], whose expression is:

$$\epsilon_C(\rho(\mathbf{r})) = \frac{A}{2} \left\{ \ln \frac{x^2}{X(x)} + \frac{2b}{Q} \tan^{-1} \frac{Q}{2x+b} - \frac{bx_0}{X(x_0)} \left[\ln \frac{(x-x_0)^2}{X(x)} + \frac{2(b+2x_0)}{Q} \tan^{-1} \frac{Q}{2x+b} \right] \right\}$$

$$x = r_s^{1/2}, \quad X(x) = x^2 + bx + c, \quad Q = (4c - b^2)^{1/2}; \quad (3.56)$$

$$A = 0.0621814, \quad x_0 = -0.409286, \quad b = 13.0720, \quad c = 42.7198$$

In addition to the energy terms for the exchange-correlation contribution (which enables the total energy to be determined) it is necessary to have corresponding terms for the potential, $V_{XC}[\rho(\mathbf{r})]$, which are used to solve the Kohn-Sham equations. These are obtained as the appropriate first derivatives using Equation (3.52).

To solve the Kohn-Sham equations a number of different approaches and strategies have been proposed. One important way in which these can differ is in the choice of basis set for expanding the Kohn-Sham orbitals. In most (but not all) DFT programs for calculating the properties of molecular systems (rather than for solid-state materials) the Kohn-Sham orbitals are expressed as a linear combination of atomic-centred basis functions:

$$\psi_i(\mathbf{r}) = \sum_{\nu=1}^K c_{\nu i} \phi_\nu \quad (3.57)$$

Several functional forms have been investigated for the basis functions ϕ_ν . Given the vast experience of using Gaussian functions in Hartree-Fock theory it will come as no surprise to learn that such functions have also been employed in density functional theory. However, these are not the only possibility: Slater type orbitals are also used, as are numerical

basis functions. We encountered Slater type orbitals in Chapter 2, but the notion of a numerical basis function is new. A numerical basis function can be generated by solving the Kohn–Sham equations for isolated atoms. This gives a set of values on a spherical polar grid centred on each atom. The variation at each grid point can be stored as a cubic spline function so enabling analytical gradients to be calculated. One advantage of a numerical basis set (if properly derived) is that it has the correct nodal behaviour close to the nucleus together with an exponential decay.

More than one function may be used to represent a particular atomic orbital. This is obviously a well-understood tactic when using Gaussian functions, but the use of basis set contractions also applies to the Slater type orbitals and the numerical basis sets. For a numerical basis set the ‘contraction’ can be derived from two functions, one corresponding to the neutral atom and the other to a positive ion.

If the basis set expansion for the Kohn–Sham orbitals in Equation (3.57) is substituted into the Kohn–Sham equations then it is possible to express them in a matrix form, identical in form to the Roothaan–Hall equations:

$$\mathbf{HC} = \mathbf{SCE} \quad (3.58)$$

In this matrix equation the elements of the Kohn–Sham matrix \mathbf{H} are given by:

$$H_{\mu\nu} = \int d\mathbf{r}_1 \phi_\mu(\mathbf{r}_1) \left\{ -\frac{\nabla_1^2}{2} - \left(\sum_{A=1}^M \frac{Z_A}{r_{1A}} \right) + \int \frac{\rho(\mathbf{r}_2)}{r_{12}} d\mathbf{r}_2 + V_{\text{XC}}[\mathbf{r}_1] \right\} \phi_\nu(\mathbf{r}_1) \quad (3.59)$$

The first two terms are straightforward and are equal to the core contribution, $H_{\mu\nu}^{\text{core}}$. The Coulomb repulsion contribution (the Hartree term) can be expanded in terms of the basis functions and the density matrix, \mathbf{P} :

$$\iint \frac{\phi_\mu(\mathbf{r}_1)\rho(\mathbf{r}_2)\phi_\nu(\mathbf{r}_1)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 = \sum_{\lambda=1}^K \sum_{\sigma=1}^K P_{\lambda\sigma} \iint \frac{\phi_\mu(\mathbf{r}_1)\phi_\nu(\mathbf{r}_1)\phi_\lambda(\mathbf{r}_2)\phi_\sigma(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \quad (3.60)$$

For a closed-shell system with N electrons the elements of the density matrix are given by:

$$P_{\mu\nu} = 2 \sum_{i=1}^{N/2} c_{\mu i} c_{\nu i} \quad (3.61)$$

This is just the same as for the Roothaan–Hall approach to Hartree–Fock theory. The overlap matrix, \mathbf{S} , is defined similarly:

$$S_{\mu\nu} = \int \phi_\mu(\mathbf{r})\phi_\nu(\mathbf{r}) d\mathbf{r} \quad (3.62)$$

The overall procedure to achieve self-consistency is very reminiscent of that used in Hartree–Fock theory, involving first an initial guess of the density by superimposing atomic densities, construction of the Kohn–Sham and overlap matrices, and diagonalisation to give the eigenfunctions and eigenvectors from which the Kohn–Sham orbitals* can be

* It is important to note that the Kohn–Sham orbitals used in density functional theory are a set of non-interacting orbitals designed to give the correct density and have no physical meaning beyond that, unlike the orbitals used in Hartree–Fock theory.

constructed and thus the density for the next iteration. This cycle continues until convergence is achieved.

The appearance of the four-centre integrals in Equation (3.60) might lead one to question the advantage of the DFT approach, at least as far as computational efficiency is concerned. Whilst these integrals can certainly be tackled using the same techniques as in Hartree–Fock theory, it is also viable in density functional theory to avoid having to calculate them by considering the left-hand side of Equation (3.60). There are two basic ways to do this. First, one can approximate the charge density by another basis set expansion:

$$\rho(\mathbf{r}) \approx \sum_k c_k \phi'_k(\mathbf{r}) \quad (3.63)$$

These auxiliary basis functions ϕ' have the same functional form as the orbital expansion and the coefficients c_k are obtained by a least-squares fitting procedure. Substituting for the density in the four-centre integrals gives a computationally less demanding three-centre, two-electron integral:

$$\iint \frac{\phi_\mu(\mathbf{r}_1)\phi_\nu(\mathbf{r}_1)\phi_\lambda(\mathbf{r}_2)\phi_\sigma(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 = \iint \frac{\phi_\mu(\mathbf{r}_1)\phi_\nu(\mathbf{r}_1)\phi'_k(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \quad (3.64)$$

The second approach focuses on the Coulomb integral and uses Poisson’s equation. Let us introduce $V_{\text{el}}(\mathbf{r}_1)$:

$$V_{\text{el}}(\mathbf{r}_1) = \int \frac{\rho(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_2 \quad (3.65)$$

Poisson’s equation relates the second derivative of the electric potential to the charge density:

$$\nabla^2 V(\mathbf{r}) = -4\pi\rho(\mathbf{r}) \quad (3.66)$$

We can thus write:

$$\nabla^2 \int \frac{\rho(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_2 = -4\pi\rho(\mathbf{r}_1) \quad (3.67)$$

This equation can be solved numerically on a grid to determine $V_{\text{el}}(\mathbf{r}_1)$. The same grid is then used to numerically integrate the four-centre, two-electron integral, Equation (3.60), as follows:

$$\iint \frac{\phi_\mu(\mathbf{r}_1)\rho(\mathbf{r}_2)\phi_\nu(\mathbf{r}_1)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \equiv \int \phi_\mu(\mathbf{r}_1)V_{\text{el}}(\mathbf{r}_1)\phi_\nu(\mathbf{r}_1) \approx \sum_{i=1}^P \phi_\mu(\mathbf{R}_i)V_{\text{el}}(\mathbf{R}_i)\phi_\nu(\mathbf{R}_i)W_i \quad (3.68)$$

In this equation the P points \mathbf{R}_i correspond to the grid used to solve the Poisson equation for V_{el} and W_i are weighting factors.

It might be wondered why these two simplifications for the four-centre, two-electron integrals can be used in density functional theory and not in Hartree–Fock theory. The reason is that the exchange contribution in Hartree–Fock theory is not a function that can be simplified (technically, it is a non-local functional), in contrast to the situation in

density functional theory. As the four-centre integrals must therefore still be determined for the exchange component in Hartree-Fock theory there is nothing to be gained from simplifying the corresponding Coulomb term.

The exchange-correlation contribution to the Kohn-Sham matrix elements (the final term in Equation (3.59)) is invariably evaluated using a grid of points. This is a consequence of the complexity of the functionals employed. The integration may then be performed using the grid directly or by fitting a further auxiliary basis set expansion with which analytical integration can be used. If a DFT program uses a basis set containing K functions and employs either a grid-based integration scheme with P points or an auxiliary basis set with P functions then the computational complexity of the calculation scales as K^2P . As P is often linearly related to K , density functional theory is often said to scale as the cube of the number of basis functions, K^3 . This contrasts with the fourth-power scaling for conventional Hartree-Fock calculations. However, many practical density functional calculations with a well-engineered computer program do not scale as the simple third power, just as practical Hartree-Fock calculations do not scale as the fourth power; these oft-quoted statements apply only to the most naïve implementations or for calculations on very small, test systems where integral neglect thresholds are not employed.

Whilst most of the programs which use density functional theory for molecular calculations employ one of the three types of basis set described thus far, there are two important alternatives to this approach. The first of these involves the solution of the Kohn-Sham equations numerically (on a grid) using what is sometimes referred to as a 'basis-set free' approach [Becke and Dickson 1990]. Such an approach is thus free from the limitations of a finite basis set expansion (provided, of course, that sufficient grid points are employed!) and can be used to evaluate different exchange-correlation functionals, as these represent the only remaining source of error. The second alternative is particularly important for the study of bulk systems such as metals and alloys and involves the use of *plane waves*. This approach will be discussed later in this chapter when we consider the general problem of using quantum mechanics to study the solid state.

3.7.3 Beyond the Local Density Approximation: Gradient-corrected Functionals

The most important feature of density functional theory is probably the way in which it directly incorporates exchange and correlation effects; the latter in particular are only truly considered in the more complex, post-Hartree-Fock approaches such as configuration interaction or many-body perturbation theory. Despite its simplicity the local density approximation performs surprisingly well. However, the local density approximation has been shown to be clearly inadequate for some problems and for this reason extensions have been developed. The most common method is to use gradient-corrected, 'non-local' functionals which depend upon the gradient of the density at each point in space and not just on its value. These gradient corrections are typically divided into separate exchange and correlation contributions. A variety of gradient corrections have been proposed in the literature. The gradient correction to the exchange functional proposed by Becke is popular [Becke 1988, 1992]; this corrects

the local spin density approximation result as follows:

$$E_X[\rho(\mathbf{r})] = E_X^{\text{LSDA}}[\rho(\mathbf{r})] - b \sum_{\sigma=\alpha,\beta} \int \rho_{\sigma}^{4/3} \frac{x_{\sigma}^2}{(1 + 6bx_{\sigma} \sinh^{-1} x_{\sigma})} d\mathbf{r}; \quad x_{\sigma} = \frac{|\nabla\rho_{\sigma}|}{\rho_{\sigma}^{4/3}} \quad (3.69)$$

$E_X^{\text{LSDA}}[\rho(\mathbf{r})]$ is the standard Slater form of the exchange energy, Equation (3.54). The form written in Equation (3.69) is for a spin-unrestricted system, from which the appropriate expression for a closed-shell system is easily derived. x_{σ} is a dimensionless parameter and b is constant with a value of 0.0042 a.u. The value of b was determined by fitting to exact exchange Hartree-Fock energies for the noble gas atoms helium to radon. Two particular features of this functional form are that in the limit $r \rightarrow \infty$ the limiting form of the exchange-correlation integral is correctly achieved and that it uses just a single parameter, b . The correlation functional of Lee, Yang and Parr is also widely used [Lee *et al.* 1988]; in its original form it was expressed as follows (for a closed-shell system):

$$E_C[\rho(\mathbf{r})] = -a \int \frac{1}{1 + d\rho^{-1/3}} \{r + b\rho^{-2/3} [C_F\rho^{5/3} - 2t_W + (\frac{1}{3}t_W + \frac{1}{18}\nabla^2\rho)e^{-cr^{-1/3}}]\} d\mathbf{r} \quad (3.70)$$

$$t_W(\mathbf{r}) = \sum_{i=1}^N \frac{|\nabla\rho_i(\mathbf{r})|^2}{\rho_i(\mathbf{r})} - \frac{1}{8}\nabla^2\rho; \quad C_F = \frac{3}{10}(3\pi^2)^{2/3}$$

a , b , c and d are constants with values 0.049, 0.132, 0.2533 and 0.349, respectively. This expression provides both local and non-local components within a single expression and the gradient contribution to second order. A combination of the standard local spin density approximation exchange result (Equation (3.54)) with the Becke gradient-exchange correction and the Lee-Yang-Parr correlation functional is currently a popular choice, commonly abbreviated to BLYP (pronounced 'blip').

3.7.4 Hybrid Hartree-Fock/Density Functional Methods

As we stated earlier, a key feature of density functional theory is the way in which correlation effects are incorporated from the beginning, unlike Hartree-Fock theory. Moreover, the incorporation of correlation into the Hartree-Fock formalism often involves significant computational overhead, as we have considered in Section 3.3. However, it is important to recognise that Hartree-Fock theory does provide an essentially exact means of treating the exchange contribution. One potentially attractive option is thus to add a correlation energy derived from DFT (e.g. the local density approximation) to the Hartree-Fock energy. In such an approach the exchange-correlation energy is written as a sum of the exact exchange term together with the correlation component from the local density approximation. This 'exact' exchange energy is obtained from the Slater determinant of the Kohn-Sham orbitals.

Unfortunately, this simple approach does not work well, but Becke has proposed a strategy which does seem to have much promise [Becke 1993a, b]. In his approach the exchange-correlation energy E_{XC} is written in the following form:

$$E_{XC} = \int_0^1 U_{XC}^{\lambda} d\lambda \quad (3.71)$$

Equation (3.71) contains a coupling parameter λ , which takes values from 0 to 1. A value of zero corresponds to a system where there is no Coulomb repulsion U_{XC} between the electrons (i.e. the Kohn-Sham non-interacting reference state). As λ increases to 1 the interelectronic Coulomb repulsion is introduced until $\lambda = 1$, which corresponds to the 'real' system with full interactions. For all values of λ the electron density is the same and equal to the density of the real system. It is not practical to perform this integral analytically and so it must be approximated. The simplest approximation is a linear interpolation:

$$E_{XC} = \frac{1}{2} (U_{XC}^0 + U_{XC}^1) \quad (3.72)$$

When $\lambda = 0$ we have U_{XC}^0 , which is the exchange-correlation potential energy of the non-interacting reference system. As there are no electronic interactions in this system there is no correlation term and so U_{XC}^0 corresponds to the pure exchange energy of the Kohn-Sham determinant and can be determined exactly. U_{XC}^1 is the exchange-correlation potential energy of the full-interacting real system. Becke proposed that this should be calculated using the local spin-density approximation. This potential energy (note that it is not the total energy, E) is available from:

$$U_{XC}^1 \approx U_{XC}^{LSDA} = \int u_{XC}[\rho_\alpha(\mathbf{r}), \rho_\beta(\mathbf{r})] d\mathbf{r} \quad (3.73)$$

u_{XC} is the exchange-correlation potential energy density of an electron gas for which appropriate expressions are available.

This so-called 'half-and-half' theory proved to be significantly better than the alternative methods based upon mixing exact exchange and correlation energies. In a refinement of the scheme, Becke recognised that there were problems with the model when $\lambda = 0$. These problems arise because the electron gas model is not appropriate near this exchange-only limit for molecular bonds. Hence a key feature of Becke's modified model is to eliminate the term U_{XC}^0 and to write the exchange-correlation energy as the following linear combination:

$$E_{XC} = E_{XC}^{LSDA} + a_0 (E_X^{\text{exact}} - E_X^{LSDA}) + a_X \Delta E_X^{GC} + a_C \Delta E_C^{GC} \quad (3.74)$$

In Equation (3.74) E_X^{exact} is the exact exchange energy (obtained from the Slater determinant of the Kohn-Sham orbitals), E_X^{LSDA} is the exchange energy under the local spin density approximation, ΔE_X^{GC} is the gradient correction for exchange and ΔE_C^{GC} is the gradient correction for correlation. a_0 , a_X and a_C are empirical coefficients obtained by least-squares fitting to experimental data (56 atomisation energies, 42 ionisation potentials, eight proton affinities and the total atomic energies of the ten first-row elements). Their values are $a_0 = 0.20$, $a_X = 0.72$ and $a_C = 0.81$. In Becke's original paper his own gradient correction for exchange was used together with a gradient correction for correlation suggested by Perdew and Wang. An alternative to this scheme is to employ the Lee-Yang-Parr correlation functional (with the gradient term) and the standard local correlation functional due to Vosko, Wilk and Nusair (VWN). This is the 'B3LYP' density functional:

$$E_{XC}^{B3LYP} = (1 - a_0) E_X^{LSDA} + a_0 E_X^{HF} + a_X \Delta E_X^{B88} + a_C E_C^{LYP} + (1 - a_C) E_C^{VWN} \quad (3.75)$$

3.7.5 Performance and Applications of Density Functional Theory

The application of density functional theory to isolated, 'organic' molecules is still in relative infancy compared with the use of Hartree-Fock methods. There continues to be a steady stream of publications designed to assess the performance of the various approaches to DFT. As we have discussed there is a plethora of ways in which density functional theory can be implemented with different functional forms for the basis set (Gaussians, Slater type orbitals, or numerical), different expressions for the exchange and correlation contributions within the local density approximation, different expressions for the gradient corrections and different ways to solve the Kohn-Sham equations to achieve self-consistency. This contrasts with the situation for Hartree-Fock calculations, which mostly use one of a series of tried and tested Gaussian basis sets and where there is a substantial body of literature to help choose the most appropriate method for incorporating post-Hartree-Fock methods, should that be desired.

A clear conclusion from such comparative studies is that density functional methods using gradient-corrected functionals can give results for a wide variety of properties that are competitive with, and in some cases superior to, *ab initio* calculations using correlation (e.g. MP2). Gradient-corrected functionals are required for the calculation of relative conformational energies and the study of intermolecular systems, particularly those involving hydrogen bonding [Sim *et al.* 1992]. As is the case with the *ab initio* methods the choice of basis set is also important in determining the results. By keeping the basis set constant (6-31G* being a popular choice) it is possible to make objective comparisons. Four examples of such comparative studies are those of Johnson and colleagues, who considered small neutral molecules [Johnson *et al.* 1993]; St-Amant *et al.*, who examined small organic molecules [St-Amant *et al.* 1995]; Stephens *et al.*, who performed a detailed study of the absorption and circular dichroism spectra of 4-methyl-2-oxetanone [Stephens *et al.* 1994]; and Frisch *et al.*, who compared a variety of density functional methods with one another and to traditional *ab initio* approaches [Frisch *et al.* 1996]. The evolution of defined sets of data such as those associated with the Gaussian-*n* series of models has also acted as a spur to those involved in developing density functional methods. For example, much of Becke's work on gradient corrections and on mixed Hartree-Fock/density function methods was evaluated using data sets originally collated for the Gaussian-1 and Gaussian-2 methods. A more recent example is a variant of the Gaussian-3 method which uses B3LYP to determine geometries and zero-point energies [Baboul *et al.* 1999].

One of the most important developments for the practical application of DFT were methods for calculating analytical gradients of the energy with respect to the nuclear coordinates. This enables molecular geometries to be optimised. Once more there are some differences between the way this is done with density functional theory compared with Gaussian-based Hartree-Fock methods. A potential problem is that the use of grid-based integration schemes makes it difficult to provide exact expressions for the gradients. However, the errors associated with the grid-based method are generally very small and do not cause problems during the optimisation.

3.8 Quantum Mechanical Methods for Studying the Solid State

3.8.1 Introduction

The quantum mechanical methods used to study the behaviour of solid-phase systems are often somewhat different to those traditionally employed for studies of individual molecules or isolated intermolecular complexes. A perfect crystalline system can be constructed by stacking copies of some repeating unit (the *unit cell*) in a systematic fashion without overlapping and without gaps. The structure of a crystal can be specified by defining the size and shape of the unit cell and the positions of the atoms within it. The unit cell is parallelepiped in shape and is characterised by three lattice vectors \mathbf{a} , \mathbf{b} and \mathbf{c} and the angles between them (Figure 3.7). It may be possible to conceive of more than one unit cell, with different unit cell parameters. In such cases a set of standard cell parameters can be obtained by the application of standardisation rules. The coordinates of the atoms in the unit cell may be expressed as fractional coordinates $(\alpha\mathbf{a}, \beta\mathbf{b}, \gamma\mathbf{c})$. Indeed, any general vector \mathbf{r} can be written in terms of these basis vectors:

$$\mathbf{r} = (\alpha\mathbf{a}, \beta\mathbf{b}, \gamma\mathbf{c}) \quad (3.76)$$

where α , β and γ are not necessarily restricted to values between 0 and 1. There are fourteen different types of basic unit cell; these are the *Bravais lattices*. Common Bravais lattices include the simple cubic, body-centred cubic and face-centred cubic (Figure 3.8). Another common structure also shown in Figure 3.8 is the hexagonal close-packed arrangement, for which the underlying Bravais lattice (called the simple hexagonal) is formed from an underlying triangular arrangement. In addition to the translational symmetry that the unit cell must possess there may be some symmetry to the arrangement of the atoms within the unit cell. The particular combination of symmetry elements in a crystal defines its *space group*. There are 230 different space groups. If there is symmetry within the unit cell then it is strictly only necessary to specify the asymmetric unit (the unique part of the structure); the positions of the other atoms can be generated using the appropriate symmetry operators.

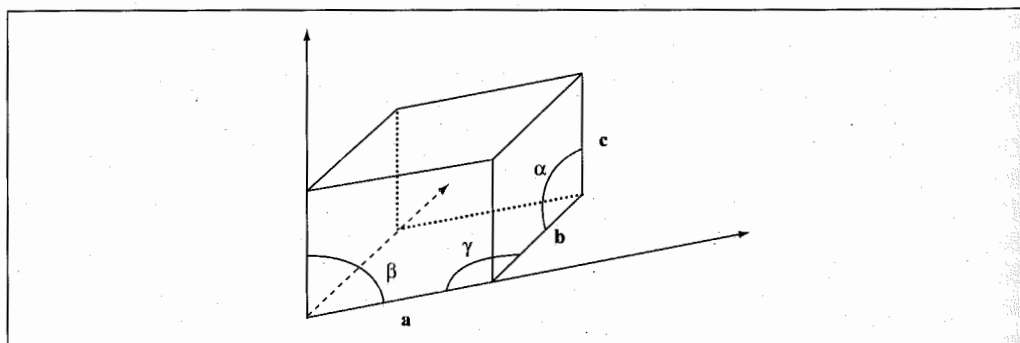


Fig. 3.7: The six parameters \mathbf{a} , \mathbf{b} , \mathbf{c} , α , β , γ which characterise the unit cell.

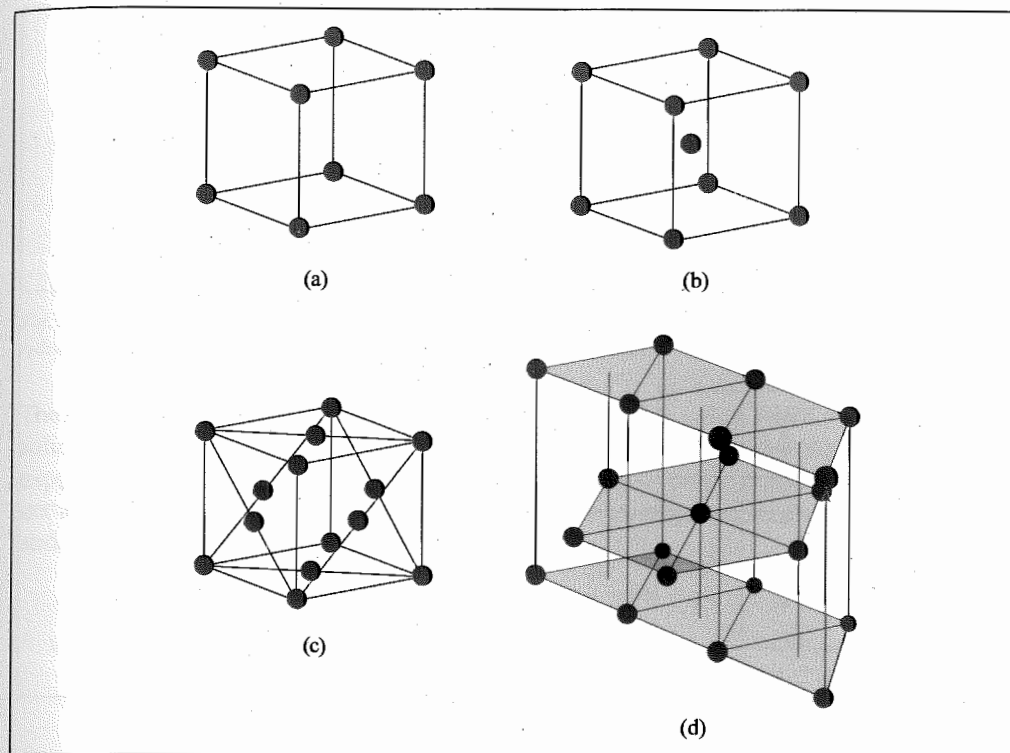


Fig. 3.8: Some basic Bravais lattices: (a) simple cubic, (b) body-centred cubic, (c) face-centred cubic and (d) simple hexagonal close-packed. (Figure adapted in part from Ashcroft N W and Mermin N D 1976. Solid State Physics. New York, Holt, Rinehart and Winston.)

Another concept that is extremely powerful when considering lattice structures is the *reciprocal lattice*. X-ray crystallographers use a reciprocal lattice defined by three vectors \mathbf{a}^* , \mathbf{b}^* and \mathbf{c}^* in which \mathbf{a}^* is perpendicular to \mathbf{b} and \mathbf{c} and is scaled so that the scalar product of \mathbf{a}^* and \mathbf{a} equals 1. \mathbf{b}^* and \mathbf{c}^* are similarly defined. In three dimensions this leads to the following definitions:

$$\mathbf{a}^* = \frac{\mathbf{b} \times \mathbf{c}}{\mathbf{a} \cdot \mathbf{b} \times \mathbf{c}}; \quad \mathbf{b}^* = \frac{\mathbf{a} \times \mathbf{c}}{\mathbf{b} \cdot \mathbf{a} \times \mathbf{c}}; \quad \mathbf{c}^* = \frac{\mathbf{a} \times \mathbf{b}}{\mathbf{c} \cdot \mathbf{a} \times \mathbf{b}} \quad (3.77)$$

Note that the denominator in each case is equal to the volume of the unit cell. The fact that \mathbf{a}^* , \mathbf{b}^* and \mathbf{c}^* have the units of 1/length gives rise to the terms 'reciprocal space' and 'reciprocal lattice'. It turns out to be convenient for our computations to work with an expanded reciprocal space that is defined by three closely related vectors \mathbf{a}^s , \mathbf{b}^s and \mathbf{c}^s , which are multiples by 2π of the X-ray crystallographic reciprocal lattice vectors:

$$\mathbf{a}^s = 2\pi\mathbf{a}^*; \quad \mathbf{b}^s = 2\pi\mathbf{b}^*; \quad \mathbf{c}^s = 2\pi\mathbf{c}^* \quad (3.78)$$

A simple illustrative example of reciprocal space is that of a 2D square lattice where the vectors \mathbf{a} and \mathbf{b} are orthogonal and of length equal to the lattice spacing, a . Here \mathbf{a}^* and \mathbf{b}^* are directed along the same directions as \mathbf{a} and \mathbf{b} respectively and have a length $1/a$

combined to give the equivalent of molecular orbitals. It is based on the assumption that the effect of orbital overlap is to modulate but not change completely the initial atomic levels. The approximation is traditionally considered most useful for describing the electronic structure of systems such as insulators and transition metals with partially filled d shells. The second approach is called the *nearly free-electron approximation*. This theory starts by considering the electrons as free particles whose motion is modulated by the presence of the lattice. The nearly free-electron approximation is traditionally considered the more suitable approach to systems such as metals where there is substantial overlap of the valence orbitals. We will outline both approaches in turn, making use of some of the fundamental principles and properties of lattices discussed earlier.

3.8.2 Band Theory and Orbital-based Approaches

Band theory is perhaps easier for chemists to understand, starting as it does from an orbital picture. We will therefore spend somewhat less space discussing this than the nearly free-electron approximation. We will start by considering the simplest problem, a 1D lattice. Initially we consider what happens if we bring together two atoms along the x axis until they are separated by a distance, a . If each atom has one s orbital, then the combined system has two molecular orbitals (one bonding and one anti-bonding). If we then add a third atom then three molecular orbitals are obtained (one bonding, one non-bonding and one anti-bonding). Four atoms give four energy levels, and so on. As more atoms are added the energy levels merge to give what is an essentially continuous *band* of energy levels (Figure 3.11). Each energy level can accommodate two electrons so if each atom contributes one electron the band will be half full. The presence of unoccupied energy levels near to the top of the filled level means that it is very easy to excite electrons from the filled to the unfilled levels. The electrons are consequently very mobile, giving rise to the special conduction and thermal properties of a metal. By contrast, if each atom contributes two electrons then the band will be completely filled. Such electrons would have to be excited to higher bands, which might, for example, be formed by the overlap of p orbitals. In an insulator the energy of this p band would typically be significantly higher than the lower s band and so excitation would require considerable energy. In a semiconductor the band gap is smaller and it may be possible to excite electrons from the top of the highest filled band (the *valence band*) to the lowest unoccupied band (the *conduction band*) at normal temperatures. These three difference scenarios are illustrated in Figure 3.11.

The periodicity of the lattice means that the values of a function (such as the electron density) will be identical at equivalent points on the lattice. Likewise there is a relationship between the wavefunction at a point (x in our 1D lattice) and at an equivalent point elsewhere on the lattice (for the 1D lattice this would be $x + na$, where n is an integer). *Bloch's theorem* provides the link; each allowed lattice wavefunction must satisfy the following relationship:

$$\psi^k(x+a) = e^{ika} \psi^k(x) \quad (3.81)$$

In this equation we have identified the wavefunction with a label, k , which for now can be considered an index; there are as many values of k as there are atoms in the 1D lattice.

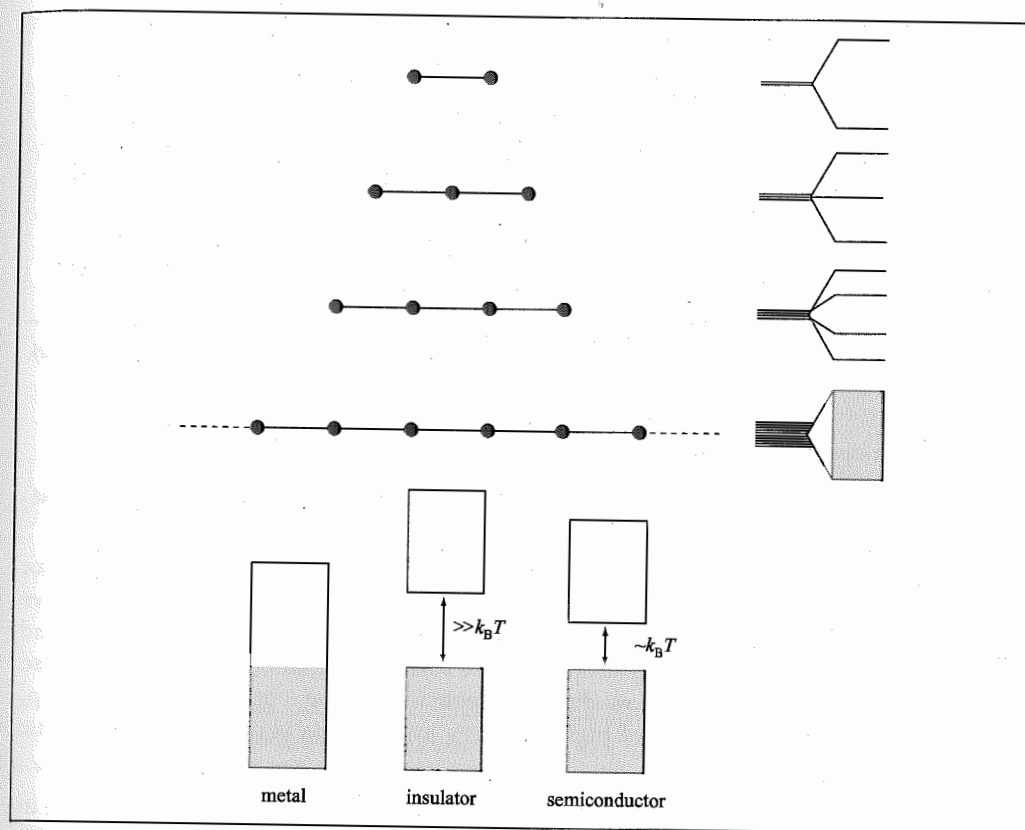


Fig. 3.11: The creation of a band of energy levels from the overlap of two, three, four, etc. atomic orbitals, which eventually gives rise to a continuum. Also shown are the conceptual differences between metals, insulators and semiconductors.

We wish to construct linear combinations of the atomic orbitals such that the overall wavefunction meets the Bloch requirement. Suppose the s orbitals in our lattice are labelled χ_n where the n th orbital is located at position $x = na$. An acceptable linear combination of these orbitals that satisfies the Bloch requirements is:

$$\psi^k = \sum_n e^{ikna} \chi_n \quad (3.82)$$

We now need to consider how the form of the wavefunction varies with k . The first situation we consider corresponds to $k = 0$, where the exponential terms are all equal to 1 and the overall wavefunction becomes a simple additive linear combination of the atomic orbitals:

$$\psi^{k=0} = \sum_n \chi_n = \chi_0 + \chi_1 + \chi_2 + \dots \quad (3.83)$$

The other situation we consider is $k = \pi/a$. Recall that $\exp(ix)$ can be written $\cos(x) + i \sin(x)$. If $k = \pi/a$ then the sine terms will all be zero, leaving just the cosine terms $\cos(n\pi)$, which can

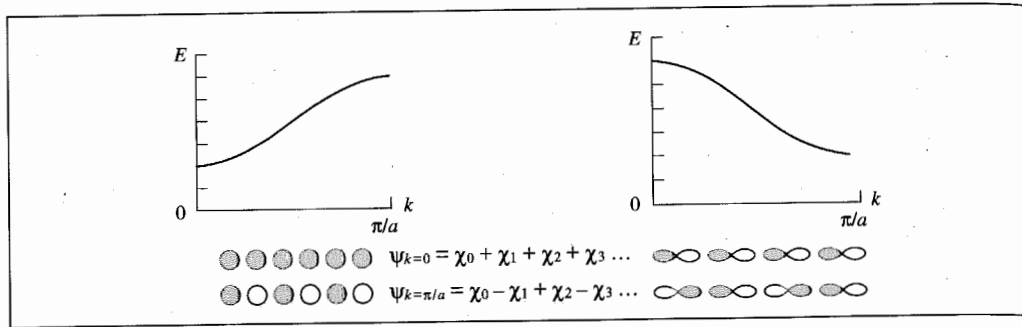


Fig. 3.12: The variation in energy with k for a 1D lattice for a set of s orbitals (left) and for a set of p_x orbitals (right). Also shown are the corresponding arrangements of orbitals.

be expressed more generally as $(-1)^n$. Hence the wavefunction is:

$$\psi^{k=\pi/a} = \sum_n (-1)^n \chi_n = \chi_0 - \chi_1 + \chi_2 - \dots \quad (3.84)$$

Equations (3.83) and (3.84) correspond to the lowest- and highest-energy wavefunctions for our simple system over this range of k . Wavefunctions for values of k between 0 and π/a have intermediate energies. The energy varies in a cosine-like manner with k between $k = 0$ and $k = \pi/a$ (Figure 3.12). Note that k can adopt negative values and that $E(-k)$ equals $E(k)$. Also worthy of note is that p orbitals show different behaviour to the s orbitals. For a set of p_x orbitals it is the $k = 0$ state that is of highest energy and $k = \pi/a$ is of lowest energy, due to their nodal behaviour.

The graph of energy versus k is called the *band structure*; the *bandwidth* is the difference in energy between the lowest and highest levels in the band. For the one-dimensional lattice the bandwidth is determined by the lattice spacing; a smaller spacing a gives a greater bandwidth in much the same way that the difference between the bonding and antibonding orbitals in H_2 increases as the atoms get closer together. As we noted above there are as many values of k (and so as many energy levels) as there are atoms in the lattice and that each energy level can accommodate two electrons.

We now move on to consider a two-dimensional square lattice in the (x, y) plane, where the inter-lattice spacing is still a . The Bloch theorem is now written in the following more general form:

$$\psi^{\mathbf{k}}(\mathbf{r} + \mathbf{T}) = e^{i\mathbf{k}\cdot\mathbf{T}} \psi^{\mathbf{k}}(\mathbf{r}) \quad (3.85)$$

In Equation (3.85) \mathbf{T} is a translation vector that maps each position into an equivalent position in a neighbouring cell, \mathbf{r} is a general positional vector and \mathbf{k} is the *wavevector* which characterises the wavefunction. \mathbf{k} has components k_x and k_y in two dimensions and is equivalent to the parameter k in the one-dimensional system. For the two-dimensional square lattice the Schrödinger equation can be expressed in terms of separate wavefunctions along the x - and y -directions. This results in various combinations of the atomic $1s$ orbitals, some of which are shown in Figure 3.13. These combinations have different energies. The lowest-energy solution corresponds to $(k_x = 0, k_y = 0)$ and is a straightforward linear

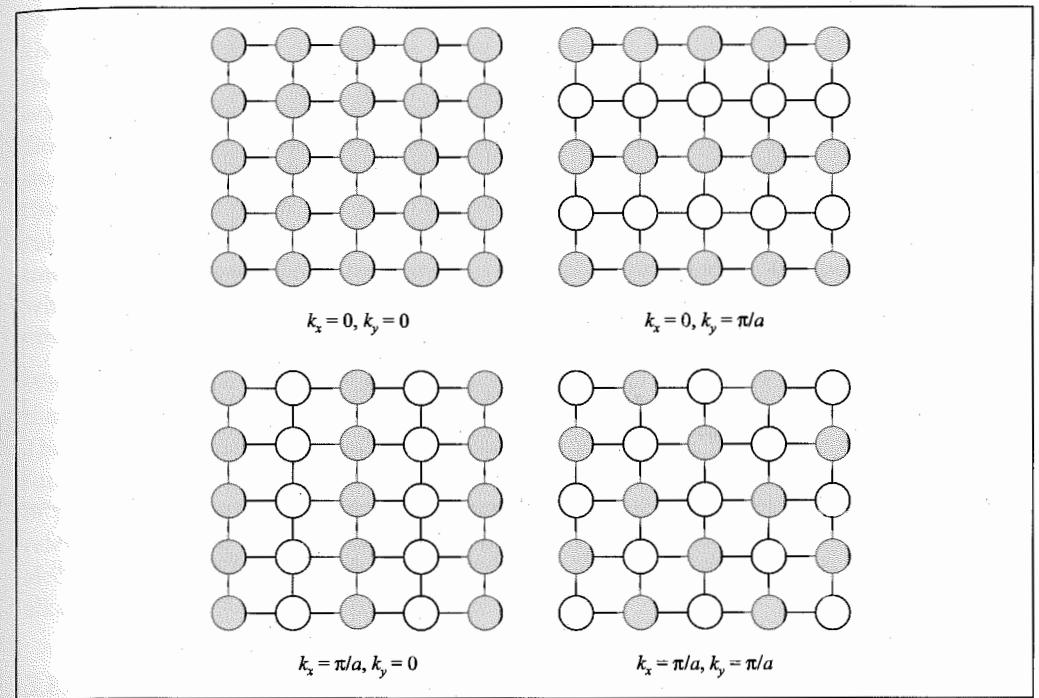


Fig. 3.13: Some of the possible combinations of atomic $1s$ orbitals for a 2D square lattice corresponding to different values of k_x and k_y . A shaded circle indicates a positive coefficient; an open circle corresponds to a negative coefficient.

combination of the atomic orbitals. The highest-energy solution corresponds to the situation where both k_x and k_y have values of π/a . The wavefunction for this high-energy solution shows a rapid variation in sign. Another important feature evident in Figure 3.13 is the wave-like nature of the various linear combinations, particularly if one imagines the lattice extending infinitely in all directions over the (x, y) plane.

The reciprocal space and the reciprocal lattice are directly related to the wavevector, \mathbf{k} ; different values of \mathbf{k} can be considered as points within the reciprocal space defined by \mathbf{a}^s , \mathbf{b}^s and \mathbf{c}^s . It turns out that, when we are calculating the wavefunction and energy levels for a solid, we need to restrict \mathbf{k} to one cell in the reciprocal lattice (typically chosen to the cell containing $\mathbf{k} = 0$, or the first Brillouin zone), otherwise there is a danger of counting some states more than once. A very common way to represent the band structure for lattice structures is to plot how the energy changes as a function of \mathbf{k} along certain lines of symmetry within the first Brillouin zone. For example, to return to our square lattice (for which the reciprocal lattice is also square) one could imagine taking a 'tour' starting at the origin ($\mathbf{k} = (0, 0)$), moving along the x axis to $\mathbf{k} = (\pi/a, 0)$ up the y axis to $\mathbf{k} = (\pi/a, \pi/a)$, and finally returning to the origin. As we undertake this tour the energy changes as shown in Figure 3.14. In this diagram we have labelled certain values of \mathbf{k} which have particular symmetry with their conventional Roman or Greek capital letters, Γ , X and M [Bradley and Cracknell 1972].

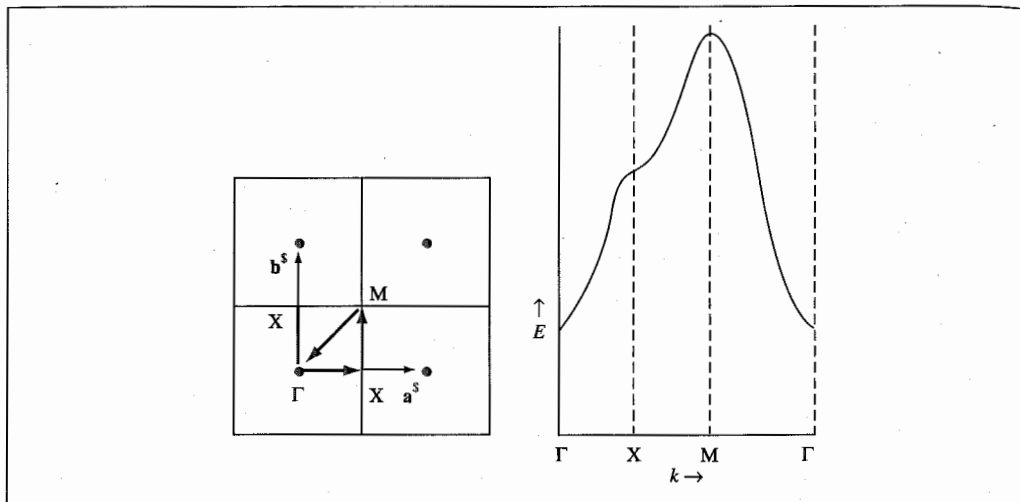


Fig. 3.14: Variation in energy for a 'tour' (Γ -X-M- Γ) of the reciprocal lattice for a 2D square lattice of hydrogen atoms. (Figure adapted in part from Hoffmann R 1988. Solids and Surfaces: A Chemist's View on Bonding in Extended Structures. New York, VCH Publishers.)

3.8.3 The Periodic Hartree-Fock Approach to Studying the Solid State

In the *periodic Hartree-Fock* approach the elements of the Fock matrix are constructed from linear combinations of so-called Bloch functions:

$$\psi_i^{\mathbf{k}}(\mathbf{r}) = \sum_{\omega} a_{\omega i}(\mathbf{k}) \varphi_{\omega}^{\mathbf{k}}(\mathbf{r}) \quad (3.86)$$

Each Bloch function is itself a linear combination of atomic orbitals:

$$\varphi_{\omega}^{\mathbf{k}}(\mathbf{r}) = \sum_{\mathbf{T}} \chi_{\omega}^{\mathbf{T}}(\mathbf{r}) \exp(i\mathbf{k} \cdot \mathbf{T}) \quad (3.87)$$

$\chi_{\omega}^{\mathbf{T}}$ is the ω th atomic orbital in the crystal cell characterised by the lattice vector \mathbf{T} . As such, this method works in real space, which contrasts with the usual implementations of the alternative plane-wave methods that we will discuss below [Dovesi *et al.* 2000]. Each atomic orbital is expressed as a linear combination of (for example) Gaussian functions, as in molecular Hartree-Fock theory. The coefficients $a_{\omega i}(\mathbf{k})$ in Equation (3.86) are obtained by solving the following matrix equation for every value of \mathbf{k} to self-consistency:

$$\mathbf{F}_{\mathbf{k}} \mathbf{A}_{\mathbf{k}} = \mathbf{S}_{\mathbf{k}} \mathbf{A}_{\mathbf{k}} \mathbf{E}_{\mathbf{k}} \quad (3.88)$$

$\mathbf{S}_{\mathbf{k}}$ is the overlap matrix for the Bloch functions for the wavevector \mathbf{k} , with $\mathbf{E}_{\mathbf{k}}$ being the energy matrix and \mathbf{A} the matrix of coefficients. $\mathbf{F}_{\mathbf{k}}$ is the Fock matrix, which consists of a sum of one- and two-electron terms. The values of \mathbf{k} are typically selected to sample from the first Brillouin zone according to a special scheme as described in Section 3.8.6. When these terms are expanded they involve infinite sums over the nuclei and electrons in the lattice. As is usual in a Hartree-Fock approach the one-electron terms involve the sum of a kinetic energy term and one due to the Coulomb interaction between the nuclei and the

electrons; the two-electron terms involve Coulomb and exchange two-electron integrals. Unfortunately, if these sums were to be evaluated individually and to completion then they would not converge to a consistent value, but would diverge. However, effective ways to determine these infinite sums have been proposed [Pisani and Dovesi 1980; Dovesi *et al.* 1983]. These involve a variety of procedures. The Coulomb interactions are divided into a series of terms corresponding to interacting and non-interacting charge distributions. The latter can then be grouped together into 'shells' and the interaction calculated using multipole expansions (see Section 4.9.1). For the shorter-range exchange interaction it is possible to truncate the integral summation at an appropriate distance without loss of accuracy. The truncation distance can depend upon the three-dimensional structure of the material and so may vary from one calculation to the next.

Within the periodic Hartree-Fock approach it is possible to incorporate many of the variants that we have discussed, such as UHF or RHF. Density functional theory can also be used. This makes it possible to compare the results obtained from these variants. Whilst density functional theory is more widely used for solid-state applications, there are certain types of problem that are currently more amenable to the Hartree-Fock method. Of particular relevance here are systems containing unpaired electrons, two recent examples being the electronic and magnetic properties of nickel oxide and alkaline earth oxides doped with alkali metal ions (Li in CaO) [Dovesi *et al.* 2000].

3.8.4 The Nearly Free-electron Approximation

Whereas the tight-binding approximation works well for certain types of solid, for other systems it is often more useful to consider the valence electrons as free particles whose motion is modulated by the presence of the lattice. Our starting point here is the Schrödinger equation for a free particle in a one-dimensional, infinitely large box:

$$\left(\frac{d^2}{dx^2}\right)\psi = -\left(\frac{2mE}{\hbar^2}\right)\psi \quad (3.89)$$

The solutions to this equation are:

$$\psi = C \exp(ikx); \quad E = (\hbar^2 k^2)/2m \quad (3.90)$$

The energy for a free particle can be related to the momentum by $E = p^2/2m$ and so the wavefunction is related to the momentum p by:

$$\psi = C \exp(\pm ipx/\hbar) \quad (3.91)$$

The wavelength of this motion is h/p and the parameter k is equal to $2\pi p/h$. Thus k has units of 1/length (i.e. reciprocal length). The energy for a free particle varies in a quadratic fashion with k and in principle any value of the energy is possible.

In two dimensions we obtain the following wavefunction:

$$\psi_{x,y} = C_x \exp(ik_x x/\hbar) C_y \exp(ik_y y/\hbar) = C \exp(i\mathbf{k} \cdot \mathbf{r}/\hbar) \quad (3.92)$$

Note that in Equation (3.92) we have expressed the wavefunction in terms of a vector, \mathbf{k} (which has components in the x and y directions of k_x and k_y) and the Cartesian vector \mathbf{r} .

The energy varies as a quadratic function of both k_x and k_y :

$$E_{x,y} = \frac{\hbar^2}{2m}(k_x^2 + k_y^2) \quad (3.93)$$

An analogous expression is obtained in three dimensions. We now need to consider periodic systems. As we have discussed, the wavefunction for a particle on a periodic lattice must satisfy Bloch's theorem, Equation (3.85). The wavevector \mathbf{k} in Bloch's theorem plays the same role in the study of periodic systems as the vector \mathbf{k} does for a free particle. One important difference is that whereas the wavevector is directly related to the momentum for a free particle (i.e. $\mathbf{k} = \mathbf{p}/\hbar$) this is not the case for the Bloch particle due to the presence of the external potential (i.e. the nuclei). However, it is very convenient to consider $\hbar\mathbf{k}$ as analogous to the momentum and it is often referred to as the *crystal momentum* for this reason. The possible values that \mathbf{k} can adopt are given by:

$$\mathbf{k} = \left(\frac{m_\alpha}{N_\alpha} \mathbf{a}^s, \frac{m_\beta}{N_\beta} \mathbf{b}^s, \frac{m_\gamma}{N_\gamma} \mathbf{c}^s \right) \quad (3.94)$$

m_α, m_β and m_γ are integers and $N_\alpha N_\beta N_\gamma = N$, the number of unit cells in the crystal. For a macroscopic system where N is very large (of the order of Avogadro's number) \mathbf{k} thus varies continuously. As we have seen before, the wavevector \mathbf{k} in the Bloch theorem (Equation (3.85)) can be considered as a point within the reciprocal lattice defined by $\mathbf{a}^s, \mathbf{b}^s$ and \mathbf{c}^s . It can also be shown (see Appendix 3.1) that a wavefunction that satisfies Bloch's theorem can be written in the following form:

$$\psi^{\mathbf{k}}(\mathbf{r}) = e^{i\mathbf{k}\cdot\mathbf{r}} u^{\mathbf{k}}(\mathbf{r}) \quad (3.95)$$

Here, $u^{\mathbf{k}}(\mathbf{r})$ is a function that is periodic on the lattice. Recall from our earlier discussions on reciprocal lattice vectors that one way to construct such a periodic function is as a Fourier series expansion of plane wavefunctions $\exp(i\mathbf{G}\cdot\mathbf{r})$:

$$u^{\mathbf{k}}(\mathbf{r}) = \sum_{\mathbf{G}} c_{\mathbf{G}}^{\mathbf{k}} \exp(i\mathbf{G}\cdot\mathbf{r}) \quad (3.96)$$

The sum runs over the reciprocal lattice vectors \mathbf{G} we considered above. A simple case is $\mathbf{G} = \mathbf{a}^s$, for which $\exp(i\mathbf{G}\cdot\mathbf{r})$ corresponds to a wave travelling perpendicular to the real-space axes \mathbf{b} and \mathbf{c} and with a wavelength such that it fits exactly into the unit cell. If $\mathbf{G} = 2\mathbf{a}^s$ then two wavelengths fit into the cell.

The external potential due to the nuclei is periodic in the lattice and it too can be written as a Fourier expansion of exponential functions of the reciprocal lattice:

$$U(\mathbf{r}) = \sum_{\mathbf{G}} U_{\mathbf{G}} \exp(i\mathbf{G}\cdot\mathbf{r}) \quad (3.97)$$

$U_{\mathbf{G}}$ is the Fourier coefficient. When this form of the potential is incorporated into the Schrödinger equation the following equation can be derived [Ashcroft and Mermin 1976]:

$$\left(\frac{\hbar^2}{2m} |\mathbf{k} + \mathbf{G}|^2 - E \right) c_{\mathbf{G}}^{\mathbf{k}} + \sum_{\mathbf{G}'} U_{\mathbf{G}'+\mathbf{G}} c_{\mathbf{G}'}^{\mathbf{k}} = 0 \quad (3.98)$$

We can recover the free-particle result (i.e. zero potential) from Equation (3.98) by setting all of the Fourier coefficients $U_{\mathbf{G}}$ to zero, in which case the equation reduces to:

$$\left(\frac{\hbar^2}{2m} |\mathbf{k} + \mathbf{G}|^2 - E \right) c_{\mathbf{G}}^{\mathbf{k}} = 0 \quad (3.99)$$

The solution of this equation requires that $E = \hbar^2 |\mathbf{k} + \mathbf{G}|^2 / 2m$ with the wavefunctions being of the form $\psi(\mathbf{r}) \propto \exp[i(\mathbf{k} + \mathbf{G})\cdot\mathbf{r}]$. Although cast in a slightly different form, this is equivalent to our earlier expression for the wavefunction of a free particle, Equation (3.92).

The summations in Equations (3.98) are over all reciprocal lattice vectors \mathbf{G} . As can be seen, for a given value of \mathbf{k} there are as many forms of this equation as there are reciprocal lattice vectors in the system. Each of these equations for the different values of \mathbf{G} gives rise to a solution which is labelled with the band index n . Thus there are as many values of n as there are reciprocal lattice vectors \mathbf{G} . Just as there are n solutions to this Schrödinger equation for a given value of \mathbf{k} , so it is also possible to consider how the energy varies with \mathbf{k} for a given value of n . To understand the entire band structure of a solid requires one to consider the variation of both \mathbf{k} and n . As we indicated above, when calculating the band structure it is usual to restrict \mathbf{k} to just the first Brillouin zone to avoid duplicate counting of states.

Let us now examine how these results can be applied to some simple one- and two-dimensional periodic systems. Initially we will consider the situation where there is no external potential and then discuss what happens when we introduce one. The first case is the one-dimensional lattice, which has reciprocal lattice vectors at $\pm 2\pi/a, \pm 4\pi/a$, etc. In order to derive the energy diagram we need to consider, for each reciprocal lattice vector \mathbf{G} , how the energy varies as we change \mathbf{k} over the first Brillouin zone (which in this case corresponds to varying \mathbf{k} from $-\pi/a$ to $+\pi/a$). The first reciprocal lattice vector is $\mathbf{G} = 0$, for which the energy simply varies quadratically with \mathbf{k} , from zero at $\mathbf{k} = 0$ to $\hbar^2(\pi/a)^2/2m$ at $\mathbf{k} = \pm 2\pi/a$. We next need to consider the two reciprocal vectors $\mathbf{G} = \pm 2\pi/a$. At the point $\mathbf{k} = 0$ the energy due to both of these reciprocal lattice vectors is $\hbar^2(2\pi/a)^2/2m$. As \mathbf{k} increases from 0 to $+\pi/a$ the value of $|\mathbf{k} + \mathbf{G}|^2$ increases for the reciprocal lattice vector $\mathbf{G} = 2\pi/a$ but it decreases for the reciprocal lattice vector $\mathbf{G} = -2\pi/a$. Conversely, as \mathbf{k} varies from 0 to $-\pi/a$ the energy increases for the reciprocal lattice vector $\mathbf{G} = -2\pi/a$ and decreases for $\mathbf{G} = 2\pi/a$. These variations in energy are shown in Figure 3.15. Two types of energy diagram are shown in this figure; one is the 'reduced-zone' scheme because the entire dependency of the energy on the wavevector is contained within the first Brillouin zone. The alternative representation is called an extended-zone scheme in which the energy levels are 'folded out' for values of \mathbf{k} beyond the first Brillouin zone.

We next need to introduce the weak potential, which acts to modulate the wavefunctions and the associated energy levels. The effects of the potential are found to be most acute where there is degeneracy of the energy levels. This arises even in the one-dimensional situation, where we have degenerate energy levels due to different reciprocal lattice vectors at $\mathbf{k} = 0$ and $\mathbf{k} = \pi/a$. The effect of the potential is to perturb these energy levels in such a way that lifts the degeneracy to create an energy gap. In the one-dimensional case the effect of the potential is to 'flatten' the energy levels in the region close to the edge of the

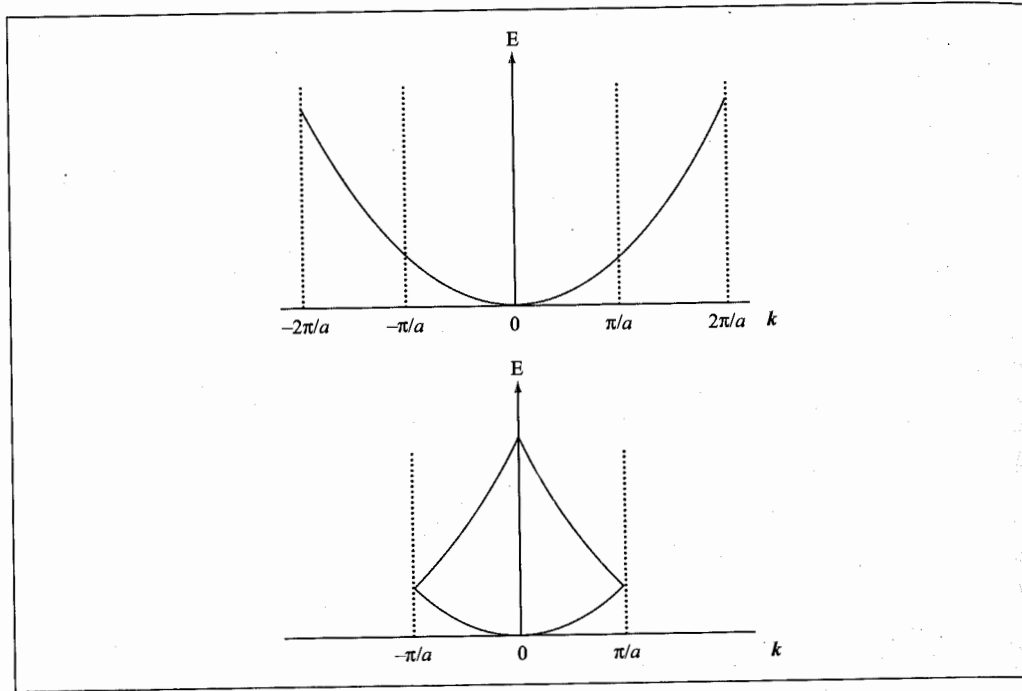


Fig. 3.15: Extended-zone and reduced-zone representations of band diagram for 1D lattice with no external potential.

Brillouin zone as shown in Figure 3.16. One way to explain the appearance of the energy gap at the edges of the Brillouin zone is to recognise that the states of a free electron are waves with a specific wavelength ($2\pi/k$ in the simple one-dimensional system). When the wavelength becomes comparable to the lattice spacing the lattice diffracts the wave and at the boundary of the Brillouin zone ($k = \pm\pi/a$) a standing wave is created. Two different standing waves are possible in a one-dimensional system, as shown in Figure 3.17. For one of the standing waves (A in Figure 3.17) the peak electron density occurs in the vicinity of the lattice points (the positive nuclei). This standing wave thus has a more favourable (i.e. lower) energy than the equivalent free travelling wave. By contrast, the peak electron density of the other standing wave (B in Figure 3.17) occurs between the nuclei and so its energy is higher. Further gaps arise at $k = \pm 2\pi/a$, and so on.

A somewhat more complex case is that of the 2D hexagonal lattice. As for the one-dimensional system we initially consider a free particle, restricting ourselves to wavevectors within the first Brillouin zone with higher-energy states being due to reciprocal lattice vectors beyond in the second, third, etc. Brillouin zones. We will consider how the energy varies as we undertake a 'tour' of the first Brillouin zone in reciprocal space starting at the origin ($\mathbf{k} = (0,0)$), then moving to one of the vertices of the hexagon (the point ($\mathbf{k} = \cos \pi/6, \sin \pi/6$)), along to the mid-point of one of the edges ($\mathbf{k} = (0, \sin \pi/6)$), and finally back to the origin (Figure 3.18). The origin, the vertex and the mid-point are all points of symmetry and are identified by the symbols Γ , K and M, respectively. For a

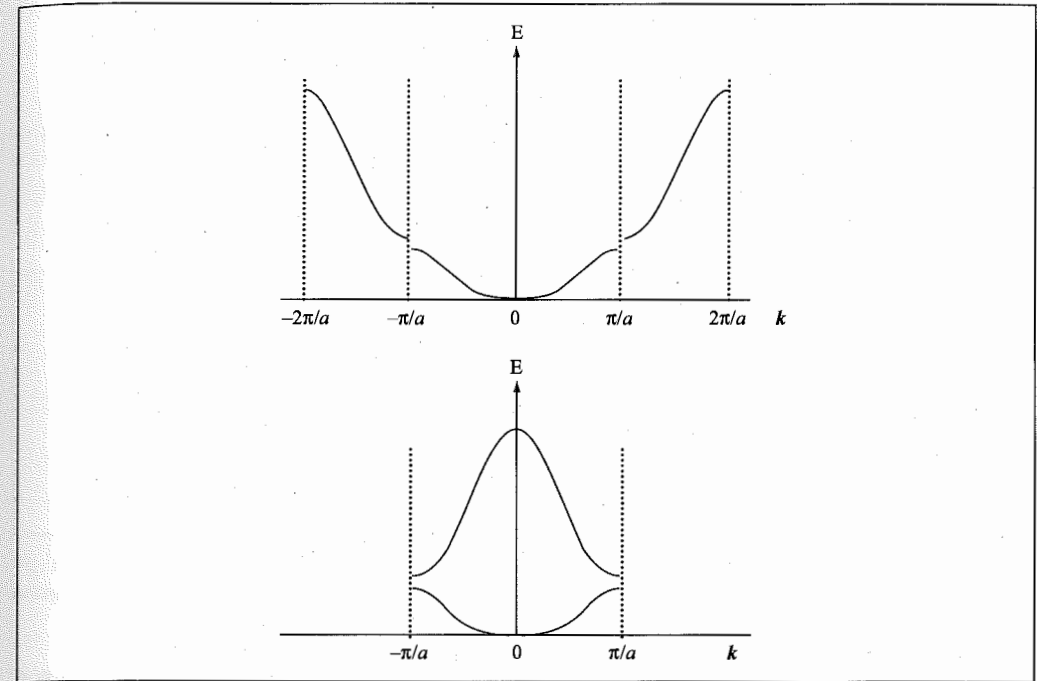


Fig. 3.16: The effect of introducing a weak potential into the 1D lattice is to lift the degeneracy of the energy levels near to the edge of the Brillouin zone (shown in both extended-zone and reduced-zone representation).

given value of k we compute the value of $|\mathbf{k} + \mathbf{G}|^2$ and thus the energy for the relevant reciprocal lattice vectors.

The simplest case is that corresponding to $\mathbf{G} = 0$. We still obtain a quadratic variation of energy with $|\mathbf{k}|$ wherever we move within the first Brillouin zone. The variation in energy for the three 'legs' of this tour can be represented in an energy band diagram as shown in Figure 3.18. As there are six nearest-neighbour cells in this system, there are six energy levels to monitor at the next stage. The distance from the origin to each of these six reciprocal lattice points is $2 \cos \pi/6$. At $k = 0$ we therefore find that all six energy levels are degenerate

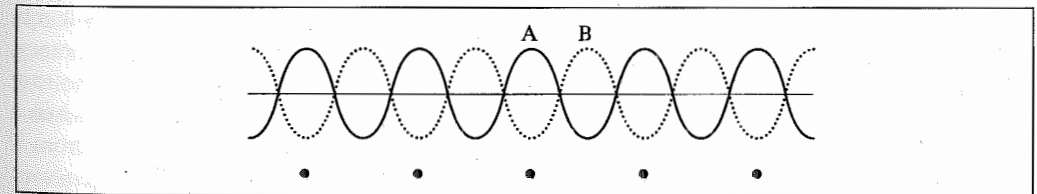


Fig. 3.17: The two possible sets of standing waves at the Brillouin zone boundary. Standing wave A concentrates electron density at the nuclei, whereas wave B concentrates electron density between the nuclei. Wave A thus has a lower energy than wave B.

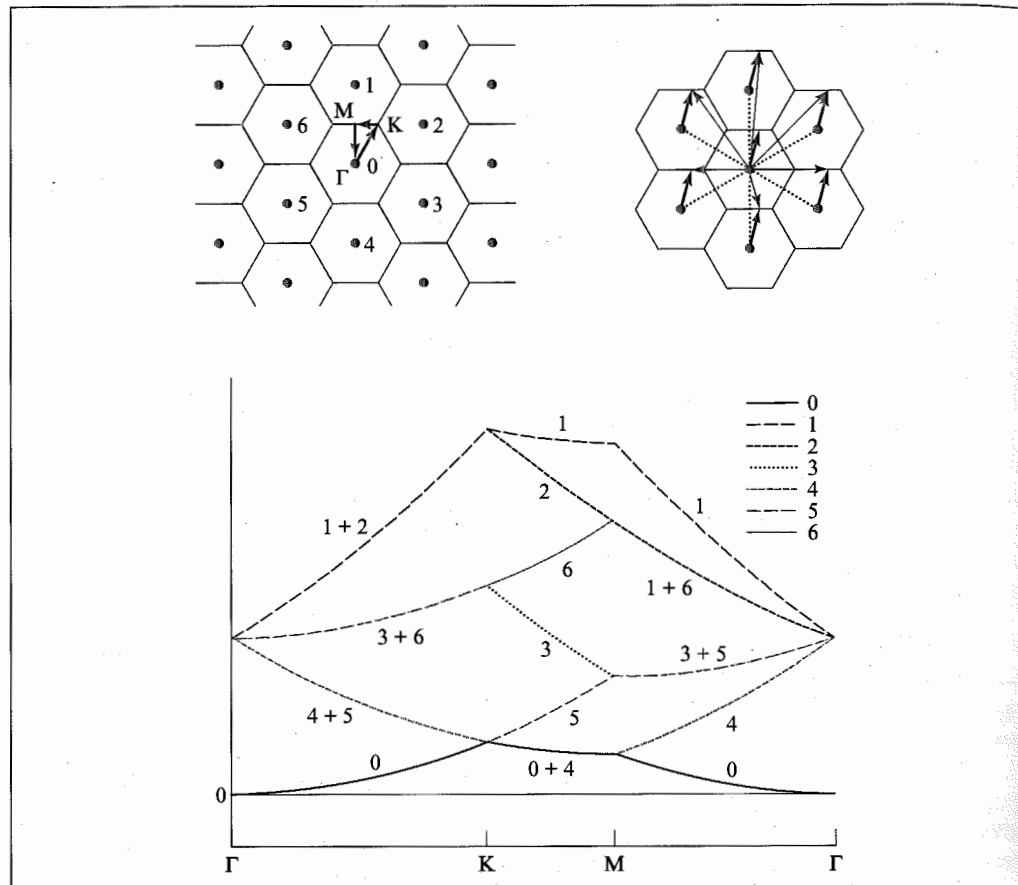


Fig. 3.18: Energy band diagram (bottom) for a free-particle 'tour' (Γ -K-M- Γ) of the reciprocal lattice for a 2D hexagonal structure (top left). A total of seven bands are shown, due to the central reciprocal lattice vector $\mathbf{G} = \mathbf{0}$ and the reciprocal lattice vectors from the six neighbouring cells. The energy varies as $|\mathbf{k} + \mathbf{G}|^2$, where the vector $\mathbf{k} + \mathbf{G}$ is computed as shown in the top right of the figure (\mathbf{k} : bold arrow; $\mathbf{k} + \mathbf{G}$: thin arrow).

and have a value of $3\hbar^2/2m$ ($(2\cos\pi/6)^2 \equiv 3$). Moving towards the point $(\cos\pi/6, \sin\pi/6)$ we find that the six vectors separate into three pairs of degenerate levels. These six reciprocal lattice points are labelled 1-6 in Figure 3.18, together with the corresponding energy levels. As the tour continues, the different energy bands show two-, three- and six-fold degeneracy, depending upon the value of \mathbf{k} . Another key feature is that along some legs of the tour certain pairs of bands are degenerate, though this degeneracy will often be lifted when a different leg is traversed. For example, the pairs 1-2, 3-6 and 4-5 are degenerate from Γ to K. Between K and M the pair 0-4 are degenerate; and on the final leg there is degeneracy between the pairs 2-6 and 3-4. When the periodic potential is introduced some, but not necessarily all, of this degeneracy will be lifted, giving rise to band gaps. The way in which this can occur is shown schematically in Figure 3.19.

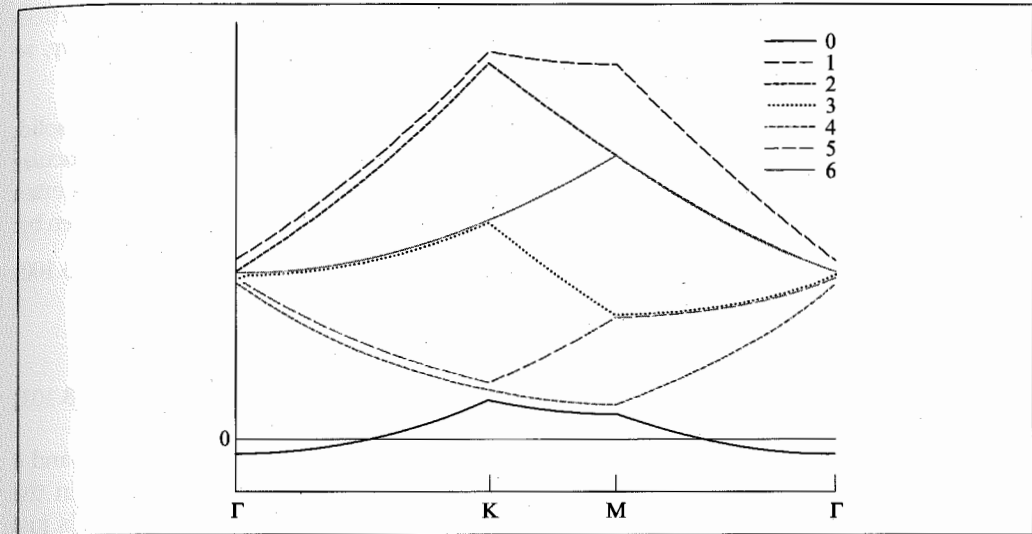


Fig. 3.19: The effect of a weak external potential is to lift degeneracy and create band gaps as illustrated for a 2D hexagonal lattice (compare with Figure 3.18).

3.8.5 The Fermi Surface and Density of States

To determine the ground state of a periodic system it is necessary to determine its band structure, by varying \mathbf{k} over the first Brillouin zone and computing at each value of \mathbf{k} the different energy bands resulting from the reciprocal lattice vectors. The number of energy levels in a band (i.e. the number of values permitted to \mathbf{k}) is equal to the number of primitive cells in the crystal, just as was the case for the orbital model in the tight-binding approximation. For each energy level corresponding to a particular value of \mathbf{k} the Pauli principle permits two electrons of opposite spin to be assigned. This process is repeated for the different bands until all the electrons have been allocated. The energy level of the highest occupied state is called the *Fermi energy* (for a metal; for an insulator, the Fermi energy is in the middle of a band gap). When all the electrons have been assigned then one of two different situations may result. In the first case all the occupied bands are completely filled. As we saw earlier, this gives rise to a band gap between the top of the highest occupied level and the bottom of the lowest empty level. The number of energy levels in each band is equal to the number of primitive cells in the crystal, so a band gap can only arise if there is an even number of electrons per primitive cell. The tight-binding approximation discussed in Section 3.8.2 may be an appropriate model to apply in this case. The second situation arises when one or more bands are partially filled. For each of these partially filled bands one can consider there to be a surface in the \mathbf{k} space that separates the occupied and the unoccupied levels, as defined by the Fermi energy. This set of surfaces is known as the *Fermi surface* and it defines a border between the occupied and unoccupied states. In many cases the Fermi surface is contained within a single band; if not, then the parts of the Fermi surface due to partially filled individual bands are known as the *branches* of the Fermi surface. The

Fermi surface will show the same underlying periodicity as the reciprocal lattice. A particularly attractive feature of the Fermi surface is that it can be measured experimentally, so providing a link between theory and experiment.

The *density of levels* is another useful way to describe the electronic structure of a solid. The density of levels indicates how many energy levels there are for a particular energy. It can thus be defined as the number of levels between E and $E + dE$. This is often normalised by volume, leading to the density of levels per unit volume $g(E)$, which is given by:

$$g(E) = \sum_n g_n(E) \quad (3.100)$$

The sum is over the bands n , with $g_n(E)$ being the density of levels in the band n :

$$g_n(E) = \frac{1}{4\pi^3} \int \delta(E - E_n(\mathbf{k})) d\mathbf{k} \quad (3.101)$$

The delta function $\delta(E - E_n(\mathbf{k}))$ has a value of 1 if $E_n(\mathbf{k})$ is in the range E to $E + dE$ and 0 otherwise. The density of states $D(E)$ is closely related to the density of levels; in the simple case where we have two electrons in each level then the density of states is just twice the density of levels. The integral of the density of states up to the Fermi level is equal to the number of electrons and the integral of the density of states multiplied by the energy is the total electronic energy:

$$N = \int D(E) dE \quad (3.102)$$

$$E_{\text{tot}} = \int D(E) E dE \quad (3.103)$$

The density of states can be usefully visualised by plotting the energy versus $D(E)$. For the simple one-dimensional situation where the energy varies in a cosine-like manner with k and the levels are equally spaced, the density of states is greatest at the top and bottom of the band (Figure 3.20). The density of states is thus inversely proportional to

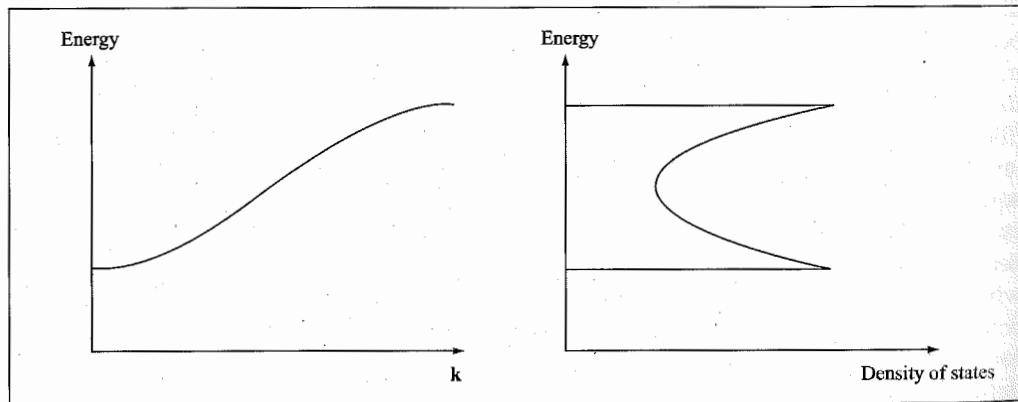


Fig. 3.20: Variation of the density of states, $D(E)$, for the simple 1D lattice, shown with the corresponding energy diagram.

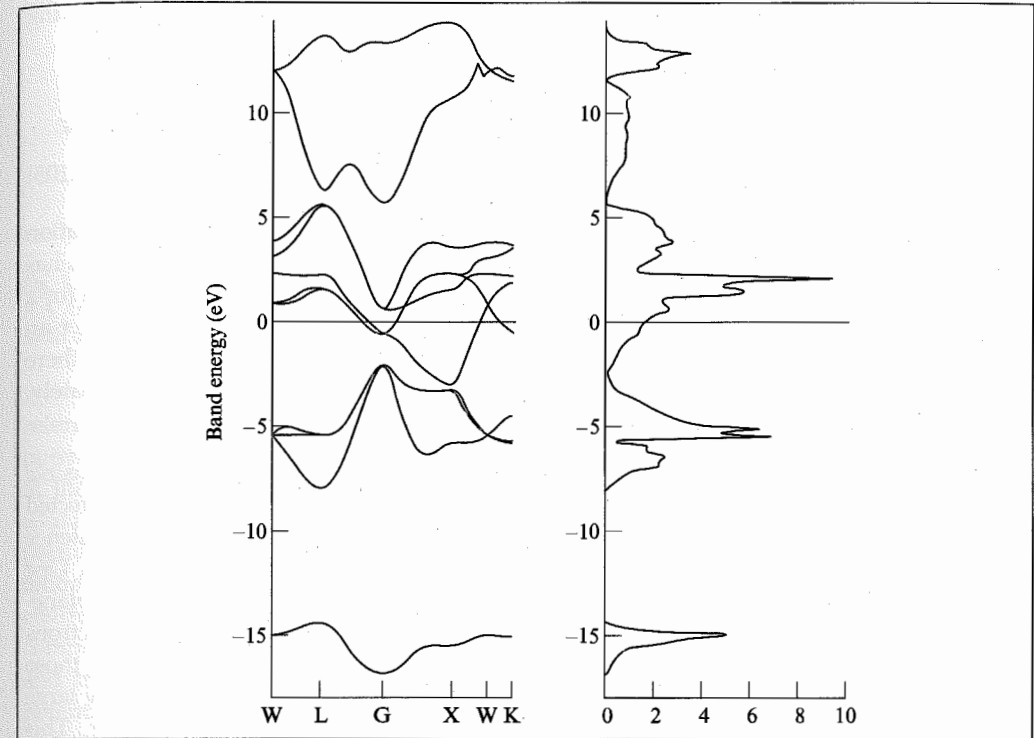


Fig. 3.21: Band structure and density of states for TiN.

the slope of the energy versus k curve; the flatter the band the greater the density of states at that energy.

The density of states is somewhat like an orbital energy diagram, but unlike the latter does not contain well-defined individual energy levels. Nevertheless, in some situations it is possible to determine from which atomic orbitals a particular energy band is largely derived. Of course, most real systems have rather more complex electronic structures than the simple cases we have used to discuss the background, as illustrated in Figure 3.21, which shows the band structure and density of states diagram for TiN.

3.8.6 Density Functional Methods for Studying the Solid State: Plane Waves and Pseudopotentials

Plane waves are often considered the most obvious basis set to use for calculations on periodic systems, not least because this representation is equivalent to a Fourier series, which itself is the natural language of periodic functions. Each orbital wavefunction is expressed as a linear combination of plane waves which differ by reciprocal lattice vectors:

$$\psi_i^{\mathbf{k}}(\mathbf{r}) = \sum_{\mathbf{G}} a_{i,\mathbf{k}+\mathbf{G}} \exp(i(\mathbf{k} + \mathbf{G}) \cdot \mathbf{r}) \quad (3.104)$$

The Kohn–Sham equations of the density functional theory then take on the following form:

$$\sum_{\mathbf{G}'} \left\{ \frac{\hbar^2}{2m} |\mathbf{k} + \mathbf{G}| \delta_{\mathbf{G}\mathbf{G}'} + V_{\text{ion}}(\mathbf{G} - \mathbf{G}') + V_{\text{elec}}(\mathbf{G} - \mathbf{G}') + V_{\text{XC}}(\mathbf{G} - \mathbf{G}') \right\} a_{i,\mathbf{k}+\mathbf{G}'} = \varepsilon_i a_{i,\mathbf{k}+\mathbf{G}'} \quad (3.105)$$

V_{ion} , V_{elec} and V_{XC} represent the electron–nuclei, electron–electron and exchange–correlation functionals, respectively. The delta function $\delta_{\mathbf{G}\mathbf{G}'}$ is zero unless $\mathbf{G} = \mathbf{G}'$, in which case it has a value of 1. There are two potential problems with the practical use of this equation for a ‘macroscopic’ lattice. First, the summation over \mathbf{G}' (a Fourier series) is in theory over an infinite number of reciprocal lattice vectors. In addition, for a macroscopic lattice there are effectively an infinite number of \mathbf{k} points within the first Brillouin zone. Fortunately, there are practical solutions to both of these problems.

We are usually interested in the valence electrons of an atom, as these are largely responsible for the chemical bonding and most physical properties. The core electrons are little affected by the atomic environment. It is therefore common only to consider explicitly the valence electrons in the calculation and to subsume the core electrons into the nuclear core. One potential drawback to the representation of valence electron wavefunctions with a plane-wave basis set is that near to the atomic nuclei the wavefunctions of the valence electrons show rapid oscillations. This is because their wavefunctions must be orthogonal to those of the core electrons. These oscillations give rise to a large kinetic energy, and a very large number of plane waves would be required to properly model this behaviour. This corresponds to taking many terms in the plane-wave expansion of the orbital, Equation (3.104). This problem is compounded by the fact that the solid systems of interest often contain elements much later in the periodic table than are usually encountered in molecular Hartree–Fock calculations. Heavy elements have many more core electrons and so an even more pronounced oscillatory behaviour. However, in this inner region the kinetic energy is largely cancelled by the high electrostatic potential energy of interaction with the nucleus. A popular way to deal with these problems is to replace the ‘true’ potential in these core regions with a much weaker one called a *pseudopotential*. This represents the way in which the valence electrons interact with the combined nucleus plus core electrons [Heine 1970]. A pseudopotential is a potential function that gives wavefunctions with the same shape as the true wavefunction outside the core region but with fewer nodes inside the core region, as illustrated in Figure 3.22. This has the effect of reducing the number of terms required for the plane wave expansion of the wavefunction, which in turn drastically reduces the scale of the computational problem.

Pseudopotentials are usually derived from all-electron atomic calculations. The valence electron pseudopotential is then required to reproduce the behaviour and properties of the valence electrons in the full calculation. For example, the energy levels with the pseudopotential should be the same as for the all-electron calculation. In addition, the pseudopotential will often depend upon the orbital angular momentum of the wavefunction (i.e. for s, p, d, etc. orbitals) and it will be desired that the total valence electron density within the core radius equals that in the all-electron situation. Such pseudopotentials are

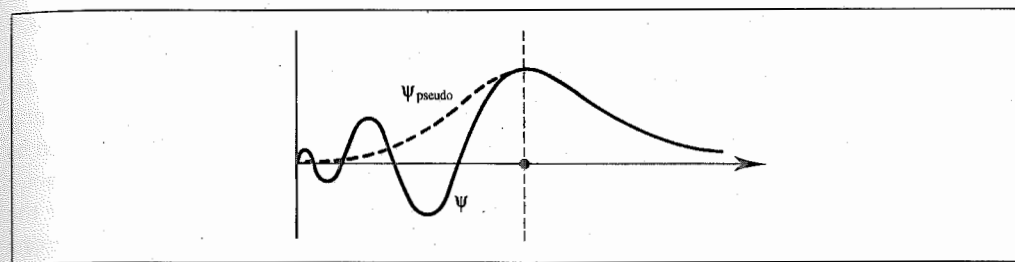


Fig. 3.22: Schematic representation of a pseudopotential. (Figure adapted from Payne M C, M P Teter, D C Allan, R A Arias and D J Joannopoulos 1992. *Iterative Minimisation Techniques for Ab initio Total-Energy Calculations: Molecular Dynamics and Conjugate Gradients*. *Reviews of Modern Physics* 64:1045–1097.)

referred to as ‘non-local norm-conserving’. An additional advantage of the use of pseudopotentials for the heavy elements is that they enable some relativistic effects to be included in the model. A number of functional forms are possible for the pseudopotentials; it is usual to assume a specific functional form and then to vary the parameters. The various pseudopotentials differ in the number of plane waves that are required for their representation and in the degree to which they can be transferred between different atomic environments. So-called ‘soft’ pseudopotentials require fewer plane waves and are therefore computationally more attractive, though there is to some extent a trade-off between softness and transferability. Subsequently developed were the ‘ultrasoft’ or ‘supersoft’ pseudopotentials, which require even fewer plane waves.

In practice, therefore, a pseudopotential is invariably employed and only plane waves with a kinetic energy ($= (\hbar^2/2m)|\mathbf{k} + \mathbf{G}|^2$) less than some cutoff are included in the calculation. The cutoff used depends on the nature of the system under investigation. For example, in the first-row elements the 2p valence orbitals approach closer to the nucleus than the comparable 3p orbitals in the second-row elements (the latter are repelled by the lower 2p states). Thus elements such as silicon or sulphur usually have softer pseudopotentials than their first-row equivalents carbon and oxygen. Everything else being equal, a higher cutoff is consequently required for the latter and hence more plane waves in the expansion (i.e. more reciprocal lattice vectors, \mathbf{G}). Note that in the plane wave expansion the basis functions are not associated with particular atoms but are defined over the whole cell (this also removes the problem of basis-set superposition errors as an additional benefit). The coefficients $a_{i,\mathbf{k}+\mathbf{G}}$ are obtained by following the usual density functional scheme: an initial guess is made of the electron density variation $\rho(\mathbf{r})$, the Kohn–Sham and overlap matrices are constructed, diagonalisation gives the eigenfunctions and eigenvectors (and thus the coefficients a) from which the Kohn–Sham orbitals can be constructed and hence the density for the next iteration.

The second important practical consideration when calculating the band structure of a material is that, in principle, the calculation needs to be performed for all \mathbf{k} vectors in the Brillouin zone. This would seem to suggest that for a macroscopic solid an infinite number of vectors \mathbf{k} would be needed to generate the band structure. However, in practice a discrete sampling over the Brillouin zone is used. This is possible because the wavefunctions at points

that are close together in k space will be almost identical and can be represented by a single representative point. Each of these discrete values is multiplied by a weight factor related to the volume of reciprocal space it represents. Obviously, the denser the set of k vectors the smaller will be the error in the calculation. Various schemes have been suggested for selecting suitable sets of k vectors which can give very accurate approximations to properties such as the charge density; the method of Monkhorst and Pack is particularly popular [Monkhorst and Pack 1976]. The selection of k vectors is also influenced by the size and shape of the system; indeed, if the unit cell is large then it may only be necessary to consider just one vector. Typically, between ten and 100 vectors are sufficient to understand the structural and electronic properties of a solid, though for certain types of problem such as calculating the optical properties of a metal many more k vectors may be required (several thousands). Ideally, one should ensure that the calculation converges both in terms of the number of wave-vectors k considered and in terms of the number of reciprocal lattice vectors G . An additional consideration is that the symmetry of the Brillouin zone itself may mean that it is not necessary for k to vary over the entire zone but that only a smaller section need be considered. For example, in our two-dimensional hexagonal close-packed case we would only have to consider the small right-angled triangle over which we undertook our 'tour'. This has an area one-twelfth that of the entire zone. This is an example of the use of the point symmetry of the Brillouin zone rather than the translational symmetry of the lattice. The small section containing the explicit k vectors required for the calculation is sometimes referred to as the irreducible part of the Brillouin zone.

3.8.7 Application of Solid-state Quantum Mechanics to the Group 14 Elements

The combination of density functional methods with pseudopotentials has been used extensively to study a wide variety of materials. Three systems that have been the subject of much interest are the group 14 elements carbon, silicon and germanium, reflecting their natural abundance, commercial importance (especially for silicon) and the large amount of experimental data available. Of particular interest is the problem of predicting the lowest-energy structure at a given volume [Cohen 1986; Mujica and Needs 1993; Needs and Mujica 1995]. In effect, this corresponds to predicting the most stable structure at a particular pressure. These elements all exist in the familiar diamond structure at normal pressures and temperatures but alternative structures can be formed by the application of pressure, at least for silicon and germanium. There has also been much speculation as to whether diamond itself could be transformed should a high enough pressure be applied. This last problem does have some practical interest as it would provide a theoretical upper limit to the pressures that could be achieved with ultra high-pressure diamond anvil cells.

There are many alternatives to the diamond structure, including body-centred cubic, face-centred cubic, hexagonal close-packed, simple hexagonal, simple cubic, β -tin, double-hexagonal close-packed and two complex tetrahedral structures: a body-centred cubic structure with eight atoms per unit cell and a simple tetragonal structure with twelve atoms per unit cell, not forgetting of course the many fullerene forms. Not all studies

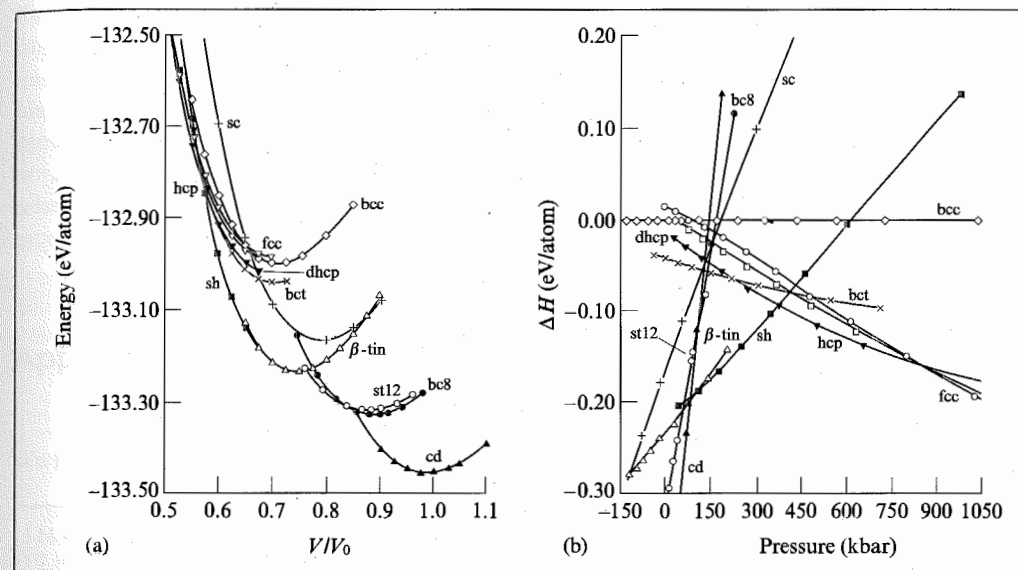


Fig. 3.23: (a) Graph of energy vs volume (scale normalised to the diamond structure) for eleven phases of silicon. (b) Enthalpy–pressure plot for the same eleven phases relative to the body-centred cubic phase. (Figures redrawn from Needs R J and A Mujica 1995. First-principles pseudopotential study of the structural phases of silicon. Physical Review B51:9652–9660.)

consider every one of these phases but by quoting the list in full we can appreciate the range of possibilities. The energy differences between many of these phases are often small and so it is particularly important to achieve an effective sampling of points in k space (recent studies suggest several thousands of such points are needed). The plane-wave cutoff can also have an effect on the results. The calculations involve minimising each structure at a number of different volumes and then fitting a polynomial to the data points. The results are usually displayed as a graph of the total energy versus the volume, as shown in Figure 3.23. Another way to display this type of data is an enthalpy–pressure plot, from which the most stable phase at any pressure is easily identified as that with the lowest enthalpy. Various bulk structural properties can also be calculated for comparison with experiment.

As we alluded, of the forms mentioned above only the diamond structure has been observed experimentally for carbon. For both silicon and germanium there is a transition to the β -tin phase around 100–130 kbar. Silicon further transforms into other structures such as the simple hexagonal with a relatively modest further increase in pressure, whereas for germanium this transition requires much more pressure. Why should this be, given that they are all in the same group? The electronic structure calculations provide some significant insights into this problem. Thus silicon has a strongly repulsive p-orbital pseudopotential due to the inner (2p) electrons, which carbon does not. This repulsion contributes to the formation of a single peak in the electron density along each Si–Si bond, whereas for carbon there are two peaks, each being near the position for the atomic p orbitals (Figure 3.24). The differences

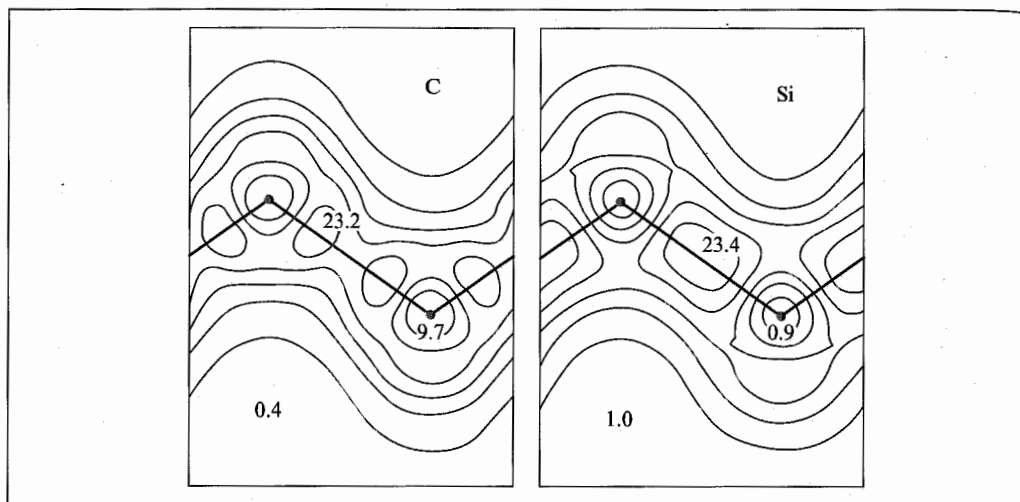


Fig. 3.24: Valence electron density for the diamond structures of carbon and silicon. (Figure redrawn from Cohen M L 1986. *Predicting New Solids and Superconductors*. Science 234:549-553.)

between silicon and germanium are ascribed to the d electron states; silicon does not have core d electrons, whereas germanium does. Certain transitions (e.g. carbon \rightarrow β -tin) do not depend upon the d character of the electronic configuration in contrast to subsequent transitions.

3.9 The Future Role of Quantum Mechanics: Theory and Experiment Working Together

Of all the methods that we will discuss in this book, quantum mechanics is probably the most widely used and the most extensively developed. The importance of the subject can be gauged in many ways, from citation counts to the number of Nobel prizes awarded. The systems studied using quantum mechanics range from the simplest molecular species (e.g. H_2^+ , HD^+ , H_3^+) to some very large and complex molecules (e.g. DNA, proteins and complex solid-state materials). Some of the most productive situations occur when experiment and theory are used in combination to tackle a problem. The methylene molecule, CH_2 , is of particular historical interest. Despite its small size, this molecule and the controversy surrounding it played an important role in establishing the role of computational quantum mechanical methods in modern-day research and the relationship between theory and experiment [Schaeffer 1986]. The early debate concentrated on the ground state of the molecule and whether its geometry was linear or bent. Early *ab initio* calculations by Foster and Boys [Foster and Boys 1960] suggested an H-C-H angle of 129° but this was refuted by spectroscopic data from Herzberg's laboratory, which were interpreted to indicate a linear geometry. Unfortunately for Foster and Boys, empirical calculations favoured by their head of department, Longuet-Higgins, also gave a linear geometry. Events came to a head when Bender and Schaeffer calculated a geometry of 135.1° and concluded that

the energy barrier between the linear and bent geometries was so large that no further improvement in the theoretical model could remove it. Soon thereafter several other experiments were undertaken, showing a bent structure. Moreover, when Herzberg re-examined his original data it was found to be consistent with a bent model. As we shall see in the remaining chapters there are many kinds of problem that can be tackled using computational chemistry methods. By no means do they always work, but there is often a synergistic relationship between experiment and theory, which means that the two combined can be much more productive than either in isolation.

Appendix 3.1 Alternative Expression for a Wavefunction Satisfying Bloch's Function

We have Equation (3.81):

$$\psi^k(x+a) = e^{ika} \psi^k(x) \quad (3.106)$$

We write $\psi(x)$ as the product of the exponential and a function $u_k(x)$:

$$u_k(x) = \psi_k(x) / \exp(ikx) \quad (3.107)$$

If we perform the same manipulation for $\psi(x+a)$ we get:

$$u_k(x+a) = \frac{\psi_k(x+a)}{e^{ik(x+a)}} = \frac{\psi_k(x) e^{ika}}{e^{ikx} e^{ika}} = \frac{\psi_k(x)}{e^{ikx}} = u_k(x) \quad (3.108)$$

Thus $u_k(x)$ is a periodic function which can be used to formulate acceptable wavefunctions:

$$\psi_k(x) = e^{ikx} u_k(x) \quad (3.109)$$

Further Reading

- Ashcroft N W and N D Mermin 1976. *Solid State Physics*. New York, Holt, Rinehart and Winston.
- Atkins P W 1991. *Quanta: A Handbook of Concepts*. Oxford, Oxford University Press.
- Atkins P W and R S Friedman 1996. *Molecular Quantum Mechanics*, 3rd edition. Oxford, Oxford University Press.
- Catlow C R A 1997. Computer Modelling as a Technique in Materials Chemistry. In Catlow C R A and A K Cheetham (Editors). *New Trends in Materials Chemistry*, NATO ASI Series C 498, Dordrecht, Kluwer.
- Catlow C R A 1998. Solids: Computer Modelling. In Schleyer, P v R, N L Allinger, T Clark, J Gasteiger, P A Kollman, H F Schaeffer III and P R Schreiner (Editors). *The Encyclopedia of Computational Chemistry*, Chichester, John Wiley & Sons.
- Gillan M J 1991. Calculating the Properties of Materials from Scratch. In Meyer M and V Pontikis (Editors). *Computer Simulation*, NATO ASI Series E 205 (Computer Simulations in Materials Science) pp. 257-281.
- Hehre W J, L Radom, P v R Schleyer and J A Pople 1986. *Ab initio Molecular Orbital Theory*. New York, John Wiley & Sons.

- Hoffmann R 1988. *Solids and Surfaces: A Chemist's View on Bonding in Extended Structures*. New York, VCH Publishers.
- Kohn W, A D Becke and R G Parr 1996. Density Functional Theory of Electronic Structure. *Journal of Chemical Physics* **100**:12974–12980.
- Kohn W and P Vashita 1983. General Density Functional Theory. In Lundquist S and N H March (Editors). *Theory of Inhomogeneous Electron Gas*, New York, Plenum, pp. 79–148.
- Kutzelnigg W and P von Herigonte 2000. Electron Correlation at the Dawn of the 21st Century. *Advances in Quantum Chemistry* **36**:185–229.
- Pisani C, R Dovesi and C Roetti 1988. Hartree-Fock *Ab initio* Treatment of Crystalline Systems. *Lecture Notes in Chemistry* Vol. 48. Berlin, Springer-Verlag.
- Pisani C, R Dovesi, C Roetti, M Cansa, R Orlando, S Casass and V R Saunders 2000. CRYSTAL and EMBED, Two Computational Tools for the *ab initio* Study of Electronic Properties of Crystals. *International Journal of Quantum Chemistry* **77**: 1032–1048.
- Schaeffer H F III (Editor) 1977. *Applications of Electronic Structure Theory*. New York, Plenum Press.
- Schaeffer H F III (Editor) 1977. *Methods of Electronic Structure Theory*. New York, Plenum Press.
- Szabo A and N S Ostlund 1982. *Modern Quantum Chemistry. Introduction to Advanced Electronic Structure Theory*. New York, McGraw-Hill.
- Wimmer E 1991 Density Functional Theory for Solids, Surface and Molecules: from Energy Bands to Molecular Bonds. In Labanowski J R and J W Andzelm (Editors). *Density Functional Methods in Chemistry*. Berlin, Springer-Verlag, pp. 7–31.

References

- Almlöf J, K Faegri Jr and K Korsell 1982. Principles for a Direct SCF Approach to LCAO-MO *Ab initio* Calculations. *Journal of Computational Chemistry* **3**:385–399.
- Ashcroft N W and N D Mermin 1976. *Solid State Physics*. New York, Holt, Rinehart & Winston.
- Baboul A G, L A Curtiss, P C Redfern and K Raghavachari 1999. Gaussian-3 Theory Using Density Functional Geometries and Zero-Point Energies. *Journal of Chemical Physics*, **110**:7650–7657.
- Becke A D 1988. Density-functional Exchange-energy Approximation with Correct Asymptotic Behaviour. *Physical Review* **A38**:3098–3100.
- Becke A D 1992. Density-functional Thermochemistry. I. The Effect of the Exchange-only Gradient Correction. *Journal of Chemical Physics* **96**:2155–2160.
- Becke A D 1993a. A New Mixing of Hartree-Fock and Local Density-functional Theories. *Journal of Chemical Physics* **98**:1372–1377.
- Becke A D 1993b. Density-functional Thermochemistry. III. The Role of Exact Exchange. *Journal of Chemical Physics* **98**:5648–5652.
- Becke A D and R M Dickson 1990. Numerical Solution of the Schroedinger Equation in Polyatomic Molecules. *Journal of Chemical Physics* **92**:3610–3612.
- Bobrowicz F W and W A Goddard III 1977. The Self-Consistent Field Equations for Generalized Valence Bond and Open-Shell Hartree-Fock Wave Functions. In Schaeffer H F III (Editor). *Modern Theoretical Chemistry III*, New York, Plenum, pp. 79–127.
- Boys S F and F Bernardi 1970. The Calculation of Small Molecular Interactions by the Differences of Separate Total Energies. Some Procedures with Reduced Errors. *Molecular Physics* **19**:553–566.
- Bradley C J and A P Cracknell 1972. *The Mathematical Theory of Symmetry in Solids*. Oxford, Clarendon Press.
- Ceperley D M and B J Alder 1980. Ground State of the Electron Gas by a Stochastic Method. *Physical Review Letters* **45**:566–569.

- Cohen M L 1986. Predicting New Solids and Superconductors. *Science* **234**:549–553.
- Cooper D L, J Gerratt and M Raimondi 1986. The Electronic Structure of the Benzene Molecule. *Nature* **323**: 699–701.
- Curtiss L A, K Raghavachari, G W Trucks and J A Pople 1991. Gaussian-2 Theory for Molecular Energies of First- and Second-row Compounds. *Journal of Chemical Physics* **94**:7221–7230.
- Curtiss L A, K Raghavachari, P C Redfern, V Rassolov and J A Pople 1998. Gaussian-3 (G3) Theory for Molecules Containing First and Second-row Atoms. *Journal of Chemical Physics* **109**:7764–7776.
- Curtiss L A, P C Redfern, K Raghavachari, V Rassolov and J A Pople 1999. Gaussian-3 Theory Using Reduced Møller-Plesset Order. *Journal of Chemical Physics* **110**:4703–4709.
- Dovesi R, R Orlando, C Roetti, C Pisani and V R Saunders 2000. The Periodic Hartree-Fock Method and Its Implementation in the CRYSTAL Code. *Physica Status Solidi* **B217**:63–88.
- Dovesi R, C Pisani, C Roetti and V R Saunders 1983. Treatment of Coulomb Interactions in Hartree-Fock Calculations of Periodic-Systems. *Physical Review* **B28**:5781–5792.
- Foster J M and S F Boys 1960. Quantum Variational Calculations for a Range of CH₂ Configurations. *Reviews in Modern Physics* **32**:305–307.
- Frisch M J, G W Trucks and J R Cheeseman 1996. Systematic Model Chemistries Based on Density Functional Theory: Comparison with Traditional Models and with Experiment. *Theoretical and Computational Chemistry (Recent Developments and Applications of Modern Density Functional Theory)* **4**:679–707.
- Gerratt J, D L Cooper, P B Karadakov and M Raimondi 1997. Modern Valence Bond Theory. *Chemical Society Reviews* pp. 87–100.
- Gunnarsson O and B I Lundqvist 1976. Exchange and Correlation in Atoms, Molecules, and Solids by the Spin-density-functional Formalism. *Physical Review* **B13**:4274–4298.
- Heine V 1970. The Pseudopotential Concept. *Solid State Physics* **24**:1–36.
- Heitler W and F London 1927. Wechselwirkung neutraler Atome und Homöopolare Bindung nach der Quantenmechanik. *Zeitschrift für Physik* **44**:455–472.
- Hohenberg P and Kohn W 1964. Inhomogeneous Electron Gas. *Physical Review* **B136**:864–871.
- Johnson B G, P M W Gill and J A Pople 1993. The performance of a family of density functional methods. *Journal of Chemical Physics* **98**:5612–5626.
- Kohn W and L J Sham 1965. Self-consistent Equations Including Exchange and Correlation Effects. *Physical Review* **A140**:1133–1138.
- Lee C, W Yang and R G Parr 1988. Development of the Colle-Salvetti Correlation Energy Formula into a Functional of the Electron Density. *Physical Review* **B37**:785–789.
- Møller C and M S Plesset 1934. Note on an Approximate Treatment for Many-Electron Systems. *Physical Review* **46**:618–622.
- Monkhorst H J and J D Pack 1976. Special Points for Brillouin-zone Integration. *Physical Review* **B13**:5188–5192.
- Morokuma K 1977. Why Do Molecules Interact? The Origin of Electron Donor-Acceptor Complexes, Hydrogen Bonding, and Proton Affinity. *Accounts of Chemical Research* **10**:294–300.
- Mujica A and R J Needs 1993. First-principles Calculations of the Structural Properties, Stability, and Band Structure of Complex Tetrahedral Phases of Germanium: ST12 and BC8. *Physical Review* **B48**:17010–17017.
- Needs R J and Mujica 1995. First-principles Pseudopotential Study of the Structural Phases of Silicon. *Physical Review* **B51**:9652–9660.
- Parr R G 1983. Density Functional Theory. *Annual Review of Physical Chemistry* **34**:631–656.
- Perdew J P and A Zunger 1981. Self-Interaction Correction to Density-Functional Approximations for Many-Electron Systems. *Physical Review* **B23**:5048–5079.
- Pisani C and R Dovesi 1980. Exact-Exchange Hartree-Fock Calculations for Periodic Systems. I. Illustration of the Method. *International Journal of Quantum Chemistry* **XVII**:501–516.

- Pople J A, M Head-Gordon and K Raghavachari 1987. Quadratic Configuration Interaction. A General Technique for Determining Electron Correlation Energies. *Journal of Chemical Physics* **87**:5968–5975.
- Pople J A, M Head-Gordon, D J Fox, K Raghavachari and L A Curtiss 1989. Gaussian-1 Theory: A General Procedure for Prediction of Molecular Energies. *Journal of Chemical Physics* **90**:5622–5629.
- Pople J A and R K Nesbet 1954. Self-consistent Orbitals for Radicals. *Journal of Chemical Physics* **22**:571–572.
- Pulay P 1977. Direct Use of the Gradient for Investigating Molecular Energy Surfaces. In Schaeffer H F III (Editor). *Applications of Electronic Structure Theory*, New York, Plenum, pp. 153–185.
- Pulay P 1980. Convergence Acceleration of Iterative Sequences. The Case of SCF Iteration. *Chemical Physics Letters* **73**:393–398.
- Pulay P 1987. Analytical Derivative Methods in quantum Chemistry. In Lawley K P (Editor). *Ab initio Methods in Quantum Chemistry - II*, New York, John Wiley & Sons, pp. 241–286.
- Roos B O, P R Taylor and E M Siegbahn 1980. A Complete Active Space SCF Method (CASSCF) Using a Density Matrix Formulated Super-CI Approach. *Chemical Physics* **48**:157–173.
- Schaeffer H F III 1986. Methylene: A Paradigm for Computational Quantum Chemistry. *Science* **231**:1100–1107.
- Sim F, St-Amant A, I Papai and D R Salahub 1992. Gaussian Density Functional Calculations on Hydrogen-Bonded Systems. *Journal of the American Chemical Society* **114**:4391–4400.
- Slater J C 1974. *Quantum Theory of Molecules and Solids Volume 4: The Self-Consistent Field for Molecules and Solids*. New York, McGraw-Hill.
- Smith B J, D J Swanton, J A Pople, H F Schaeffer III and L Radom 1990. Transition Structures for the Interchange of Hydrogen Atoms within the Water Dimer. *Journal of Chemical Physics* **92**:1240–1247.
- St-Amant A, W D Cornell, P A Kollman and T A Halgren 1995. Calculation of Molecular Geometries, Relative Conformational Energies, Dipole Moments and Molecular Electrostatic Potential Fitted Charges of Small Organic Molecules of Biochemical Interest by Density Functional Theory. *Journal of Computational Chemistry* **16**:1483–1506.
- Stephens P J, F J Devlin, C F Chabalowski and M J Frisch 1994. *Ab Initio* Calculation of Vibrational Absorption and Circular Dichroism Spectra Using Density Functional Force Fields. *Journal of Physical Chemistry* **98**:11623–11627.
- Umeyama H and K Morokuma 1977. The Origin of Hydrogen Bonding. An Energy Decomposition Study. *Journal of the American Chemical Society* **99**:1316–1332.
- Vosko S H, L Wilk and M Nusair 1980. Accurate Spin-dependent Electron Liquid Correlation Energies for Local Spin Density Calculations: A Critical Analysis. *Canadian Journal of Physics* **58**:1200–1211.
- Wimmer E 1997. Electronic Structure Methods. In Catlow C R A and A K Cheetham (Editors). *New Trends in Materials Chemistry*, NATO ASI Series C 498. Dordrecht, Kluwer.

CHAPTER FOUR

Empirical Force Field Models: Molecular Mechanics

4.1 Introduction

Many of the problems that we would like to tackle in molecular modelling are unfortunately too large to be considered by quantum mechanics. Quantum mechanical methods deal with the electrons in a system, so that even if some of the electrons are ignored (as in the semi-empirical schemes) a large number of particles must still be considered, and the calculations are time-consuming. Force field methods (also known as molecular mechanics) ignore the electronic motions and calculate the energy of a system as a function of the nuclear positions only. Molecular mechanics is thus invariably used to perform calculations on systems containing significant numbers of atoms. In some cases force fields can provide answers that are as accurate as even the highest-level quantum mechanical calculations, in a fraction of the computer time. However, molecular mechanics cannot of course provide properties that depend upon the electronic distribution in a molecule.

That molecular mechanics works at all is due to the validity of several assumptions. The first of these is the Born–Oppenheimer approximation, without which it would be impossible to contemplate writing the energy as a function of the nuclear coordinates at all. Molecular mechanics is based upon a rather simple model of the interactions within a system with contributions from processes such as the stretching of bonds, the opening and closing of angles and the rotations about single bonds. Even when simple functions (e.g. Hooke's law) are used to describe these contributions the force field can perform quite acceptably. Transferability is a key attribute of a force field, for it enables a set of parameters developed and tested on a relatively small number of cases to be applied to a much wider range of problems. Moreover, parameters developed from data on small molecules can be used to study much larger molecules such as polymers.

4.1.1 A Simple Molecular Mechanics Force Field

Many of the molecular modelling force fields in use today for molecular systems can be interpreted in terms of a relatively simple four-component picture of the intra- and inter-molecular forces within the system. Energetic penalties are associated with the deviation of bonds and angles away from their 'reference' or 'equilibrium' values, there is a function

that describes how the energy changes as bonds are rotated, and finally the force field contains terms that describe the interaction between non-bonded parts of the system. More sophisticated force fields may have additional terms, but they invariably contain these four components. An attractive feature of this representation is that the various terms can be ascribed to changes in specific internal coordinates such as bond lengths, angles, the rotation of bonds or movements of atoms relative to each other. This makes it easier to understand how changes in the force field parameters affect its performance, and also helps in the parametrisation process. One functional form for such a force field that can be used to model single molecules or assemblies of atoms and/or molecules is:

$$\begin{aligned} \mathcal{V}(\mathbf{r}^N) = & \sum_{\text{bonds}} \frac{k_i}{2} (l_i - l_{i,0})^2 + \sum_{\text{angles}} \frac{k_i}{2} (\theta_i - \theta_{i,0})^2 + \sum_{\text{torsions}} \frac{V_n}{2} (1 + \cos(n\omega - \gamma)) \\ & + \sum_{i=1}^N \sum_{j=i+1}^N \left(4\epsilon_{ij} \left[\left(\frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] + \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} \right) \end{aligned} \quad (4.1)$$

$\mathcal{V}(\mathbf{r}^N)$ denotes the potential energy, which is a function of the positions (\mathbf{r}) of N particles (usually atoms). The various contributions are schematically represented in Figure 4.1. The first term in Equation (4.1) models the interaction between pairs of bonded atoms, modelled here by a harmonic potential that gives the increase in energy as the bond length l_i deviates from the reference value $l_{i,0}$. The second term is a summation over all valence angles in the molecule, again modelled using a harmonic potential (a valence angle is the angle formed between three atoms A–B–C in which A and C are both bonded to B). The third term in Equation (4.1) is a torsional potential that models how the energy changes as a bond rotates. The fourth contribution is the non-bonded term. This is calculated between all pairs of atoms (i and j) that are in different molecules or that are in

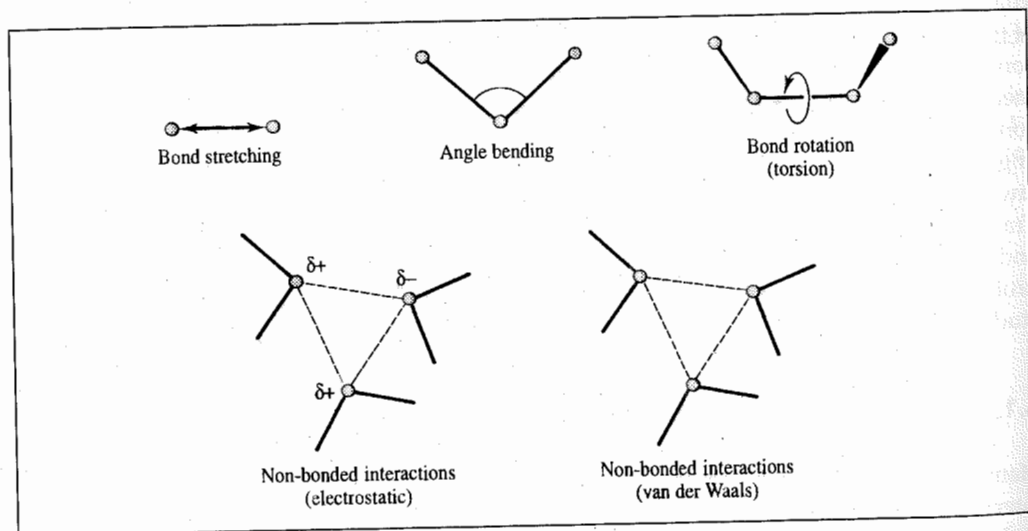


Fig. 4.1: Schematic representation of the four key contributions to a molecular mechanics force field: bond stretching, angle bending and torsional terms and non-bonded interactions.

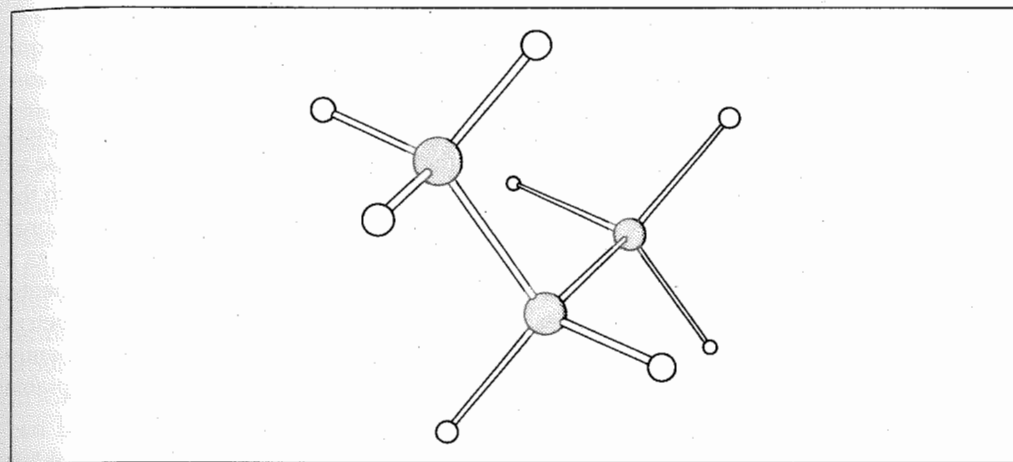


Fig. 4.2: A typical force field model for propane contains ten bond-stretching terms, eighteen angle-bending terms, eighteen torsional terms and 27 non-bonded interactions.

the same molecule but separated by at least three bonds (i.e. have a $1, n$ relationship where $n \geq 4$). In a simple force field the non-bonded term is usually modelled using a Coulomb potential term for electrostatic interactions and a Lennard-Jones potential for van der Waals interactions.

We shall discuss the nature of these different contributions in more detail in Sections 4.3–4.10, but here we consider how the simple force field of Equation (4.1) would be used to calculate the energy of a conformation of propane (Figure 4.2). Propane has ten bonds: two C–C bonds and eight C–H bonds. The C–C bonds are symmetrically equivalent but the C–H bonds fall into two classes, one group corresponding to the two hydrogens bonded to the central methylene (CH_2) carbon and one group corresponding to the six hydrogens bonded to the methyl carbons. In some sophisticated force fields different parameters would be used for these two different types of C–H bond, but in most force fields the same bonding parameters (i.e. k_i and $l_{i,0}$) would be used for each of the eight C–H bonds. This is an example of the way in which the same parameters can be used for a wide variety of molecules. There are 18 different valence angles in propane, comprising one C–C–C angle, ten C–C–H angles and seven H–C–H angles. Note that all angles are included in the force field model even though some of them may not be independent of the others. There are 18 torsional terms: twelve H–C–C–H torsions and six H–C–C–C torsions. Each of these is modelled with a cosine series expansion that has minima at the *trans* and *gauche* conformations. Finally, there are 27 non-bonded terms to calculate, comprising 21 H–H interactions and six H–C interactions. The electrostatic contribution would be calculated using Coulomb's law from partial atomic charges associated with each atom and the van der Waals contribution as a Lennard-Jones potential with appropriate ϵ_{ij} and σ_{ij} parameters. A sizeable number of terms are thus included in the force field model, even for a molecule as simple as propane. Even so, the number of terms (73) is many fewer than the number of integrals that would be involved in an equivalent *ab initio* quantum mechanical calculation.

an atom type that is different from the carbon atom in an isolated five-membered ring such as histidine, which in turn is different from the atom type of a carbon atom in a benzene ring. Indeed, the AMBER force field uses different atom types for a histidine amino acid depending upon its protonation state (Figure 4.3). Other, more general, force fields would assign these atoms to the same generic 'sp² carbon' atom type. It is often found that force fields which are designed for modelling specific classes of molecule (such as proteins and nucleic acids, in the case of AMBER) use more specific atom types than force fields designed for general-purpose use.

We now discuss in some detail the individual contributions to a molecular mechanics force field, giving a selection of the various functional forms that are in common use. We shall then consider the important task of parametrisation, in which values for the many force constants are derived. Our discussion will be illuminated by examples chosen from contemporary force fields in widespread use and the MM2/MM3/MM4 and AMBER force fields in particular.

4.3 Bond Stretching

The potential energy curve for a typical bond has the form shown in Figure 4.4. Of the many functional forms used to model this curve, that suggested by Morse is particularly useful. The Morse potential has the form:

$$v(l) = D_e \{1 - \exp[-a(l - l_0)]\}^2 \quad (4.2)$$

D_e is the depth of the potential energy minimum and $a = \omega \sqrt{\mu/2D_e}$, where μ is the reduced mass and ω is the frequency of the bond vibration. ω is related to the stretching constant of the bond, k , by $\omega = \sqrt{k/\mu}$. l_0 is the reference value of the bond. The Morse potential is not usually used in molecular mechanics force fields. In part this is because it is not particularly amenable to efficient computation but also because it requires three parameters to be specified for each bond. Moreover, it is rare in molecular mechanics calculations for

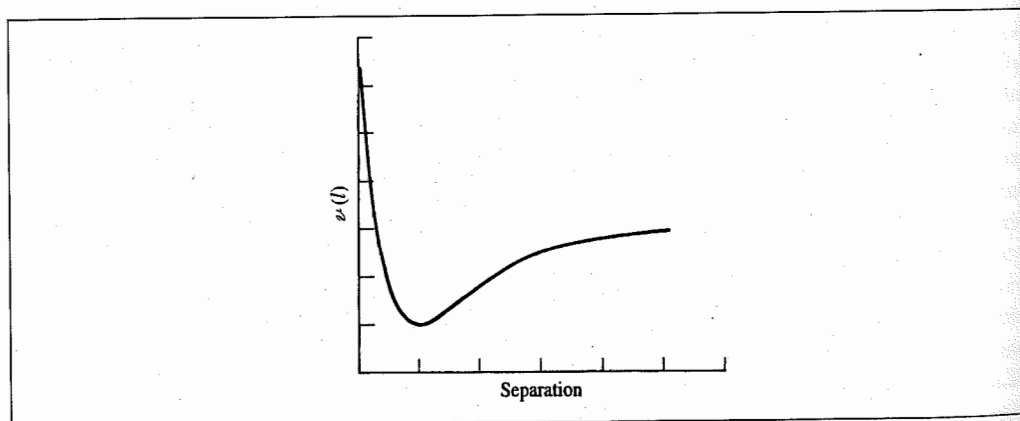


Fig. 4.4: Variation in bond energy with interatomic separation.

bonds to deviate significantly from their equilibrium values; the Morse curve describes a wide range of behaviour from the strong equilibrium behaviour to dissociation. Consequently, simpler expressions are often used. The most elementary approach is to use a Hooke's law formula in which the energy varies with the square of the displacement from the reference bond length l_0 :

$$v(l) = \frac{k}{2}(l - l_0)^2 \quad (4.3)$$

The astute reader will have noticed our use of the term 'reference bond length' (sometimes called the 'natural bond length') for the parameter l_0 . This parameter is commonly called the 'equilibrium' bond length, but to do so can be misleading. The reference bond length is the value that the bond adopts when all other terms in the force field are set to zero. The equilibrium bond length, by contrast, is the value that is adopted in a minimum energy structure, when all other terms in the force field contribute. The complex interplay between the various components in the force field means that the bond may well deviate slightly from its reference value in order to compensate for other contributions to the energy. It is also important to recognise that 'real' molecules undergo vibrational motion (even at absolute zero, there is a zero-point energy due to vibrational motion). A true bond-stretching potential is not harmonic but has a shape similar to that in Figure 4.4, which means that the 'average' length of the bond in a vibrating molecule will deviate from the equilibrium value for the hypothetical motionless state. The effects are usually small, but they are significant if one wishes to predict bond lengths to thousandths of an ångström. When comparing the results of calculations with experimental data, one must also remember that different experimental techniques measure different 'equilibrium' values, especially when the experiments are performed at different temperatures. The errors in experimentally determined bond lengths can be quite large; for example, libration of a molecule in a crystal means that the bond lengths determined by X-ray methods at room temperature may have errors as large as 0.015 Å. MM2 was parametrised to fit the values obtained by electron diffraction, which give the mean distances between atoms averaged over the vibrational motion at room temperature.

The forces between bonded atoms are very strong and considerable energy is required to cause a bond to deviate significantly from its equilibrium value. This is reflected in the magnitude of the force constants for bond stretching; some typical values from the MM2 force field are shown in Table 4.1, where it can be seen that those bonds one would

Bond	l_0 (Å)	k (kcal mol ⁻¹ Å ⁻²)
Csp ³ -Csp ³	1.523	317
Csp ³ -Csp ²	1.497	317
Csp ² -Csp ²	1.337	690
Csp ² =O	1.208	777
Csp ³ -Nsp ³	1.438	367
C-N (amide)	1.345	719

Table 4.1 Force constants and reference bond lengths for selected bonds [Allinger 1977].

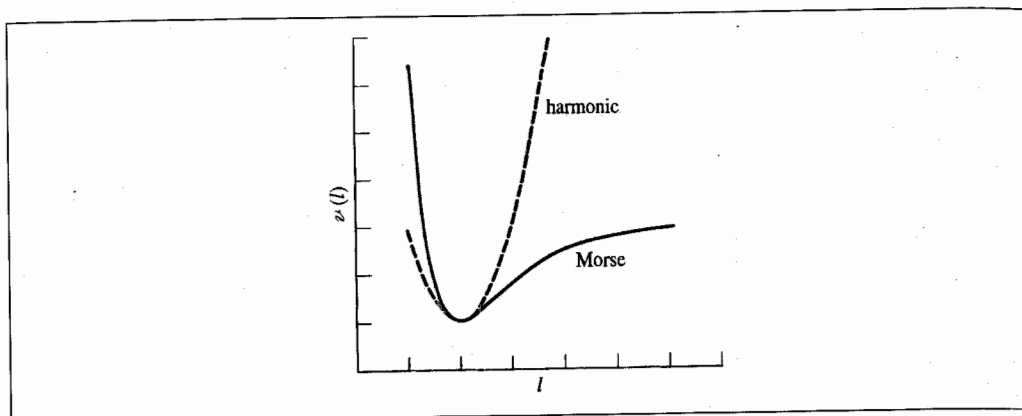


Fig. 4.5: Comparison of the simple harmonic potential (Hooke's law) with the Morse curve.

intuitively expect to be stronger have large force constants (contrast C–C with C=C and N≡N). A deviation of just 0.2 Å from the reference value l_0 with a force constant of $300 \text{ kcal mol}^{-1} \text{ \AA}^{-2}$ would cause the energy of the system to rise by 12 kcal/mol.

The Hooke's law functional form is a reasonable approximation to the shape of the potential energy curve at the bottom of the potential well, at distances that correspond to bonding in ground-state molecules. It is less accurate away from equilibrium (Figure 4.5). To model the Morse curve more accurately, cubic and higher terms can be included and the bond-stretching potential can be written as follows:

$$v(l) = \frac{k}{2}(l - l_0)^2 [1 - k'(l - l_0) - k''(l - l_0)^2 - k'''(l - l_0)^3 \dots] \quad (4.4)$$

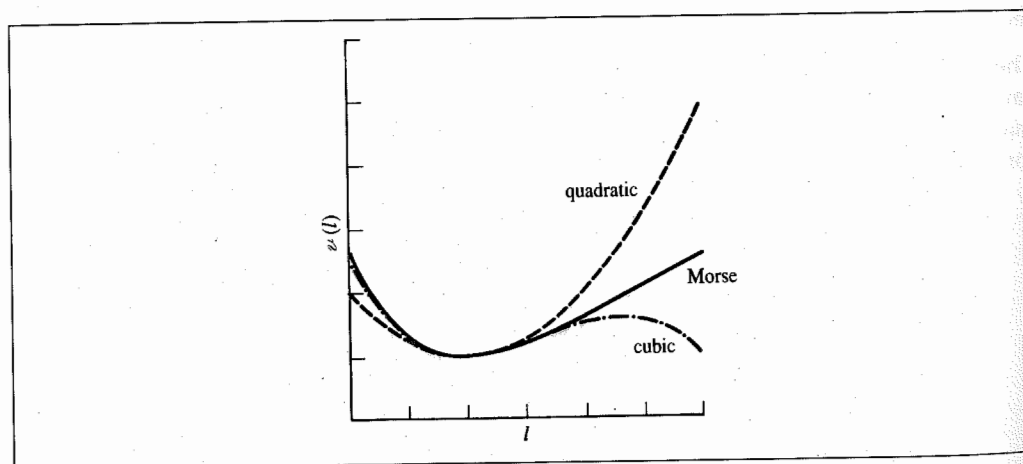


Fig. 4.6: A cubic bond-stretching potential passes through a maximum but gives a better approximation to the Morse curve close to the equilibrium structure than the quadratic form.

An undesirable side-effect of an expansion that includes just a quadratic and a cubic term (as is employed in MM2) is that, far from the reference value, the cubic function passes through a maximum. This can lead to a catastrophic lengthening of bonds (Figure 4.6). One way to accommodate this problem is to use the cubic contribution only when the structure is sufficiently close to its equilibrium geometry and is well inside the 'true' potential well. MM3 also includes a quartic term; this eliminates the inversion problem and leads to an even better description of the Morse curve.

4.4 Angle Bending

The deviation of angles from their reference values is also frequently described using a Hooke's law or harmonic potential:

$$v(\theta) = \frac{k}{2}(\theta - \theta_0)^2 \quad (4.5)$$

The contribution of each angle is characterised by a force constant and a reference value. Rather less energy is required to distort an angle away from equilibrium than to stretch or compress a bond, and the force constants are proportionately smaller, as can be observed in Table 4.2.

Angle	θ_0	k (kcal mol ⁻¹ deg ⁻¹)
Csp ³ –Csp ³ –Csp ³	109.47	0.0099
Csp ³ –Csp ³ –H	109.47	0.0079
H–Csp ³ –H	109.47	0.0070
Csp ³ –Csp ² –Csp ³	117.2	0.0099
Csp ³ –Csp ² =Csp ²	121.4	0.0121
Csp ³ –Csp ² =O	122.5	0.0101

Table 4.2 Force constants and reference angles for selected angles [Allinger 1977].

As with the bond-stretching terms, the accuracy of the force field can be improved by the incorporation of higher-order terms. MM2 contains a quartic term in addition to the quadratic term. Higher-order terms have also been included to treat certain pathological cases such as very highly strained molecules. The general form of the angle-bending term then becomes:

$$v(\theta) = \frac{k}{2}(\theta - \theta_0)^2 [1 - k'(\theta - \theta_0) - k''(\theta - \theta_0)^2 - k'''(\theta - \theta_0)^3 \dots] \quad (4.6)$$

4.5 Torsional Terms

The bond-stretching and angle-bending terms are often regarded as 'hard' degrees of freedom, in that quite substantial energies are required to cause significant deformations from

their reference values. Most of the variation in structure and relative energies is due to the complex interplay between the torsional and non-bonded contributions.

The existence of barriers to rotation about chemical bonds is fundamental to understanding the structural properties of molecules and conformational analysis. The three minimum-energy staggered conformations and three maximum-energy eclipsed structures of ethane are a classic example of the way in which the energy changes with a bond rotation. Quantum mechanical calculations suggest that this barrier to rotation can be considered to arise from antibonding interactions between the hydrogen atoms on opposite ends of the molecule; the antibonding interactions are minimised when the conformation is staggered and are at a maximum when the conformation is eclipsed. Many force fields are used for modelling flexible molecules where the major changes in conformation are due to rotations about bonds; in order to simulate this it is essential that the force field properly represents the energy profiles of such changes.

Not all molecular mechanics force fields use torsional potentials; it may be possible to rely upon non-bonded interactions between the atoms at the end of each torsion angle (the 1,4 atoms) to achieve the desired energy profile. However, most force fields for 'organic' molecules do use explicit torsional potentials with a contribution from each bonded quartet of atoms A-B-C-D in the system. Thus there would be nine individual torsional terms for ethane and 24 for benzene ($6 \times \text{C-C-C-C}$, $12 \times \text{C-C-C-H}$ and $6 \times \text{H-C-C-H}$). Torsional potentials are almost always expressed as a cosine series expansion. One functional form is:

$$v(\omega) = \sum_{n=0}^N \frac{V_n}{2} [1 + \cos(n\omega - \gamma)] \quad (4.7)$$

ω is the torsion angle.

An alternative but equivalent expression is:

$$v(\omega) = \sum_{n=0}^N C_n \cos(\omega)^n \quad (4.8)$$

V_n in Equation (4.7) is often referred to as the 'barrier' height, but to do so is misleading, obviously so when more than one term is present in the expansion. Moreover, other terms in the force field equation contribute to the barrier height as a bond is rotated, especially the non-bonded interactions between the 1,4 atoms. The value of V_n does, however, give a qualitative indication of the relative barriers to rotation; for example, V_n for an amide bond will be larger than for a bond between two sp^3 carbon atoms. n in Equation (4.7) is the *multiplicity*; its value gives the number of minimum points in the function as the bond is rotated through 360° . γ (the phase factor) determines where the torsion angle passes through its minimum value. For example, the energy profile for rotation about the single bond between two sp^3 carbon atoms could be represented by a single torsional term with $n = 3$ and $\gamma = 0^\circ$. This would give a threefold rotational profile with minima at torsion angles of $+60^\circ$, -60° and 180° and maxima at $\pm 120^\circ$ and 0° . A double bond between two sp^2 carbon atoms would have $n = 2$ and $\gamma = 180^\circ$, giving minima at 0° and 180° . The value of V_n would also be significantly larger for the double bond than for the single

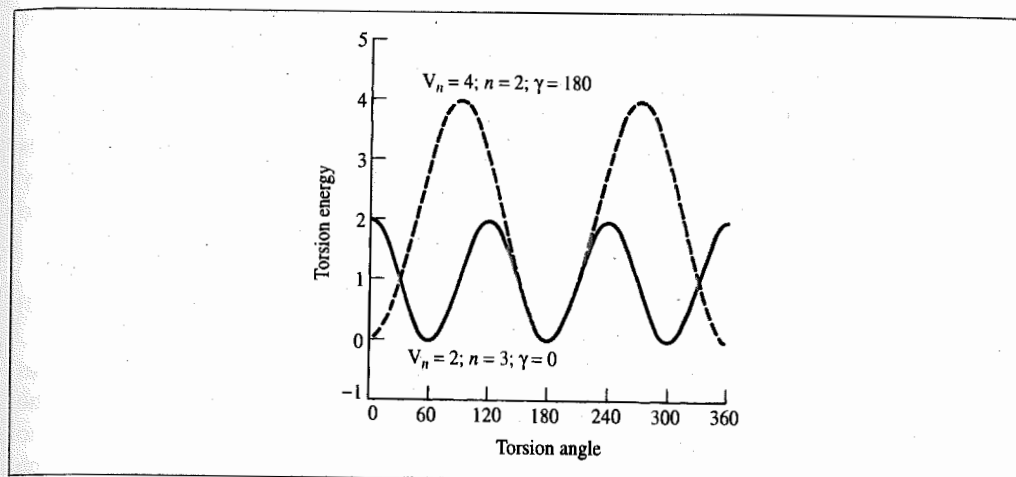


Fig. 4.7: Torsional potential varies as shown for different values of V_n , n and γ .

bond. The effects of varying V_n , n and γ are illustrated in Figure 4.7 for commonly occurring torsional potentials.

Many of the torsional terms in the AMBER force field contain just one term from the cosine series expansion, but for some bonds it was found necessary to include more than one term. For example, to correctly model the tendency of O-C-C-O bonds to adopt a *gauche* conformation, a torsional potential with two terms was used for the O-C-C-O contribution:

$$v(\omega_{\text{C-O-O-C}}) = 0.25(1 + \cos 3\omega) + 0.25(1 + \cos 2\omega) \quad (4.9)$$

The torsional energy for a $\text{OCH}_2\text{-CH}_2\text{O}$ fragment (found in the sugars in DNA) varies with the torsion angle ω as shown in Figure 4.8. Another feature of the AMBER force field is its use of general torsional parameters. The energy profile for rotation about a bond that is described by a general torsional potential depends solely upon the atom types of the two

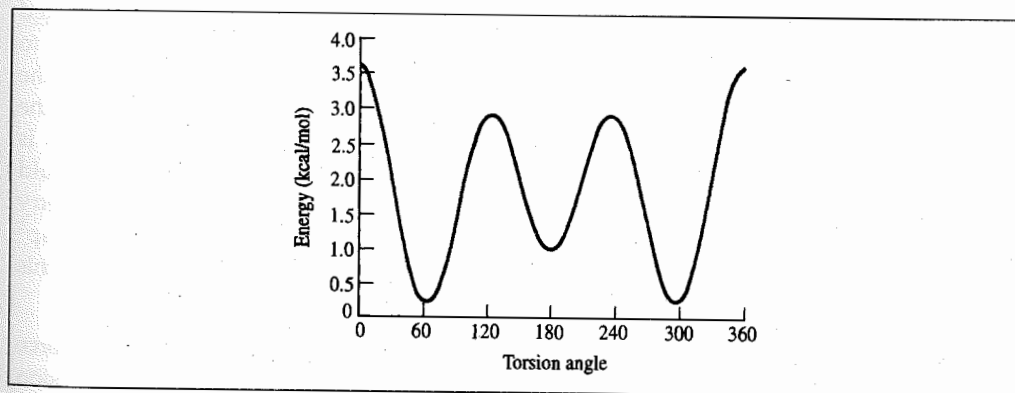


Fig. 4.8: Variation in torsional energy (AMBER force field) with O-C-C-O torsion angle (ω) for $\text{OCH}_2\text{-CH}_2\text{O}$ fragment. The minimum energy conformations arise for $\omega = 60^\circ$ and 300° .

atoms that comprise the central bond and not upon the atom types of the terminal atoms. For example, all torsion angles in which the central bond is between two sp^3 -hybridised carbon atoms (e.g. H-C-C-H, C-C-C-C, H-C-C-C) are assigned the same torsional parameters, unless the torsion is a special case such as O-C-C-O. In its treatment of the torsional contribution, AMBER takes a position intermediate between those force fields which only ever use a single term in the torsional expansion and those which consistently use more terms for all torsions. MM2 falls into the latter category; it uses three terms in the expansion:

$$v(\omega) = \frac{V_1}{2} (1 + \cos \omega) + \frac{V_2}{2} (1 - \cos 2\omega) + \frac{V_3}{2} (1 + \cos 3\omega) \quad (4.10)$$

A physical interpretation has been ascribed to each of the three terms in the MM2 torsional expansion from an analysis of *ab initio* calculations on simple fluorinated hydrocarbons. The first, onefold term corresponds to interactions between bond dipoles, which are due to differences in electronegativity between bonded atoms. The twofold term is due to the effects of hyperconjugation (in alkanes) and conjugation effects (in alkenes), which provide 'double bond' character to the bond. The threefold term corresponds to steric interactions between the 1,4 atoms. It was found that the additional terms in the torsional potential were especially important for systems containing heteroatoms, such as the halogenated hydrocarbons and molecules containing CCOC and CCNC fragments.

With careful parametrisation a force field which uses more than one term in the torsional expansion will be more successful than a force field that uses only a single term (and this is borne out by the MM2 force field). The major drawback is that many parameters are required to model even a modest range of molecules.

4.6 Improper Torsions and Out-of-plane Bending Motions

Let us consider how cyclobutanone would be modelled using a force field containing just standard bond-stretching and angle-bending terms of the type in Equation (4.1). The equilibrium structure obtained with such a force field would have the oxygen atom located out of the plane formed by the adjoining carbon atom and the two carbon atoms bonded to it, as shown in Figure 4.9. In this structure, the angles to the oxygen adopt values close to the reference value of 120° . Experimentally, it is found that the oxygen atom remains in the plane of the cyclobutane ring, even though the C-C=O angles are large (133°). This is because the π -bonding energy, which is maximised in the coplanar arrangement, would be much reduced if the oxygen were bent out of the plane. To achieve the desired geometry it is necessary to incorporate an additional term (or terms) in the force field that keeps the sp^2 carbon and the three atoms bonded to it in the same plane. The simplest way to achieve this is to use an *out-of-plane* bending term.

There are several ways in which out-of-plane bending terms can be incorporated into a force field. One approach is to treat the four atoms as an 'improper' torsion angle (i.e. a torsion angle in which the four atoms are not bonded in the sequence 1-2-3-4). One way to define an improper torsion for cyclobutane would involve the atoms 1-5-3-2 in Figure 4.9.

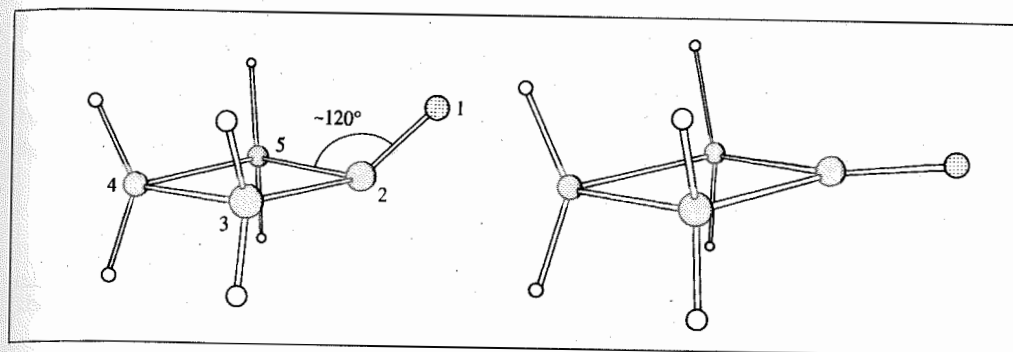


Fig. 4.9: Without an out-of-plane term, the oxygen atom in cyclobutane is predicted to lie out of the plane of the ring (left) rather than in the plane.

A torsional potential of the following form is then used to maintain the improper torsion angle at 0° or 180° :

$$v(\omega) = k(1 - \cos 2\omega) \quad (4.11)$$

Various other ways to incorporate the out-of-plane bending contribution are possible. For example, one definition that is closer to the notion of an 'out-of-plane bend' involves a calculation of the angle between a bond from the central atom and the plane defined by the other three atoms (Figure 4.10). A value of 0° corresponds to all four atoms being coplanar. A third approach is to calculate the height of the central atom above a plane defined by the other three atoms (Figure 4.10). With these two definitions the deviation of the out-of-plane coordinate (be it an angle or a distance) can be modelled using a harmonic potential of the form

$$v(\theta) = \frac{k}{2}\theta^2; \quad v(h) = \frac{k}{2}h^2 \quad (4.12)$$

Of these three functional forms, the improper torsion definition is most widely used as it can then be easily included with the 'proper' torsional terms in the force field. However, the

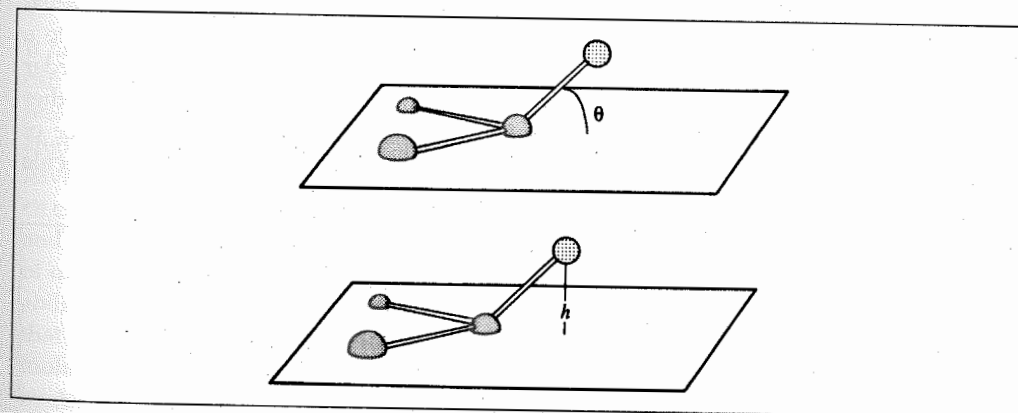


Fig. 4.10: Two ways to model the out-of-plane bending contributions.

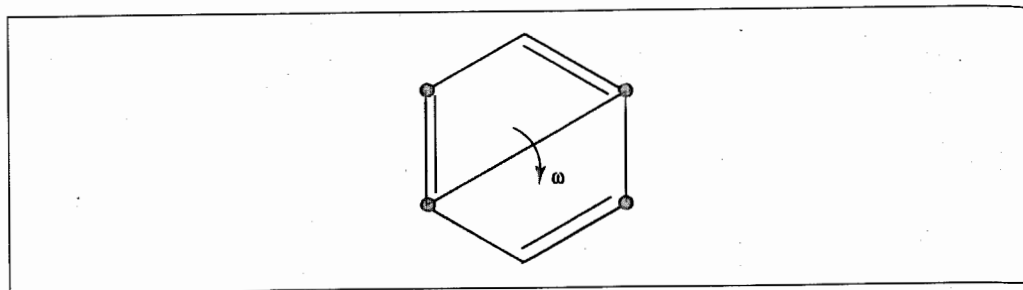


Fig. 4.11: Improper torsional terms can be used to keep a benzene ring planar.

other two functional forms may be better ways to implement out-of-plane bending in the force field. Out-of-plane terms may also be used to achieve a particular geometry. For example, if it is desired to ensure that an aromatic ring such as benzene maintains an approximately planar structure then this can be achieved using a suitable set of out-of-plane bending terms involving atoms on opposite sides of the ring (Figure 4.11). Improper torsional terms are commonly used in the so-called united atom force fields to maintain stereochemistry at chiral centres (see Section 4.14). It is important to remember that out-of-plane terms may not always be necessary, and that to include such terms may have a deleterious effect on the performance of the force field. Vibrational frequencies in particular are often rather sensitive to the presence of out-of-plane terms.

4.7 Cross Terms: Class 1, 2 and 3 Force Fields

The presence of *cross terms* in a force field reflects coupling between the internal coordinates. For example, as a bond angle is decreased it is found that the adjacent bonds stretch to reduce the interaction between the 1,3 atoms, as illustrated in Figure 4.12. Cross terms were found to be important in force fields designed to predict vibrational spectra that were the forerunners of molecular mechanics force fields, and so it is not surprising that cross terms must often be included in a molecular mechanics force field to achieve optimal performance. One should in principle include cross terms between all contributions to a force field. However, only a few cross terms are generally found to be necessary in order to reproduce structural properties accurately; more may be needed to reproduce other properties such as vibrational frequencies, which are more sensitive to the presence of such terms. In general, any interactions involving motions that are far apart in a molecule

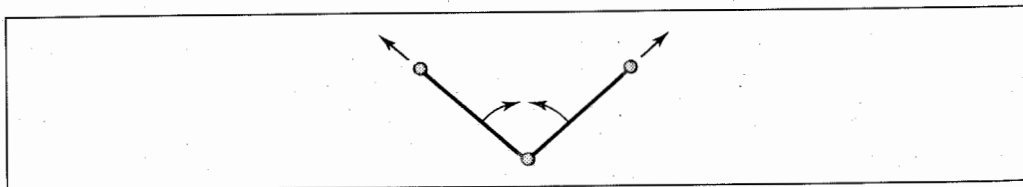


Fig. 4.12: Coupling between the stretching of the bonds as an angle closes.

can usually be set to zero. Most cross terms are functions of two internal coordinates, such as stretch–stretch, stretch–bend and stretch–torsion terms, but cross terms involving more than two internal coordinates such as the bend–bend–torsion have also been used. Various functional forms are possible for the cross terms. For example, the stretch–stretch cross term between two bonds 1 and 2 can be modelled as:

$$v(l_1, l_2) = \frac{k_{1,2}}{2} [(l_1 - l_{1,0})(l_2 - l_{2,0})] \quad (4.13)$$

The stretching of the two bonds adjoining an angle could be modelled using an equation of the following form (as in MM2, MM3 and MM4):

$$v(l_1, l_2, \theta) = \frac{k_{l_1, l_2, \theta}}{2} [(l_1 - l_{1,0}) + (l_2 - l_{2,0})](\theta - \theta_0) \quad (4.14)$$

In a *Urey–Bradley* force field, angle bending is achieved using 1,3 non-bonded interactions rather than an explicit angle-bending potential. The stretch–bond term in such a force field would be modelled by a harmonic function of the distance between the 1,3 atoms:

$$v(r_{1,3}) = \frac{k_{r_{1,3}}}{2} (r_{1,3} - r_{1,3}^0)^2 \quad (4.15)$$

A stretch–torsion cross term can be used to model the stretching of a bond that occurs in an eclipsed conformation. Two possible functional forms are:

$$v(l, \omega) = k(l - l_0) \cos n\omega \quad (4.16)$$

$$v(l, \omega) = k(l - l_0)[1 + \cos n\omega] \quad (4.17)$$

n is the periodicity of the rotation about the bond ($n = 3$ for sp^3 – sp^3 bonds).

Torsion–bend and torsion–bend–bend terms may also be included; the latter, for example, would couple two angles A–B–C and B–C–D to a torsion angle A–B–C–D. Maple, Dinur and Hagler used quantum mechanics calculations to investigate which of the cross terms are most important and suggested that the stretch–stretch, stretch–bend, bend–bend, stretch–torsion and bend–bend–torsion were most important [Dinur and Hagler 1991] (schematically illustrated in Figure 4.13).

It has been suggested that the presence of cross terms (together with some other features) can provide a general way to classify force fields [Hwang *et al.* 1994]. A class I force field was considered one which is restricted to harmonic terms (e.g. for bond stretching and angle bending) and which does not have any cross terms. A class II force field would have anharmonic terms (e.g. through the use of Morse potentials or quartic terms) and explicit cross terms to account for the coupling between coordinates. The presence of these higher and cross terms would tend to improve the ability of the force field to predict the properties of more unusual systems (such as those which are highly strained) and also to enhance its ability to reproduce vibrational spectra. Another characteristic of a class II force field was that it could be used without modification to model the properties of isolated small molecules, condensed phases and macromolecular systems. It was subsequently suggested by Allinger [Allinger *et al.* 1996b] that a class III force field would also take account of chemical effects and other features such as electronegativity and hyperconjugation. A classic

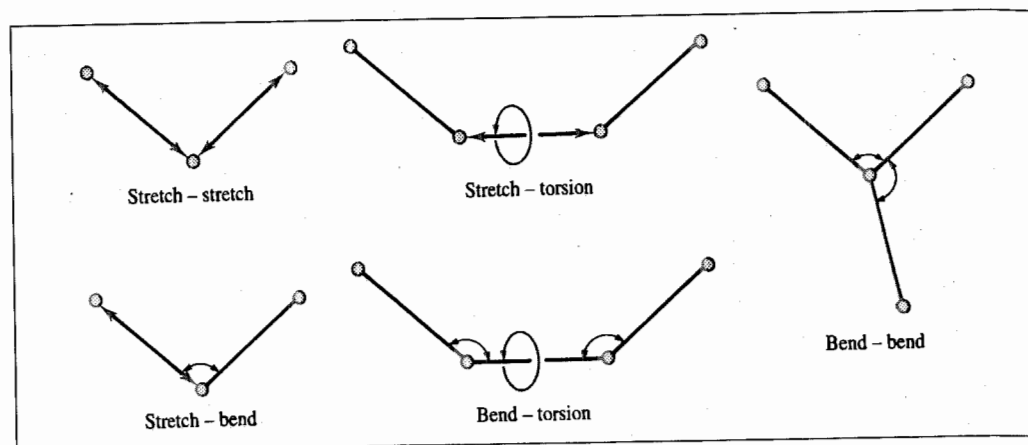


Fig. 4.13: Schematic illustration of the cross terms believed to be most important in force fields. (Adapted from Dinur U and A T Hagler 1991. *New Approaches to Empirical Force Fields*. In *Reviews in Computational Chemistry*, Lipkowitz K B and D B Boyd (Editors). New York, VCH Publishers, pp. 99-164.)

example of the latter effect (hyperconjugation) is the change in the length of the C–H bond in acetaldehyde with rotation about the C–C bond. When the C–H bond is perpendicular to the plane of the carbonyl group there is maximum overlap between the σ orbital of the C–H bond and the π^* orbital of the carbonyl carbon. Donation of electron density from the C–H bond to this π^* orbital is accompanied by a lengthening of the bond and a greater contribution from the charged resonance structure (Figure 4.14). When the bond to the hydrogen atom is in the plane the overlap is minimal. *Ab initio* calculations suggested that the bond length changed by 0.006 Å between the two forms. This effect was incorporated within MM4 by a term of the following form:

$$\Delta l = k(1 - \cos 2\omega) \quad (4.18)$$

This is a kind of torsion–stretch cross term but different from the one where the central bond changes with torsion angle. There has been some considerable debate about the existence and origin of the hyperconjugative effects, but low-temperature X-ray crystallographic experiments on appropriate compounds together with *ab initio* calculations certainly reveal a detectable effect.

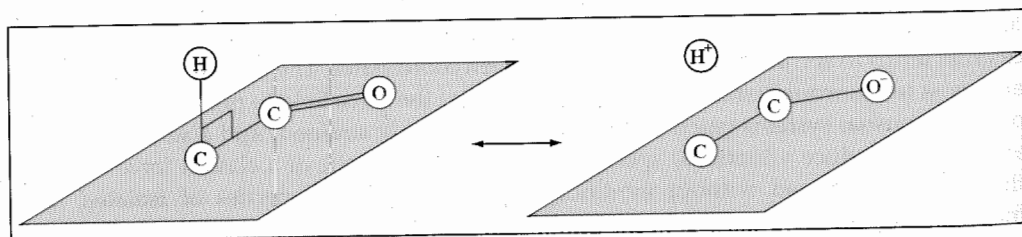


Fig. 4.14: Valence bond representation of the hyperconjugation effect which leads to a lengthening of the C–H bond in acetaldehyde.

4.8 Introduction to Non-bonded Interactions

Independent molecules and atoms interact through non-bonded forces, which also play an important role in determining the structure of individual molecular species. The non-bonded interactions do not depend upon a specific bonding relationship between atoms. They are 'through-space' interactions and are usually modelled as a function of some inverse power of the distance. The non-bonded terms in a force field are usually considered in two groups, one comprising electrostatic interactions and the other van der Waals interactions.

4.9 Electrostatic Interactions

4.9.1 The Central Multipole Expansion

Electronegative elements attract electrons more than less electronegative elements, giving rise to an unequal distribution of charge in a molecule. This charge distribution can be represented in a number of ways, one common approach being an arrangement of fractional point charges throughout the molecule. These charges are designed to reproduce the electrostatic properties of the molecule. If the charges are restricted to the nuclear centres they are often referred to as *partial atomic charges* or *net atomic charges*. The electrostatic interaction between two molecules (or between different parts of the same molecule) is then calculated as a sum of interactions between pairs of point charges, using Coulomb's law:

$$\mathcal{V} = \sum_{i=1}^{N_A} \sum_{j=1}^{N_B} \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} \quad (4.19)$$

N_A and N_B are the numbers of point charges in the two molecules. This approach to the representation and calculation of electrostatic interactions will be considered in more detail in Section 4.9.2. First, we shall consider an alternative approach to the calculation of electrostatic interactions which treats a molecule as a single entity and is (in principle at least) capable of providing a very efficient way to calculate electrostatic intermolecular interactions. This is the *central multipole expansion*, which is based upon the electric moments or multipoles: the charge, dipole, quadrupole, octopole, and so on introduced in Section 2.7.3. These moments are usually represented by the following symbols: q (charge), μ (dipole), Θ (quadrupole) and Φ (octopole). We are often interested in the lowest non-zero electric moment. Thus species such as Na^+ , Cl^- , NH_4^+ or CH_3CO_2^- have the charge as their lowest non-zero moment. For many uncharged molecules the dipole is the lowest non-zero moment. Molecules such as N_2 and CO_2 have the quadrupole as their lowest non-zero moment. The lowest non-zero moment for methane and tetrafluoromethane is the octopole. Each of these multipole moments can be represented by an appropriate distribution of charges. Thus a dipole can be represented using two charges placed an appropriate distance apart. A quadrupole can be represented using four charges and an octopole by eight charges. A complete description of the charge distribution around a molecule requires all of the non-zero electric moments to the specified. For some molecules,

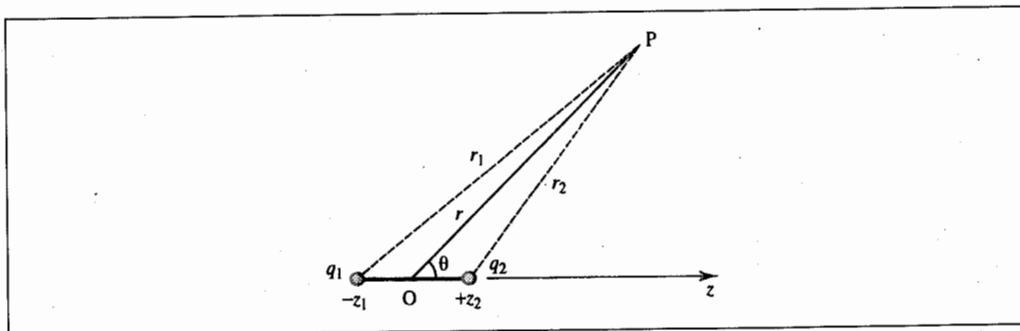


Fig. 4.15: The electrostatic potential due to two point charges.

the lowest non-zero moment may not be the most significant and it may therefore be unwise to ignore the higher-order terms in the expansion without first checking their values.

To illustrate how the multipolar expansion is related to a distribution of charges in a system, let us consider the simple case of a molecule with two charges q_1 and q_2 , positioned at $-z_1$ and z_2 , respectively (Figure 4.15). The electrostatic potential at point P (a distance r from the origin, r_1 from charge q_1 and r_2 from charge q_2) is then given by:

$$\phi(r) = \frac{1}{4\pi\epsilon_0} \left(\frac{q_1}{r_1} + \frac{q_2}{r_2} \right) \quad (4.20)$$

By applying the cosine rule this can be written as follows (see Figure 4.15):

$$\phi(r) = \frac{1}{4\pi\epsilon_0} \left(\frac{q_1}{\sqrt{r^2 + z_1^2 + 2rz_1 \cos \theta}} + \frac{q_2}{\sqrt{r^2 + z_2^2 - 2rz_2 \cos \theta}} \right) \quad (4.21)$$

If $r \gg z_1$ and $r \gg z_2$ then this expression can be expanded as follows:

$$\phi(r) = \frac{1}{4\pi\epsilon_0} \left(\frac{q_1 + q_2}{r} + \frac{(q_2 z_2 - q_1 z_1) \cos \theta}{r^2} + \frac{(q_1 z_1^2 + q_2 z_2^2)(3 \cos^2 \theta - 1)}{2r^3} + \dots \right) \quad (4.22)$$

We can now associate the appropriate terms in the expansion with the various electric moments:

$$\phi(r) = \frac{1}{4\pi\epsilon_0} \left(\frac{q}{r} + \frac{\mu \cos \theta}{r^2} + \frac{\Theta(3 \cos^2 \theta - 1)}{2r^3} + \dots \right) \quad (4.23)$$

Thus $(q_1 + q_2)$ is the charge; $(q_2 z_2 - q_1 z_1)$ is the dipole; $(q_1 z_1^2 + q_2 z_2^2)$ is the quadrupole, and so on. One interesting feature about a charge distribution is that only the first non-zero moment is independent of the choice of origin. Thus, if a molecule is electrically neutral (i.e. $q_1 + q_2 = 0$) then its dipole moment is independent of the choice of origin. This can be demonstrated for our two-charge system as follows. If the position of the origin is now moved to a point $-z'$, then the dipole moment relative to this new origin is given by:

$$\mu' = q_2(z_2 + z') - q_1(z_1 - z') = \mu + qz' \quad (4.24)$$

Only if the total charge on the system (q) equals zero will the dipole moment be unchanged. Similar arguments can be used to show that if both the charge and the dipole moment are zero then the quadrupole moment is independent of the choice of origin. For convenience, the origin is often taken to be the centre of mass of the charge distribution.

The electric moments are examples of *tensor properties*: the charge is a rank 0 tensor (which is the same as a scalar quantity); the dipole is a rank 1 tensor (which is the same as a vector, with three components along the x , y and z axes); the quadrupole is a rank 2 tensor with nine components, which can be represented as a 3×3 matrix. In general, a tensor of rank n has 3^n components.

For a distribution of charges (one not restricted to lie along one of the Cartesian axes), the dipole moment is given by:

$$\mu = \sum q_i \mathbf{r}_i \quad (4.25)$$

The components of the dipole moment along the x , y and z axes are $\sum q_i x_i$, $\sum q_i y_i$ and $\sum q_i z_i$. The analogous way to define the quadrupole moment is as follows:

$$\Theta = \begin{pmatrix} \sum q_i x_i^2 & \sum q_i x_i y_i & \sum q_i x_i z_i \\ \sum q_i y_i x_i & \sum q_i y_i^2 & \sum q_i y_i z_i \\ \sum q_i z_i x_i & \sum q_i z_i y_i & \sum q_i z_i^2 \end{pmatrix} \quad (4.26)$$

This definition of the quadrupole is obviously dependent upon the orientation of the charge distribution within the coordinate frame. Transformation of the axes can lead to alternative definitions that may be more informative. Thus the quadrupole moment is commonly defined as follows:

$$\Theta = \frac{1}{2} \begin{pmatrix} \sum_i q_i (3x_i^2 - r_i^2) & 3 \sum_i q_i x_i y_i & 3 \sum_i q_i x_i z_i \\ 3 \sum_i q_i x_i z_i & \sum_i q_i (3y_i^2 - r_i^2) & 3 \sum_i q_i y_i z_i \\ 3 \sum_i q_i x_i z_i & 3 \sum_i q_i y_i z_i & \sum_i q_i (3z_i^2 - r_i^2) \end{pmatrix} \quad (4.27)$$

In Equation (4.27) $r_i^2 = x_i^2 + y_i^2 + z_i^2$. This definition enables one to assess the deviation from spherical symmetry as a spherically symmetric charge distribution will have

$$\sum_i q_i x_i^2 = \sum_i q_i y_i^2 = \sum_i q_i z_i^2 = \frac{1}{3} \sum_i q_i r_i^2 \quad (4.28)$$

and so the diagonal elements of the tensor will be zero. Quadrupoles are also reported in terms of the *principal axes*; these are three mutually perpendicular axes α , β and γ , which are linear combinations of x , y and z such that the quadrupole tensor is diagonal (i.e. off-diagonal elements are zero):

$$\Theta = \begin{pmatrix} \Theta_{\alpha\alpha} & 0 & 0 \\ 0 & \Theta_{\beta\beta} & 0 \\ 0 & 0 & \Theta_{\gamma\gamma} \end{pmatrix} \quad (4.29)$$

Let us now consider the effect of placing another molecule with a linear charge distribution (charges q'_1 and q'_2) with its centre of mass at the point P. The relative orientation of the two

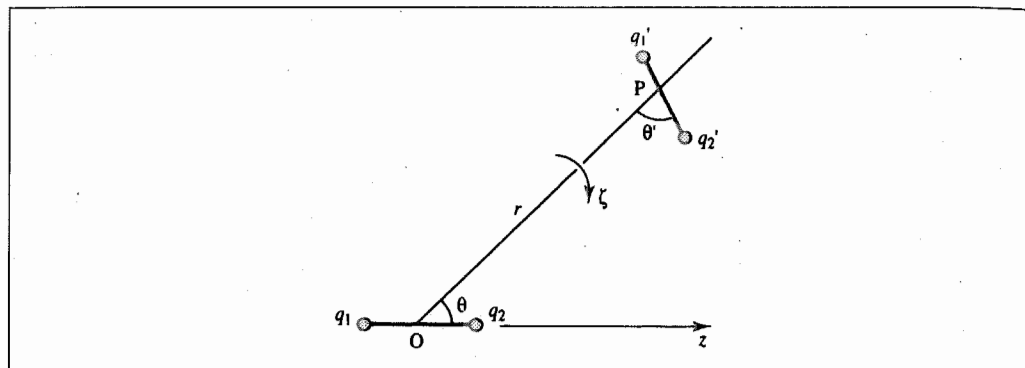
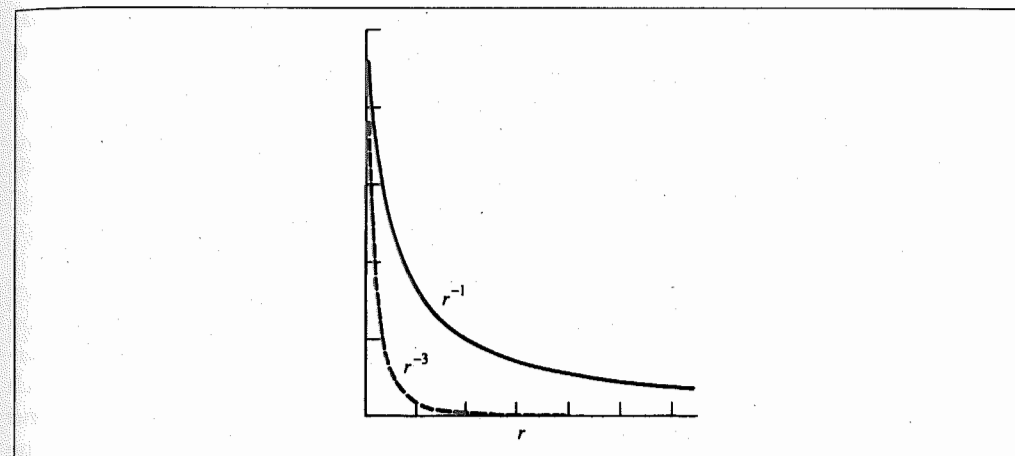


Fig. 4.16: The relative orientation of two dipoles.

molecules can be described in terms of four parameters (the distance joining their centres of mass and three angles as shown in Figure 4.16). The electrostatic interaction between the two molecules is calculated by multiplying each charge by the potential at that point and adding the result for each charge. The following expression is the result [Buckingham 1959]:

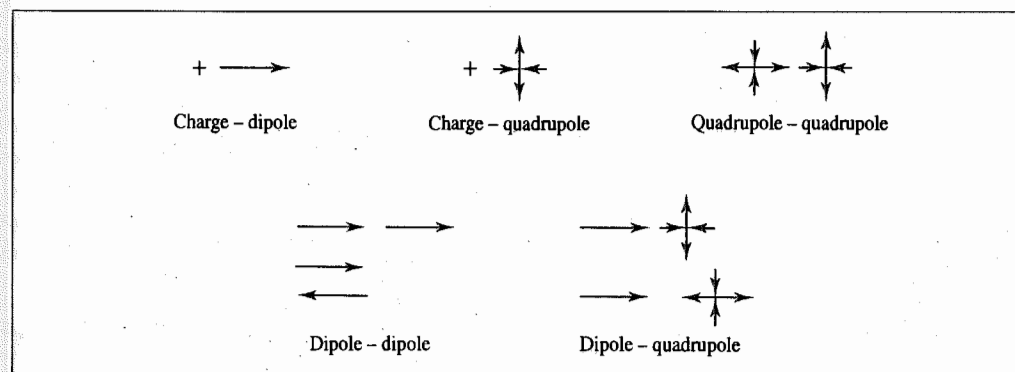
$$V(q, q') = \frac{1}{4\pi\epsilon_0} \left\{ \begin{aligned} & \frac{qq'}{r} \\ & + \frac{1}{r^2} (q\mu' \cos \theta + q'\mu \cos \theta') \\ & + \frac{\mu\mu'}{r^3} (2 \cos \theta \cos \theta' + \sin \theta \sin \theta' \cos \zeta) \\ & + \frac{1}{2r^3} [q\Theta' (3 \cos^2 \theta' - 1) + q'\Theta (3 \cos^2 \theta - 1)] \\ & + \frac{3}{2r^4} [\mu\Theta' \{ \cos \theta (3 \cos^2 \theta' - 1) + 2 \sin \theta \sin \theta' \cos \theta' \cos \zeta \} \\ & \quad + \mu'\Theta \{ \cos \theta' (3 \cos^2 \theta - 1) + 2 \sin \theta' \sin \theta \cos \theta \cos \zeta \}] \\ & + \frac{3\Theta\Theta'}{4r^5} [1 - 5 \cos^2 \theta - 5 \cos^2 \theta' + 17 \cos^2 \theta \cos^2 \theta' \\ & \quad + 2 \sin^2 \theta \sin^2 \theta' \cos^2 \zeta + 16 \sin \theta \sin \theta' \cos \theta \cos \theta' \cos \zeta] \\ & + \dots \end{aligned} \right\} \quad (4.30)$$

The energy of interaction between two charge distributions is thus an infinite series that includes charge-charge, charge-dipole, dipole-dipole, charge-quadrupole, dipole-quadrupole interactions, quadrupole-quadrupole terms, and so on. These terms depend on different inverse powers of the separation r . If the molecules are neutral (i.e. $q = q' = 0$) then the leading term in the expansion is that due to the dipole-dipole interaction, which varies as r^{-3} . This is a key result, for the range of the dipole-dipole interaction (r^{-3}) is much less than that of the Coulomb interaction (r^{-1}), Figure 4.17. This will be important in later chapters, where we shall collect atoms together into neutral groups. The electrostatic interaction

Fig. 4.17: The charge-charge energy decays much more slowly ($\propto r^{-1}$) than the dipole-dipole energy ($\propto r^{-3}$).

between these groups then decays as r^{-3} rather than the r^{-1} dependence of each individual charge-charge interaction. This can be seen in Figure 4.17, in which the functions r^{-1} and r^{-3} have been plotted as a function of distance. Even when the dipole-dipole interaction energy has fallen off almost to zero the charge-charge interaction energy is still significant. In general, the interaction energy between two multipoles of order n and m decreases as $r^{-(n+m+1)}$. It should be emphasised again that these expressions are only valid when the separation of the two molecules, r , is much larger than the internal dimensions of the molecules. The favourable arrangements for the various multipoles are shown in Figure 4.18.

A central multipole expansion therefore provides a way to calculate the electrostatic interaction between two molecules. The multipole moments can be obtained from the wavefunction and can therefore be calculated using quantum mechanics (see Section 2.7.3) or can be determined from experiment. One example of the use of a multipole expansion is

Fig. 4.18: The most favourable orientations of various multipoles. (Figure adapted from Buckingham A D 1959. *Molecular Quadrupole Moments*. Quarterly Reviews of the Chemical Society 13:183-214.)

the benzene model of Claessens, Ferrario and Ryckaert [Claessens *et al.* 1983]. Benzene has no charge and no dipole moment, but it does have a sizeable quadrupole. The inclusion of the quadrupole was found to give clearly superior results in molecular dynamics simulations of the liquid state over models that lacked any electronic contribution.

The main advantage of the multipolar description for calculating the electrostatic interactions between molecules is its efficiency. For example, the charge-charge interaction energy between two benzene molecules would require 144 individual charge-charge interactions with a partial atomic charge model rather than the single quadrupole-quadrupole term. Unfortunately, the multipole expansion is not applicable when the molecules are separated by distances comparable with the molecular dimensions. The formal condition for convergence of the multipolar interaction energy is that the distance between two interacting molecules should be larger than the sum of the distances from the centre of each molecule to the furthest part of its charge distribution. If a sphere is constructed around each molecule, positioned on its centre of mass, with a radius that encompasses all of the charge distribution, then the multipole expansion for the interaction between two molecules will converge if these spheres do not intersect. Even if one requires the sphere to encompass just the nuclei in a molecule (i.e. ignoring the fact that the charge distribution around a molecule extends to infinity) there may still be problems. For example, the convergence sphere for a molecule such as butane would extend beyond the van der Waals radii in some directions, enabling other molecules to penetrate the convergence sphere, as illustrated in Figure 4.19. Another problem is that the multipolar expansion may be slow to converge. The multipolar expansion is often located at the centre of mass, but this may not be the best choice to achieve the most rapid convergence.

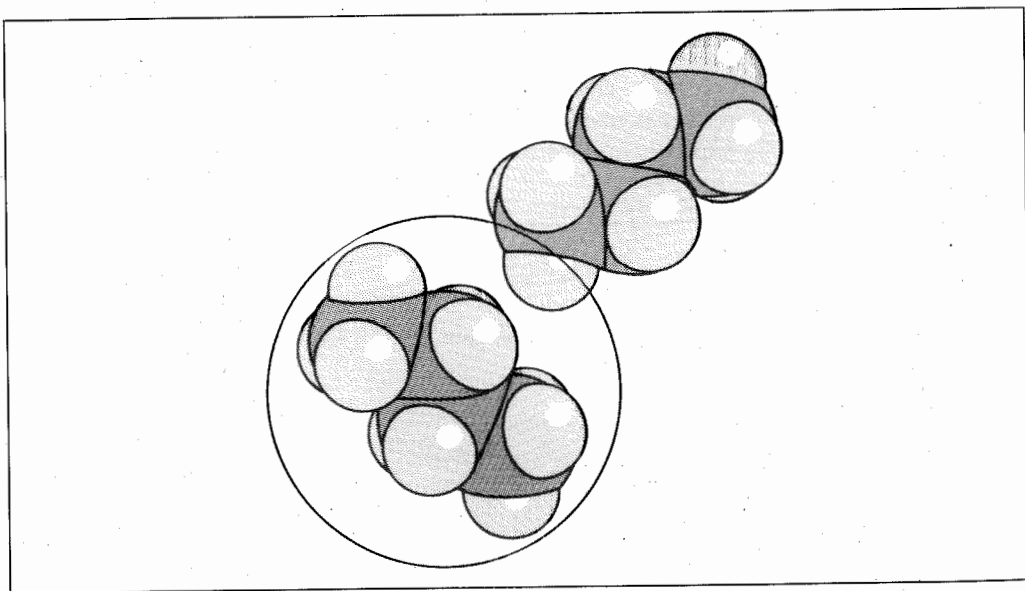


Fig. 4.19: The convergence sphere of the multipole expansion for a molecule such as butane may be penetrated by another molecule.

There are other difficulties with the central multipole expansion. The multipole moments are properties of the entire molecule and so cannot be used to determine intramolecular interactions. The central multipole model thus tends to be restricted to calculations involving small molecules that are kept fixed in conformation during the calculation, and where the interactions between molecules act at their centres of mass. It can be a complicated procedure to calculate the forces acting on a molecule with a multipole model. The interaction between multipoles of zero order (i.e. charges) gives rise to a simple translational force. Multipoles of a higher order have directionality, and interactions between these produce a torque, or twisting force. Moreover, whereas the charge-charge forces are equal and opposite, the torque acting on molecule *i* due to another molecule *j* is not necessarily equal and opposite to the torque on molecule *j* due to molecule *i*.

4.9.2 Point-charge Electrostatic Models

We therefore return to the point-charge model for calculating electrostatic interactions. If sufficient point charges are used then all of the electric moments can be reproduced and the multipole interaction energy, Equation (4.30), is exactly equal to that calculated from the Coulomb summation, Equation (4.19).

An accurate representation of a molecule's electrostatic properties may require charges to be placed at locations other than at the atomic nuclei. A simple example of this is molecular nitrogen, which has a dipole moment of zero. The total charge on nitrogen is zero, and so an atomic partial charge model would put zero charge on each nucleus. However, nitrogen does have a quadrupole moment and this significantly affects its properties. The simplest way to model this is to place three partial charges along the bond: a charge of $-q$ at each nucleus and $+2q$ at the centre of mass. The quadrupole-quadrupole interaction between two nitrogen molecules can then be calculated by summing nine pairs of charge-charge interactions. The value of q can be calculated using the following relationship between the quadrupole moment and the partial charge:

$$\Theta = 2q(l/2)^2 \quad (4.31)$$

l is the bond length. The experimental quadrupole moment is consistent with a charge, q , of approximately $0.5e$. In fact, a better representation of the electrostatic potential around the nitrogen molecule is obtained using the five-charge model shown in Figure 4.20.

An alternative to the point charge model is to assign dipoles to the bonds in the molecule. The electrostatic energy is then given as a sum of dipole-dipole interaction energies. This approach (which is adopted in MM2/MM3/MM4) can be unwieldy for molecules that have a formal charge and which require charge-charge and charge-dipole terms to be included in the energy expression. Charged species are dealt with more naturally using the point charge model.

4.9.3 Calculating Partial Atomic Charges

Given the widespread use of the partial atomic charge model, it is important to consider how the charges are obtained. For simple species the atomic charges required to reproduce the

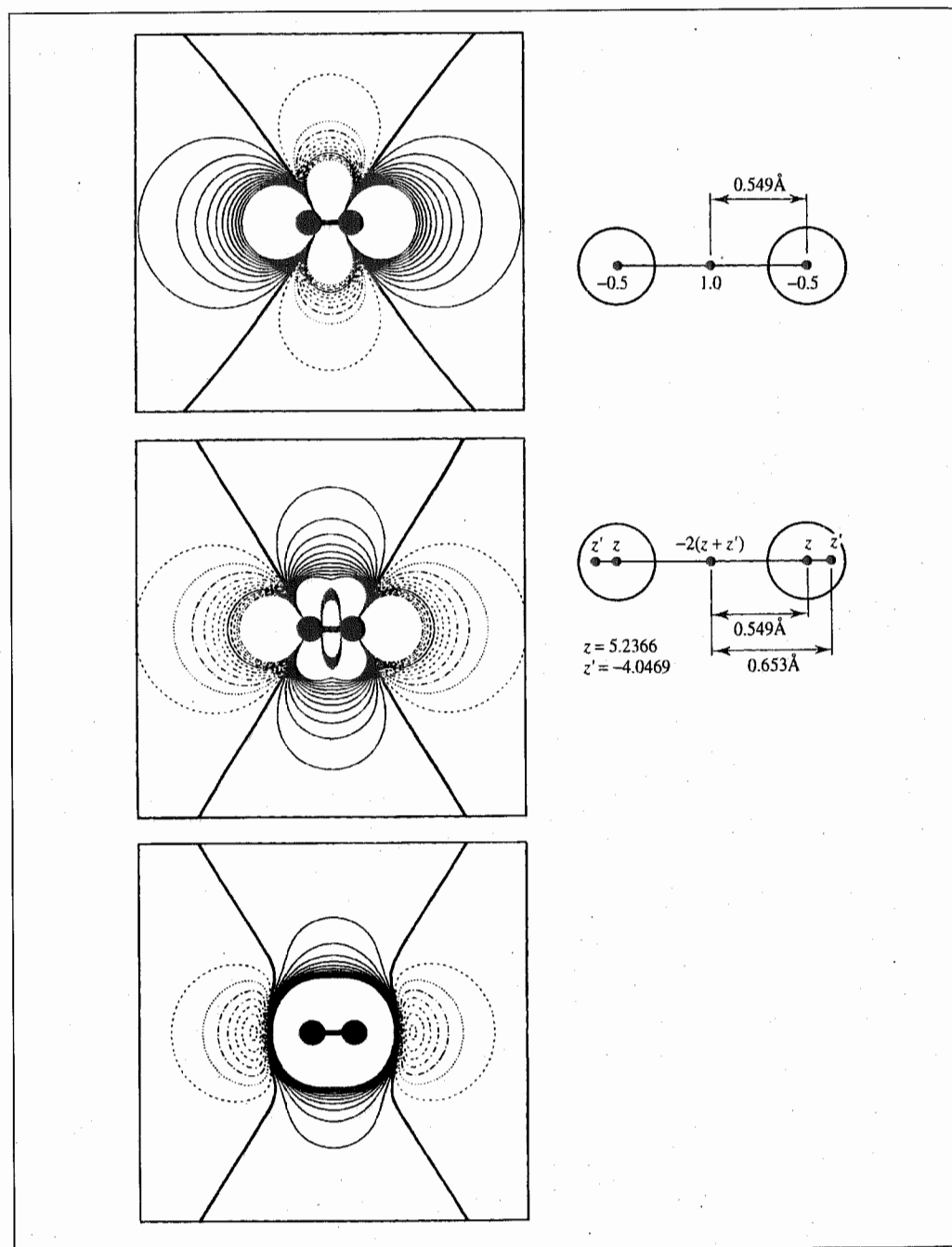


Fig. 4.20: Two charge models for N_2 with the electrostatic potentials that they generate. Also shown is the electrostatic potential calculated using ab initio quantum mechanics (6-31G* basis set.) Negative contours are dashed and the zero contour is bold.

electric moments can be calculated exactly if the geometry is known. For example, the experimentally determined dipole moment of HF (1.82 D) can be reproduced by placing equal but opposite charges of $0.413e$ on the two atomic nuclei (assuming a bond length of 0.917 \AA). The tetrahedral arrangement of the hydrogens about the carbon in methane means that each hydrogen atom has an identical charge equal to one quarter the charge on the carbon. The molecule is electrically neutral with zero dipole and quadrupole moments but a non-zero octopole moment, which can be reproduced using a hydrogen charge of approximately $0.14e$.

In some cases the atomic charges are chosen to reproduce thermodynamic properties calculated using a molecular dynamics or Monte Carlo simulation. A series of simulations is performed and the charge model is modified until satisfactory agreement with experiment is obtained. This approach can be quite powerful despite its apparent simplicity, but it is only really practical for small molecules or simple models.

The electrostatic properties of a molecule are a consequence of the distribution of the electrons and the nuclei and thus it is reasonable to assume that one should be able to obtain a set of partial atomic charges using quantum mechanics. Unfortunately, the partial atomic charge is not an experimentally observable quantity and cannot be unambiguously calculated from the wavefunction. This explains why numerous ways to determine partial atomic charges have been proposed, and why there is still considerable debate as to the 'best' method to derive them. Indirect comparisons of the various methods are possible, usually by calculating appropriate quantities from the charge model and then comparing the results with either experiment or quantum mechanics. For example, one might examine how well the charge model reproduces the experimental or quantum mechanical multipole moments or the electrostatic potential around the molecule.

We have already encountered in Section 2.7.5 the population analysis method for calculating partial atomic charges. Such sets of charges (commonly referred to as *Mulliken charges* when obtained from that particular partitioning scheme) are often considered to be inappropriate for accurately representing the interactions between molecules. This is because Mulliken charges are primarily dependent upon the constitution of the molecule - how the atoms are bonded together - rather than being designed to reproduce the properties that determine how molecules interact with each other, such as the electrostatic potential. The importance of the electrostatic potential in intermolecular interactions has resulted in much interest in schemes that calculate charges consistent with this particular property.

4.9.4 Charges Derived from the Molecular Electrostatic Potential

The electrostatic potential at a point is the force acting on a unit positive charge placed at that point. The nuclei give rise to a positive (i.e. repulsive) force, whereas the electrons give rise to a negative potential. The electrostatic potential is an observable quantity that can be determined from a wavefunction using Equations (2.222) and (2.223):

$$\phi(\mathbf{r}) = \phi_{\text{nucl}}(\mathbf{r}) + \phi_{\text{elec}}(\mathbf{r}) = \sum_{A=1}^M \frac{Z_A}{|\mathbf{r} - \mathbf{R}_A|} - \int \frac{d\mathbf{r}' \rho(\mathbf{r}')}{|\mathbf{r}' - \mathbf{r}|} \quad (4.32)$$

The electrostatic potential is a continuous property and is not easily represented by an analytical function. Consequently, it is necessary to derive a discrete representation for use in numerical analysis. The objective is to derive the set of partial charges (usually partial atomic charges) that best reproduces the quantum mechanical electrostatic potential at a series of points surrounding the molecule. A solution to this problem was suggested by Cox and Williams [Cox and Williams 1981]. The electrostatic potential at each of the chosen points is calculated from the wavefunction. A least-squares fitting procedure is then employed to determine the set of partial atomic charges that best reproduces the electrostatic potential at the points, subject to the constraint that the sum of the charges should be equal to the net charge on the molecule. Symmetry conditions may also be imposed to ensure that the charges on symmetrically equivalent atoms are equal. It is also possible to require the atomic charges to reproduce other electrostatic properties of the molecules such as the dipole moment. The fitting procedure minimises the sum of squares of the differences in the electrostatic potential. Thus, if the electrostatic potential at a point is ϕ_i^0 and if the value from the charge model is ϕ_i^{calc} , then the objective is to minimise the following function:

$$R = \sum_{i=1}^{N_{\text{points}}} w_i (\phi_i^0 - \phi_i^{\text{calc}})^2 \quad (4.33)$$

N_{points} is the number of points and w_i is a weighting factor that enables different points to be given different degrees of 'importance' in the fitting process. One of the charges is dependent on the values of the others (because the sum must equal Z , the molecular charge). This N th charge has a value given by:

$$q_N = Z - \sum_{j=1}^{N-1} q_j \quad (4.34)$$

The electrostatic potential due to the charges q_j at the point i is given by Coulomb's law:

$$\phi_i^{\text{calc}} = \sum_{j=1}^{N-1} \frac{q_j}{4\pi\epsilon_0 r_{ij}} + \frac{Z - \sum_{j=1}^{N-1} q_j}{4\pi\epsilon_0 r_{iN}} \quad (4.35)$$

r_{ij} is the distance from the charge j to the point i . At a minimum value of the error function, R , the first derivative is equal to zero with respect to all charges q_k :

$$\frac{\partial R}{\partial q_k} = -2 \sum_{i=1}^{N_{\text{points}}} w_i (\phi_i^0 - \phi_i^{\text{calc}}) \left(\frac{\partial \phi_i^{\text{calc}}}{\partial q_k} \right) = 0 \quad (4.36)$$

This equation can be written in the following form:

$$\sum_{i=1}^{N_{\text{points}}} w_i \left(\phi_i^0 - \frac{Z}{r_{iN}} \right) \left(\frac{1}{r_{ik}} - \frac{1}{r_{iN}} \right) = \sum_{j=1}^{N-1} \left[\sum_{i=1}^{N_{\text{points}}} w_i \left(\frac{1}{r_{ik}} - \frac{1}{r_{iN}} \right) \left(\frac{1}{r_{ij}} - \frac{1}{r_{iN}} \right) \right] \frac{q_j}{4\pi\epsilon_0} \quad (4.37)$$

When expressed in this way, then the set of equations can be recast as a matrix equation of the form $\mathbf{A}\mathbf{q} = \mathbf{a}$. The charges \mathbf{q} are then determined using standard matrix methods via $\mathbf{q} = \mathbf{A}^{-1}\mathbf{a}$.

The points i ($1, 2, \dots, N_{\text{points}}$) where the potential is fitted can be chosen in a variety of ways but should be taken from the region where it is most important to model intermolecular interactions correctly. This region is just beyond the van der Waals radii of the atoms involved. Cox and Williams selected points from a regular grid in a shell defined by two surfaces, one corresponding to the union of the van der Waals radii plus 1.2 Å and the others approximately 1 Å beyond that. The CHELP procedure of Chirlian and Francl [Chirlian and Francl 1987] uses spherical shells, 1 Å apart, centred on each atom with points symmetrically distributed on the surface. Any points within the van der Waals radius of any atom in the system are discarded and the shells extend to 3 Å from the van der Waals surface of the molecule. The CHELP method employs a Lagrange multiplier method to find the atomic charges, rather than an iterative least-squares procedure. This minimises the error function R (Equation (4.33)) subject to the constraint that the charges sum to the total molecular charge. Such an analysis yields a set of $N + 1$ equations in $N + 1$ unknowns and can be solved using standard matrix methods. The CHELPG algorithm of Breneman and Wiberg [Breneman and Wiberg 1990] combines the regular grid of points of Cox and Williams with the Lagrange multiplier method of Chirlian and Francl as the results from CHELP were found to change if the molecule was reoriented in the coordinate system. In CHELPG a cubic grid of points (spaced 0.3–0.8 Å apart) is used and all grid points that lie within the van der Waals radius of any atom are discarded, together with all points that lie further than 2.8 Å away from any atom.

The algorithm of Singh and Kollman used to derive the charges in the 1984 AMBER force field uses points on a series of molecular surfaces, constructed using gradually increasing van der Waals radii for the atoms [Singh and Kollman 1984]. The points at which the potential was fitted were located on these shells. For the 1995 AMBER force field a modified version of this electrostatic potential method was employed (termed 'restrained electrostatic potential fit', or RESP [Bayly *et al.* 1993]). The RESP algorithm uses hyperbolic restraints on non-hydrogen atoms. These restraints have the effect of reducing the charges on some atoms, particularly buried carbon atoms, which can be assigned artificially high charges in standard electrostatic potential fitting methods. The RESP charges also vary less with the molecular conformation.

4.9.5 Deriving Charge Models for Large Systems

Molecular mechanics is used to model systems containing thousands of atoms such as polymers. How then can charges be derived for such species? Clearly one cannot routinely perform quantum mechanical calculations on a molecule with so many atoms and so it must be broken into fragments of a suitable size. In some cases the fragments might appear relatively easy to define; for example, many polymeric systems are constructed by connecting together chemically defined monomeric units. The atomic charges for each monomer should be obtained from calculations on suitable fragments that recreate the immediate local environment of the fragment in the larger molecule. For example, partial atomic charges for amino acids are often obtained from calculations on a 'dipeptide' fragment (see Figure 4.21), which is more akin to the environment within a protein than in an isolated amino acid.

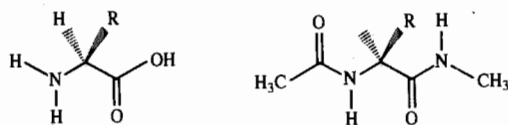


Fig. 4.21: The charges used for calculations on proteins are best derived using a suitable fragment for each amino acid that reflects the environment within the protein (right), rather than the isolated amino acid (left).

The charge sets obtained from electrostatic potential fitting can be highly dependent upon the basis set used to derive the wavefunction. Moreover, the charges do not always improve if a larger basis set is used. It is generally considered that the 6-31G⁺ basis set gives reasonable results for calculations relevant to condensed phases. In many cases it is possible to scale the results of a calculation using a small basis set or even a lower level of theory (such as a semi-empirical calculation) to obtain results comparable with those of a high-level calculation. Of the various semi-empirical methods available, MNDO appears to give the best correspondence with the charges derived from *ab initio* calculations, and scaling factors have been determined by several research groups [Ferenczy *et al.* 1990; Luque *et al.* 1990; Bezler *et al.* 1990]. An additional complicating factor is that the charges obtained from electrostatic potential fitting will often depend upon the conformation for which the quantum mechanical calculation was performed [Williams 1990]. One solution is to perform a series of charge calculations for different conformations and then use a charge model in which each charge is weighted according to the relative population of that particular conformation as calculated from the Boltzmann distribution [Reynolds *et al.* 1992]. In a few charge models the charges vary continuously with the conformation [Rappé and Goddard 1991; Dinur and Hagler 1995].

4.9.6 Rapid Methods for Calculating Atomic Charges

Some methods calculate atomic charges solely from information about the atoms present in the molecule and the way in which the atoms are connected. The great advantage of such methods is that they are very fast and can be used to calculate the charge distributions for large numbers of molecules (e.g. in a database). We will consider the Gasteiger and Marsili method [Gasteiger and Marsili 1980] as an example.

The Gasteiger–Marsili approach uses the concept of the *partial equalisation of orbital electronegativity*. Electronegativity is a concept well known to chemists, being defined by Pauling as ‘the power of an atom to attract electrons to itself’. Mulliken subsequently defined the electronegativity of an atom A as the average of its ionisation potential I_A and its electron affinity E_A :

$$\chi_A = \frac{1}{2}(I_A + E_A) \quad (4.38)$$

As Mulliken pointed out, the ionisation potential and electron affinity are specific to a given valence state of an atom, and therefore the electronegativities of an atom’s valence states would not be expected to be the same. This idea can be extended to the concept of orbital

electronegativity, which is the electronegativity of a specific orbital in a given valence state. For example, an sp orbital has a higher electronegativity than an sp³ orbital. The orbital electronegativity will also depend on the occupancy of the orbital; an empty orbital will be better able to attract an electron than an orbital with a single electron, which in turn will be better than an orbital with two electrons. The electronegativity of an orbital will also be affected by the charges in other orbitals. Gasteiger and Marsili assumed a polynomial relationship between the orbital electronegativity $\chi_{\mu A}$ of an orbital ϕ_{μ} in atom A and the charge Q_A on the atom A:

$$\chi_{\mu A} = a_{\mu} + b_{\mu A} Q_A + \chi_{\mu A} Q_A^2 \quad (4.39)$$

Values of the coefficients a , b and c were derived for common elements in their usual valence states (for example, for carbon there are different values for sp³, sp²π and spπ² valence states).

Electrons flow from the less electronegative elements to the more electronegative ones. This flow of electrons results in a positive charge on the less electronegative atoms and a negative charge on the more electronegative atoms, and as such the flow acts to equalise the electronegativities. Total equalisation of electronegativity does not, however, lead to chemically sensible results. This effect is modelled in the Gasteiger and Marsili approach by an iterative procedure, in which less and less charge is transferred between bonded atoms at each step. The electron charge transferred from an atom A to an atom B (where B is more electronegative than A) in iteration k is given by:

$$Q^{(k)} = \frac{\chi_B^{(k)} - \chi_A^{(k)}}{\chi_A^+} \alpha^k \quad (4.40)$$

In Equation (4.40), $Q^{(k)}$ is the charge (in electrons) transferred; $\chi_A^{(k)}$ and $\chi_B^{(k)}$ are the electronegativities of the atoms A and B; χ_A^+ is the electronegativity of the cation of the less electronegative atom and α is a damping factor which is raised to the power k . Gasteiger and Marsili set α to $\frac{1}{2}$. The charge on each atom is initially assigned its formal charge. In each iteration, the electronegativities are calculated using Equation (4.39) and hence the charge to be transferred. The total charge on an atom at the end of each iteration is thus obtained by adding the charge transferred from all bonds to the atom to the value of the charge from the previous iteration. The damping factor α^k reduces the influence of the more electronegative atoms. This influence decreases with each iteration. With a damping factor of $\frac{1}{2}$ rapid convergence is achieved, usually within four or five steps.

A somewhat related method is the charge equilibration method of Rappé and Goddard [Rappé and Goddard 1991]. This is employed in the ‘Universal Force Field’ (UFF) [Rappé *et al.* 1992] as a general method for calculating charge distributions over a very wide range of molecules (in principle, the entire periodic table). An additional feature of the method is that the charges are dependent upon the molecular geometry and so can change during the course of a calculation such as a molecular dynamics simulation. The starting point for this approach is a series expansion of the energy of an isolated atom in terms of the charge:

$$v_A(q) = v_{A0} + q_A \left(\frac{\partial v}{\partial q} \right)_{A0} + \frac{1}{2} q_A^2 \left(\frac{\partial^2 v}{\partial q^2} \right)_{A0} + \dots \quad (4.41)$$

Truncating this expansion after second-order terms and considering three specific states (for charges of 0, +1 and -1) leads to:

$$v_A(0) = v_{A0} \quad (4.42)$$

$$v_A(+1) = v_{A0} + q_A \left(\frac{\partial v}{\partial q} \right)_{A0} + \frac{1}{2} q_A^2 \left(\frac{\partial^2 v}{\partial q^2} \right)_{A0} \quad (4.43)$$

$$v_A(-1) = v_{A0} - q_A \left(\frac{\partial v}{\partial q} \right)_{A0} + \frac{1}{2} q_A^2 \left(\frac{\partial^2 v}{\partial q^2} \right)_{A0} \quad (4.44)$$

Now the energy of the positive species is the ionisation potential (*IP*) and the energy of the negative species is minus the electron affinity (*EA*). Combining these results gives:

$$\left(\frac{\partial v}{\partial q} \right)_{A0} = \frac{1}{2} (IP + EA) = \chi_A^0 \quad (4.45)$$

$$\left(\frac{\partial^2 v}{\partial q^2} \right)_{A0} = IP - EA \quad (4.46)$$

As usual, χ_A is the electronegativity. Rappé and Goddard suggested that for a neutral atom with a singly occupied orbital the difference between the ionisation potential and the electron affinity would correspond to the Coulomb repulsion between two electrons placed in that orbital (the orbital would be unoccupied in the positive ion and doubly occupied in the negative species). Writing this difference as J_{AA}^0 (referred to as the *idempotential*) leads to:

$$v_A(q) = v_{A0} + \chi_A^0 q_A + \frac{1}{2} J_{AA}^0 q_A^2 \quad (4.47)$$

Both the electronegativity and the idempotential can be derived from atomic data, though such atomic data generally need to be corrected for use in molecular systems. In order to use these equations to derive a set of charges for a molecule we first consider the total electrostatic energy of the system:

$$\mathcal{V}(q_1 \cdots q_N) = \sum_{i=1}^N (v_{A0} + \chi_A^0 q_A + \frac{1}{2} q_A^2 J_{AA}^0) + \sum_{A=1}^N \sum_{B=A+1}^N q_A q_B J_{AB} \quad (4.48)$$

In this equation J_{AB} represents a formulation of the Coulomb energy between charges q_A and q_B . For well-separated atoms a simple $1/r$ dependency is used. However, this simple Coulomb law is not appropriate for atoms whose charge distributions overlap. In such circumstances (which particularly arise for bonded atoms) there is a significant shielding correction. This shielding correction is a Coulomb integral (Equation (2.107)), with the atomic density being described using a single Slater type orbital whose precise form depends on the nature (ns, np or nd) of the outer valence orbital together with the covalent radius.

In order to derive the actual charges we first incorporate the factors J_{AA}^0 (the limiting value of J_{AA} as the distance tends to zero) into the double summation in Equation (4.48):

$$\mathcal{V}(q_1 \cdots q_N) = \sum_{A=1}^N (v_{A0} + \chi_A^0 q_A) + \frac{1}{2} \sum_{A=1}^N \sum_{B=1}^N q_A q_B J_{AB} \quad (4.49)$$

We can then take the derivative of the energy with respect to q_A , which leads to:

$$\frac{\partial \mathcal{V}}{\partial q_A} = \chi_A^0 + \sum_{B=1}^N q_B J_{AB} = \chi_A^0 + J_{AA}^0 q_A + \sum_{B=1; B \neq A}^N q_B J_{AB} \quad (4.50)$$

The derivative of the energy with respect to the charge is an atomic chemical potential; at equilibrium these chemical potentials will all be equal. The electrons move from regions of low electronegativity (high electrochemical potential) to regions of high electronegativity (low electrochemical potential). A further constraint is that the sum of the atomic charges must sum to the total charge on the molecule. These conditions enable a set of simultaneous equations to be written (subject to per-element limits on the charge on any given atom).

The presence of the $q_A q_B$ term with its implied distance dependency means that the charges depend upon the molecular geometry. Thus, should the conformation of a molecule change the atomic charges will also change. Just three parameters are required for each atom in the system (the electronegativity, the idempotential and the covalent radius).

4.9.7 Beyond Partial Atomic Charge Models

Most of the charge models that we have considered so far place the charge on the nuclear centres. Atom-centred charges have many advantages. For example, the electrostatic forces due to charge-charge interactions then act directly on the nuclei. This is important if one wishes to calculate the forces on the nuclei as is required for energy minimisation or a molecular dynamics simulation. Nuclear-centred charges do nevertheless suffer from some drawbacks. In particular, they assume that the charge density about each atom is spherically symmetrical. However, an atom's valence electrons are often distributed in a far from spherical manner, especially in molecules that contain features such as lone pairs and π electron clouds above aromatic ring systems.

4.9.8 Distributed Multipole Models

One way to represent the anisotropy of a molecular charge distribution is to use *distributed multipoles*. In this model, point charges, dipoles, quadrupoles and higher multipoles are distributed throughout the molecule. These distributed multipoles can be determined in various ways but the distributed multipole analysis (DMA) model of A J Stone [Stone 1981; Stone and Alderton 1985] is probably the best-known example. The DMA method calculates the multipoles from a quantum mechanics wavefunction defined in terms of Gaussian basis functions. As we saw in Section 2.6, the overlap between two Gaussian functions can be represented by another Gaussian located at a point (P) along the line that connects them. Each product of basis functions $\phi_\mu \phi_\nu$ thus corresponds to a charge density at P. This density can be expressed as a multipole expansion about P. The highest multipole moment in the local expansion depends upon the basis set used; no multipole moment higher than the sum of the angular quantum numbers of the basis set is possible. Thus, when using a basis set that contains just s and p functions there will be local multipoles no higher than the quadrupole. The crucial feature is that the local multipole expansion

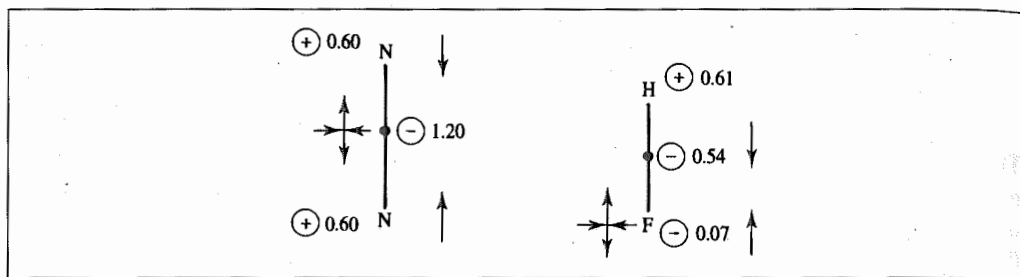


Fig. 4.22: Distributed multipole models for N_2 and HF. (Figure adapted from Stone A J and M Alderton 1985. *Distributed Multipole Analysis Methods and Applications*. Molecular Physics 56:1047-1064.)

about P can be represented as a multipole expansion about another nearby point S. In the distributed multipole approach, a set of site points is chosen and then the local multipole expansion for each pair of basis functions is 'moved' from the relevant point P to one of the sites S.

There are no limitations on the number or location of the multipole sites S; a natural set to use is obtained by placing a site point on each atomic nucleus. In some applications (especially for small molecules) additional sites are defined at the centres of bonds. For example, Stone derived a distributed multipole model for nitrogen from a Dunning [5s4p2d] basis set with two polarisation functions. This model contains charges of +0.60 on the nuclei and a charge of -1.20 at the centre of the bond, together with a dipole on each of the two nuclei and a quadrupole located at the centre of the bond (see Figure 4.22). For HF charges are placed on the two nuclei and at the centre of the bond with a dipole and a quadrupole on the fluorine and a small dipole at the centre of the bond (Figure 4.22). In larger molecules not every atom may be given a site, such as hydrogen atoms bonded to apolar atoms. It is also possible to restrict the order of the multipole expansion at a given atom so that, for example, only a charge component would be present on a polar hydrogen with the higher moments being represented by multipoles on the atom to which it is bonded. An important consideration when choosing the multipole sites is that, when a local multipole expansion is moved, the resulting multipole expansion is no longer a truncated series. However, the smaller the distance between P and the corresponding site point S, the quicker the series converges. In practice, therefore, each local multipole moment expansion is either moved to the nearest site point or is divided between the two nearest site points when they are equally close. With a basis set that contains just s and p functions and multipole sites at the atomic nuclei, it is usually found that the distributed multipole series converges rapidly after the quadrupole term. The multipoles themselves can vary considerably with the basis set used to perform the *ab initio* calculation, but the various electronic properties derived from them usually do not change much.

The distributed multipole model automatically includes non-spherical, anisotropic effects due to features such as lone pairs or π electrons. The original applications of the DMA approach were to small molecules such as diatomics and triatomics. The method has since been used to develop models for nuclei acids and for peptides and has even been applied to the undecapeptide cyclosporin [Price *et al.* 1989], which contains 199 atoms (the

quantum mechanical calculation on this molecule used 1000 basis functions). However, distributed multipole models have not yet been widely incorporated into force fields, not least because of the additional computational effort required. It can be complicated to calculate the atomic forces with the distributed multipole model; in particular, multipoles that are not located on atoms generate torques, which must be analysed further to determine the forces on the nuclei.

4.9.9 Using Charge Schemes to Study Aromatic-Aromatic Interactions

The attractive interactions between molecules containing π systems have long been studied by theoreticians and experimentalists. Such systems are involved in a variety of phenomena, including the stacking of the nucleic acid bases in DNA, the packing of aromatic molecules in crystals and interactions between amino acid side chains in proteins. A variety of orientations are observed for aromatic dimers, ranging from edge-on, T-shaped structures to face-to-face structures (Figure 4.23). Within these two families the molecules can move relative to each other, so that, for example, in a face-to-face arrangement the atoms are overlaid or are staggered. In the T-shaped structure the large quadrupole moments of the benzene molecules adopt their most favourable orientation.

One very simple model of the interactions in such systems was devised by Hunter and Saunders [Hunter and Saunders 1990], who wanted to explain the stacking behaviour of aromatic systems such as the porphyrins shown in Figure 4.24. It is experimentally observed that these molecules adopt a cofacial arrangement with their centres offset as shown. Hunter and Saunders placed point charges not only at the nuclei but also at locations above and below each atom, perpendicular to the plane of the ring. Thus in benzene each carbon atom was given a charge of +1 and also had two associated charges of $-\frac{1}{2}$ above and below the ring (Figure 4.25). The electrostatic interaction between two ring systems is

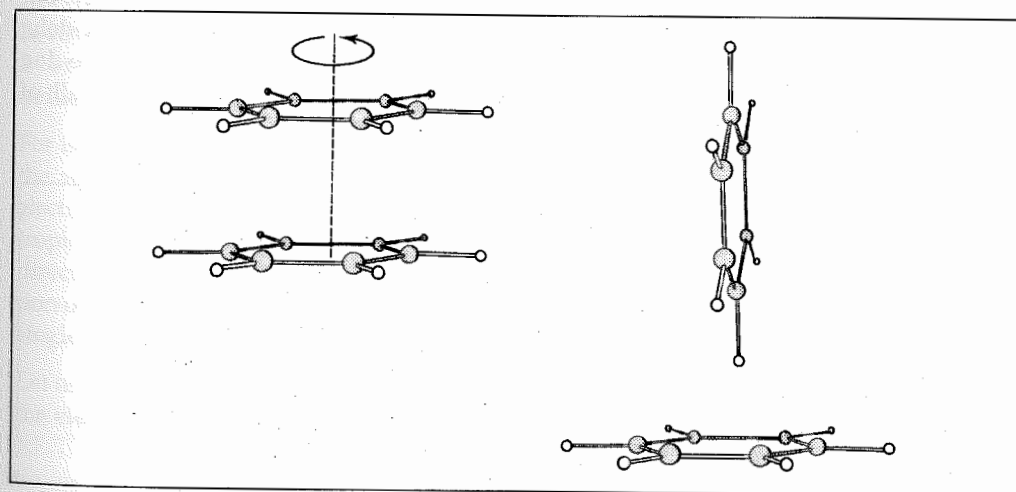


Fig. 4.23: Face-to-face (left) and T-shaped (right) orientations of the benzene dimer.

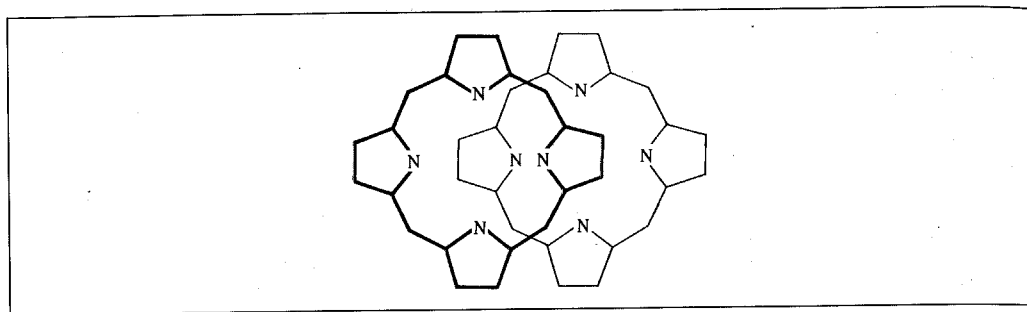


Fig. 4.24: Porphyrin system typical of those studied by Hunter and Saunders [Hunter and Saunders 1990].

calculated in the usual way by summing the charge-charge interactions using Coulomb's law. A major advantage of the Hunter-Saunders approach is its computational simplicity. Moreover, it can be extended to cover a wide range of atom types and so applied to many systems [Vinter 1994] with particular emphasis on simulating DNA [Hunter 1993, Packer *et al.* 2000]. Hunter and Saunders summarised the results of their investigations on porphyrins in three rules:

1. π - π repulsion dominates in a face-to-face geometry;
2. π - σ attraction dominates in an edge-on geometry;
3. π - σ attraction dominates in an offset π -stacked geometry.

The interactions between aromatic systems have also been studied using point charge models, central multipoles and distributed multipoles. Fowler and Buckingham examined homodimers of *sym*-triazine and 1,3,5-trifluorobenzene (Figure 4.26) [Fowler and Buckingham 1991]. They were particularly keen to calculate how the electrostatic energy changed as the rings were twisted in the face-to-face geometry. All but one of the energy models suggested that the staggered orientations were the arrangements of minimum energy, but the energy difference between the eclipsed and staggered structures varied widely, depending upon the model. The central multipole model was found to be ineffective due to convergence problems. Three different point-charge models were considered, all of

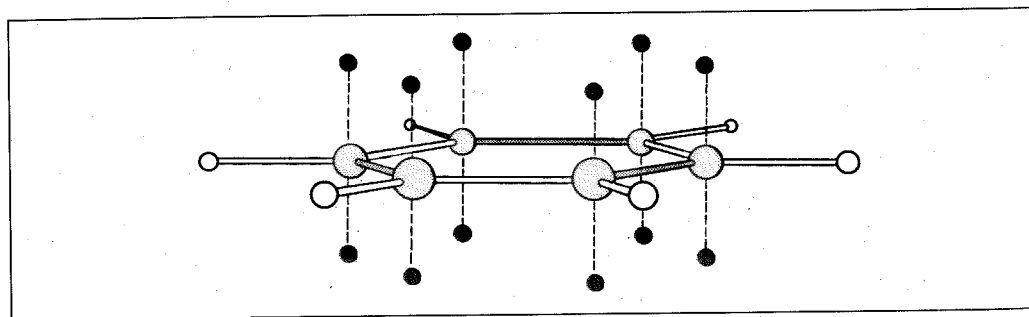


Fig. 4.25: Anisotropic model of benzene developed by Hunter and Saunders [Hunter and Saunders 1990].

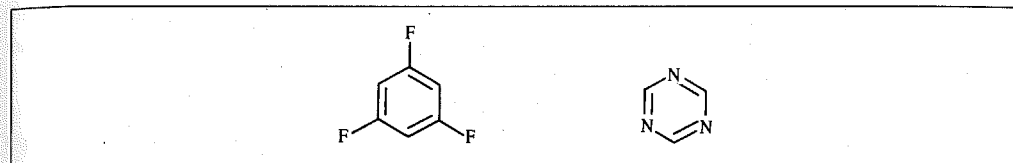


Fig. 4.26: *Sym*-triazine and 1,3,5-trifluorobenzene.

which gave acceptable energy curves. The distributed multipole model also performed well, being comparable to the most accurate of the point-charge models.

4.9.10 Polarisation

Our discussion of electronic effects has concentrated so far on 'permanent' features of the charge distribution. Electrostatic interactions also arise from changes in the charge distribution of a molecule or atom caused by an external field, a process called *polarisation*. The primary effect of the external electric field (which in our case will be caused by neighbouring molecules) is to induce a dipole in the molecule. The magnitude of the induced dipole moment μ_{ind} is proportional to the electric field E , with the constant of proportionality being the polarisability α :

$$\mu_{\text{ind}} = \alpha E \quad (4.51)$$

The energy of interaction between a dipole μ_{ind} and an electric field E (the induction energy) is determined by calculating the work done in charging the field from zero to E , using the following integral:

$$v(\alpha, E) = - \int_0^E dE \mu_{\text{ind}} = - \int_0^E dE \alpha E = -\frac{1}{2} \alpha E^2 \quad (4.52)$$

In strong electric fields contributions to the induced dipole moment that are proportional to E^2 or E^3 can also be important, and higher-order moments such as quadrupoles can also be induced. We will not be concerned with such contributions.

For isolated atoms, the polarisability is isotropic - it does not depend on the orientation of the atom with respect to the applied field, and the induced dipole is in the direction of the electric field, as in Equation (4.51). However, the polarisability of a molecule is often anisotropic. This means that the orientation of the induced dipole is not necessarily in the same direction as the electric field. The polarisability of a molecule is often modelled as a collection of isotropically polarisable atoms. A small molecule may alternatively be modelled as a single isotropic polarisable centre.

Let us consider the electric field due to a dipole μ aligned along the z axis. The magnitude of the electric field at a point P due to the dipole (see Figure 4.27) is:

$$E(r, \theta) = \frac{\mu \sqrt{1 + 3 \cos^2 \theta}}{4\pi\epsilon_0 r^3} \quad (4.53)$$

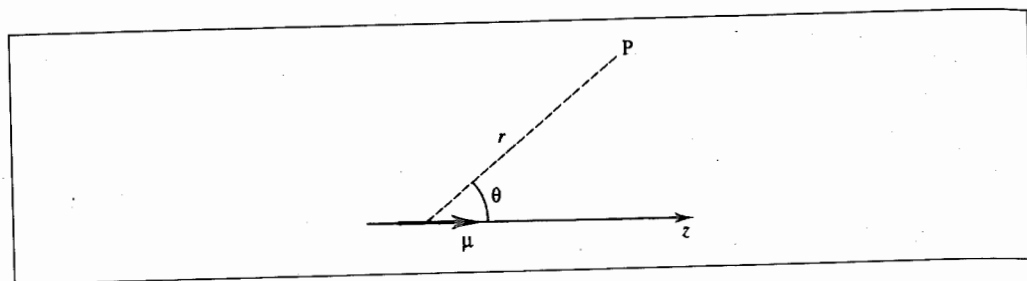


Fig. 4.27: Electric field at point P due to dipole at the origin.

The induction energy with another molecule of polarisability α placed at P is therefore

$$v(r, \theta) = -\alpha \mu^2 \frac{1 + 3 \cos^2 \theta}{(4\pi\epsilon_0 r^3)^2} \quad (4.54)$$

The interaction between a dipole and an induced dipole is independent of the disorienting effect of thermal motion, whereas the dipole-dipole interaction between two permanent dipoles does vary with the relative orientation of the two dipoles. This is because the induced dipole follows the direction of the permanent dipole even as the molecules change their orientations as a consequence of molecular collisions.

An important consideration when modelling polarisation effects is that the dipole induced on a molecule (A) will affect the charge distribution of another molecule (B). The electric field at A due to the dipole(s) on B will in turn be affected. The presence of other molecules can also influence the interaction. Consider the polarisation interaction between a polar molecule and a neighbour (Figure 4.28). A third molecule may reduce the size of the electric field on the second molecule and so lower the induction energy. This type of three-body effect will be particularly significant when polarisable atoms are close to polar groups. Polarisation is a cooperative effect and, as such, is modelled using a set of coupled equations which are typically solved iteratively. Initially, the induced dipoles are set to zero. An initial approximation to each induced dipole is then calculated from the permanent charges (i.e. partial atomic charges). The electric field due to these induced dipoles is then added to the electric field from the permanent charges. This gives a refined value of the electric field from which a new induced dipole can be determined. The calculation continues until the induced dipoles do not change significantly between iterations.

A variety of schemes for including polarisation into molecular mechanics force fields have been devised. One approach is to model the polarisation effects at the atomic level, with

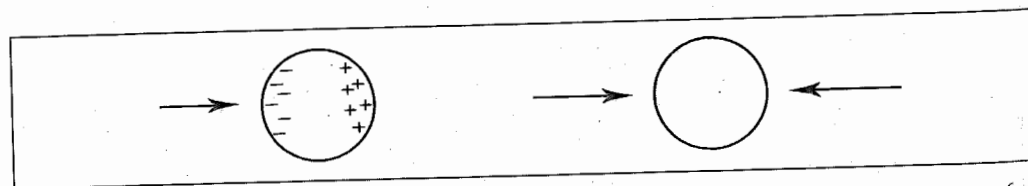


Fig. 4.28: The polarisation interaction between a dipole and a polarisable molecule can be affected by the presence of a second dipole (right) and is therefore a many-body effect.

dipoles being induced on each atom [Dang *et al.* 1991]. The magnitude of the dipole induced on an atom i is given by:

$$\mu_{\text{ind},i} = \alpha_i E_i \quad (4.55)$$

α_i is the atomic polarisability, assumed to be isotropic. Appropriate values of α_i have been determined for various systems. The electric field, E_i , at atom i is the vector sum of the field due to the permanent and induced dipoles of the other atoms in the system:

$$E_i = \sum_{j \neq i} \frac{q_j \mathbf{r}_{ij}}{r_{ij}^3} + \sum_{j \neq i} \frac{\mu_j}{r_{ij}^3} \left(3 \frac{\mathbf{r}_{ij} \mathbf{r}_{ij}^T}{r_{ij}^2} - 1 \right) \quad (4.56)$$

\mathbf{r}_i and \mathbf{r}_j are the position vectors of the atoms i and j . Convergence of these equations in procedures such as molecular dynamics, where successive configurations are generated, can be accelerated if the induced dipoles obtained at each current step are used as the starting points for the next configuration.

An alternative way to model polarisation effects is exemplified by the water model of Sprik and Klein [Sprik and Klein 1988], where the polarisation centre is represented as a collection of closely spaced charges whose values are permitted to vary but whose total sums to zero. In the water model, shown in Figure 4.29, four tetrahedrally arranged charges are used to model the polarisation centre. These charges endow the molecule with an induced dipole moment of any magnitude and direction. The charges are determined iteratively for each configuration of the system. The isotropic polarisability of a simple ion can similarly be treated using two charges of equal magnitude but opposite sign placed either side of the ion. The direction of the 'bond' linking the two polarisation charges and the ion can reorient to change the direction of the induced dipole. In a subsequent refinement of this model Sprik and Klein replaced the point charges by Gaussian charge distributions at the polarisation sites; these were better at modelling features such as hydrogen bonding.

One appealing approach is the dynamically fluctuating charge model of Berne and colleagues [Rick *et al.* 1994]. This method has much in common with the charge equilibration scheme of Rappé and Goddard (see Section 4.9.6) in its use of the electronegativity equalisation approach, which ensures that the atomic chemical potentials are equal in the molecule. The charges are considered as dynamically fluctuating variables, along with the atomic nuclei in a molecular dynamics simulation. This means that the charges evolve in a natural

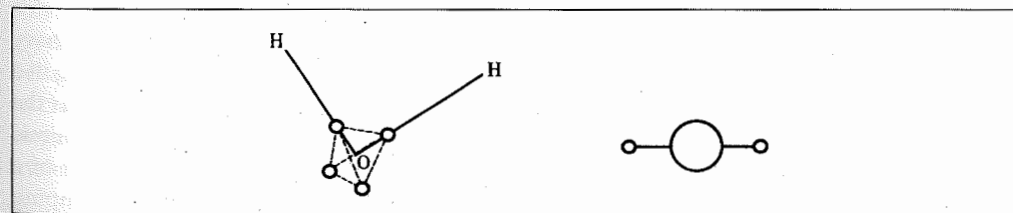


Fig. 4.29: Polarisable models of water and ions developed by Sprik and Klein. (Figure adapted from Sprik M 1993. *Effective Pair Potentials and Beyond*. In *Computer Simulation in Chemical Physics*, Allen M P, D J Tildesley (Editors). Dordrecht, Kluwer).

manner during the simulation rather than having to determine a new set of charges at each iteration of the procedure. This fluctuating charge model includes intramolecular interactions and so the traditional Coulombic $1/r$ expression is not appropriate. Rather, the charges are replaced by charge distributions (formulated as Slater s orbitals) whose interaction is calculated using a Coulomb integral expression. This interaction is effectively identical to the standard Coulomb expression for intermolecular interactions, only differing for the intramolecular contribution.

One feature of this oscillating charge model is that it requires rather less computational effort than traditional polarisation models. It also implicitly preserves the higher-order multipole terms, which need to be explicitly incorporated in some of the alternative approaches. Ions are represented by two partial charges (which sum to the required integral ionic charge) which are connected by a harmonic spring. The mass of one of these two species is made much greater than the other so that the heavier site remains near the centre of mass as the spring oscillates. This particular model has been used for simulations of pure liquid water [Rick *et al.* 1994], the solvation of amides [Rick and Berne 1996] and to investigate the effects of polarisability on the hydration of the chloride ion in water clusters [Stuart and Berne 1996]. These calculations predicted that the chloride ions were located on the outside of the clusters, even when they contained more than 100 water molecules. This was in contrast to equivalent calculations using a non-polarisable model, the difference being attributed to the presence of fluctuations in the dipole strengths of the water molecules in the cluster, which are, as a consequence, more mobile.

Due to the computational expense, polarisation effects are often included in a calculation only when their effect is likely to be significant, such as simulations of ionic solutions. These systems usually contain atoms or ions and small molecules only. It is important to be aware of the following problem when using atomic polarisabilities. Consider a diatomic molecule. The application of an external field will induce dipoles on both atoms. The dipole on one atom will also contribute to the electric field at the other atom, and thereby influence its induced dipole, but the model takes no account of the fact that the charge distributions on the two atoms are inherently linked. For this reason (and for reasons of computational efficiency) it is common to treat small molecules such as water as single polarisable centres when calculating polarisation effects.

4.9.11 Solvent Dielectric Models

All of the formulae that we have written for electrostatic energies, potentials and forces have included the permittivity of free space, ϵ_0 . This is as one would expect for species acting in a vacuum. However, under some circumstances a different dielectric model is used in the equations for the electrostatic interactions. This is often done when it is desired to mimic solvent effects, without actually including any explicit solvent molecules. One effect of a solvent is to dampen the electrostatic interactions. A very simple way to model this damping effect is to increase the permittivity, most easily by using an appropriate value for the relative permittivity in the Coulomb's law equation (i.e. $\epsilon = \epsilon_0 \epsilon_r$). An alternative approach is to make the dielectric dependent upon the separation of the charged species; this gives rise

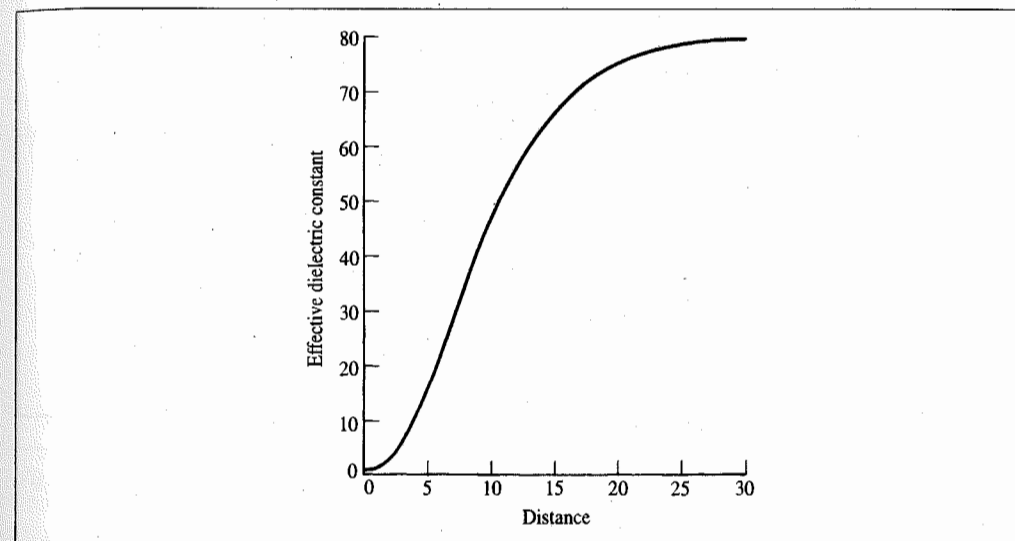


Fig. 4.30: A sigmoidal dielectric model smoothly varies the effective permittivity from 80 to 1 as shown.

to the so-called distance-dependent dielectric models. The simplest implementation of a distance-dependent dielectric is to make the relative permittivity proportional to the distance. The interaction energy between two charges q_i and q_j then becomes:

$$v(r) = \frac{1}{4\pi\epsilon_0} \frac{q_i q_j}{r^2} \quad (4.57)$$

The simple distance-dependent dielectric has no physical basis and so it is not generally recommended, except when no alternative is possible. More sophisticated distance-dependent functions can also be employed. Many of these have an approximately sigmoidal shape in which the relative permittivity is low at short distances and then rises towards the bulk value at long distances. One example of such a function is [Smith and Pettit 1994]:

$$\epsilon_{\text{eff}}(r) = \epsilon_r - \frac{\epsilon_r - 1}{2} [(rS)^2 + 2rS + 2] e^{-rS} \quad (4.58)$$

The value of ϵ_{eff} varies from a value of 1 at zero separation to ϵ_r (the bulk permittivity of the solvent) at large distances, in a manner determined by the parameter S (which is typically given a value between 0.15 \AA^{-1} and 0.3 \AA^{-1} ; Figure 4.30). Sigmoidal functions give better behaviour than the simple distance-dependent dielectric model. However, it may be difficult to choose the appropriate value for the bulk dielectric ϵ_r , when performing calculations on large solutes, as the shortest distance between two charges may be through the solute molecule rather than through the solvent (Figure 4.31).

The polarisation term can be a major contributor to the free energy of solvation of a solute, and a variety of schemes have been devised to incorporate such effects where the solvent is modelled as a continuum. We shall discuss these methods in more detail in Sections 11.9–11.12.

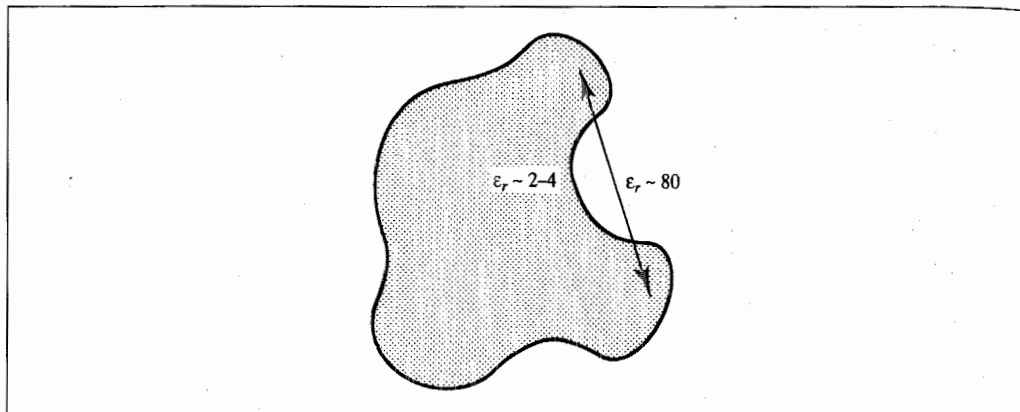


Fig. 4.31: A line joining two points may pass through regions of different permittivity.

4.10 Van der Waals Interactions

Electrostatic interactions cannot account for all of the non-bonded interactions in a system. The rare gas atoms are an obvious example; all of the multipole moments of a rare gas atom are zero and so there can be no dipole-dipole or dipole-induced dipole interactions. But there clearly must be interactions between the atoms; how else could rare gases have liquid and solid phases or show deviations from ideal gas behaviour? Deviations from ideal gas behaviour were famously quantitated by van der Waals, thus the forces that give rise to such deviations are often referred to as van der Waals forces.

If we were to study the interaction between two isolated argon atoms using a molecular beam experiment then we would find that the interaction energy varies with the separation in a manner as shown in Figure 4.32. The other rare gases show a similar behaviour. The essential features of this curve are as follows. The interaction energy is zero at infinite distance (and indeed is negligible even at relatively short distances). As the separation is reduced, the energy decreases, passing through a minimum at a distance of approximately 3.8 \AA for argon. The energy then rapidly increases as the separation decreases further. The force between the atoms, which equals minus the first derivative of the potential energy with respect to distance, is also shown in Figure 4.32. A variety of experiments have been used to provide evidence for the nature of the van der Waals interactions, including gas imperfections, molecular beams, spectroscopic studies and measurements of transport properties.

4.10.1 Dispersive Interactions

The curve in Figure 4.32 is usually considered to arise from a balance between attractive and repulsive forces. The attractive forces are long-range, whereas the repulsive forces act at short distances. The attractive contribution is due to *dispersive forces*. London first showed how the dispersive force could be explained using quantum mechanics [London 1930] and so this interaction is sometimes referred to as the London force. The dispersive force

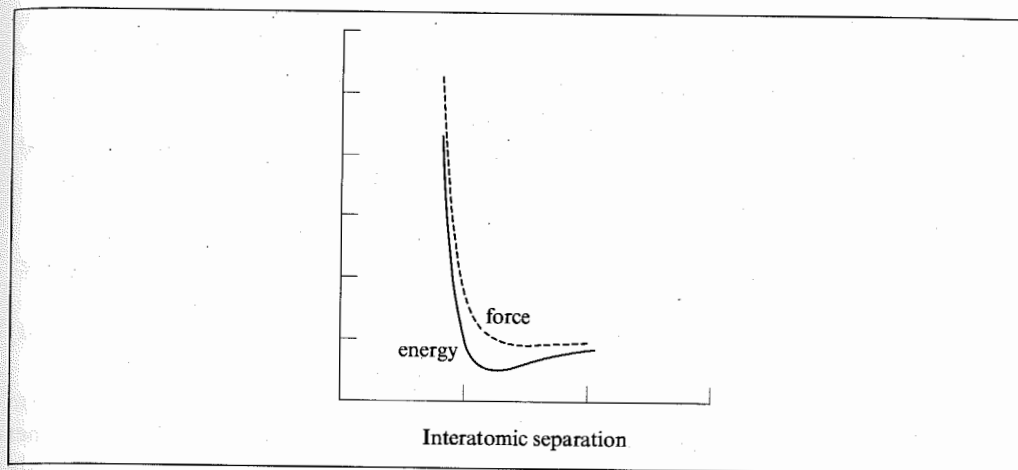


Fig. 4.32: The interaction energy and the force between two argon atoms.

is due to instantaneous dipoles which arise during the fluctuations in the electron clouds. An instantaneous dipole in a molecule can in turn induce a dipole in neighbouring atoms, giving rise to an attractive inductive effect.

A simple model to explain the dispersive interaction was proposed by Drude. This model consists of 'molecules' with two charges, $+q$ and $-q$, separated by a distance r . The negative charge performs simple harmonic motion with angular frequency ω along the z axis about the stationary positive charge (Figure 4.33). If the force constant for the oscillator is k and if the mass of the oscillating charge is m , then the potential energy of an isolated Drude molecule is $\frac{1}{2}kz^2$, where z is the separation of the two charges. ω is related to the force constant by $\omega = \sqrt{k/m}$. The Schrödinger equation for a Drude molecule is:

$$-\frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial z^2} + \frac{1}{2}kz^2 \psi = E\psi \quad (4.59)$$

This is the Schrödinger equation for a simple harmonic oscillator. The energies of the system are given by $E_\nu = (\nu + \frac{1}{2}) \times \hbar\omega$ and the zero-point energy is $\frac{1}{2}\hbar\omega$.

We now introduce a second Drude molecule, identical to the first, with the positive charge also located on the z axis and an oscillating negative charge (Figure 4.33). When the two molecules are infinitely separated, they do not interact and the total ground-state energy of the system is

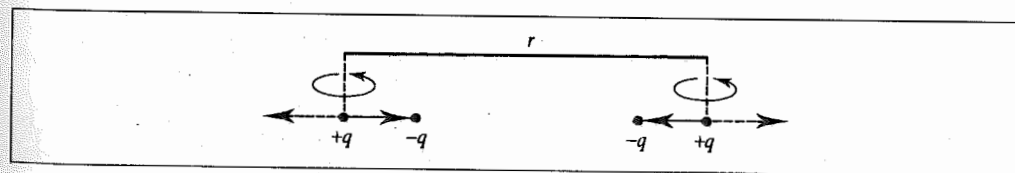


Fig. 4.33: The Drude model for dispersive interactions. (Figure adapted from Rigby M, E B Smith, W A Wakeham and G C Maitland 1986. *The Forces Between Molecules*. Oxford, Clarendon Press.)

just twice the zero-point energy of a single molecule, $\hbar\omega/2\pi$. As the molecules approach (along the z axis) there are interactions between the two dipoles, and the interaction energy between the two 'molecules' can be shown to be approximately given by (see Appendix 4.1):

$$v(r) = -\frac{\alpha^4 \hbar\omega}{2(4\pi\epsilon_0)^2 r^6} \quad (4.60)$$

The Drude model thus predicts that the dispersion interaction varies as $1/r^6$.

The two-dimensional Drude model can be extended to three dimensions, the result being:

$$v(r) = -\frac{3\alpha^4 \hbar\omega}{4(4\pi\epsilon_0)^2 r^6} \quad (4.61)$$

The Drude model only considers the dipole-dipole interaction; if higher-order terms, due to dipole-quadrupole, quadrupole-quadrupole, etc., interactions are included as well as other terms in the binomial expansion, then the energy of the Drude model is more properly written as a series expansion:

$$v(r) = \frac{C_6}{r^6} + \frac{C_8}{r^8} + \frac{C_{10}}{r^{10}} + \dots \quad (4.62)$$

All of the coefficients C_n are negative, implying an attractive interaction. Despite its simplicity, the Drude model gives quite reasonable results; if just the C_6 term is included then for argon the resulting dispersion energy is only about 25% too small.

4.10.2 The Repulsive Contribution

Below about 3 \AA , even a small decrease in the separation between a pair of argon atoms causes a large increase in the energy. This increase has a quantum mechanical origin and can be understood in terms of the Pauli principle, which formally prohibits any two electrons in a system from having the same set of quantum numbers. The interaction is due to electrons with the same spin, therefore the short-range repulsive forces are often referred to as *exchange forces*. They are also known as overlap forces. The effect of exchange is to reduce the electrostatic repulsion between pairs of electrons by forbidding them to occupy the same region of space (i.e. the internuclear region). The reduced electron density in the internuclear region leads to repulsion between the incompletely shielded nuclei. At very short internuclear separations, the interaction energy varies as $1/r$ due to this nuclear repulsion, but at larger separations the energy decays exponentially, as $\exp(-2r/a_0)$, where a_0 is the Bohr radius.

4.10.3 Modelling Van der Waals Interactions

The dispersive and exchange-repulsive interactions between atoms and molecules can be calculated using quantum mechanics, though such calculations are far from trivial, requiring electron correlation and large basis sets. For a force field we require a means to model the interatomic potential curve accurately (Figure 4.32), using a simple empirical

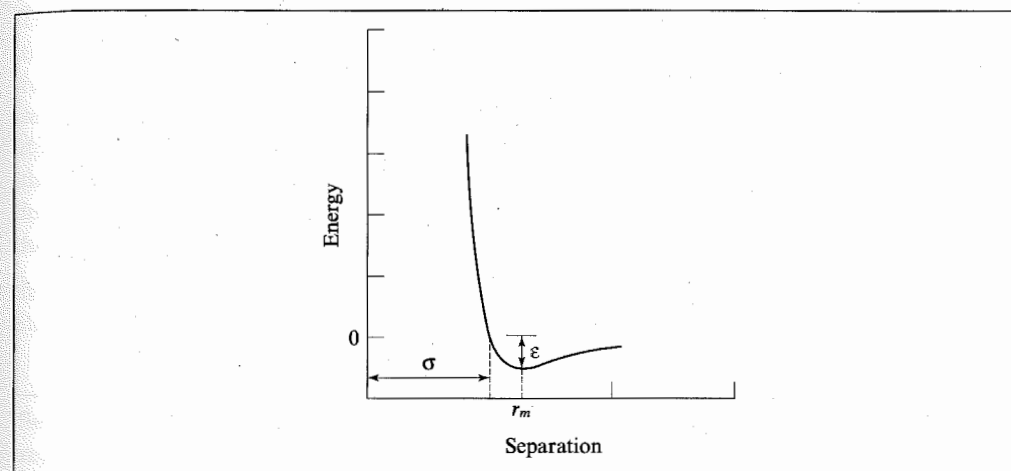


Fig. 4.34: The Lennard-Jones potential.

expression that can be rapidly calculated. The need for a function that can be rapidly evaluated is a consequence of the large number of van der Waals interactions that must be determined in many of the systems that we would like to model. The best known of the van der Waals potential functions is the *Lennard-Jones 12-6 function*, which takes the following form for the interaction between two atoms:

$$v(r) = 4\epsilon \left[\left(\frac{\sigma}{r} \right)^{12} - \left(\frac{\sigma}{r} \right)^6 \right] \quad (4.63)$$

The Lennard-Jones 12-6 potential contains just two adjustable parameters: the collision diameter σ (the separation for which the energy is zero) and the well depth ϵ . These parameters are graphically illustrated in Figure 4.34. The Lennard-Jones equation may also be expressed in terms of the separation at which the energy passes through a minimum, r_m (also written r^*). At this separation, the first derivative of the energy with respect to the internuclear distance is zero (i.e. $\partial v/\partial r = 0$), from which it can easily be shown that $r_m = 2^{1/6}\sigma$. We can thus also write the Lennard-Jones 12-6 potential function as follows:

$$v(r) = \epsilon \{ (r_m/r)^{12} - 2(r_m/r)^6 \} \quad (4.64)$$

or

$$v(r) = A/r^{12} - C/r^6 \quad (4.65)$$

A is equal to ϵr_m^{12} (or $4\epsilon\sigma^{12}$) and C is equal to $2\epsilon r_m^6$ (or $4\epsilon\sigma^6$).

The Lennard-Jones potential is characterised by an attractive part that varies as r^{-6} and a repulsive part that varies as r^{-12} . These two components are drawn in Figure 4.35. The r^{-6} variation is of course the same power-law relationship found for the leading term in theoretical treatments of the dispersion energy such as the Drude model. There are no strong theoretical arguments in favour of the repulsive r^{-12} , especially as quantum

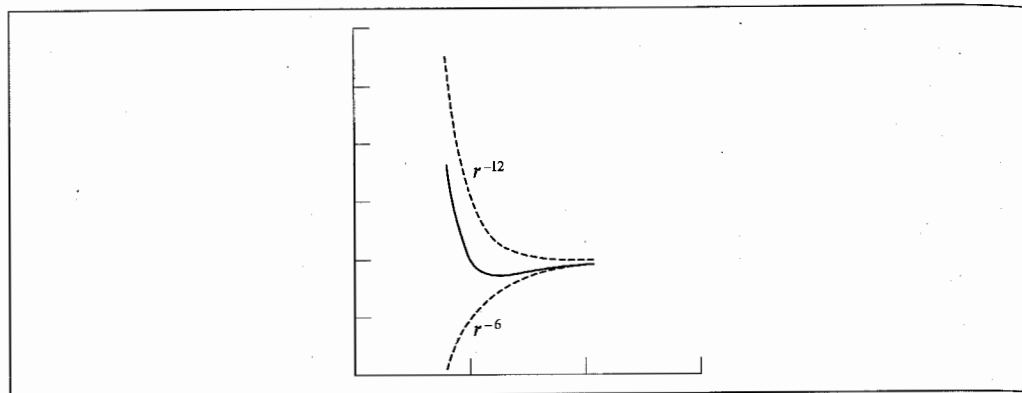


Fig. 4.35: The Lennard-Jones potential is constructed from a repulsive component (αr^{-12}) and an attractive component (αr^{-6}).

mechanics calculations suggest an exponential form. The twelfth power term is found to be quite reasonable for rare gases but is rather too steep for other systems such as hydrocarbons. However, the 6-12 potential is widely used, particularly for calculations on large systems, as r^{-12} can be rapidly calculated by squaring the r^{-6} term. The r^{-6} term can also be calculated from the square of the distance without having to perform a computationally expensive square root calculation. Different powers have also been used for the repulsive part of the potential; values of 9 or 10 give a less steep curve and are used in some force fields. Lennard-Jones' original potential has been written in the following general form:

$$v(r) = k\varepsilon \left[\left(\frac{\sigma}{r} \right)^n - \left(\frac{\sigma}{r} \right)^m \right]; \quad k = \frac{n}{n-m} \left(\frac{n}{m} \right)^{m/(n-m)} \quad (4.66)$$

Equation (4.66) returns the Lennard-Jones potential for $n = 12$ and $m = 6$.

Halgren has proposed an alternative functional form designed to be simple enough to be easily incorporated into molecular mechanics calculations whilst also improving the ability to reproduce experimental data [Halgren 1992, 1996a, b]. In this sense it is an attempt to improve on the Lennard-Jones potential without introducing the complexity of some of the potentials employed by spectroscopists. This potential has the general form:

$$v(r) = \varepsilon_{ij} \left(\frac{1 + \delta}{\rho_{ij} + \delta} \right)^{(n-m)} \left(\frac{1 + \gamma}{\rho_{ij}^m + \gamma} - 2 \right) \quad (4.67)$$

In this equation $\rho_{ij} = r_{ij}/r_{ij}^*$. The constants δ and γ apply to all interactions between the atoms i and j . This potential reduces to the standard Lennard-Jones 12-6 potential if the following choice of parameters is used: $n = 12$, $m = 6$, $\delta = \gamma = 0$. Halgren proposed a 'buffered 14-7' potential in which $n = 14$, $m = 7$, $\delta = 0.07$ and $\gamma = 0.12$, giving the following equation:

$$v(r) = \varepsilon_{ij} \left(\frac{1.07r_{ij}^*}{r_{ij} + 0.07r_{ij}^*} \right)^7 \left(\frac{1.12r_{ij}^{*7}}{r_{ij}^7 + 0.12r_{ij}^{*7}} \right) \quad (4.68)$$

There were several reasons for developing this functional form. First was the desire to keep the potential finite as the interatomic potential approaches zero (unlike the Lennard-Jones function, which becomes infinite). Second, it gives a more accurate reproduction of the series expansion for the dispersion interaction, Equation (4.62). Third, if a larger value of d is used then the repulsive component is greatly reduced without significantly changing the distance at which the potential crosses zero or the depth of the energy minimum. This feature is useful for optimising structures with crude initial geometries; other functional forms can have significant problems with such situations.

In the buffered 14-7 potential the minimum-energy separation r_{ii}^* for an atom i depends on its atomic polarisability:

$$r_{ii}^* = A_i \alpha_i^{1/4} \quad (4.69)$$

Several formulations in which the r^{-12} term in the standard Lennard-Jones formulation is replaced by a theoretically more realistic exponential expression have been proposed. These include the *Buckingham potential*:

$$v(r) = \varepsilon \left[\frac{6}{\alpha - 6} \exp[-\alpha(r/r_m - 1)] - \frac{\alpha}{\alpha - 6} \left(\frac{r_m}{r} \right)^6 \right] \quad (4.70)$$

There are three adjustable parameters in the Buckingham potential (ε , r_m and α). A value of α between approximately 14 and 15 gives a potential that closely corresponds to the Lennard-Jones 12-6 potential in the minimum-energy region. When using the Buckingham potential it is important to remember that at very short distances the potential becomes strongly attractive, as shown in Figure 4.36. This could lead to nuclei being fused together during a calculation, and so the program must check that atoms are not becoming too close. The

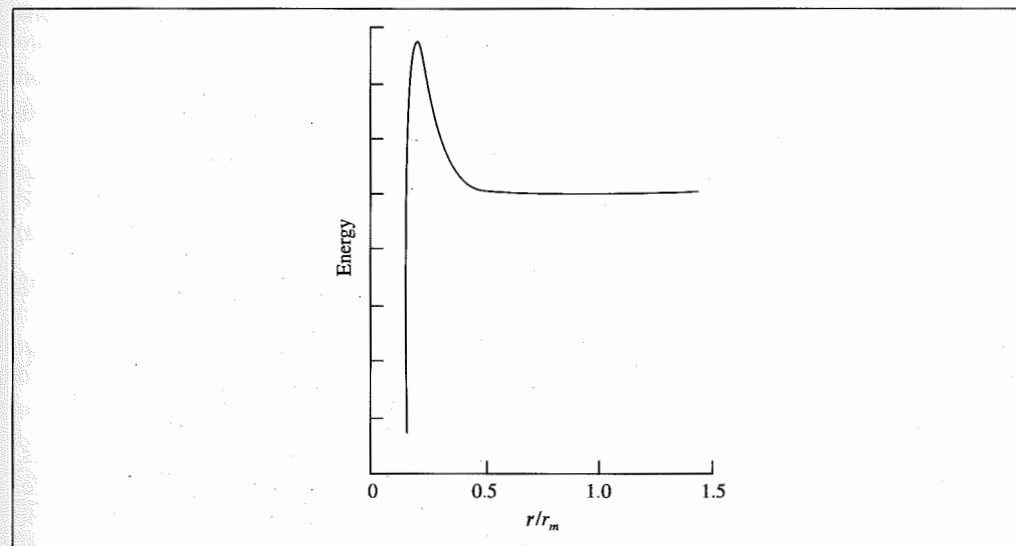


Fig. 4.36: A drawback of the Buckingham potential is that it becomes steeply attractive at short distances.

atoms come close together some redistribution of the charge along the connecting bonds would be expected that would act to reduce the interaction. Such a charge redistribution would not be possible for two atoms at a similar distance apart if they were in different molecules.

The parameters for the van der Waals interactions can be obtained in a variety of ways. In the early force fields, such parameters were often determined from an analysis of crystal packing. The objective of such studies was to produce a set of van der Waals parameters which enabled the experimental geometries and thermodynamic properties such as the heat of sublimation to be reproduced as accurately as possible. More recent force fields derive their van der Waals parameters using liquid simulations in which the parameters are optimised to reproduce a range of thermodynamic properties such as the densities and enthalpies of vaporisation for appropriate liquids.

4.10.5 Reduced Units

The Lennard-Jones potential is completely specified by the two parameters ϵ and σ . This means that the results of a calculation performed on (say) liquid argon can be easily converted to give equivalent results for another noble gas. For this reason it is common to simulate the rare gases in terms of reduced units with ϵ and σ both set to 1. The results can then be converted to any system as appropriate. For example, the reduced density ρ^* is related to the real density by $\rho^* = \rho\sigma^3$; the reduced energy E^* is given by $E^* = E/\epsilon$, and so on. Electrostatic interactions given by Coulomb's law are also often written in terms of a reduced unit of charge, which corresponds to each charge being divided by $\sqrt{4\pi\epsilon_0}$. This means that Coulomb's law takes the less cumbersome form:

$$v(q_1, q_2) = q_1 q_2 / r_{12} \quad \text{or} \quad v(q_1, q_2) = q_1 q_2 / \epsilon_r r_{12} \quad (4.79)$$

4.11 Many-body Effects in Empirical Potentials

The electrostatic and van der Waals energies that we have considered so far are calculated between pairs of interaction sites. The total non-bonded interaction energy is thus determined by adding together the interactions between all pairs of sites in the system. However, the interaction between two molecules can be affected by the presence of a third, fourth or more molecules. For example, the interaction energy between three molecules A, B and C is not in general given by the sum of the pairwise interaction energies: $v(A, B, C) \neq v(A, B) + v(A, C) + v(B, C)$. We have already seen an example of a non-pairwise contribution, namely the polarisation interaction, which is determined using a self-consistent procedure.

Three-body effects can significantly affect the dispersion interaction. For example, it is believed that three-body interactions account for approximately 10% of the lattice energy of crystalline argon. For very precise work, interactions involving more than three atoms may have to be taken into account, but they are usually small enough to be ignored. A potential that includes both two- and three-body interactions would be written in the following

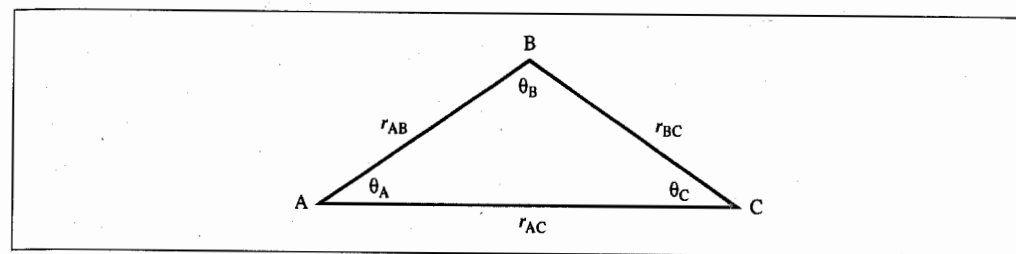


Fig. 4.37: Calculating the three-body Axilrod-Teller contribution.

general form:

$$\mathcal{V}(\mathbf{r}^N) = \sum_{i=1}^N \sum_{j=i+1}^N v^{(2)}(r_{ij}) + \sum_{i=1}^N \sum_{j=i+1}^N \sum_{k=j+1}^N v^{(3)}(r_{ij}, r_{ik}, r_{jk}) \quad (4.80)$$

Axilrod and Teller investigated the three-body dispersion contribution and showed that the leading term is:

$$v^{(3)}(r_{AB}, r_{AC}, r_{BC}) = \nu_{A,B,C} \frac{3 \cos \theta_A \cos \theta_B \cos \theta_C}{(r_{AB} r_{AC} r_{BC})^3} \quad (4.81)$$

θ_A , θ_B and θ_C are the internal angles of the triangle with sides of length r_{AB} , r_{AC} and r_{BC} (Figure 4.37). $\nu_{A,B,C}$ is a constant characteristic of the three species A, B and C. If A, B and C are identical then $\nu_{A,B,C}$ is approximately related to the Lennard-Jones coefficient C_6 and the polarisability by

$$\nu_{A,B,C} = -\frac{3\alpha C_6}{4(4\pi\epsilon_0)} \quad (4.82)$$

The effect of the Axilrod-Teller term (also known as the triple-dipole correction) is to make the interaction energy more negative when three molecules are linear but to weaken it when the molecules form an equilateral triangle. This is because the linear arrangement enhances the correlations of the motions of the electrons, whereas the equilateral arrangement reduces it.

The three-body contribution may also be modelled using a term of the form $v^{(3)}(r_{AB}, r_{AC}, r_{BC}) = K_{A,B,C} \{\exp(-\alpha r_{AB}) \exp(-\beta r_{AC}) \exp(-\gamma r_{BC})\}$ where K , α , β and γ are constants describing the interaction between the atoms A, B and C. Such a functional form has been used in simulations of ion-water systems, where polarisation alone does not exactly model configurations when there are two water molecules close to an ion [Lybrand and Kollman 1985]. The three-body exchange repulsion term is thus only calculated for ion-water-water trimers when the species are close together.

The computational effort is significantly increased if three-body terms are included in the model. Even with a simple pairwise model, the non-bonded interactions usually require by far the greatest amount of computational effort. The number of bond, angle and torsional terms increases approximately with the number of atoms (N) in the system, but the number of non-bonded interactions increases with N^2 . There are $N(N-1)/2$ distinct pairs of

interactions to evaluate for a pairwise potential. If three-body effects are included then there are $N(N-1)(N-2)/6$ unique three-body interactions. A system with 1000 atoms has 499 500 pairwise interactions and 166 167 000 three-body interactions. In general, there are approximately $N/3$ times more three-body terms than two-body terms and so it is clear why it is often considered preferable to avoid calculating the three-body interactions.

4.12 Effective Pair Potentials

Fortunately, it is found that a significant proportion of the many-body effects can be incorporated into a pairwise model, if properly parametrised. The pair potentials most commonly used in molecular modelling are thus 'effective' pairwise potentials; they do not represent the true interaction energy between two isolated particles but are parametrised to include many-body effects in the pairwise energy. Similarly, polarisation effects can be implicitly included in a force field by the simple expedient of enhancing the electrostatic interaction. This can be done by using larger partial charges than those for an isolated molecule. This is most obviously manifested in larger multipole moments; the dipole moment of a single water molecule is 1.85 D, whereas the dipole moment of many simple water models designed to simulate liquid water are significantly larger (closer to the experimental value for liquid water of 2.6 D).

A notable example of a potential that does include many-body terms is the Barker-Fisher-Watts potential for argon, which combines a pairwise potential with an Axilrod-Teller triple

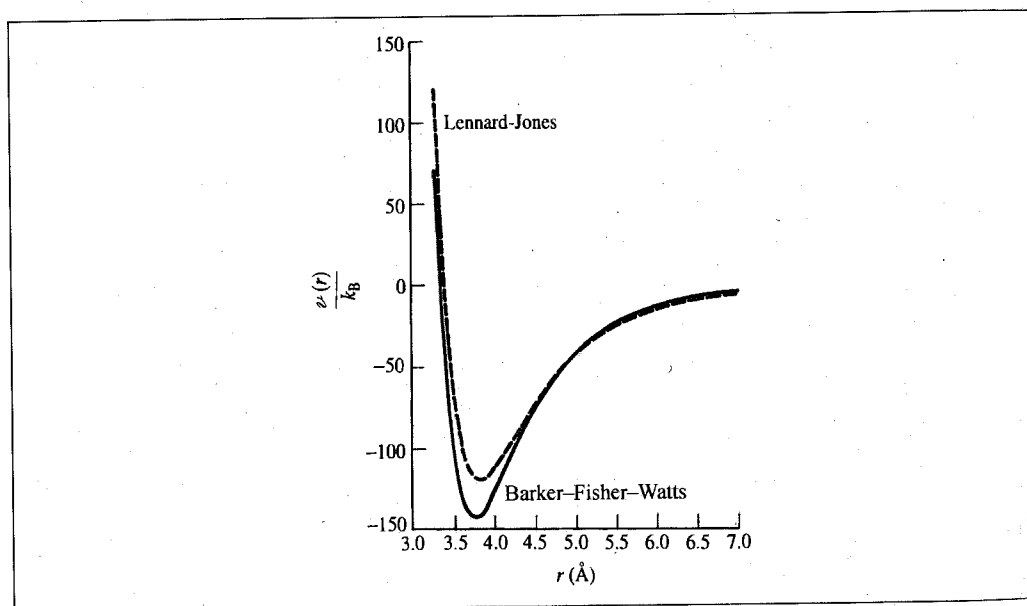


Fig. 4.38: Comparison of the Lennard-Jones potential for argon with the Barker-Fisher-Watts pair potential; k_B is Boltzmann's constant.

potential [Barker *et al.* 1971]. The pair potential is a linear combination of two potentials that each take the following form:

$$v^*(r^*) = e^{\alpha(1-r^*)} [A_0 + A_1(r^* - 1) + A_2(r^* - 1)^2 + A_3(r^* - 1)^3 + A_4(r^* - 1)^4 + A_5(r^* - 1)^5] + \frac{C_6}{\delta + R^{*6}} + \frac{C_8}{\delta + R^{*8}} + \frac{C_{10}}{\delta + R^{*10}} \quad (4.83)$$

This potential function contains eleven constants: $\alpha, A_0 \dots A_5, C_6, C_8, C_{10}$ and δ . The function is expressed in terms of r^* , which is given by $r^* = r/r_m$, where r_m is the separation at the minimum in the potential. The 'true' interaction energy as a function of the separation, r , is then obtained by multiplying $v^*(r^*)$ by the depth of the potential well, ε :

$$v(r) = \varepsilon v^*(r^*) \quad (4.84)$$

A comparison of the pairwise contribution to the Barker-Fisher-Watts potential with the Lennard-Jones potential for argon is shown in Figure 4.38.

4.13 Hydrogen Bonding in Molecular Mechanics

Some force fields replace the Lennard-Jones 6-12 term between hydrogen-bonding atoms by an explicit hydrogen-bonding term, which is often described using a 10-12 Lennard-Jones potential:

$$v(r) = \frac{A}{r^{12}} - \frac{C}{r^{10}} \quad (4.85)$$

This function is used to model the interaction between the donor hydrogen atom and the heteroatom acceptor atom. Its use is intended to improve the accuracy with which the geometry of hydrogen-bonding systems is predicted. Other force fields incorporate a more complicated hydrogen-bonding function that takes into account deviations from the geometry of the hydrogen bond and is thus dependent upon the coordinates of the donor and acceptor atoms as well as the hydrogen atom. For example, the YETI force field [Vedani 1988] uses the following form for its hydrogen bonding term:

$$v_{\text{HB}} = \left(\frac{A}{r_{\text{H}\dots\text{Acc}}^{12}} - \frac{C}{r_{\text{H}\dots\text{Acc}}^{10}} \right) \cos^2 \theta_{\text{Don}\dots\text{H}\dots\text{Acc}} \cos^4 \omega_{\text{H}\dots\text{Acc-LP}} \quad (4.86)$$

The energy in Equation (4.86) depends upon the distance from the hydrogen to the acceptor, the angle subtended at the hydrogen by the bonds to the donor and the acceptor, and the deviation of the hydrogen bond from the closest lone-pair direction at the acceptor atom ($\omega_{\text{H}\dots\text{Acc-LP}}$ in Equation (4.86), Figure 4.39).

The GRID program [Goodford 1985] that is used for finding energetically favourable regions in protein binding sites uses a direction-dependent 6-4 function:

$$v_{\text{HB}} = \left(\frac{C}{d^6} - \frac{D}{d^4} \right) \cos^m \theta \quad (4.87)$$

θ is the angle subtended at the hydrogen and m is usually set to 4.

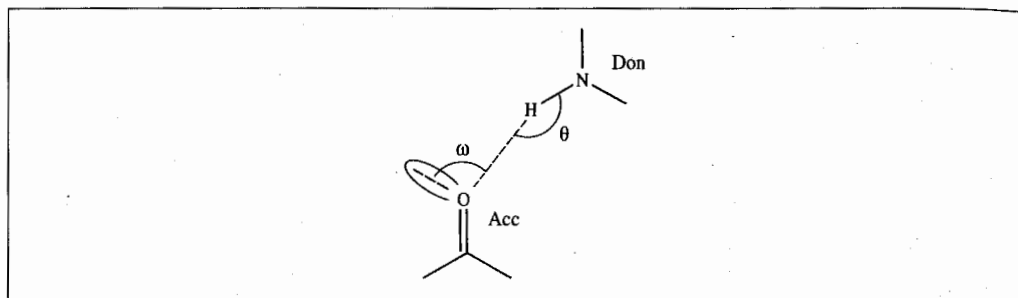


Fig. 4.39: Definition of hydrogen-bond geometry used in YETI force field.

By no means do all force fields contain explicit hydrogen-bonding terms; most rely upon electrostatic and van der Waals interactions to reproduce hydrogen bonding.

4.14 Force Field Models for the Simulation of Liquid Water

Many of the concepts that we have considered so far can be illustrated by examining some of the empirical models that have been developed to study water. Despite its small size, water acts as a paradigm for the different force field models that we have discussed. Moreover, many of its properties can be easily determined using computer simulation methods and so readily compared with experiment. It is also one of the most challenging systems to model accurately. A wide range of water models have been proposed. The computational efficiency with which the energy can be calculated using a given model is often an important factor as there may be a very large number of water molecules present, together with a solute; most of the force fields used to simulate liquid water thus use effective pairwise potentials with no explicit three-body terms or polarisation effects.

Water models can be conveniently divided into three types. In the simple interaction-site models each water molecule is maintained in a rigid geometry and the interaction between molecules is described using pairwise Coulombic and Lennard-Jones expressions. Flexible models permit internal changes in conformation of the molecule. Finally, models have been developed that explicitly include the effects of polarisation and many-body effects.

4.14.1 Simple Water Models

The 'simple' water models use between three and five interaction sites and a rigid water geometry. The TIP3P [Jorgensen *et al.* 1983] and SPC [Berendsen *et al.* 1981] models use a total of three sites for the electrostatic interactions; the partial positive charges on the hydrogen atoms are exactly balanced by an appropriate negative charge located on the oxygen atom. The van der Waals interaction between two water molecules is computed using a Lennard-Jones function with just a single interaction point per molecule centred on the oxygen atom; no van der Waals interactions involving the hydrogen atoms are calculated. The TIP3P and SPC models differ slightly in the geometry of each water molecule, in the

	SPC	SPC/E	TIP3P	BF	TIP4P	ST2
$r(\text{OH}), \text{\AA}$	1.0	1.0	0.9572	0.96	0.9572	1.0
HOH, deg	109.47	109.47	104.52	105.7	104.52	109.47
$A \times 10^{-3}, \text{kcal } \text{\AA}^{12}/\text{mol}$	629.4	629.4	582.0	560.4	600.0	238.7
$C, \text{kcal } \text{\AA}^6/\text{mol}$	625.5	625.5	595.0	837.0	610.0	268.9
$q(\text{O})$	-0.82	-0.8472	-0.834	0.0	0.0	0.0
$q(\text{H})$	0.41	0.4238	0.417	0.49	0.52	0.2375
$q(\text{M})$	0.0	0.0	0.0	-0.98	-1.04	-0.2375
$r(\text{OM}), \text{\AA}$	0.0	0.0	0.0	0.15	0.15	0.8

Table 4.3 A comparison of various water models [Jorgensen *et al.* 1983]. For the ST2 potential, $q(\text{M})$ is the charge on the 'lone pairs', which are a distance 0.8 \AA from the oxygen atom (see Figure 4.40).

hydrogen charges and in the Lennard-Jones parameters. These differences are indicated in Table 4.3, which also includes data for the SPC/E model [Berendsen *et al.* 1987], which is an updated version of the SPC model. The four-site models such as that of Bernal and Fowler [Bernal and Fowler 1933] (which is now relatively little used but is important for historical reasons as it dates from 1933) and Jorgensen's TIP4P model [Jorgensen *et al.* 1983] shift the negative charge from the oxygen atom to a point along the bisector of the HOH angle towards the hydrogens (Figure 4.40). The parameters for these two models are also given in the table. The most commonly used five-site model is the ST2 potential of Stillinger and Rahman [Stillinger and Rahman 1974]. Here, charges are placed on the hydrogen atoms and on two lone-pair sites on the oxygen. The electrostatic contribution is modulated so that for oxygen-oxygen distances below 2.016 \AA it is zero and for distances greater than 3.1287 \AA it takes its full value. Between these two distances the electrostatic contribution is modulated using a function that smoothly varies from 0.0 at the shorter distance to 1.0 at the longer distance (see Section 6.7.3).

The experimentally determined dipole moment of a water molecule in the gas phase is 1.85 D. The dipole moment of an individual water molecule calculated with any of these simple models is significantly higher; for example, the SPC dipole moment is 2.27 D and that for TIP4P is 2.18 D. These values are much closer to the effective dipole moment of liquid water, which is approximately 2.6 D. These models are thus all effective pairwise models. The simple water models are usually parametrised by calculating various properties using molecular dynamics or Monte Carlo simulations and then modifying the

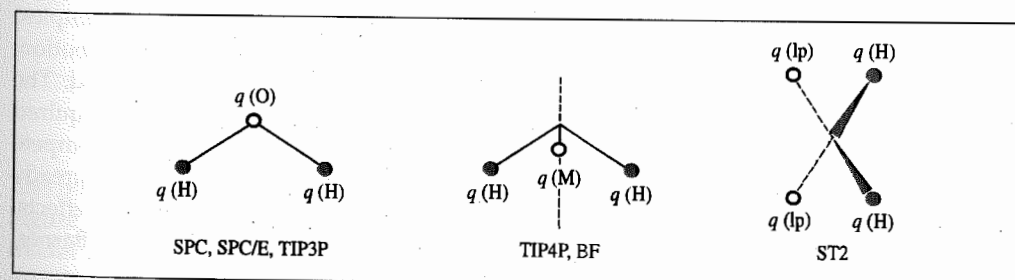


Fig. 4.40: Some 'simple' water models (Table 4.3) [Jorgensen *et al.* 1983].

parameters until the desired level of agreement between experiment and theory is achieved. Thermodynamic and structural properties are usually used in the parametrisation, such as the density, radial distribution function, enthalpy of vaporisation, heat capacity, diffusion coefficient and dielectric constant.* It is found that some properties such as the density and the enthalpy of vaporisation are predicted rather well by all of the models, but there is significant variation in the values for other properties such as the dielectric constant [Jorgensen *et al.* 1983]. When comparing the different models, it is also important to take account of the computational effort each requires. Thus, nine site-site distances must be calculated for each water dimer using a three-site model; ten are required for a four-site model, and seventeen for the ST2 model.

The use of a rigid model for water is obviously an approximation, and it means that some properties cannot be determined at all. For example, only when internal flexibility is included can the vibrational spectrum be calculated and compared with experiment. Flexibility is most easily incorporated by 'grafting' bond-stretching and angle-bending terms onto the potential function for a rigid model. Such an approach needs to be done with care. For example, Ferguson has developed a flexible model for water that is based upon the SPC model [Ferguson 1995]. The partial charges and van der Waals parameters in this model were slightly different from those in the rigid model, and flexibility was achieved using cubic and harmonic bond-stretching terms and a harmonic angle-bending term. The calculated values compared well with experimental results for a wide range of thermodynamic and structural properties, including the dielectric constant and self-diffusion coefficient.

4.14.2 Polarisable Water Models

The simple models give very good results for a wide range of properties of pure liquid water. However, there is some concern that they are not appropriate models to use for the most accurate work. This is especially the case for inhomogeneous systems where there are strong electric field gradients due to the presence of ions, and at the solute-solvent interface. Under such circumstances models that explicitly include polarisation effects and three-body terms are considered to be more appropriate. The inclusion of an explicit polarisation term should also enhance the ability of the model to reproduce the behaviour of water in other phases (e.g. solid and vapour) and at the interface between different phases. The dipole moment of an isolated water molecule in such a model should thus be closer to the gas-phase value rather than to the 'effective' value in liquid water. The simplest way to include polarisation is to use an isotropic molecular polarisability contribution; an alternative is to use atom-centred polarisabilities or the variable charge method. The incorporation of polarisability may significantly increase the computational effort required for a liquid simulation, and even then only the best polarisable models currently compete with the well-established models that use effective pairwise potentials. We have already considered some of the polarisable water models in our discussion of polarisation effects. One early attempt to incorporate such effects into a water model was made by Barnes,

* A discussion of the calculation of these properties from computer simulation is given in Section 6.2.

Finney, Nicholas and Quinn [Barnes *et al.* 1979]. Their polarisable electropole water model represented the charge distribution by a multipole expansion comprising a dipole of 1.855 D and a quadrupole moment that was determined from quantum mechanical calculations on an isolated molecule. Polarisation effects were calculated using an isotropic molecular polarisability from the electric fields being produced by the dipoles and quadrupoles of surrounding molecules. The model also used a spherically symmetric Lennard-Jones function. A more recent study used the fluctuating charge model with both the TIP4P and SPC geometries [Rick *et al.* 1994]. The charges were assigned to reproduce the correct dipole moment of the gas-phase molecule (in contrast to the equivalent non-polarisable models). Of the two geometries, the TIP4P model gave the better results for various properties. The dielectric properties were considered particularly well reproduced, including features in the dielectric spectrum arising from the translational motion of a water molecule in the cage of its neighbours. This feature is not present in fixed-charge models. Moreover, the computational cost with this particular model was only about 1.1 times that of the fixed-charge equivalent.

4.14.3 *Ab initio* Potentials for Water

The final category of water model that we shall consider are the '*ab initio*' potentials. These are based upon *ab initio* quantum mechanical calculations on small clusters of water molecules. One example of this type is the NCC model of Nieser, Corongiu and Clementi, which combines a two-molecule potential with a polarisation term [Nieser *et al.* 1990]. They had previously tried to explicitly include both three- and four-body effects but found this model computationally too expensive. The two-body model uses partial charges on the hydrogen atoms and a compensating negative charge on a site located along the bisector of the HOH angle, as in the TIP4P model. The equation used is:

$$\begin{aligned} \mathcal{V}_{\text{two-body}} = & q^2 \left(\frac{1}{R_{13}} + \frac{1}{R_{14}} + \frac{1}{R_{23}} + \frac{1}{R_{24}} \right) \\ & + \frac{4q^2}{R_{78}} - 2q^2 \left(\frac{1}{R_{81}} + \frac{1}{R_{82}} + \frac{1}{R_{73}} + \frac{1}{R_{74}} \right) \\ & + A_{\text{OO}} e^{-B_{\text{OO}}R_{56}} + A_{\text{HH}} (e^{-B_{\text{HH}}R_{13}} + e^{-B_{\text{HH}}R_{14}} + e^{-B_{\text{HH}}R_{23}} + e^{-B_{\text{HH}}R_{24}}) \\ & + A_{\text{OH}} (e^{-B_{\text{OH}}R_{53}} + e^{-B_{\text{OH}}R_{54}} + e^{-B_{\text{OH}}R_{61}} + e^{-B_{\text{OH}}R_{62}}) \\ & - A'_{\text{OH}} (e^{-B'_{\text{OH}}R_{53}} + e^{-B'_{\text{OH}}R_{54}} + e^{-B'_{\text{OH}}R_{61}} + e^{-B'_{\text{OH}}R_{62}}) \\ & + A_{\text{PH}} (e^{-B_{\text{PH}}R_{73}} + e^{-B_{\text{PH}}R_{74}} + e^{-B_{\text{PH}}R_{81}} + e^{-B_{\text{PH}}R_{82}}) \\ & + A_{\text{PO}} (e^{-B_{\text{PO}}R_{76}} + e^{-B_{\text{PO}}R_{85}}) \end{aligned} \quad (4.88)$$

The points P are the locations where the negative charge is placed (numbered 7 and 8 in Figure 4.41) and the terms A_{PH} and A_{PO} are used to enhance the performance of the model at short distances. q is the charge on each hydrogen. The polarisation term is calculated in an iterative manner using induced dipoles along each O-H bond. The NCC model was parametrised by fitting to the energies and other properties of 250 trimer and

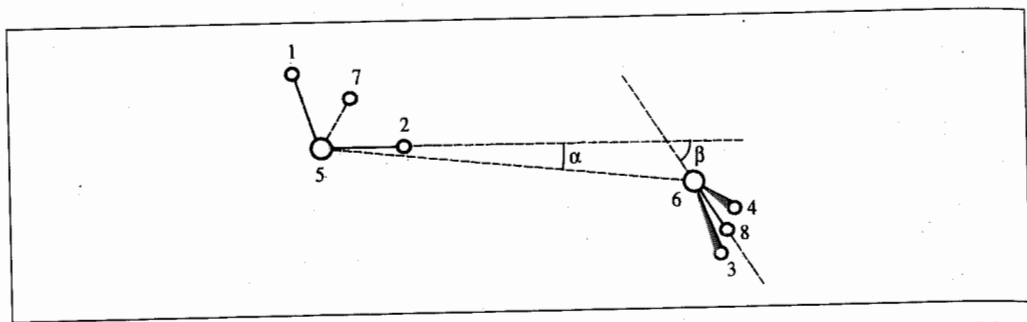


Fig. 4.41: The NCC water model. (After Corongiu G 1992. *Molecular Dynamics Simulation for Liquid Water Using a Polarisable and Flexible Potential*. *International Journal of Quantum Chemistry* 42:1209-1235.)

350 dimer configurations determined with high-level *ab initio* methods and large basis sets. The water trimer data was used to fit the many-body parameters (i.e. the locations of the induced dipole moments and the point charges, together with the polarisability and the value of the hydrogen charge). The dimer data were then used to fit the remaining terms in the potential.

The original NCC potential was designed as a rigid water model and performed well in tests of its ability to reproduce experimental data for both water dimers and liquid water. A flexible version has also been developed [Corongiu 1992], with the energy being expressed as a function of the three internal coordinates (two bond lengths and one angle) with terms up to quartics:

$$\begin{aligned}
 \mathcal{V}_{\text{intra}} = & \frac{1}{2} f_{RR} (\delta_1^2 + \delta_2^2) + \frac{1}{2} f_{\theta\theta} (\delta_3^2) + f_{RR'} \delta_1 \delta_2 + f_{R\theta} (\delta_1 + \delta_2) \delta_3 \\
 & + \frac{1}{R_e} [f_{RRR} (\delta_1^3 + \delta_2^3) + f_{\theta\theta\theta} \delta_3^3 + f_{RRR'} (\delta_1 + \delta_2) \delta_1 \delta_2 \\
 & + f_{RR\theta} (\delta_1^2 + \delta_2^2) \delta_3 + f_{RR'\theta} \delta_1 \delta_2 \delta_3 + f_{R\theta\theta} (\delta_1 + \delta_2) \delta_3^2] \\
 & + \frac{1}{R_e^2} [f_{RRRR} (\delta_1^4 + \delta_2^4) + f_{\theta\theta\theta\theta} \delta_3^4 + f_{RRR'} (\delta_1^2 + \delta_2^2) \delta_1 \delta_2 \\
 & + f_{RRR'\theta} \delta_1^2 \delta_2^2 + f_{RRR\theta} (\delta_1^3 + \delta_2^3) \delta_3] \\
 & + \frac{1}{R_e^2} [f_{RRR'\theta} (\delta_1 + \delta_2) \delta_1 \delta_2 \delta_3 + f_{RR\theta\theta} (\delta_1^2 + \delta_2^2) \delta_3^2 \\
 & + f_{RR'\theta\theta} \delta_1 \delta_2 \delta_3 + f_{R\theta\theta\theta} (\delta_1 + \delta_2) \delta_3^3]
 \end{aligned} \quad (4.89)$$

where $\delta_1 = R_1 - R_e$, $\delta_2 = R_2 - R_e$ and $\delta_3 = R_e(\theta - \theta_e)$.

The functional form of the NCC model demonstrates the complexity of some empirical models (and this for a molecule that contains only three atoms!). We should also note that the development of empirical models from *ab initio* quantum mechanical data is an approach that is already well established and looks likely to be a method that is more widely used in the future.

4.15 United Atom Force Fields and Reduced Representations

In our discussion so far, we have assumed that all of the atoms in the system are explicitly represented in the model. However, as the number of non-bonded interactions scales with the square of the number of interaction sites present, there are clear advantages if the number of interaction sites can be reduced. The simplest way to do this is to subsume some or all of the atoms (usually just the hydrogen atoms) into the atoms to which they are bonded. A methyl group would then be modelled as a single 'pseudo-atom' or 'united atom'. The van der Waals and electrostatic parameters would be modified to take account of the adjoining hydrogen atoms. Considerable computational savings are possible; for example, if butane is modelled as a four-site model rather than one with twelve atoms then the van der Waals interaction between two butane molecules involves the calculation of sixteen terms rather than 144. Other hydrocarbons are often represented using united atom models. Many of the earliest calculations on proteins used united atom representations. In this case, not all of the hydrogen atoms in the protein are subsumed into their adjacent atoms, but just those that are bonded to carbon atoms. Hydrogen atoms bonded to polar atoms such as nitrogen and oxygen are able to participate in hydrogen-bonding interactions, which are modelled much better if these hydrogens are explicitly represented.

One drawback with a united atom force field is that chiral centres may be able to invert during a calculation. This was found to be a problem with the united atom force fields for proteins. The alpha carbon in the peptide unit (C_α in Figure 4.42) is bonded to a hydrogen atom and to the side chain (glycine and proline are slightly different; see Section 10.1). A united atom force field model would not explicitly include the alpha hydrogen. Unfortunately, the stereochemistry at the alpha carbon can then invert during a calculation. This should be avoided as the naturally occurring amino acids have a defined stereochemistry (as shown in Figure 4.42). This inversion may be prevented through the use of an improper torsion term (e.g. N-C-C_α-R) to keep the side chain in the correct relative position.

In a united atom force field the van der Waals centre of the united atom is usually associated with the position of the heavy (i.e. non-hydrogen) atom. Thus, for a united CH₃ or CH₂ group the van der Waals centre would be located at the carbon atom. It would be more accurate to associate the van der Waals centre with a position that was offset slightly from the carbon position, in order to reflect the presence of the hydrogen atoms. Toxvaerd has developed such a model that gives superior performance for alkanes than do the simple united atom models, particularly for simulations at high pressures [Toxvaerd 1990]. In

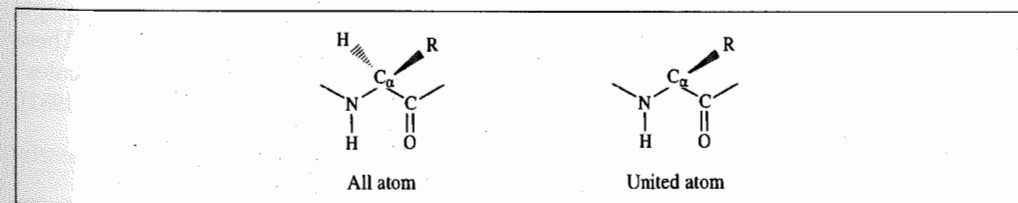


Fig. 4.42: Representations of the naturally occurring amino acids.

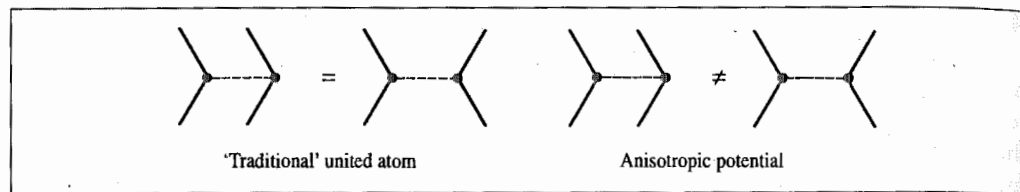


Fig. 4.43: The interaction energy between the two arrangements shown is equal in a 'traditional' united atom force field but different in the Toxvaerd anisotropic model. (Figure adapted from Toxvaerd S 1990. *Molecular Dynamics Calculations of the Equation of State of Alkanes*. The Journal of Chemical Physics 93:4290-4295.)

Toxvaerd's model the interaction sites are located at the geometrical centres of the CH_2 or CH_3 groups. The forces between these sites act on the united atom mass centre, which remains located on the carbon atom (with a mass of 14 for a CH_2 group and 15 for a CH_3 group). As the interaction site is no longer located at an atomic nucleus the forces acting on the masses are more complicated to calculate, but little additional computational expense is required. The effect of using such an anisotropic potential is nicely illustrated by the two arrangements of methylene units shown schematically in Figure 4.43. In the united atom model both arrangements would have the same energies and forces, but this is not so with the Toxvaerd anisotropic potential.

4.15.1 Other Simplified Models

In some force field models, even simpler representations are used than the united atom approach, with entire groups of atoms being modelled as single interaction points. For example, a benzene ring might be modelled as a single site with appropriately chosen parameters.

Yet other models have no obvious relationship to any 'real' molecule but are useful because their simplicity enables larger or more extensive calculations to be performed than would otherwise be possible. The polymer field is full of such models, as we shall discuss in Section 8.6. Another area where such models have been widely applied is in the study of liquid crystals. Liquid crystals are able to form phases that are characterised by a long-range order of the molecular orientations in at least one dimension. Many of the molecules that exhibit liquid crystalline behaviour are rod-shaped, but disc-like molecules can also form liquid crystalline phases. Some typical examples of molecules that can show such behaviour are shown in Figure 4.44. In the liquid crystalline state the rod-shaped molecules are aligned with their long axes pointing in approximately the same direction. Some very simple computer models have been used to investigate the behaviour of liquid crystals. These simple models enable large simulations to be performed on assemblies of many 'molecules'. One example of such a simplified model is the Gay-Berne potential [Gay and Berne 1981], which models the anisotropic interaction between two particles as:

$$v(r_{ij}) = 4\epsilon(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}}) \left\{ \left[\frac{\sigma_0}{r_{ij} - \sigma(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}}) + \sigma_0} \right]^{12} - \left[\frac{\sigma_s}{r_{ij} - \sigma(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}}) + \sigma_s} \right]^6 \right\} \quad (4.90)$$

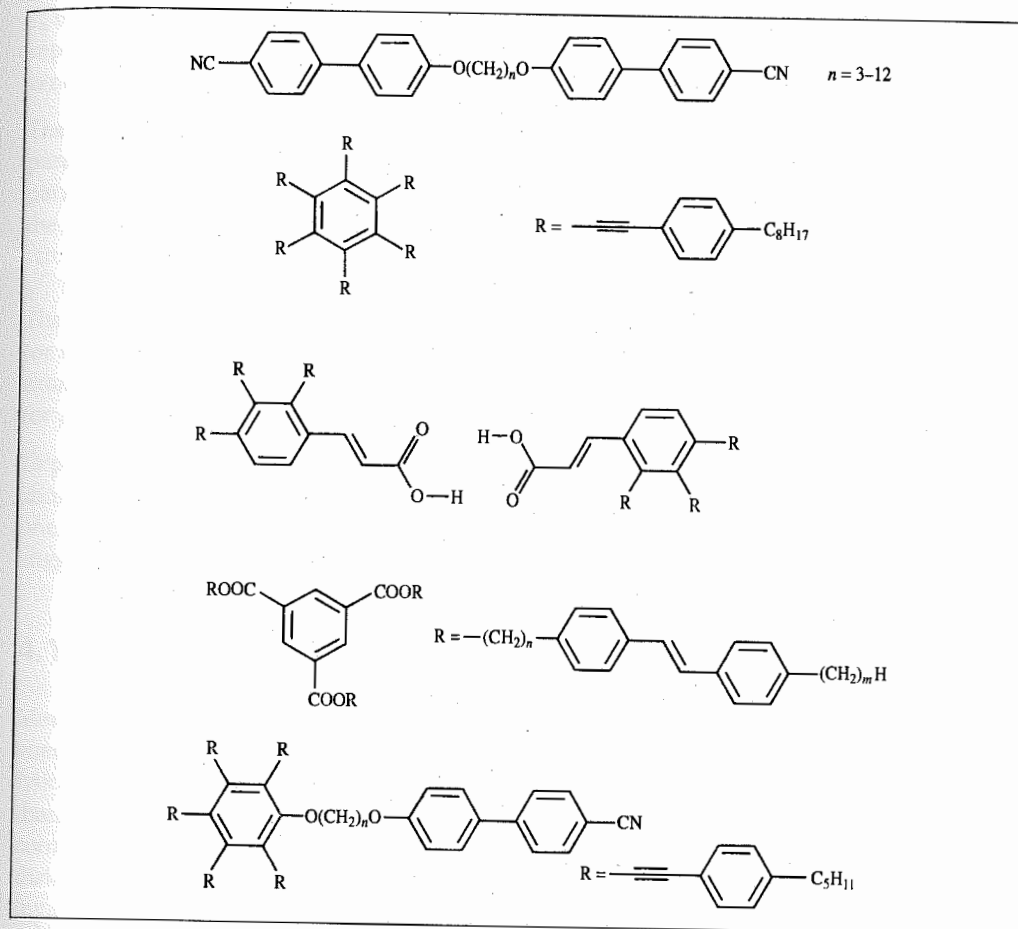


Fig. 4.44: Some typical liquid crystal molecules.

$\hat{\mathbf{u}}_i$ and $\hat{\mathbf{u}}_j$ are unit vectors that describe the orientations of the two molecules i and j and $\hat{\mathbf{r}}$ is a unit vector along the line connecting their centres (Figure 4.45). The molecules can be considered as ellipsoids which have a shape that is reflected in two size parameters, σ_s and σ_e , which are the separations at which the attractive and repulsive terms in the potential cancel for end-to-end and side-by-side arrangements respectively. These are incorporated into the potential via the parameter σ :

$$\sigma(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}}) = \sigma_0 \left\{ 1 - \frac{\chi}{2} \left[\frac{(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{r}} + \hat{\mathbf{u}}_j \cdot \hat{\mathbf{r}})^2}{1 + \chi(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{u}}_j)} + \frac{(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{r}} - \hat{\mathbf{u}}_j \cdot \hat{\mathbf{r}})^2}{1 - \chi(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{u}}_j)} \right] \right\}^{-1/2} \quad (4.91)$$

where

$$\chi = \frac{(\sigma_e/\sigma_s)^2 - 1}{(\sigma_e/\sigma_s)^2 + 1} \quad (4.92)$$

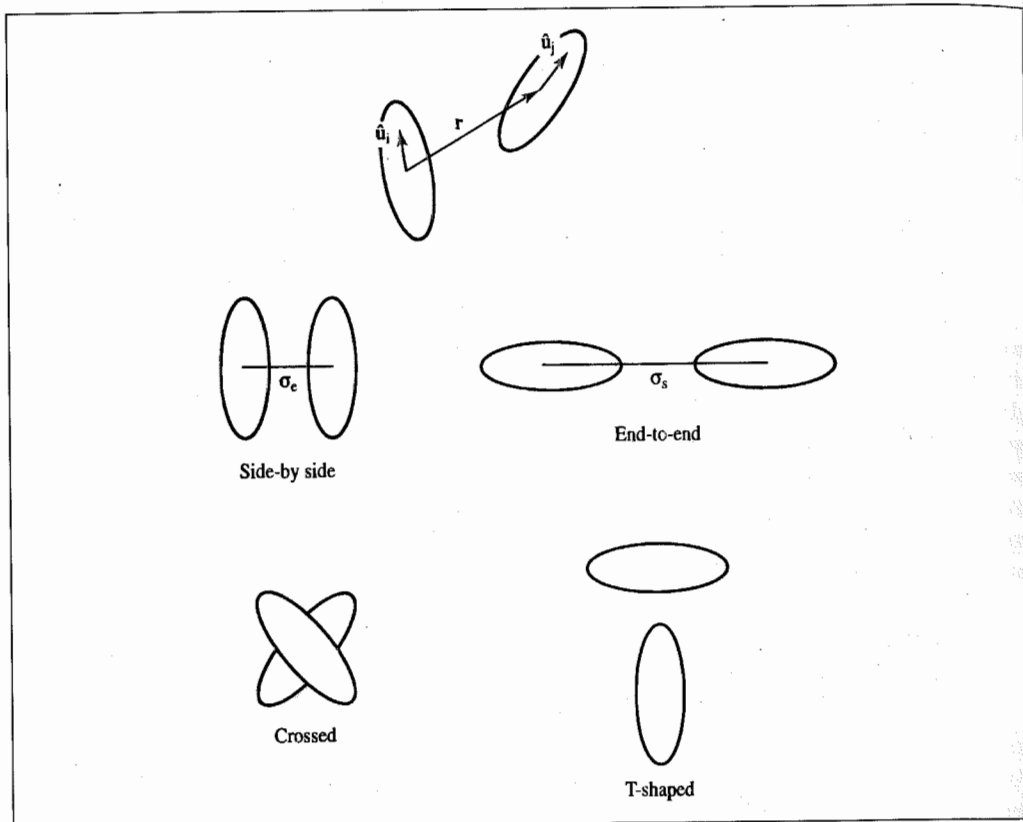


Fig. 4.45: The Gay-Berne model for liquid crystal systems and some typical arrangements.

χ is the *shape anisotropy parameter*; it is zero for spherical particles and is equal to 1 for infinitely long rods and -1 for infinitely thin discs; σ_0 is typically set equal to σ_s .

The energy term is also orientation-dependent and is written as follows:

$$\varepsilon(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}}) = \varepsilon_0 \varepsilon'^{\mu}(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}}) \varepsilon^{\nu}(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j) \quad (4.93)$$

where

$$\varepsilon(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j) = [1 - \chi^2(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{u}}_j)^2]^{-1/2} \quad (4.94)$$

$$\varepsilon'(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}}) = \left\{ 1 - \frac{\chi'}{2} \left[\frac{(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{r}} + \hat{\mathbf{u}}_j \cdot \hat{\mathbf{r}})^2}{1 + \chi'(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{u}}_j)} + \frac{(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{r}} - \hat{\mathbf{u}}_j \cdot \hat{\mathbf{r}})^2}{1 - \chi'(\hat{\mathbf{u}}_i \cdot \hat{\mathbf{u}}_j)} \right] \right\}$$

χ' measures the anisotropy of the attractive forces:

$$\chi' = \frac{1 - (\varepsilon_e/\varepsilon_s)^{1/\mu}}{(\varepsilon_e/\varepsilon_s)^{1/\mu} + 1} \quad (4.95)$$

ε_e is the well depth for an end-to-end arrangement of the ellipsoids when the attractive and repulsive contributions cancel, and ε_s is the corresponding well depth for the side-by-side arrangement (Figure 4.45).

The Gay-Berne potential is rather complex but is governed by a relatively small number of parameters, some of which have readily interpretable meanings. The effect of changing the parameters can be most clearly understood by considering certain orientations, such as the side-by-side, end-to-end, crossed and T-shaped structures (Figure 4.45). In the crossed structure the well depth $\varepsilon(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}})$ and the separation $\sigma(\hat{\mathbf{u}}_i, \hat{\mathbf{u}}_j, \hat{\mathbf{r}})$ are independent of χ and χ' . The ratio of the well depths for the end-to-end and side-by-side arrangements is $\varepsilon_e/\varepsilon_s$. The exponents μ and ν are considered adjustable parameters. One way to obtain values for these is to fit the Gay-Berne function to arrangements of Lennard-Jones particles. For example, Luckhurst, Stevens and Phippen determined a value of 1 for ν and a value of 2 for μ by fitting to a linear array of four Lennard-Jones centres [Luckhurst *et al.* 1990].

Depending upon the parameters chosen, simulations performed using the Gay-Berne potential show behaviour typical of liquid crystalline materials. Moreover, by modifying the potential one can determine what contributions affect the liquid crystalline properties and so help to suggest what types of molecule should be made in order to attain certain properties.

4.16 Derivatives of the Molecular Mechanics Energy Function

Many molecular modelling techniques that use force-field models require the derivatives of the energy (i.e. the force) to be calculated with respect to the coordinates. It is preferable that analytical expressions for these derivatives are available because they are more accurate and faster than numerical derivatives. A molecular mechanics energy is usually expressed in terms of a combination of internal coordinates of the system (bonds, angles, torsions, etc.) and interatomic distances (for the non-bonded interactions). The atomic positions in molecular mechanics are invariably expressed in terms of Cartesian coordinates (unlike quantum mechanics, where internal coordinates are often used). The calculation of derivatives with respect to the atomic coordinates usually requires the chain rule to be applied. For example, for an energy function that depends upon the separation between two atoms (such as the Lennard-Jones potential, Coulomb electrostatic interaction or bond-stretching term) we can write:

$$r_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2} \quad (4.96)$$

$$\frac{\partial v}{\partial x_i} = \frac{\partial v}{\partial r_{ij}} \frac{\partial r_{ij}}{\partial x_i} \quad (4.97)$$

$$\frac{\partial r_{ij}}{\partial x_i} = \frac{(x_i - x_j)}{r_{ij}} \quad (4.98)$$

Thus, for the Lennard-Jones potential:

$$\frac{\partial v}{\partial r_{ij}} = \frac{24\varepsilon}{r_{ij}} \left[-2 \left(\frac{\sigma}{r_{ij}} \right)^{12} + \left(\frac{\sigma}{r_{ij}} \right)^6 \right] \quad (4.99)$$

The force in the x direction acting on atom i due to its interaction with atom j is given by:

$$\mathbf{f}_{x_i} = (\mathbf{x}_i - \mathbf{x}_j) \frac{24\epsilon}{r_{ij}^2} \left[2 \left(\frac{\sigma}{r_{ij}} \right)^{12} - \left(\frac{\sigma}{r_{ij}} \right)^6 \right] \quad (4.100)$$

Analytical expressions for the derivatives of the other terms that are commonly found in force fields are also available [Niketic and Rasmussen 1977]. Similar expressions must be derived from scratch when new functional forms are developed.

4.17 Calculating Thermodynamic Properties Using a Force Field

A molecular mechanics program will return an 'energy value' for any configuration or conformation of the system. This value is properly described as a 'steric energy' and is the energy of the system relative to a zero point that corresponds to a hypothetical molecule in which all of the bond lengths, valence angles, torsions and non-bonded separations are set to their strainless values. It is not necessary to know the actual value of the zero point to calculate the *relative* energies of different configurations or different conformations of the system.

Molecular mechanics can be used to calculate heats of formation. To do so requires the energy to form the bonds in the molecule to be added to the steric energy. These bond energies are typically obtained by fitting to experimentally determined heats of formation and are stored as empirical parameters within the force field. The accuracy with which heats of formation can be predicted with molecular mechanics is, in appropriate cases, comparable with experiment. Thus, the steric energy of a given structure may vary considerably from one force field to another, but its heat of formation should be much closer (if the force fields have been properly parametrised).

A third type of 'energy' that can be obtained from a molecular mechanics calculation is the 'strain energy'. Differences in steric energy are only valid for different conformations or configurations of the same system. Strain energies enable different molecules to be compared. To determine the strain energy it is usual to define some 'strainless' reference point. The reference points can be chosen in many ways and so many different definitions of strain energy have been proposed in the literature. For example, Allinger and co-workers defined the reference point using a set of 'strainless' compounds such as the all-*trans* conformations of the straight-chain alkanes from methane to hexane. From this set of compounds it was possible to derive a set of strainless energy parameters for constituent parts of the molecules. The inherent strain energy of a hydrocarbon is then obtained by subtracting the reference 'strainless' energy from the actual steric energy calculated using the force field. One interesting conclusion of this study was that chair cyclohexane has an inherent strain energy due to the presence of 1,4 van der Waals interactions between the carbon atoms within the ring.

The sources of strain are often quantified by examining the different components (bonds, angles, etc.) of the force field. Such analyses can provide useful information, especially for cases such as highly strained rings. However, in many molecules the strain is distributed

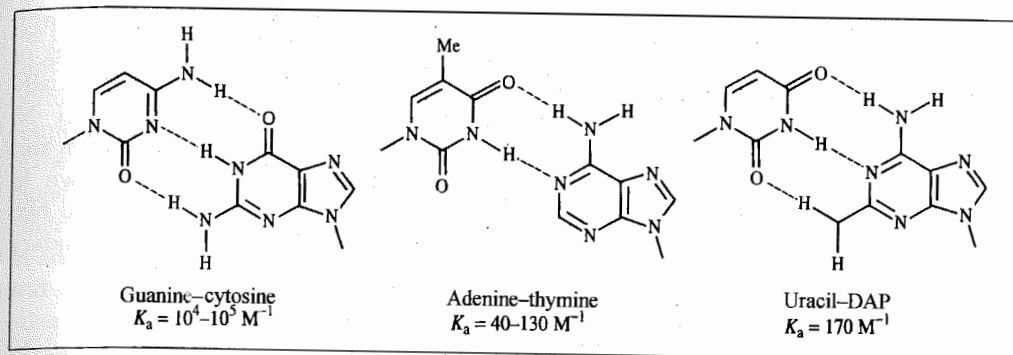


Fig. 4.46: The DNA base pairs guanine (G), cytosine (C), adenine (A) and thymine (T). The uracil-2,6-diaminopyridine pair can also form three hydrogen bonds but has a much lower association constant than G-C.

among a variety of internal parameters (and in any case is force-field-dependent). For intermolecular interactions the interpretation can be easier, for the 'interaction energy' is simply equal to the difference between the energies of the two isolated species and the energy of the intermolecular complex. A good example of this type of calculation and the conclusions that can be drawn from it is the study by Jorgensen and Pranata [Jorgensen and Pranata 1990] of the interaction between analogues of the DNA base pairs. In the double helical structure of DNA the bases pair up adenine (A) with thymine (T) and guanine (G) with cytosine (C) (Figure 4.46).

The association constant of the G-C base pair in chloroform is between 10^4 M^{-1} and 10^5 M^{-1} whereas the association between the A-T base pair is significantly weaker, at $40\text{--}130 \text{ M}^{-1}$. One obvious reason for this difference is that there are three hydrogen bonds in the G-C base pair and only two in the A-T base pair. However, a simple hydrogen-bond count

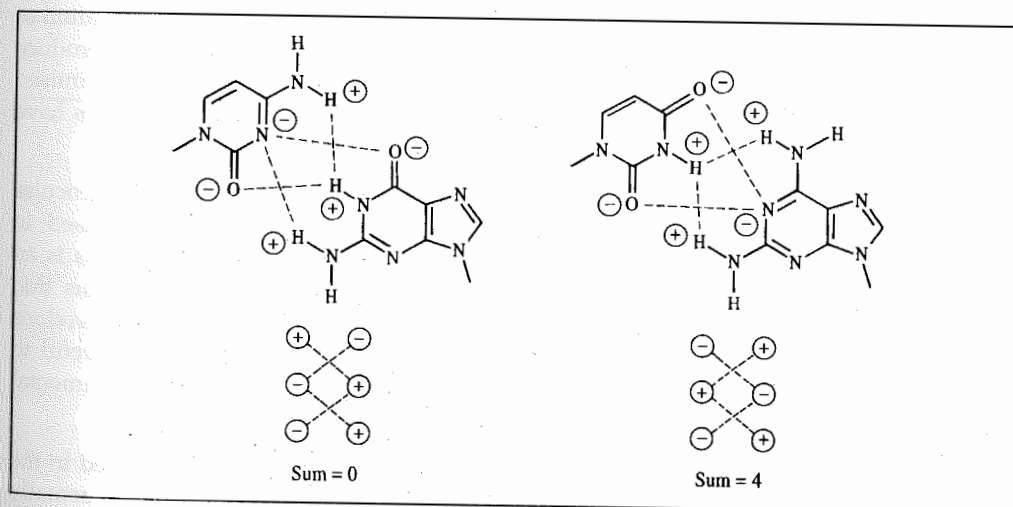


Fig. 4.47: Secondary interactions in guanine-cytosine and uracil-DAP.

does not explain all of the data, for synthetic analogues show a significant variation in their association constants, despite having three hydrogen bonds. The weak binding of the uracil-2,6-diaminopyridine (DAP) system (Figure 4.46) could be considered especially anomalous as it contains the same types of hydrogen bond as in G-C ($\text{NH}_2 \cdots \text{O}$, $\text{NH} \cdots \text{N}$, $\text{NH}_2 \cdots \text{O}$). A qualitative explanation for this phenomenon was proposed by Jorgensen and Pranata who examined the secondary interactions in these complexes. As shown in Figure 4.47, the G-C system contains two unfavourable secondary interactions and two favourable ones, an overall sum of zero. In the uracil-DAP system, all four secondary interactions are unfavourable.

4.18 Force Field Parametrisation

A force field can contain a large number of parameters, even if it is intended for calculations on only a small set of molecules. Parametrisation of a force field is not a trivial task. A significant amount of effort is required to create a new force field entirely from scratch, and even the addition of a few parameters to an existing force field in order to model a new class of molecules can be a complicated and time-consuming procedure. The performance of a force field is often particularly sensitive to just a few of the parameters (usually the non-bonded and torsional terms), so it is often sensible to spend more time optimising these parameters rather than others (such as the bond-stretching and angle-bending terms), the values of which do not greatly affect the results.

The first step is to select the data that are going to be used to guide the parametrisation process. Molecular mechanics force fields may be used to determine a variety of structurally related properties and the parametrisation data should be chosen accordingly. The geometries and relative conformational energies of certain key molecules are usually included in the data set. It is increasingly common to include vibrational frequencies in the parametrisation; these are usually more difficult to reproduce but the incorporation of appropriate cross terms can often help. Some force fields are parametrised to reproduce thermodynamic properties using computer simulation techniques. The OPLS (optimised parameters for liquid simulations [Jorgensen and Tirado-Reeves 1988]) parameters have been obtained in this way.

Unfortunately, experimental data may be non-existent or difficult to obtain for particular classes of molecules. Quantum mechanics calculations are thus increasingly used to provide the data for the parametrisation of molecular mechanics force fields. This is an important development because it greatly extends the range of chemical systems that can be treated using the force-field approach. *Ab initio* calculations are able to reproduce experimental results for small representative systems. Clearly, one should be careful to properly validate a force field derived in such a way by testing against experimental data if at all possible.

Once a functional form for the force field has been chosen and the data to be used in the parametrisation identified, there are then two basic methods that can be used to actually obtain the parameters. The first approach is 'parametrisation by trial and error', in which

the parameters are gradually refined to give better and better fits to the data. It is difficult to simultaneously modify a large number of parameters in such a strategy and so it is usual to perform the parametrisation in stages. It is important to remember that there is some coupling between all of the degrees of freedom and so for the most sensitive work none of the parameters can truly be taken in isolation. Parameters for the hard degrees of freedom (bond stretching and angle bending) can, however, often be treated separately from the others (indeed the bond and angle parameters are often transferred from one force field to another without modification). By contrast, the soft degrees of freedom (non-bonded and torsional contributions) are closely coupled and can significantly influence each other. One protocol that can be quite successful is to first establish a series of van der Waals parameters. The electrostatic model is then determined (e.g. by electrostatic potential fitting). Finally, the torsional potentials are determined by ensuring that the torsional barriers are reproduced together with the relative energies of the different conformations. Of course, it may be necessary to modify any of the parameters at any stage should the results be inadequate and so parametrisation is invariably an iterative procedure.

As experimental information on torsional barriers is often sparse or non-existent, quantum mechanical calculations are widely used to determine torsional potentials. The general strategy is as follows. First, a molecular fragment that adequately represents the rotatable bond of interest and its immediate environment is chosen. A series of structures are then generated by rotating about the bond and their energies determined using quantum mechanics. The torsional potential is then fitted to reproduce the energy curve, in conjunction with the van der Waals potential and partial charges. This procedure can be illustrated using the study of Pranata and Jorgensen who wanted to perform some calculations on FK506, a potent immunosuppressant (Figure 4.48) [Pranata and Jorgensen 1991]. FK506 contains a ketoamide functionality that has a *trans* conformation when the molecule is bound to its receptor but which is *cis* in the crystal structure of isolated FK506. NMR experiments suggested that the molecule adopts both *cis* and *trans* conformations in solution. This part of the molecule is clearly implicated in its function and so it was considered important

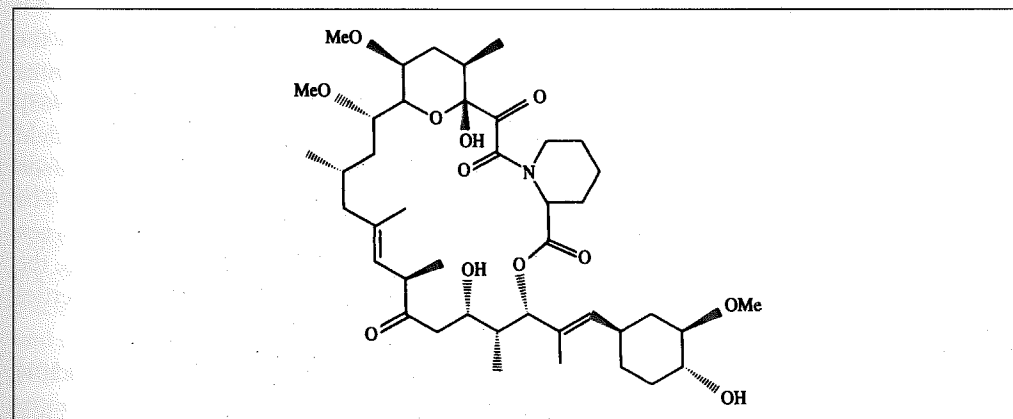


Fig. 4.48: The immunosuppressant FK506.

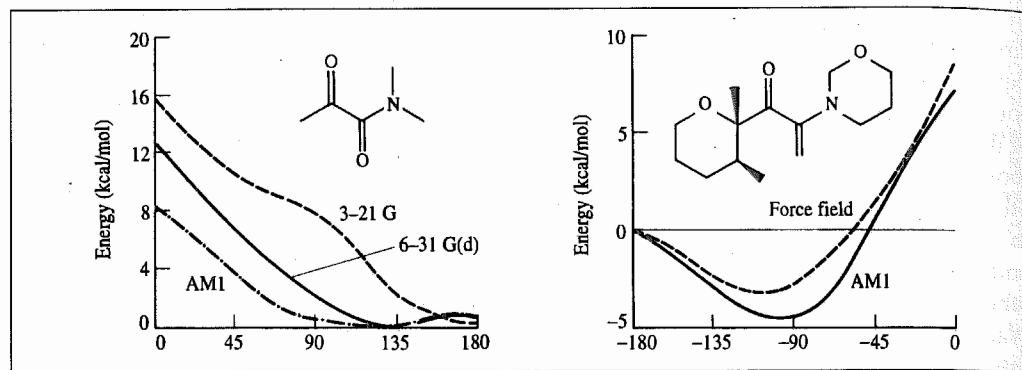


Fig. 4.49: Fragments used to derive and evaluate parameters for the ketoamide functionality in FK506. (Figure redrawn from J Pranata and W L Jorgensen 1991. *Computational Studies on FK506: Conformational Search and Molecular Dynamics Simulations in Water*. The Journal of the American Chemical Society 113:9483-9493.)

to correctly model the torsional potential about this bond. Pranata and Jorgensen intended to use the AMBER force field for their calculations but the force field contained no parameters for this link.

Molecular orbital calculations were performed on *N,N*-dimethyl- α -ketopropanamide (Figure 4.49, left), which was chosen as an appropriate model system. Semi-empirical calculations using AM1 and *ab initio* calculations using a 6-31G(d) basis set suggested that the minimum energy conformation corresponded to a torsion angle of 124° and 135°, respectively, with the *anti* conformation being slightly higher in energy (~0.7 kcal/mol). However, an analogous calculation using the 3-21G basis set did predict that the *anti* conformation was at a minimum (Figure 4.49). Crystal structures of compounds containing this fragment revealed that an orthogonal structure was commonly encountered. Torsional parameters were then fitted to the 6-31G(d) potential and evaluated by calculating an energetic profile for rotation in a larger fragment of the FK506 molecule using the force field and comparing it with that obtained using AM1 (Figure 4.49, right).

An alternative approach to parametrisation, pioneered by Lifson and co-workers in the development of their 'consistent' force fields, is to use least-squares fitting to determine the set of parameters that gives the optimal fit to the data [Lifson and Warshel 1968]. Again, the first step is to choose a set of experimental data that one wishes the force field to reproduce (or calculate using quantum mechanics, if appropriate). Warshel and Lifson used thermodynamic data, equilibrium conformations and vibrational frequencies. The 'error' for a given set of parameters equals the sum of squares of the differences between the observed and calculated values for the set of properties. The objective is to change the force field parameters to minimise the error. This is done by assuming that the properties can be related to the force field by a Taylor series expansion:

$$\Delta y(\mathbf{x} + \delta \mathbf{x}) = \Delta y(\mathbf{x}) + \mathbf{Z}\delta \mathbf{x} + \dots \quad (4.101)$$

Δy is a vector of the differences between the calculated and experimental data and is a vector whose components are the force field parameters. \mathbf{Z} is a matrix whose elements are the

derivatives of each property with respect to each of the parameters, $\partial x/\partial y$. An iterative procedure is used to minimise the sum of squares of the differences, Δy^2 . The method is easily modified to enable various weighting factors to be assigned to the different pieces of experimental data, so that (for example) the thermodynamic data could be given greater importance than the vibrational frequencies.

A well-known application of the least-squares approach to the optimisation of a force field was performed by Hagler, Huler and Lifson, who derived a force field for peptides by fitting to crystal data of a variety of appropriate compounds [Hagler *et al.* 1977; Hagler and Lifson 1974]. A key result of their work was that no explicit hydrogen bond term was required to model the hydrogen-bonding interactions, but that a combination of appropriate electrostatic and van der Waals models was sufficient. A group led by Hagler more recently developed a force field based upon the results of *ab initio* quantum mechanics calculations on small molecules, again using least-squares fitting [Maple *et al.* 1988]. The quantum mechanics calculations were performed not only on small molecules at equilibrium geometries but also on structures that were distorted from equilibrium. For each geometry the energy was calculated together with the first and second derivatives of the energy. This provided a wealth of data for the subsequent fitting procedure. This research has resulted in many new algorithms for the derivation of force-field parameters and has also challenged some of the assumptions about the development and functional form of force fields. One feature of the resulting force field, named CFF (standing for consistent force field), is that it contains rather more cross terms than other force fields. This can be ascribed to the objective of accurately reproducing vibrational spectra.

4.19 Transferability of Force Field Parameters

The range of systems that have been studied by force field methods is extremely varied. Some force fields have been developed to study just one atomic or molecular species under a wider range of conditions. For example, the chlorine model of Rodger, Stone and Tildesley [Rodger *et al.* 1988] can be used to study the solid, liquid and gaseous phases. This is an anisotropic site model, in which the interaction between a pair of sites on two molecules depends not only upon the separation between the sites (as in an isotropic model such as the Lennard-Jones model) but also upon the orientation of the site-site vector with respect to the bond vectors of the two molecules. The model includes an electrostatic component which contains dipole-dipole, dipole-quadrupole and quadrupole-quadrupole terms, and the van der Waals contribution is modelled using a Buckingham-like function.

Other force fields are designed for use with specific classes of molecules; we have already encountered the AMBER force field, which is designed for calculations on proteins and nucleic acids. Yet other force fields are intended to be applied to a wide range of molecules, and indeed some force fields are designed to model the entire periodic table. Intuitively, one might expect a 'specialised' force field to perform better than a 'general' force field, and while this is certainly true for the best of the specialised force fields, a good general force field can often outperform a poor specific force field.

The ability to transfer parameters from one molecule to another is crucial for any force field. Without it, the task of parametrisation would be impossible, because so many parameters would be required, and the force field would have no predictive ability. Transferability has a number of important consequences for the development and application of force fields. The problem of transferability is often first encountered when a molecular mechanics program fails to run because parameters are missing for the molecule being studied. One must somehow find values for the missing parameters. Some programs automatically 'guess' force field parameters; it is wise to check these assignments as they may be suspect. For the developer of a force field, a compromise must often be found between a complex functional form and a large number of atom types. It is also important to try to ensure that the errors in the force field are balanced, in the sense that it would be silly to spend a lot of time getting (say) the bond-stretching terms just right, if the van der Waals parameters give rise to large errors.

An alternative to 'guessing' parameters (which, if done properly, can sometimes give quite reasonable results) is to construct the force field in such a way that the parameters can be derived from atomic properties. This is particularly pertinent to those force fields which are designed to be used on a very wide range of elements and atom types, such as the Universal Force Field [Rappé *et al.* 1992]. This force field is claimed to model the entire periodic table and as such it would probably be impossible to derive individual parameters for each of the terms; indeed, the data required for such an exercise does not exist for many cases. Thus the UFF has a set of atom types which are characterised by atomic number, hybridisation and formal oxidation state. Reference bond lengths are initially set equal to the sum of the two relevant atomic bond radii and then corrected for bond order and the relative electronegativities of the two atoms. Bond force constants are obtained from Badger's rules, under which the force constant is proportional to the product of the 'effective atomic charges' for the two atoms and inversely proportional to the cube of the interatomic distance:

$$k_{ij} \propto \frac{q_i^* q_j^*}{r_{ij}^3} \quad (4.102)$$

The effective atomic charges are either obtained by fitting to data on diatomic molecules (where it exists) or by interpolation or extrapolation from this fit.

Transferability can be helped by using the same parameters for as wide a range of situations as possible. The non-bonded terms are particularly problematic in this regard; it would, in principle, be necessary to have parameters for the non-bonded interactions between all possible pairs of atom types. This would give rise to a very large number of parameters. It is therefore commonly assumed that the same set of van der Waals parameters can be used for most, if not all, atoms of the same element. For example, all carbon atoms (sp^3 , sp^2 , sp , etc.) would be treated with the same set of van der Waals parameters, all nitrogens by a common set, and so on. The torsional terms may also be generalised, so that the torsional parameters depend solely upon the atom types of the two atoms that form the central bond, rather than on all four atoms that comprise the torsion angle, as described in Section 4.5 for the AMBER force field.

4.20 The Treatment of Delocalised π Systems

The bonds in conjugated π systems are often of different lengths. For example, the central bond in butadiene is approximately 1.47 Å long, but the two terminal $CH=CH_2$ bonds are approximately 1.34 Å. If butadiene is modelled using a force field in which all four carbon atoms are assigned the same atom type (e.g. 'carbon sp^2 ') then each bond will be assigned the same bonding parameters and in the equilibrium structure all carbon-carbon bonds will be almost identical in length. A similar situation arises for aromatic systems. For example, not all the bonds in naphthalene are of equal length (unlike benzene). The bond lengths in a delocalised π system depend upon the bond orders; the higher the bond order, the shorter the bond.

In some cases it may be possible to circumvent this problem by creating a model specific to the conjugated system. For butadiene the central carbon-carbon bond of the π system could be treated in a different manner to the two terminal bonds, for example by using one atom type for the $-CH=$ carbon atoms and one for the $=CH_2$ carbon atoms in butadiene. This approach might be acceptable if we wanted to perform an extensive series of calculations on substituted butadienes, but it does compromise the transferability of the force field parameters. An alternative is to incorporate a molecular orbital calculation into the force field. Two variants on this theme have been developed. In one approach, the π and σ systems are treated separately [Warshel and Karplus 1972; Warshel and Lippicirella 1981]. For a given geometry, a self-consistent field quantum mechanical calculation is performed on the π system, typically with an appropriate semi-empirical theory. Molecular mechanics is simultaneously applied to the σ system. The energies of the quantum mechanical and molecular mechanical calculations are added together, and the geometry is modified to minimise this combined energy. A obvious assumption inherent in this approach is that the π and σ systems can be separated, which may be difficult to justify when deviations from planarity are present. Nevertheless, the approach has been extended to include those containing conjugated nitrogen and oxygen atoms, which has enabled the study of the properties of not only the ground states of some important biological chromophores (such as porphyrins) but also their excited states [Warshel and Lippicirella 1981].

An alternative approach is exemplified by the MM2/MM3/MM4 family of programs. First, a molecular orbital calculation is performed on the π system. If the initial conformation of the system is non-planar the calculation is performed on the equivalent planar system. The force field parameters are then modified according to the quantum mechanical bond orders. In MMP2 (the name given to the special version of MM2 which incorporated these features) these parameters are the force constant for the bonds in the π system, the reference bond lengths and the torsional barriers [Sprague *et al.* 1987; Allinger and Sprague 1973]. The system is then subjected to the usual molecular mechanics treatment using the new force field parameters. A linear relationship between the stretching constants and the bond orders, and between the reference bond lengths and the bond orders was found to give good results. Initially, the torsional barriers were assumed to be proportional to the square of the bond orders, but this relationship was modified slightly in subsequent versions

of the program. Thus in MM4 the V_2 and V_3 terms become:

$$V_2 = [A + p_{ij}^{\omega=0} \beta_{ij}] V_2^0 \quad (4.103)$$

$$V_3 = K_{V_3} [1 - p_{ij}(\omega)] V_3^0 \quad (4.104)$$

In Equation (4.103) p_{ij} is the bond order about the central bond $i-j$ of the torsion angle calculated for a torsion angle of zero and β_{ij} is the resonance integral from the molecular orbital calculation. The parameter A has a value of -0.09 and so the V_2 term is lower for those conjugated bonds with a lower bond order. In Equation (4.104) p_{ij} is now the bond order for the bond $i-j$ calculated for the torsion angle ω . K_{V_3} equals 1.25 and so V_3 increases with decreasing bond order. A bond with a lower bond order (and so a lower V_2 and a higher V_3) is thus more likely to deviate from planarity.

4.21 Force Fields for Inorganic Molecules

It may come as a surprise to many readers to learn that the earliest force field calculations on inorganic molecules were reported at much the same time as the first calculations on organic systems. For example, Corey and Bailar described the use of empirical force field calculations on octahedral complexes of cobalt in 1959 [Corey and Bailar 1959]. The range of metal-containing systems that can be considered by force field methods has steadily expanded since then. Moreover, many systems of commercial interest contain metals or other elements not usually found in 'organic' or 'biochemical' systems.

Some inorganic systems (such as certain coordination complexes) are little different to organic systems from a force field point of view; the bonding can be represented in a similar way and many of the force field parameters originally developed for organic systems can be transferred without modification. However, inorganic molecules do have certain properties which makes them more difficult to model than their organic counterparts. Perhaps the two most striking properties are the much wider range of geometries and the presence of highly delocalised bonds. Thus inorganic molecules include square planar and sawhorse (e.g. SF_4) shapes for four coordination and T-shaped for three coordination. Coordination numbers higher than four are also possible, with five (square pyramidal, trigonal bipyramidal) and six (octahedral and trigonal prismatic) being particularly common. To model such systems using conventional organic force fields would often be problematic because their geometries do not have a high degree of symmetry. For example, in a trigonal bipyramid there are in principle three different types of bond angle subtended at the central atom (90° , 120° and 180°). Moreover, in such systems the atoms are often equivalent (interchanging them gives the same structure back). However, if these atoms are assigned different force field parameters then this equivalence is not reproduced by the calculation. At least in these cases there is an obvious localised bonding scheme that can be applied; this is often not possible with organometallic molecules. For example, how should the bonding in ferrocene be represented in a force field calculation? Is there a bond between the iron and each of the carbon atoms in the two cyclopentadienyl rings? Is there a 'bond' from the iron to the centre of each of the rings? A yet further complication is that significant deviations from ideal geometries are often observed due to electronic effects such as the Jahn-Teller effect.

Whilst there is no universal solution to these problems within the context of a single force field similar to those used in organic chemistry, for certain situations it is possible to use an organic-like force field with only relatively small modifications. For obvious reasons those complexes with a high degree of symmetry are most amenable to such a treatment. Thus octahedral and square planar complexes are the simplest to model because of their symmetry (in addition to the geometries common in organic chemistry). However, even these have two types of equilibrium angle (180° and 90°). The situation can be much more complicated for the other geometries or for structures where the geometry about the metal is a distortion of a regular arrangement. A Urey-Bradley treatment of the bonding about the metal can often be quite successful in achieving the correct geometries. Here, there are no angle-bending terms at the metal but terms due to pairs of atoms bonded to the metal.

It is much more difficult to use such a force field to model metal π systems, where the bonding between the metal and the ligand is not easily represented by a conventional bonding picture. As we have discussed, metal atoms can adopt a wide range of geometries in π complexes, which are often significantly distorted from regular structures. Nevertheless, force fields have been developed which can cope with such systems, as well as being able to model more traditional systems such as organic compounds. These force fields often use a rather different functional form from Equation (4.1) and the parameters are obtained in a different way. One distinctive feature of both the Universal Force Field and the SHAPES force field developed by Landis and co-workers [Allured *et al.* 1991; Cleveland and Landis 1996] is the way in which angle bending is treated. The harmonic potential that is commonly employed in standard force fields is inappropriate to model the distortion of systems as the angle approaches 180° . UFF [Rappé *et al.* 1993] uses a cosine Fourier series for each angle ABC:

$$v(\theta) = K_{\text{ABC}} \sum_{n=0}^m C_n \cos n\theta \quad (4.105)$$

The coefficients C_n are chosen to ensure that the function has a minimum at the appropriate reference bond angle. For linear, trigonal, square planar and octahedral coordination, Fourier series with just two terms are used with a C_0 term and a term for $n = 1, 2, 3$ or 4, respectively:

$$v(\theta) = K_{\text{ABC}} [1 - \cos(n\theta)] \quad (4.106)$$

Thus, for example, if $n = 4$ then the function has minima at both 90° and 180° as required for octahedral geometries. The general case is exemplified by the H-O-H angle in water, where it is desired to have a minimum in the energy at an angle of 104.5° . Moreover, at this angle (θ_0) the second derivative of the energy equals the force constant. If in addition it is required that the energy is a maximum at 180° the following expression results:

$$v(\theta) = K_{\text{ABC}} [C_0 + C_1 \cos(\theta) + C_2 (\cos 2\theta)] \quad (4.107)$$

The three coefficients are defined as:

$$C_2 = \frac{1}{4 \sin^2(\theta_0)}; \quad C_1 = -4C_2 \cos(\theta_0); \quad C_0 = C_2 [2 \cos^2(\theta_0) + 1] \quad (4.108)$$

The SHAPES angle-bending term is very similar:

$$v(\theta) = K_{ABC} \sum_{n=0}^m [1 + \cos(n\theta - \delta)] \quad (4.109)$$

δ is the phase shift. Landis subsequently developed a formulation (called VALBOND) for the angle-bending term that is based on valence bond theory and which can produce results that compare well with *ab initio* calculations [Landis *et al.* 1995, 1998]. For example, using just one set of C–H parameters the H–C–H bond angles in ethene, formaldehyde and both singlet and triplet carbene match closely those found experimentally. One key practical advantage of this method is that it is not necessary to define equilibrium bond angles.

4.22 Force Fields for Solid-state Systems

Empirical potential models are widely used to study the solid state, complementing the quantum mechanical approaches we discussed in Chapter 3. One important difference between solid-state materials and 'organic' molecules (and indeed, some inorganic complexes) is that whilst the latter can generally be described using a localised bond model this is not always the case for the former. As a consequence, molecular mechanics approaches of the kind we have discussed so far in this chapter can be applied successfully only to certain types of material. Ionic and metallic systems especially require an alternative approach. Perhaps the key difference between solid-state materials and isolated molecules is the way in which the electrostatic terms are considered. As we shall see in Sections 6.7 and 6.8 it is common to truncate such interactions at some cutoff distance. However, solid-state modelling is concerned with materials that have long-range order; moreover, they often contain highly charged species. This means that the use of cutoffs can have a particularly detrimental effect, necessitating the use of special techniques such as the Ewald summation that enable more accurate interaction energies to be calculated. First, however, we shall consider the treatment of covalent systems which are amenable to the 'organic' style of molecular mechanics force field treatment, as exemplified by the study of zeolites.

4.22.1 Covalent Solids: Zeolites

Zeolites are materials generally composed of silicon, aluminium, oxygen and a metal cation or proton. They have a multitude of commercial uses including catalysis and separation (e.g. they are used in oil refining to separate linear and branched alkanes). Many of these important properties are a consequence of the presence within the zeolite of channels of molecular dimensions. It is therefore natural that molecular modelling techniques should be used to investigate the intrinsic properties of such materials and the way in which they interact with adsorbates.

The size of many zeolite systems means that considerable computational resources may be required for the calculation. In some cases therefore, such as the study of adsorption

processes, the zeolite is kept rigid and attention is concentrated on the intermolecular interactions between the zeolite and the adsorbate. This is often done using a combination of van der Waals and electrostatic terms; a Lennard-Jones potential may be used for the van der Waals component, but a Buckingham-like potential is often preferred. Electrostatic interactions can be very important for zeolites. However, the partial charges used in the various published force fields can vary enormously (from $0.4e$ to as much as $1.9e$ for the silicon atoms in silicates).

It is obviously an approximation to keep the zeolite rigid, and in more complex models the structure can vary. Many of the force fields that have been developed to model zeolites are very similar to the valence force fields used for organic and biological molecules, typically containing bond-stretching, angle-bending and torsional terms in addition to the non-bonded interactions. One important consideration when modelling zeolites is that very little energy is required to deform the Si–O–Si bond over an extremely wide range (at least 120° to 180°). This is shown in Figure 4.50, which shows the results of *ab initio* calculations using a 3-21G* basis set for $\text{H}_3\text{SiOSiH}_3$. The Fourier series expansions used by the UFF and SHAPES force fields for the angle-bending terms are designed to cope with such angular variation; Nicholas, Hopfinger, Trouw and Iton suggested the following quartic potential as an alternative specifically for the Si–O–Si angle [Nicholas *et al.* 1991]:

$$v(\theta) = \frac{k_1}{2}(\theta - \theta_0)^2 + \frac{k_2}{2}(\theta - \theta_0)^3 + \frac{k_3}{2}(\theta - \theta_0)^4 \quad (4.110)$$

With the correct choice of the parameters k_i and θ_0 the *ab initio* data in Figure 4.50 could be reproduced very well. In this force field a Urey–Bradley term was also included between the silicon atoms in such angles to model the lengthening of the Si–O bond as the angle decreased.

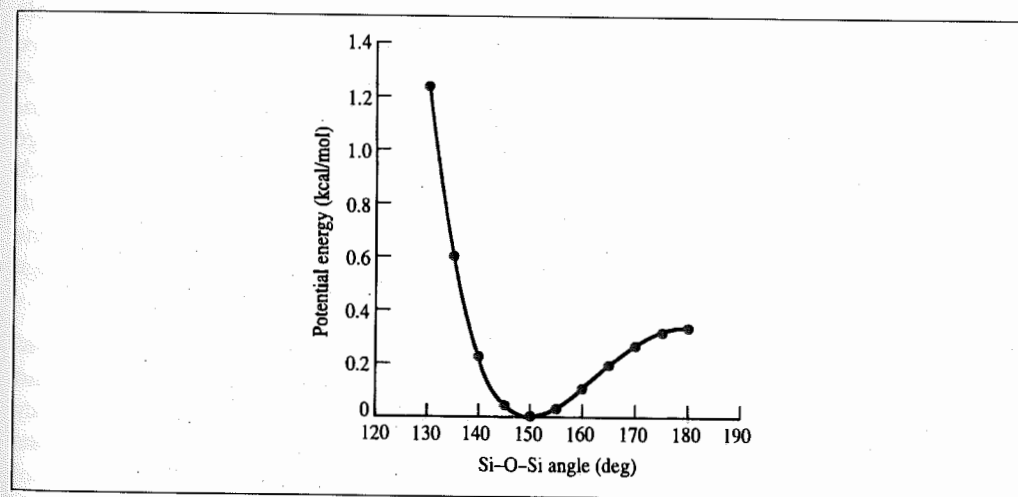


Fig. 4.50: Variation in energy with the Si–O–Si angle. (Figure redrawn from Grigoras S and T H Lane 1988. *Molecular Parameters for Organosilicon Compounds Calculated from Ab Initio Computations*. Journal of Computational Chemistry 9:25–39.)

4.22.2 Ionic Solids

The covalent approach is rarely appropriate for ionic and polar solids such as oxides and halides. The usual starting point for studying such systems is to write the potential as a series expansion of pairwise, three-body, etc., terms:

$$\mathcal{V} = \mathcal{V}_0 + \sum_{i=1}^N \sum_{j=i+1}^N v_{ij}(r) + \sum_{i=1}^N \sum_{j=i+1}^N \sum_{k=j+1}^N v_{ijk}(r) + \dots \quad (4.111)$$

One of the oldest of such models is due to Born [Born 1920], who restricted the series to pairwise terms, which were in turn divided into long-range Coulomb interactions and short-range repulsive forces. If an inverse power law is used for the repulsive term the potential energy is thus:

$$\mathcal{V} = \sum_{i=1}^N \sum_{j=i+1}^N \left(\frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} + \frac{A}{r_{ij}^n} \right) \quad (4.112)$$

The simplest way to apply such an equation is to assume that the charges q are equal to the oxidation states of the relevant species and that the repulsive potential only acts between nearest neighbours (though in common with many solid-state calculations the long-range ionic interaction is generally calculated for all possible interactions using an approach such as the Ewald sum, Section 6.8). This only leaves the two parameters A and n whose determination in principle requires only two pieces of experimental data (though the values obtained may vary quite considerably depending upon which data is chosen). An obvious extension of the simple form of Equation (4.112) is to model the short-range interactions by an alternative functional form; the Buckingham potential is commonly employed.

For a simple material such as sodium chloride the oxidation state assumption is a reasonable one. However, for other systems this is not necessarily the case. Various methods have been proposed for determining appropriate sets of non-integral charges. One strategy is to examine the distribution of charge within the material, as can be obtained from high-resolution X-ray experiments. However, there is no unique way to partition the charge unless there is zero bonding overlap between the ions. The atoms-in-molecules approach (see Section 2.7.7) may be a good way to do this but this is not the only option. It is worth mentioning that one advantage of the formal charge approach is that it can facilitate the transferability of potentials from one material to another whilst still maintaining charge neutrality.

The Born model with integral or partial charges assumes that the ions have zero polarisability. This is reasonable for small cations such as Li^+ or Mg^{2+} but can introduce significant errors for other systems. One property that clearly demonstrates this is the high-frequency dielectric constant. At a suitably high frequency only the electrons can keep up with the external field and the dielectric constant is given by the Clausius-Mosotti relationship:

$$\frac{(\epsilon_r - 1)}{(\epsilon_r + 2)} = \frac{4\pi}{3V_m} \sum_{i=1}^N \alpha_i \quad (4.113)$$

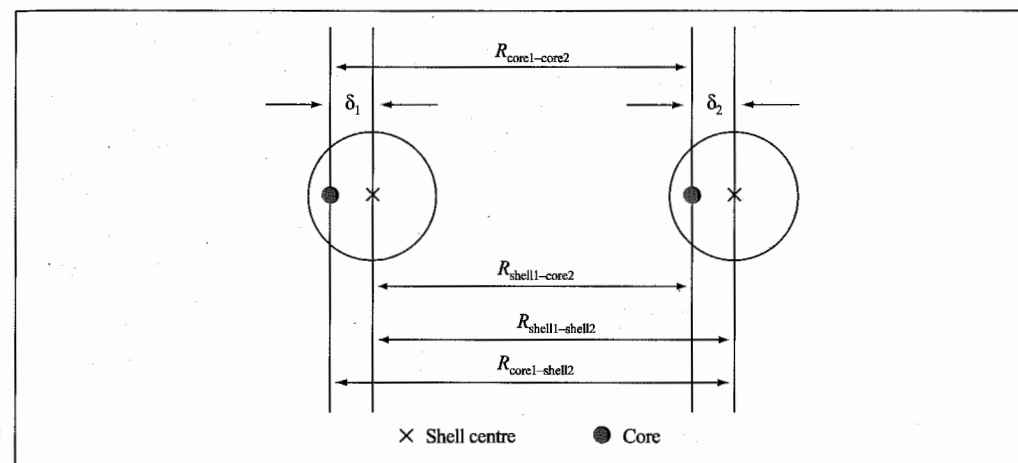


Fig. 4.51: The Dick-Overhauser shell model.

ϵ_r is the relative permittivity, V_m is the molar volume and α_i is the polarisability of the i th ion with the sum being over the N ions. If the ions were not polarisable then ϵ_r would have a value of 1. As we have seen, one way to incorporate polarisation is to assign a point polarisability to each ion. However, this model does not often give good results, at least for certain properties. This is because it fails to account for the coupling between polarisation and short-range repulsion effects. Thus polarisation causes distortions in the distribution of the valence electrons, and short-range repulsion is itself a consequence of the overlap between such electrons. The overall effect of short-range repulsion is to reduce polarisation effects. One model that can take this coupling into account is the shell model of Dick and Overhauser [Dick and Overhauser 1958] (Figure 4.51). In this model the ion is represented by a massive core linked to a massless shell by a harmonic spring. Both the core and the shell have charges associated with them. In an electric field the shell retains its charge but moves with respect to the core. The polarisability of an isolated ion in this model is proportional to Y^2/k where k is the spring constant of the harmonic spring and Y is the charge on the shell. The electrostatic interaction energy equals the sum over all ions and shells, not counting any interaction between an ion and its own shell. Although it is appealing to assume that the shells somehow play the role of the valence electrons this is probably an over-interpretation if only due to the fact that the shell charges, Y , do not necessarily assume small negative values.

Three-body and higher terms are sometimes incorporated into solid-state potentials. The Axilrod-Teller term is the most obvious way to achieve this. For systems such as the alkali halides this makes a small contribution to the total energy. Other approaches involve the use of terms equivalent to the harmonic angle-bending terms in valence force fields; these have the advantage of simplicity but, as we have already discussed, are only really appropriate for small deviations from the equilibrium bond angle. Nevertheless, it can make a significant difference to the quality of the results in some cases.

As for molecular systems, the parameters used to study the solid state can be derived using both experimental and theoretical data. There is a long tradition of using quantum mechanical calculations to extract such potentials. Whereas it is now common for the sophisticated Hartree-Fock and density functional theory approaches to be used for such parameter derivations, an approach called electron gas theory (a crude version of density functional theory) played a significant historical role and is still used [Allan and Mackrodt 1994]. One example of the way in which *ab initio* quantum mechanical calculations can play a role in this process is provided by the derivation of a potential model for α -Al₂O₃ [Gale *et al.* 1992]. Previous attempts to derive empirical potentials for this material (using a shell model combined with a Buckingham potential) were not entirely successful; in particular these did not correctly predict that the corundum structure should have the lowest energy. One interesting feature of these earlier parameterisations was the great variation in the core and shell charges; for example, in one of the models the aluminium core and shell charges were 1.617 and 1.383 respectively; in another they were 10.6063 and -8.0563. A feature of the periodic Hartree-Fock calculations (see Section 3.8.3) was the use of distorted structures to provide more information on the nature of the energy surface, which was found to give better results.

4.23 Empirical Potentials for Metals and Semiconductors

Perhaps the most important consideration when discussing the development and use of empirical potentials for studying atomic solids is that pairwise potential models are often not very suitable. The performance of pairwise potential models can be bad for transition metals and even worse for semiconductors! There are a number of reasons why this is so, many of which are due to the fundamental behaviour of pairwise potentials for certain experimental properties. The most oft-quoted properties are as follows:

1. The ratio between the cohesive energy and the melting temperature, $E_c/k_B T$. The cohesive energy is the energy cost of removing an atom from within the solid matrix. This ratio is observed to be approximately 30 in metals but about 10 in pairwise systems.
2. The ratio between the vacancy formation energy and the cohesive energy, E_v/E_c . This ratio is between $\frac{1}{4}$ and $\frac{1}{3}$ in metals but closer to unity in two-body systems (exactly 1 if the structure is not permitted to relax). This can be understood as follows. Suppose each atom in a solid has Z neighbours. If one of the atoms is removed then the coordination of the surrounding Z atoms will fall to $Z - 1$. Using a pairwise energy model the vacancy formation energy is thus Z times the atom-atom bond energy. The cohesive energy is the energy to reduce the coordination of an atom from Z to zero and so would also equal Z times the atom-atom bond energy. The energy change for both of these processes is thus equal for the pairwise model.
3. The ratio between the elastic constants C_{12}/C_{44} . Elastic constants will be discussed in Section 5.10; for a cubic solid there are three distinct values, which are labelled C_{11} , C_{12} and C_{44} . For a two-body system the ratio is exactly 1 (this is known as the Cauchy relationship). For metals and oxides deviation from unity is common; gold has a particularly high value, which is indicative of its high malleability.

4. The surface properties of metals are such that the surface tends to relax inwards but systems described by two-body interactions tend to relax outwards.

The main reason for the failure of pairwise potentials is that they are unable to deal simultaneously with both surface and bulk environments. Thus on the surface there are generally fewer bonds, but these tend to be stronger than in the bulk, where there are more, but weaker, bonds. Several many-body potentials have been devised to try to address this problem. Many of these potentials have a similar, sometimes mathematically equivalent, functional form. This reflects their common origins in some form of quantum mechanical description of bonding. However, they differ in their underlying approach, the degree to which they conform to these quantum mechanical origins and the way in which they are parametrised. Here we will outline various models: the Finnis-Sinclair model (and the Sutton-Chen extension), the embedded-atom model, the Stillinger-Weber model and the Tersoff model.

The origins of the Finnis-Sinclair potential [Finnis and Sinclair 1984] lie in the density of states and the *moments theorem*. Recall that the density of states $D(E)$ (see Section 3.8.5) describes the distribution of electronic states in the system. $D(E)$ gives the number of states between E and $E + \delta E$. Such a distribution can be described in terms of its *moments*. The moments are usually defined relative to the energy of the atomic orbital from which the molecular orbitals are formed. The m th moment, μ^m , is given by:

$$\mu^m = \sum_n (E - E_{\text{atomic}})^m D(E) \quad (4.114)$$

The summation runs over the molecular orbitals or bonds. The first moment is the mean of the distribution. If the moments are defined relative to the atomic orbital energy then this first moment will be zero. The second moment (the sum of the squares of the deviations) is the width of the distribution (the variance). The third moment describes how skewed the distribution is about the mean. If all the moments are known then the distribution can be completely characterised. Of these various moments one would expect the second to be most related to the binding energy, as this indicates how much the energy levels in the solid differ from those in the atom. Indeed, a high correlation is found to exist between the binding energy and the square root of the second moment. Armed with this relationship it would be possible to predict the binding energy for perfect lattices where the atomic environments were identical. However, a more useful model is one based on a local atomic environment ('real' materials contain features such as surfaces and defects). This requires a local density of states to be defined for each atom, $d_i(E)$, where the contribution of each molecular orbital is weighted by the amount of the orbital on the atom. In a linear combination of atomic orbitals (LCAO) model this weight is the sum of the squares of the basis set coefficients for those atomic orbitals centred on the atom. The global density of states is equal to the sum of the local densities of states over all atoms and the electronic binding energy for each atom equals the integral of $d_i(E)E$:

$$E_i^{\text{el}} = \int d_i(E)E dE \quad (4.115)$$

Thus, if we knew the second moment of the local density of states we should be able to determine the atomic binding energy via the square root relationship. However, as quantum

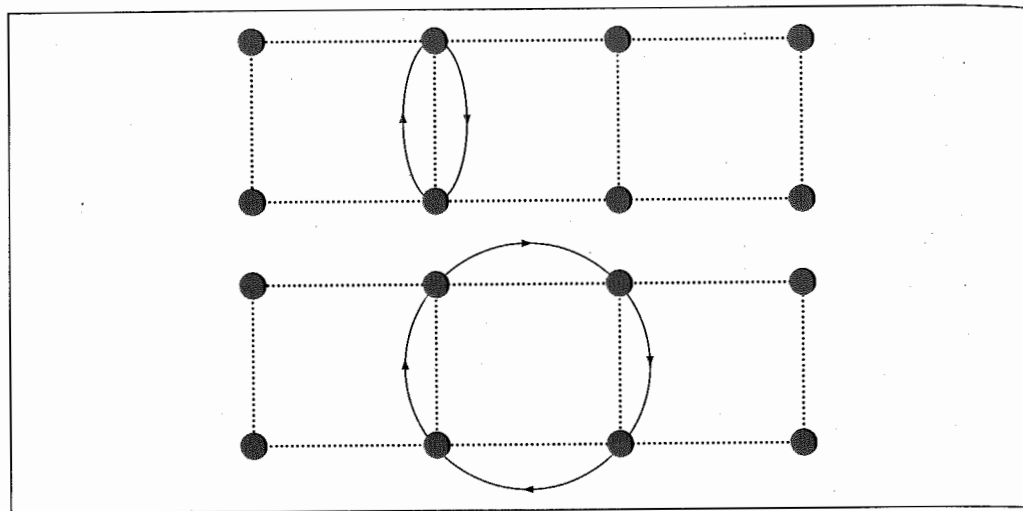


Fig. 4.52: Calculating paths using the moments theorem. Illustrated are paths of lengths 2 and 4.

mechanics is the only way we currently know of to determine the density of states, this might seem rather self-defeating. This is the role of the *moments theorem*, which relates the bonding topology to the moments of the local density of states without requiring an explicit calculation of the electronic energy levels.

The moments theorem states that the m th moment of the local density of states on an atom i is determined by the sum of all paths of length m over neighbouring atoms that start and end at i . For the second moment these paths involve just two 'hops', from the atom in question to a neighbour and back again (Figure 4.52). For the higher moments, the number of possible paths increases dramatically and becomes a challenging calculation. However, for the second moment the number of paths of length 2 is simply equal to the number of nearest neighbours, Z . Consequently, the local electronic binding energy for each atom is approximately equal to the square root of the number of neighbours. This is the *second-moment approximation*:

$$E_i^{\text{el}} \propto \sqrt{Z_i} \quad (4.116)$$

As an aside, we can easily show how this satisfies the ratio E_v/E_c (property 2, page 240). The energy E_v associated with Z atoms having their coordination reduced from Z to $Z - 1$ will be $Z[\sqrt{Z} - \sqrt{Z-1}]$. The cohesive energy E_c is proportional to \sqrt{Z} . For typical values of Z this gives E_v/E_c as approximately $\frac{1}{2}$.

In the Finnis-Sinclair potential a pairwise contribution is added to the many-body term to give the following form:

$$\mathcal{V} = \sum_{i=1}^N \sum_{j=i+1}^N P(r_{ij}) + \sum_{i=1}^N A\sqrt{\rho_i} \quad (4.117)$$

$P(r_{ij})$ is the pairwise potential, which, depending upon the model, can be considered to include electrostatic and repulsive contributions. The second term is a function of the electron density, ρ_i , and varies with the square root, in keeping with the second-moment approximation. The electron density for an atom includes contributions from the neighbouring atoms as follows:

$$\rho_i = \sum_{j=1, j \neq i}^N \phi_{ij}(r_{ij}) \quad (4.118)$$

$\phi_{ij}(r_{ij})$ is a short-range, decreasing function of the distance between the two atoms i and j . In the original Finnis-Sinclair model the function $\phi_{ij}(r_{ij})$ was written as a parabolic function of the interatomic distance, $(r_{ij} - r_c)^2$, where r_c is a cutoff distance chosen to lie between the second and third neighbouring shells. ϕ_{ij} is zero beyond this cutoff and zero beyond.

The Finnis-Sinclair potential can be written in a more general form by replacing the number of neighbouring atoms by an exponential function of the distance between atoms. This is necessary because the number of neighbours is not always straightforward to define, especially in disordered systems and near defects. An exponential function also reflects the fact that electron densities decay exponentially from the nucleus. Moreover, the pairwise potential can also be written as an exponential function of distance to give the following general equation:

$$\mathcal{V} = \sum_{i=1}^N \left\{ \sum_{j=1, j \neq i}^N A e^{-\alpha r_{ij}} - B \left[\sum_{j=1, j \neq i}^N e^{-\beta r_{ij}} \right]^{1/2} \right\} \quad (4.119)$$

Sutton and Chen extended the potential to longer range to enable the study of certain problems such as the interactions between clusters of atoms [Sutton and Chen 1990]. Their objective was to combine the superior Finnis-Sinclair description of short-range interactions with a van der Waals tail to model the long-range interactions. The form of the Sutton-Chen potential is:

$$\mathcal{V} = \varepsilon \left\{ \sum_{i=1}^N \sum_{j=i+1}^N \left(\frac{a}{r_{ij}} \right)^n - c \sum_{i=1}^N \left[\sum_{j=1, j \neq i}^N \left(\frac{a}{r_{ij}} \right)^m \right]^{1/2} \right\} \quad (4.120)$$

In this equation, ε and a are parameters with dimensions of energy and length respectively, c is a dimensionless (positive) parameter, and m and n are integers such that n is greater than m . The use of power-law relationships in the Sutton-Chen potential has a number of useful consequences, analogous to the scaling properties of the Lennard-Jones potential. For example, for a given crystal structure (e.g. hexagonal close-packed, face-centred cubic, body-centred cubic, etc.) the value of c is fixed. Moreover, if two metals are described by the same values of m and n then the results for one system may be converted directly to the other by rescaling the energy and length parameters ε and a . Typical values for m are between 6 and 8 and for n between 9 and 12.

The embedded-atom method [Daw and Baskes 1984] is an empirical embodiment of a simplified quantum mechanical model for bonding in solids called *effective medium*

theory. The key feature of effective medium theory is the replacement of the complex environment around each atom by a simplified model known as jellium. The jellium environment corresponds to a homogeneous electron gas with a positive background. Each atom is considered to be surrounded by a sphere with a radius such that the electronic charge within each sphere due to the background jellium is equal and opposite to the charge on the atom. In the embedded-atom method the background electron density is replaced by a sum of electron densities from the neighbouring atoms. The many-body term is known as an *embedding function*; this gives the energy of each atom as a function of the electron density, ρ_i . In the embedded-atom method the electron density ρ_i equals the sum of the electron densities ϕ_{ij} from neighbouring atoms (Equation (4.118)). In the Daw and Baskes model a Coulomb potential was used for the pairwise potential but with an effective charge $Z(r)$ that decreases gradually with internuclear distance. The embedding function was represented with a cubic spline equation that has a single minimum and goes to zero at vanishing density. The densities were obtained from quantum mechanical calculations.

Both the Finnis–Sinclair and the embedded-atom potentials (together with others that we have not considered here) can be represented using a very similar functional form. However, it is important to realise that they differ in the way that they connect to the first-principles, quantum mechanical model of bonding. They also differ in the procedures used to parametrise the models, so that different parametrisations may be reported for the same material.

The construction of empirical potentials for semiconductors is considered to be an even greater challenge than for metals. In our earlier discussion of the use of density functional methods to determine the electronic structure of the group 14 elements carbon, silicon and germanium we referred to the fact that, whilst the most stable form of silicon is the diamond structure, as pressure is applied so new structures can be obtained. That such a variety of structures can be achieved indicates that they are rather close in energy. Another interesting property of silicon is that in the liquid form it is a metal and the liquid is more dense than the solid. Two of the potentials that have been applied to these systems are the Stillinger–Weber and the Tersoff potentials. The Stillinger–Weber potential [Stillinger and Weber 1985] uses a two-body and three-body term:

$$\mathcal{V} = \sum_{i=1}^N \sum_{j=i+1}^N f_2(r_{ij}) + \sum_{i=1}^N \sum_{j=i+1}^N \sum_{k=j+1}^N [h(r_{ij}, r_{ik}, \theta_{ijk}) + h(r_{ji}, r_{jk}, \theta_{ijk}) + h(r_{ki}, r_{kj}, \theta_{ijk})] \quad (4.121)$$

$$f_2(r_{ij}) = A(Br_{ij}^{-p} - r_{ij}^{-q}) \exp[(r_i - a)^{-1}] \quad (4.122)$$

$$h(r_{ij}, r_{ik}, \theta_{ijk}) = \lambda \exp[\gamma(r_{ij} - a)^{-1} + \gamma(r_{ik} - a)^{-1}] (\cos \theta_{ijk} + \frac{1}{3})^2 \quad (4.123)$$

These equations all use distances and energies in reduced units and the functional form is designed to go to zero without discontinuities at the cutoff distance $r = a$. There are seven parameters ($A, B, p, q, a, \lambda, \gamma$), which were determined by a search procedure, with care being taken to ensure that the diamond structure was the most stable periodic arrangement and that the melting point and liquid structure (as determined by molecular dynamics simulations) were in reasonable agreement with experiment. The three-body term is

designed to favour the tetrahedral geometry found in the diamond structure, which is why it works reasonably well for this form of crystalline silicon. However, it does not perform so well for the other solid forms, which have a different atomic geometry, or for other properties such as the liquid structure.

The Tersoff potential [Tersoff 1988] is based on a model known as the *empirical bond-order potential*. This potential can be written in a form very similar to the Finnis–Sinclair potential:

$$\mathcal{V} = \sum_{i=1}^N \left\{ \sum_{j=1, j \neq i}^N A e^{\alpha r_{ij}} - b_{ij} B e^{-\beta r_{ij}} \right\} \quad (4.124)$$

The key term is b_{ij} , which is the bond order between the atoms i and j . This parameter depends upon the number of bonds to the atom i ; the strength of the 'bond' between i and j decreases as the number of bonds to the atom i increases. The original bond-order potential [Abell 1985] is mathematically equivalent to the Finnis–Sinclair model if the bond order b_{ij} is given by:

$$b_{ij} = \left(1 + \sum_{k=1; k \neq i, k \neq j}^N e^{-\beta(r_{ik} - r_{ij})} \right)^{-1/2} \quad (4.125)$$

It can be readily confirmed that b_{ij} decreases as the number of bonds N increases and/or their length (r_{ik}) decreases. This relationship between the bond strength and the number of neighbours provides a useful way to rationalise the structure of solids. Thus the high coordination of metals suggests that it is more effective for them to form more bonds, even though each individual bond is weakened as a consequence. Materials such as silicon achieve the balance for an intermediate number of neighbours and molecular solids have the smallest atomic coordination numbers.

The Tersoff potential was designed specifically for the group 14 elements and extends the basic empirical bond-order model by including an angular term. The interaction energy between two atoms i and j using this potential is:

$$v_{ij} = f_C(r_{ij}) [A e^{-\lambda_1 r_{ij}} - b_{ij} B e^{-\lambda_2 r_{ij}}]$$

where

$$b_{ij} = (1 + \beta^n \zeta_{ij}^n)^{-1/2n}; \quad \zeta_{ij} = \sum_{k \neq i, j} f_C(r_{ik}) g(\theta_{ijk}) \exp[\lambda_3^3 (r_{ij} - r_{ik})^3] \quad (4.126)$$

$$g(\theta) = 1 + \frac{c^2}{d^2} - \frac{c^2}{[d^2 + (h - \cos \theta)^2]}$$

The function f_C is a smoothing function with the value 1 up to some distance r_{ij} (typically chosen to include just the first neighbour shell) and then smoothly tapers to zero at the cutoff distance. b_{ij} is the bond-order term, which incorporates an angular term dependent upon the bond angle θ_{ijk} . The Tersoff potential is more broadly applicable than the Stillinger–Weber potential, but does contain more parameters.

Appendix 4.1 The Interaction Between Two Drude Molecules

In the system comprising two Drude molecules (see Section 4.9.1), an additional term must be included in the Hamiltonian [Rigby *et al.* 1986]. This additional term arises from the interactions between the two dipoles. The instantaneous dipole of each molecule is $qz(t)$, where $z(t)$ is the separation of the charges. Thus, if we label the molecules 1 and 2, we can write the dipole-dipole interaction energy as:

$$v(\mu_1, \mu_2) = -\frac{2\mu_1\mu_2}{4\pi\epsilon_0 r^3} = -\frac{2z_1z_2q^2}{4\pi\epsilon_0 r^3} \quad (4.127)$$

r is the separation of the two molecules. The Schrödinger equation for this system is thus:

$$-\frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial z_1^2} - \frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial z_2^2} + \left[\frac{1}{2}kz_1^2 + \frac{1}{2}kz_2^2 - \frac{2z_1z_2q^2}{4\pi\epsilon_0 r^3} \right] \psi = E\psi \quad (4.128)$$

This equation can be solved by making the following substitutions:

$$a_1 = \frac{z_1 + z_2}{\sqrt{2}}; \quad a_2 = \frac{z_1 - z_2}{\sqrt{2}}; \quad k_1 = k - \frac{2q^2}{4\pi\epsilon_0 r^3}; \quad k_2 = k + \frac{2q^2}{4\pi\epsilon_0 r^3} \quad (4.129)$$

These reduce Equation (4.128) to

$$-\frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial a_1^2} - \frac{\hbar^2}{2m} \frac{\partial^2 \psi}{\partial a_2^2} + \left[\frac{1}{2}k_1a_1^2 + \frac{1}{2}k_2a_2^2 \right] \psi = E\psi \quad (4.130)$$

This is the Schrödinger equation for two independent (i.e. non-interacting) oscillators with frequencies given as follows:

$$\omega_1 = \omega \sqrt{1 - \frac{2q^2}{4\pi\epsilon_0 r^3 k}}; \quad \omega_2 = \omega \sqrt{1 + \frac{2q^2}{4\pi\epsilon_0 r^3 k}} \quad (4.131)$$

$\omega/2\pi$ is the frequency of an isolated Drude molecule. The ground state energy of the system is therefore just the sum of the zero-point energies of the two oscillators: $E_0 = \frac{1}{2}\hbar(\omega_1 + \omega_2)$.

If we now substitute for ω_1 and ω_2 and expand the square roots using the binomial theorem, then we obtain the following:

$$E_0(r) = \hbar\omega - \frac{q^4\hbar\omega}{2(4\pi\epsilon_0)^2 r^6 k^2} - \dots \quad (4.132)$$

The interaction energy of the two oscillators is the difference between this zero-point energy and the energy of the system when the oscillators are infinitely separated and so:

$$v(r) = -\frac{q^4\hbar\omega}{2(4\pi\epsilon_0)^2 r^6 k^2} \quad (4.133)$$

The force constant, k , is related to the polarisability of the molecule, α as follows. Suppose a single Drude molecule is exposed to an external electric field \mathbf{E} . In the electric field, a force qE acts on each charge (in opposite directions as the charges are of opposite sign). This force causes the charges to separate and equilibrium is reached when the restoring force due to the stretching of the bond (kz) is equal to the electrostatic force: $qE = kz$. This separation

of the charges is equivalent to a static dipole given by $\mu_{ind} = qz = q^2E/k$. However, the induced dipole is also related to the polarisability by $\mu_{ind} = \alpha E$. Thus the polarisability can be written in terms of the force constant k : $\alpha = q^2/k$. With this substitution the result for the Drude model in two dimensions is:

$$v(r) = -\frac{\alpha^4\hbar\omega}{2(4\pi\epsilon_0)^2 r^6} \quad (4.134)$$

In three dimensions the equivalent result is:

$$v(r) = -\frac{3\alpha^4\hbar\omega}{4(4\pi\epsilon_0)^2 r^6} \quad (4.135)$$

Further Reading

- Bowen J P and N L Allinger 1991. Molecular Mechanics: The Art and Science of Parameterisation. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 2. New York, VCH Publishers, pp. 81-97.
- Brenner D W, O A Shendreuva and D A Areshkin 1998. Quantum-Based Analytic Interatomic Forces and Materials Simulation. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 12. New York, VCH Publishers, pp. 207-239.
- Burkert U and N L Allinger 1982. *Molecular Mechanics*. ACS Monograph 177. Washington D.C., American Chemical Society.
- Dykstra C E 1993. Electrostatic Interaction Potentials in Molecular Force Fields. *Chemical Reviews* 93:2339-2353.
- Landis C R, D M Root and T Cleveland 1995. Molecular Mechanics Force Fields for Modeling Inorganic and Organometallic Compounds. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 6. New York, VCH Publishers, pp. 73-148.
- Niketic S R and K Rasmussen 1977. *The Consistent Force Field: A Documentation*. Berlin, Springer-Verlag.
- Price S L 2000. Towards More Accurate Model Intermolecular Potentials for Organic Molecules. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 14. New York, VCH Publishers, pp. 225-289.
- Rigby M, E B Smith, W A Wakeham and G C Maitland 1981. *Intermolecular Forces: Their Origin and Determination*. Oxford, Clarendon Press.
- Rigby M, E B Smith, W A Wakeham and G C Maitland 1986. *The Forces Between Molecules*. Oxford, Clarendon Press.
- Van der Graaf B, S L Njo and K S Smirnov 2000. Introduction to Zeolite Modeling. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 14. New York, VCH Publishers, pp. 137-223.
- Williams D E 1991. Net Atomic Charge and Multipole Models for the *Ab Initio* Molecular Electric Potential. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 2. New York, VCH Publishers, pp. 219-271.

References

- Abell G C 1985. Empirical Chemical Pseudopotential Theory of Molecular and Metallic Bonding. *Physical Review* B31:6184-6196.

- Allan, N L and W C Mackrodt 1994. Density Functional Theory and Interionic Potentials. *Philosophical Magazine* **B69**:871-878.
- Allinger N L 1977. Conformational Analysis 130. MM2. A Hydrocarbon Force Field Utilizing V_1 and V_2 Torsional Terms. *Journal of the American Chemical Society* **99**:8127-8134.
- Allinger N L, K Chen and J-H Lii 1996a. An Improved Force Field (MM4) for Saturated Hydrocarbons. *Journal of Computational Chemistry* **17**:642-668.
- Allinger N L, K Chen, J A Katzenelenbogen, S R Wilson and G M Anstead 1996b. Hyperconjugative Effects on Carbon-Carbon Bond Lengths in Molecular Mechanics (MM4). *Journal of Computational Chemistry* **17**:747-755.
- Allinger N L, F Li and L Yan 1990a. Molecular Mechanics. The MM3 Force Field for Alkenes. *Journal of Computational Chemistry* **11**:848-867.
- Allinger N L, F Li, L Yan and J C Tai 1990b. Molecular Mechanics (MM3) Calculations on Conjugated Hydrocarbons. *Journal of Computational Chemistry* **11**:868-895.
- Allinger N L and J T Sprague 1973. Calculation of the Structures of Hydrocarbons Containing Delocalised Electronic Systems by the Molecular Mechanics Method. *Journal of the American Chemical Society* **95**:3893-3907.
- Allinger N L, Y H Yuh and J-J Lii 1989. Molecular Mechanics. The MM3 Force Field for Hydrocarbons I. *Journal of the American Chemical Society* **111**:8551-9556.
- Allured V S, C M Kelly and C R Landis 1991. SHAPES Empirical Force-Field - New Treatment of Angular Potentials and Its Application to Square-Planar Transition-Metal Complexes. *Journal of the American Chemical Society* **113**:1-12.
- Barker J A, R A Fisher and R O Watts 1971. Liquid Argon: Monte Carlo and Molecular Dynamics Calculations. *Molecular Physics* **21**:657-673.
- Barnes P, J L Finney, J D Nicholas and J E Quinn 1979. Cooperative Effects in Simulated Water. *Nature* **282**:459-464.
- Bayly C I, P Cieplak, W D Cornell and P A Kollman 1993. A Well-Behaved Electrostatic Potential Based Method for Deriving Atomic Charges - The RESP Model. *Journal of Physical Chemistry* **97**:10269-10280.
- Berendsen H C, J P M Postma, W F van Gunsteren and J Hermans 1981. Interaction Models for Water in Relation to Protein Hydration. In Pullman B (Editor). *Intermolecular Forces*. Dordrecht, Reidel, pp. 331-342.
- Berendsen H J C, J R Grigera and T P Straatsma 1987. The Missing Term in Effective Pair Potentials. *Journal of Physical Chemistry* **91**:6269-6271.
- Bernal J D and R H Fowler 1933. A Theory of Water and Ionic Solution, with Particular Reference to Hydrogen and Hydroxyl Ions. *Journal of Chemical Physics* **1**:515-548.
- Bezler B H, K M Merz Jr and P A Kollman 1990. Atomic Charges Derived from Semi-Empirical Methods. *Journal of Computational Chemistry* **11**:431-439.
- Born M 1920. Volumen and Hydratationswärme der Ionen. *Zeitschrift für Physik* **1**:45-48.
- Breneman C M and K B Wiberg 1990. Determining Atom-Centred Monopoles from Molecular Electrostatic Potentials. The Need for High Sampling Density in Formamide Conformational Analysis. *Journal of Computational Chemistry* **11**:361-373.
- Buckingham A D 1959. Molecular Quadrupole Moments. *Quarterly Reviews of the Chemical Society* **13**:183-214.
- Chirlian L E and M M Francl 1987. Atomic Charges Derived from Electrostatic Potentials: A Detailed Study. *Journal of Computational Chemistry* **8**:894-905.
- Claessens M, M Ferrario and J-P Ryckaert 1983. The Structure of Liquid Benzene. *Molecular Physics* **50**:217-227.
- Cleveland T and C R Landis 1996. Valence Bond Concepts Applied to the Molecular Mechanics Description of Molecular Shapes. 2. Applications to Hypervalent Molecules of the P-Block. *Journal of the American Chemical Society* **118**:6020-6030.

- Corey E J and J C Bailar Jr 1959. The Stereochemistry of Complex Inorganic Compounds. XXII. Stereospecific Effects in Complex Ions. *Journal of the American Chemical Society* **81**:2620-2629.
- Cornell W D, P Cieplak, C I Bayly, I R Gould, K M Merz Jr, D M Ferguson, D C Spellmeyer, T Fox, J W Caldwell and P A Kollman 1995. A Second Generation Force Field for the Simulation of Proteins, Nucleic Acids and Organic Molecules. *Journal of the American Chemical Society* **117**:5179-5197.
- Corongiu G 1992. Molecular Dynamics Simulation for Liquid Water Using a Polarisable and Flexible Potential. *International Journal of Quantum Chemistry* **42**:1209-1235.
- Cox S R and D E Williams 1981. Representation of the Molecular Electrostatic Potential by a New Atomic Charge Model. *Journal of Computational Chemistry* **2**:304-323.
- Dang L X, J E Rice, J Caldwell and P A Kollman 1991. Ion Solvation in Polarisable Water: Molecular Dynamics Simulations. *Journal of the American Chemical Society* **113**:2481-2486.
- Daw M S and M I Baskes 1984. Embedded-atom Method: Derivation and Application to Impurities, Surfaces, and Other Defects in Metals. *Physical Review* **B29**:6443-6453.
- Dick B G and A W Overhauser 1958. Theory of the Dielectric Constants of Alkali Halide Crystals. *Physical Review* **112**:90-103.
- Dinur U and A T Hagler 1991. New Approaches to Empirical Force Fields. In K B Lipkowitz and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 2. New York, VCH Publishers, pp. 99-164.
- Dinur U and A T Hagler 1995. Geometry-Dependent Atomic Charges: Methodology and Application to Alkanes, Aldehydes, Ketones and Amides. *Journal of Computational Chemistry* **16**:154-170.
- Ferency G G, C A Reynolds and W G Richards 1990. Semi-Empirical AM1 Electrostatic Potentials and AM1 Electrostatic Potential Derived Charges - A Comparison with *Ab Initio* Values. *Journal of Computational Chemistry* **11**:159-169.
- Ferguson D M 1995. Parameterisation and Evaluation of a Flexible Water Model. *Journal of Computational Chemistry* **16**:501-511.
- Finnis M W and J E Sinclair 1984. A Simple Empirical *N*-body Potential for Transition Metals. *Philosophical Magazine* **A50**:45-55.
- Fowler P W and A D Buckingham 1991. Central or Distributed Multipole Moments? Electrostatic Models of Aromatic Dimers. *Chemical Physics Letters* **176**:11-18.
- Gale J D, C R A Catlow and W C Mackrodt 1992. Periodic *Ab Initio* Determination of Interatomic Potentials for Alumina. *Modelling and Simulation in Materials Science and Engineering* **1**:73-81.
- Gasteiger J and M Marsili 1980. Iterative Partial Equalization of Orbital Electronegativity - Rapid Access to Atomic Charges. *Tetrahedron* **36**:3219-3288.
- Gay J G and B J Berne 1981. Modification of the Overlap Potential to Mimic a Linear Site-Site Potential. *Journal of Chemical Physics* **74**:3316-3319.
- Goodford P J 1985. A Computational Procedure for Determining Energetically Favorable Binding Sites on Biologically Important Macromolecules. *Journal of Medicinal Chemistry* **28**:849-857.
- Hagler A T, E Huler and S Lifson 1977. Energy Functions for Peptides and Proteins. I. Derivation of a Consistent Force Field Including the Hydrogen Bond from Amide Crystals. *Journal of the American Chemical Society* **96**:5319-5327.
- Hagler A T and S Lifson 1974. Energy Functions for Peptides and Proteins. II. The Amide Hydrogen Bond and Calculation of Amide Crystal Properties. *Journal of the American Chemical Society* **96**:5327-5335.
- Halgren T A 1992. Representation of van der Waals (vdW) Interactions in Molecular Mechanics Force Fields: Potential Form, Combination Rules, and vdW Parameters. *Journal of the American Chemical Society* **114**:7827-7843.
- Halgren T A 1996a. Merck Molecular Force Field I. Basis, Form, Scope, Parameterisation and Performance of MMFF94. *Journal of Computational Chemistry* **17**:490-519.
- Halgren T A 1996b. Merck Molecular Force Field II: MMFF94 van der Waals and Electrostatic Parameters for Intermolecular Interactions. *Journal of Computational Chemistry* **17**:520-552.

- Hill T L 1948. Steric Effects. I. Van der Waals Potential Energy Curves. *Journal of Chemical Physics* **16**:399-404.
- Hunter C A 1993. Sequence-dependent DNA structure. The role of base stacking interactions. *Journal of Molecular Biology* **230**:1024-1054.
- Hunter C A and J K M Saunders 1990. The Nature of π - π Interactions. *The Journal of the American Chemical Society* **112**:5525-5534.
- Hwang M J, T P Stockfish and A T Hagler 1994. Derivation of Class II Force Fields. 2. Derivation and Characterisation of a Class II Force Field, CFF93, for the Alkyl Functional Group and Alkane Molecules. *Journal of the American Chemical Society* **116**:2515-2525.
- Jorgensen W L, J Chandrasekhar, J D Madura, R W Impey and M L Klein 1983. Comparison of Simple Potential Functions for Simulating Liquid Water. *Journal of Chemical Physics* **79**:926-935.
- Jorgensen W L and J Pranata 1990. Importance of Secondary Interactions in Triply Hydrogen Bonded Complexes: Guanine-Cytosine vs Uracil-2,6-Diaminopyridine. *Journal of the American Chemical Society* **112**:2008-2010.
- Jorgensen W L and J Tirado-Rives 1988. The OPLS Potential Functions for Proteins - Energy Minimizations for Crystals of Cyclic-Peptides and Crambin. *Journal of the American Chemical Society* **110**:1666-1671.
- Landis C R, T Cleveland and T K Firman 1995. Making Sense of the Shapes of Simple Metal Hydrides. *Journal of the American Chemical Society* **117**:1859-1860.
- Landis C R, T K Firman, D M Root and T Cleveland 1998. A Valence Bond Perspective on the Molecular Shapes of Simple Metal Alkyls and Hydrides. *Journal of the American Chemical Society* **120**:1842-1854.
- Lifson S and A Warshel 1968. Consistent Force Field for Calculations of Conformations, Vibrational Spectra and Enthalpies of Cycloalkane and *n*-Alkane Molecules. *Journal of Chemical Physics* **49**:5116-5129.
- Lii J-H and N L Allinger 1989. Molecular Mechanics. The MM3 Force Field for Hydrocarbons. 2. Vibrational Frequencies and Thermodynamics. *Journal of the American Chemical Society* **111**:8566-8582.
- London F 1930. Zur Theori und Systematik der Molekularkräfte. *Zeitschrift für Physik* **63**:245-279.
- Luckhurst G R, R A Stephens and R W Phippen 1990. Computer Simulation Studies of Anisotropic Systems XIX. Mesophases Formed by the Gay-Berne Model Mesogen. *Liquid Crystals* **8**:451-464.
- Luque F J, F Ilas and M Orozco 1990. Comparative Study of the Molecular Electrostatic Potential Obtained from Different Wavefunctions - Reliability of the Semi-Empirical MNDO Wavefunction. *Journal of Computational Chemistry* **11**:416-430.
- Lybrand T P and P A Kollman 1985. Water-Water and Water-Ion Potential Functions Including Terms for Many Body Effects. *Journal of Chemical Physics* **83**:2923-2933.
- Maple J R, U Dinur and A T Hagler 1988. Derivation of Force Fields for Molecular Mechanics and Molecular Dynamics from *Ab Initio* Energy Surfaces. *Proceedings of the National Academy of Sciences USA* **85**:5350-5354.
- Nevins N, K Chen and N L Allinger 1996a. Molecular Mechanics (MM4) Calculations on Alkenes. *Journal of Computational Chemistry* **17**:669-694.
- Nevins N, K Chen and N L Allinger 1996b. Molecular Mechanics (MM4) Calculations on Conjugated Hydrocarbons. *Journal of Computational Chemistry* **17**:695-729.
- Nevins N, K Chen and N L Allinger 1996c. Molecular Mechanics (MM4) Vibrational Frequency Calculations for Alkenes and Conjugated Hydrocarbons. *Journal of Computational Chemistry* **17**:730-746.
- Nicholas J B, A J Hopfinger, F R Trouw and L E Iton 1991. Molecular Modelling of Zeolite Structure. 2. Structure and Dynamics of Silica Sodalite and Silicate Force Field. *The Journal of the American Chemical Society* **113**:4792-4800.
- Niesar U, G Corongiu, E Clementi, G R Keller and D K Bhattacharya 1990. Molecular Dynamics Simulations of Liquid Water Using the NCC *Ab Initio* Potential. *Journal of Physical Chemistry* **94**:7949-7956.

- Niketic S R and K Rasmussen 1977. *The Consistent Force Field: A Documentation*. Berlin, Springer-Verlag.
- Packer M J, M P Dauncey and C A Hunter 2000. Sequence-dependent DNA Structure: Dinucleotide Conformational Maps. *Journal of Molecular Biology* **295**:71-83.
- Pranata J and W L Jorgensen 1991. Computational Studies on FK506: Computational Search and Molecular Dynamics Simulations in Water. *Journal of the American Chemical Society* **113**:9483-9493.
- Price S L, R J Harrison and M F Guest 1989. An *Ab Initio* Distributed Multipole Study of the Electrostatic Potential Around an Undecapeptide Cyclosporin Derivative and a Comparison with Point Charge Electrostatic Models. *Journal of Computational Chemistry* **10**:552-567.
- Rappé A K, C J Casewit, K S Colwell, W A Goddard III and W M Skiff 1992. UFF, a Full Periodic Table Force Field for Molecular Mechanics and Molecular Dynamics Simulations. *Journal of the American Chemical Society* **114**:10024-10035.
- Rappé A K, K S Colwell and C J Casewit 1993. Application of a Universal Force Field to Metal Complexes. *Inorganic Chemistry* **32**:3438-3450.
- Rappé A K and W A Goddard III 1991. Charge Equilibration for Molecular Dynamics Simulations. *Journal of Physical Chemistry* **95**:3358-3363.
- Reynolds C A, J W Essex and W G Richards 1992. Atomic Charges for Variable Molecular Conformations. *Journal of the American Chemical Society* **114**:9075-9079.
- Rick S W and B J Berne 1996. Dynamical Fluctuating Charge Force Fields: The Aqueous Solvation of Amides. *Journal of the American Chemical Society* **118**:672-679.
- Rick S W, S J Stuart and B J Berne 1994. Dynamical Fluctuating Charge Force Fields: Application to Liquid Water. *Journal of Chemical Physics* **101**:6141-6156.
- Rigby M, E B Smith, W A Wakeham and G C Maitland 1986. *The Forces Between Molecules*. Oxford, Clarendon Press.
- Rodger P M, A J Stone and D J Tildesley 1988. The Intermolecular Potential of Chlorine. A Three Phase Study. *Molecular Physics* **63**:173-188.
- Singh U C and P A Kollman 1984. An Approach to Computing Electrostatic Charges for Molecules. *Journal of Computational Chemistry* **5**:129-145.
- Smith P E and B M Pettitt 1994. Modelling Solvent in Biomolecular Systems. *Journal of Physical Chemistry* **98**:9700-9711.
- Sprague J T, J C Tai, Y Yuh and N L Allinger 1987. The MMP2 Computational Method. *Journal of Computational Chemistry* **8**:581-603.
- Sprink M and M L Klein 1988. A Polarizable Model for Water Using Distributed Charge Sites. *Journal of Chemical Physics* **89**:7556-7560.
- Stillinger F H and A Rahman 1974. Improved Simulation of Liquid Water by Molecular Dynamics. *Journal of Chemical Physics* **60**:1545-1557.
- Stillinger F H and T A Weber 1985. Computer Simulation of Local Order in Condensed Phases of Silicon. *Physical Review* **B31**:5262-5271.
- Stone A J 1981. Distributed Multipole Analysis, or How to Describe a Molecular Charge Distribution. *Chemical Physics Letters* **83**:233-239.
- Stone A J and M Alderton 1985. Distributed Multipole Analysis Methods and Applications. *Molecular Physics* **56**:1047-1064.
- Stuart S J and B J Berne 1996. Effects of Polarizability on the Hydration of the Chloride Ion. *Journal of Physical Chemistry* **100**:11934-11943.
- Sutton A P and J Chen 1990. Long-range Finnis-Sinclair Potentials. *Philosophical Magazine Letters* **61**:139-146.
- Tersoff J 1988. New Empirical Approach for the Structure and Energy of Covalent Systems. *Physical Review* **B37**:6991-7000.
- Toxvaerd S 1990. Molecular Dynamics Calculation of the Equation of State of Alkanes. *Journal of Chemical Physics* **93**:4290-4295.

- Vedani A 1988. YETI: An Interactive Molecular Mechanics Program for Small-Molecular Protein Complexes. *Journal of Computational Chemistry* **9**:269–280.
- Vinter J G 1994. Extended Electron Distributions Applied to the Molecular Mechanics of Some Intermolecular Interactions. *Journal of Computer-Aided Molecular Design* **8**:653–668.
- Warshel A and M Karplus 1972. Calculation of Ground and Excited State Potential Surfaces of Conjugated Molecules. I. Formulation and Parameterisation. *Journal of the American Chemical Society* **94**:5612–5622.
- Warshel A and A Lippicirella 1981. Calculations for Ground- and Excited-State Potential Surfaces for Conjugated Heteroatomic Molecules. *Journal of the American Chemical Society* **103**:4664–4673.
- Weiner SJ, P A Kollman, D A Case, U C Singh, C Ghio, G Alagona, S Profeta and P Weiner 1984. A New Force Field for Molecular Mechanical Simulation of Nucleic Acids and Proteins. *Journal of the American Chemical Society* **106**:765–784.
- Williams D E 1990. Alanyl Dipeptide Potential-Derived Net Atomic Charges and Bond Dipoles, and Their Variation with Molecular Conformation. *Biopolymers* **29**:1367–1386.

CHAPTER FIVE

Energy Minimisation and Related Methods for Exploring the Energy Surface

5.1 Introduction

For all except the very simplest systems the potential energy is a complicated, multi-dimensional function of the coordinates. For example, the energy of a conformation of ethane is a function of the 18 internal coordinates or 24 Cartesian coordinates that are required to completely specify the structure. As we discussed in Section 1.3, the way in which the energy varies with the coordinates is usually referred to as the *potential energy surface* (sometimes called the *hypersurface*). In the interests of brevity all references to 'energy' should be taken to mean 'potential energy' for the rest of this chapter, except where explicitly stated otherwise. For a system with N atoms the energy is thus a function of $3N - 6$ internal or $3N$ Cartesian coordinates. It is therefore impossible to visualise the entire energy surface except for some simple cases where the energy is a function of just one or two coordinates. For example, the van der Waals energy of two argon atoms (as might be modelled using the Lennard-Jones potential function) depends upon just one coordinate: the interatomic distance. Sometimes we may wish to visualise just a part of the energy surface. For example, suppose we take an extended conformation of pentane and rotate the two central carbon-carbon bonds so that the torsion angles vary from 0° to 360° , calculating the energy of each structure generated. The energy in this case is a function of just two variables and can be plotted as a contour diagram or as an isometric plot, as shown in Figure 5.1.

We will use the term 'energy surface' to refer not only to systems in which the bonding remains unchanged, as in these two examples, but also where bonds are broken and/or formed. Our discussion will be appropriate to both quantum mechanics and molecular mechanics, except where otherwise stated.

In molecular modelling we are especially interested in minimum points on the energy surface. Minimum energy arrangements of the atoms correspond to stable states of the system; any movement away from a minimum gives a configuration with a higher energy. There may be a very large number of minima on the energy surface. The minimum with the very lowest energy is known as the *global energy minimum*. To identify those geometries of

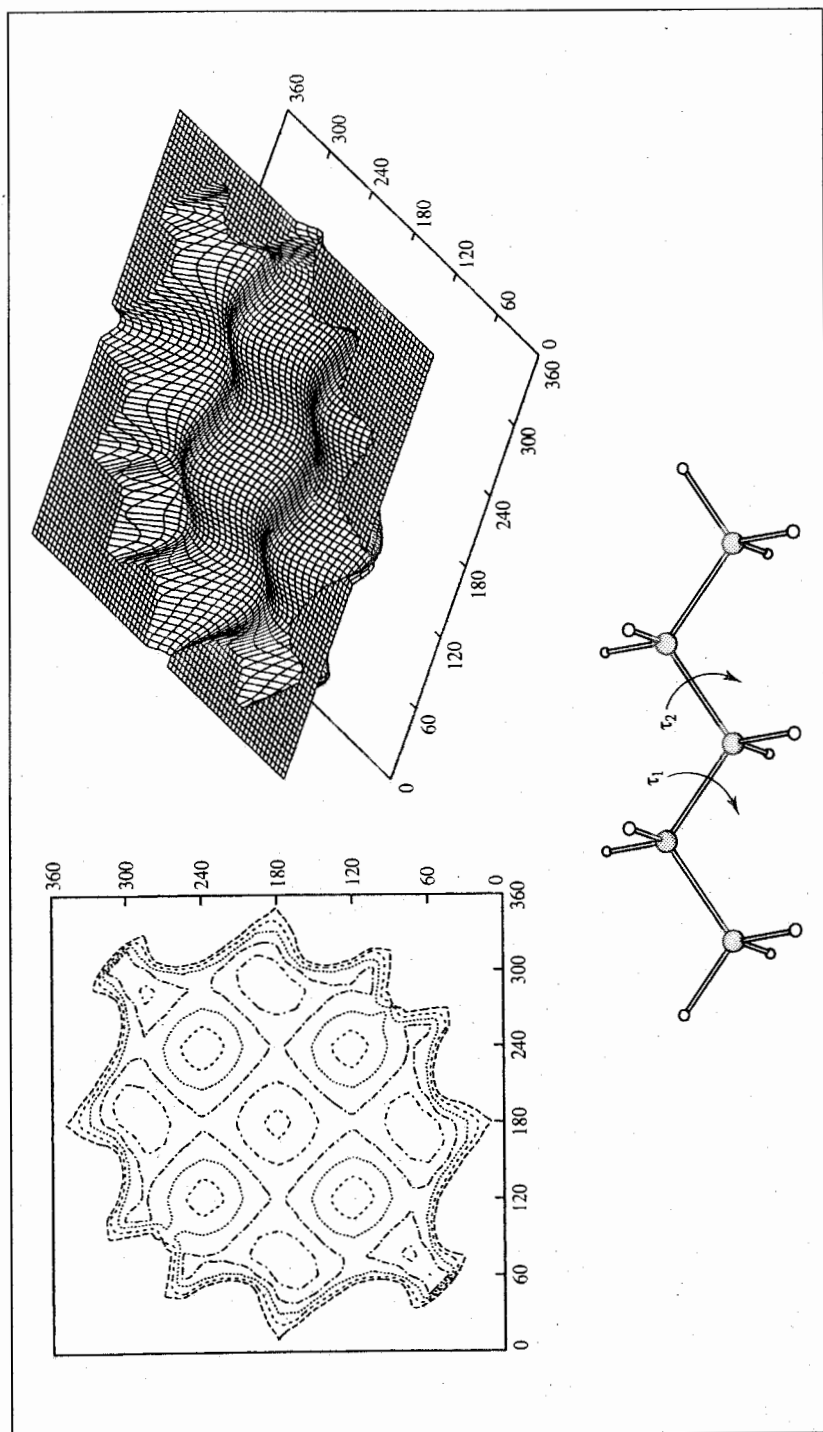


Fig. 5.1: Variation in the energy of pentane with the two torsion angles indicated and represented as a contour diagram and isometric plot. Only the lowest-energy regions are shown.

the system that correspond to minimum points on the energy surface we use a *minimisation algorithm*. There is a vast literature on such methods and so we will concentrate on those approaches that are most commonly used in molecular modelling. We may also be interested to know how the system changes from one minimum energy structure to another. For example, how do the relative positions of the atoms vary during a reaction? What structural changes occur as a molecule changes its conformation? The highest point on the pathway between two minima is of especial interest and is known as the *saddle point*, with the arrangement of the atoms being the *transition structure*. Both minima and saddle points are stationary points on the energy surface, where the first derivative of the energy function is zero with respect to all the coordinates.

A geographical analogy can be a helpful way to illustrate many of the concepts we shall encounter in this chapter. In this analogy minimum points correspond to the bottom of valleys. A minimum may be described as being in a 'long and narrow valley' or 'a flat and featureless plain'. Saddle points correspond to mountain passes. We refer to algorithms taking steps 'uphill' or 'downhill'.

5.1.1 Energy Minimisation: Statement of the Problem

The minimisation problem can be formally stated as follows: given a function f which depends on one or more independent variables x_1, x_2, \dots, x_i , find the values of those variables where f has a minimum value. At a minimum point the first derivative of the function with respect to each of the variables is zero and the second derivatives are all positive:

$$\frac{\partial f}{\partial x_i} = 0; \quad \frac{\partial^2 f}{\partial x_i^2} > 0 \quad (5.1)$$

The functions of most interest to us will be the quantum mechanics or molecular mechanics energy with the variables x_i being the Cartesian or the internal coordinates of the atoms. Molecular mechanics minimisations are nearly always performed in Cartesian coordinates, where the energy is a function of $3N$ variables; it is more common to use internal coordinates (as defined in the Z-matrix) with quantum mechanics. For analytical functions, the minimum of a function can be found using standard calculus methods. However, this is not generally possible for molecular systems due to the complicated way in which the energy varies with the coordinates. Rather, minima are located using numerical methods, which gradually change the coordinates to produce configurations with lower and lower energies until the minimum is reached. To illustrate how the various minimisation algorithms operate, we shall consider a simple function of two variables: $f(x, y) = x^2 + 2y^2$. This function is represented as a contour diagram in Figure 5.2. The function has one minimum point, located at the origin. In our examples we will attempt to locate the minimum from the point (9.0, 9.0). Although this is a function of just two variables for the purposes of illustration, all of the methods that we shall consider can be applied to functions of many more variables.

We can classify minimisation algorithms into two groups: those which use derivatives of the energy with respect to the coordinates and those which do not. Derivatives can be useful

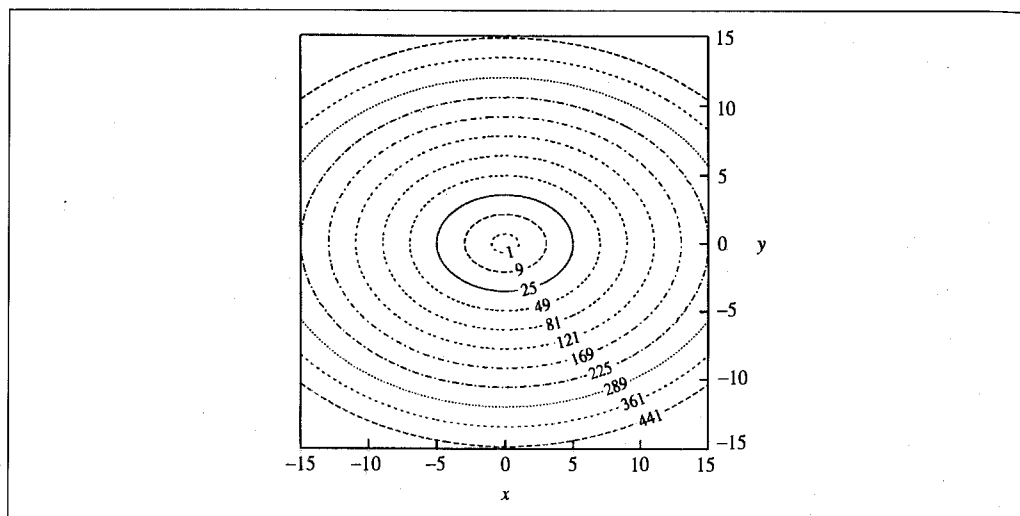


Fig. 5.2: The function $x^2 + 2y^2$.

because they provide information about the shape of the energy surface, and, if used properly, they can significantly enhance the efficiency with which the minimum is located. There are many factors that must be taken into account when choosing the most appropriate algorithm (or combination of algorithms) for a given problem; the ideal minimisation algorithm is the one that provides the answer as quickly as possible, using the least amount of memory. No single minimisation method has yet proved to be the best for all molecular modelling problems and so most software packages offer a choice of methods. In particular, a method that works well with quantum mechanics may not be the most suitable for use with molecular mechanics. This is partly because quantum mechanics is usually used to model systems with fewer atoms than molecular mechanics; some operations that are integral to certain minimisation procedures (such as matrix inversion) are trivial for small systems but formidable for systems containing thousands of atoms. Quantum mechanics and molecular mechanics also require different amounts of computational effort to calculate the energies and the derivatives of the various configurations. Thus an algorithm that takes many steps may be appropriate for molecular mechanics but inappropriate for quantum mechanics.

Most minimisation algorithms can only go downhill on the energy surface and so they can only locate the minimum that is nearest (in a downhill sense) to the starting point. Thus, Figure 5.3 shows a schematic energy surface and the minima that would be obtained starting from three points A, B and C. The minima can be considered to correspond to the locations where a ball rolling on the energy surface under the influence of gravity would come to rest. To locate more than one minimum or to locate the global energy minimum we therefore usually require a means of generating different starting points, each of which is then minimised. Some specialised minimisation methods can make uphill moves to seek out minima lower in energy than the nearest one, but no algorithm has yet proved capable of locating the

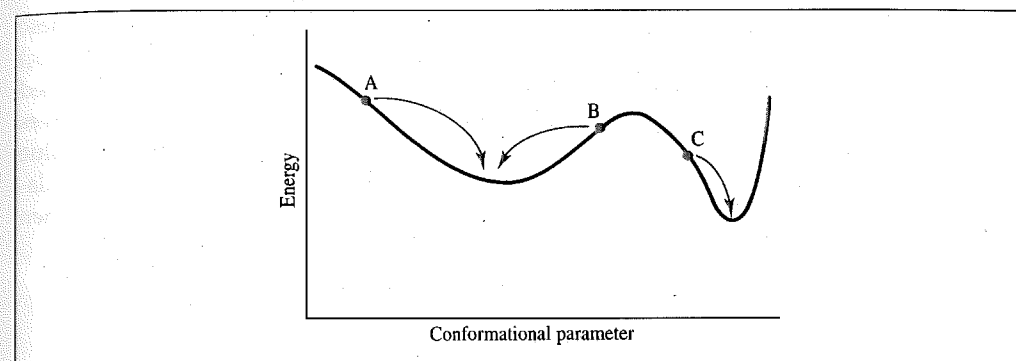


Fig. 5.3: A schematic one-dimensional energy surface. Minimisation methods move downhill to the nearest minimum. The statistical weight of the narrow, deep minimum may be less than a broad minimum which is higher in energy.

global energy minimum from an arbitrary starting position. The shape of the energy surface may be important if one wishes to calculate the relative populations of the various minimum energy structures. For example, a deep and narrow minimum may be less highly populated than a broad minimum that is higher in energy as the vibrational energy levels will be more widely spaced in the deeper minimum and so less accessible. For this reason the global energy minimum may not be the most highly populated minimum. In any case, the 'active' structure (e.g. the biologically active conformation of a drug molecule) may not correspond to the global minimum, or to the most highly populated conformation, or even to a minimum energy structure at all.

The input to a minimisation program consists of a set of initial coordinates for the system. The initial coordinates may come from a variety of sources. They may be obtained from an experimental technique, such as X-ray crystallography or NMR. In other cases a theoretical method is employed, such as a conformational search algorithm. A combination of experimental and theoretical approaches may also be used. For example, to study the behaviour of a protein in water one may take an X-ray structure of the protein and immerse it in a solvent 'bath', where the coordinates of the solvent molecules have been obtained from a Monte Carlo or molecular dynamics simulation.

5.1.2 Derivatives

In order to use a derivative minimisation method it is obviously necessary to be able to calculate the derivatives of the energy with respect to the variables (i.e. the Cartesian or internal coordinates, as appropriate). Derivatives may be obtained either analytically or numerically. The use of analytical derivatives is preferable as they are exact, and because they can be calculated more quickly; if only numerical derivatives are available then it may be more effective to use a non-derivative minimisation algorithm. The problems of calculating analytical derivatives with quantum mechanics and molecular mechanics were discussed in Sections 3.4.3 and 4.16, respectively.

Nevertheless, under some circumstances it is necessary to use numerical derivatives. These can be calculated as follows. If one of the coordinates x_i is changed by a small change (δx_i) and the energy for the new arrangement is computed then the derivative $\partial E/\partial x_i$ is obtained by dividing the change in energy (δE) by the change in coordinate ($\delta E/\delta x_i$). This strictly gives the derivative at the mid-point between the two points x_i and $x_i + \delta x_i$. A more accurate value of the derivative at the point x_i may be obtained (at the cost of an additional energy calculation) by evaluating the energy at two points, $x_i + \delta x_i$ and $x_i - \delta x_i$. The derivative is then obtained by dividing the difference in the energies by $2\delta x_i$.

5.2 Non-derivative Minimisation Methods

5.2.1 The Simplex Method

A *simplex* is a geometrical figure with $M + 1$ interconnected vertices, where M is the dimensionality of the energy function. For a function of two variables the simplex is thus triangular in shape. A tetrahedral simplex is used for a function of three variables and so for an energy function of $3N$ Cartesian coordinates the simplex will have $3N + 1$ vertices; if internal coordinates are used then the simplex will have $3N - 5$ vertices. Each vertex corresponds to a specific set of coordinates for which an energy can be calculated. For our function $f(x, y) = x^2 + 2y^2$ the simplex method would use a triangular simplex.

The simplex algorithm locates a minimum by moving around on the potential energy surface in a fashion that has been likened to the motion of an amoeba. Three basic kinds of move are possible. The most common type of move is a reflection of the vertex with the highest value through the opposite face of the simplex, in an attempt to generate a new point that has a lower value. If this new point is lower in energy than any of the other points in the simplex then a 'reflection and expansion' move may be applied. When a 'valley floor' is reached then a reflection move will fail to produce a better point. Under such circumstances the simplex contracts along one dimension from the highest point. If this fails to reduce the energy then a third type of move is possible, in which the simplex contracts in all directions, pulling around the lowest point. These three moves are illustrated in Figure 5.4.

To implement the simplex algorithm it is first necessary to generate the vertices of the initial simplex. The initial configuration of the system corresponds to just one of these vertices. The remaining points can be obtained in a variety of ways, but one simple method is to add a constant increment to each coordinate in turn. The energy of the system is calculated at the new point, giving the function value for the relevant vertex.

The simplex method is most useful where the initial configuration of the system is very high in energy, because it rarely fails to find a better solution. However, it can be rather expensive in terms of computer time due to the large number of energy evaluations which are required (merely to generate the initial simplex requires $3N + 1$ energy evaluations). For this reason the simplex method is often used in combination with a different minimisation algorithm: a few steps of the simplex method are used to refine the initial structure and then a more efficient method can take over.

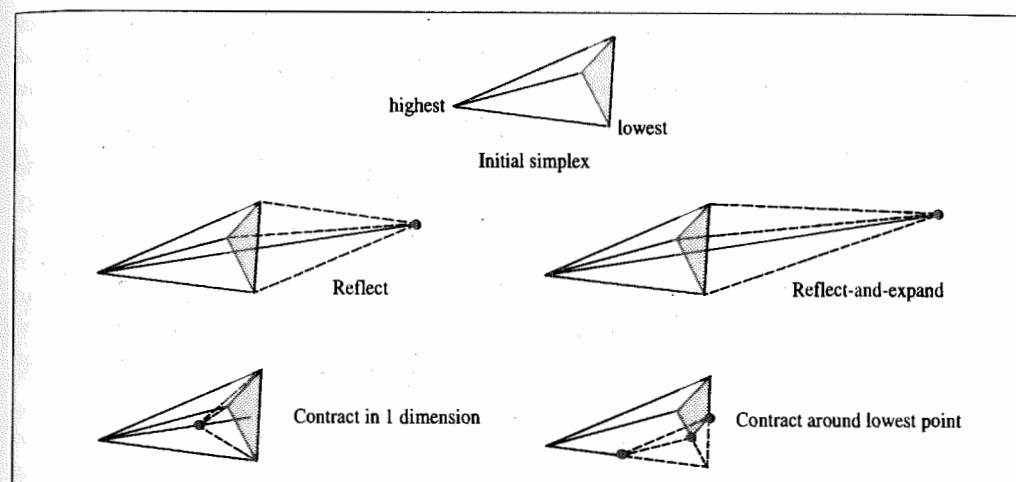


Fig. 5.4: The three basic moves permitted to the simplex algorithm (reflection, and its close relation reflect-and-expand; contract in one dimension and contract around the lowest point). (Figure adapted from Press W H, B P Flannery, S A Teukolsky and W T Vetterling 1992. Numerical Recipes in Fortran. Cambridge, Cambridge University Press.)

Let us consider the application of the simplex method to our quadratic function, $f = x^2 + 2y^2$ (Figure 5.5). Suppose our initial simplex contains vertices located at the points (9, 9), (11, 9) and (9, 11), which have been generated by adding a constant factor 2 to each of the variables in turn. The values of the function at these points are 243, 283 and 323, respectively. The vertex with the highest function value is at (9, 11) and so in the first iteration this point is reflected through the opposite face of the triangle to generate a point with coordinates (11, 7) and a function value of 219 (we do not use the reflect-and-expand move in our

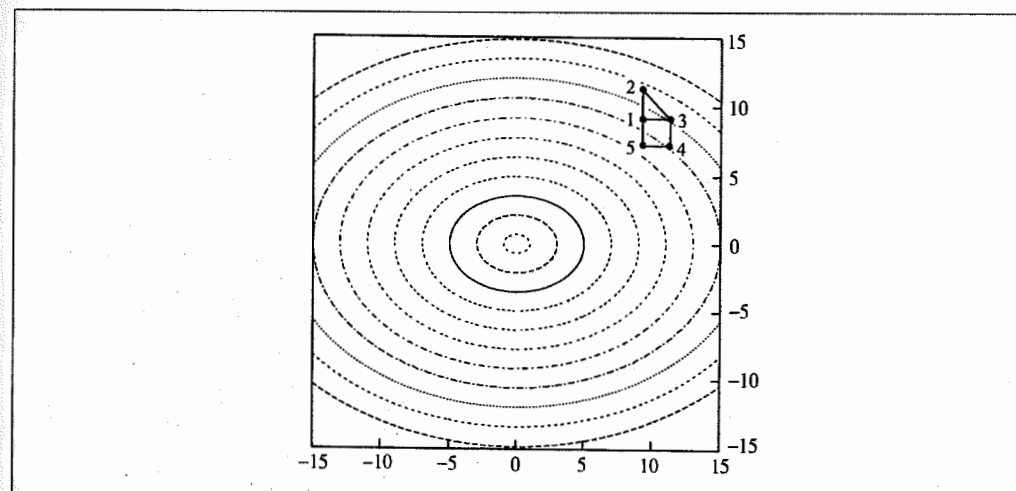


Fig. 5.5: The first few steps of the simplex algorithm with the function $x^2 + 2y^2$. The initial simplex corresponds to the triangle 123. Point 2 has the largest value of the function and the next simplex is the triangle 134. The simplex for the third step is 145.

illustration). The highest vertex is now at (11, 9), which is reflected through the opposite face of the simplex to give the point (9, 7), where the function has a value of 179. In fact, for this admittedly artificial problem the simplex algorithm takes more than 30 steps to find a point where the function has a value less than 0.1.

Why does the simplex contain one more vertex than the number of degrees of freedom? The reason is that with fewer than $M + 1$ vertices the algorithm cannot explore the whole energy surface. Suppose we use only a two-vertex simplex to explore our quadratic energy surface. A simplex with just two vertices is a straight line. The only moves that would be possible in this case would be to other points that lie on this line; none of the energy surface away from the line would be explored. Similarly, if we have a function of three variables and restrict the simplex to a triangle then we will only be able to explore the region of space that lies in the same plane as the triangle, whereas the minimum may not lie in this plane.

5.2.2 The Sequential Univariate Method

The simplex method is rarely considered suitable for quantum mechanical calculations, due to the number of energy evaluations that must be performed. The sequential univariate method is a non-derivative method that is considered more appropriate in this case. This method systematically cycles through the coordinates in turn. For each coordinate, two new structures are generated by changing the current coordinate (i.e. $x_i + \delta x_i$ and $x_i + 2\delta x_i$). The energies of these two structures are calculated. A parabola is then fitted through the three points corresponding to the two distorted structures and the original structure. The minimum point in this quadratic function is determined and the coordinate is then changed to the position of the minimum. The procedure is illustrated in Figure 5.6. When the changes in all the coordinates are

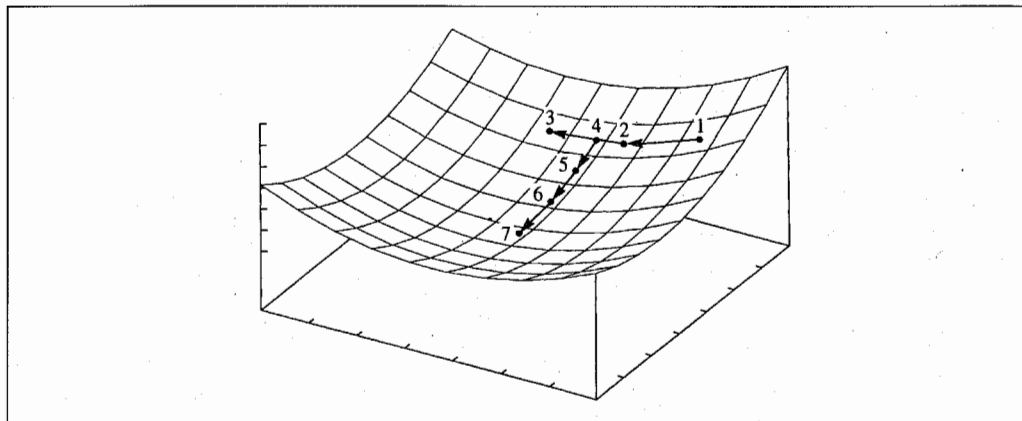


Fig. 5.6: The sequential univariate method. Starting at the point labelled 1 two steps are made along one of the coordinates to give points 2 and 3. A parabola is fitted to these three points and the minimum located (point 4). The same procedure is then repeated along the next coordinate (points 5, 6 and 7). (Figure adapted from Schlegel H B 1987. *Optimization of Equilibrium Geometries and Transition Structures*. In Lawley K P (Editor). *Ab Initio Methods in Quantum Chemistry - I*. New York, John Wiley, pp. 249-286.)

sufficiently small then the minimum is deemed to have been reached, otherwise a new iteration is performed. The sequential univariate method usually requires fewer function evaluations than the simplex method but it can be slow to converge especially if there is strong coupling between two or more of the coordinates or when the energy surface is analogous to a long narrow valley.

5.3 Introduction to Derivative Minimisation Methods

Derivatives provide information that can be very useful in energy minimisation, and derivatives are used by most popular minimisation methods. The direction of the first derivative of the energy (the gradient) indicates where the minimum lies, and the magnitude of the gradient indicates the steepness of the local slope. The energy of the system can be lowered by moving each atom in response to the force acting on it; the force is equal to minus the gradient. Second derivatives indicate the curvature of the function, information that can be used to predict where the function will change direction (i.e. pass through a minimum or some other stationary point).

When discussing derivative methods it is useful to write the function as a Taylor series expansion about the point x_k :

$$\mathcal{V}(x) = \mathcal{V}(x_k) + (x - x_k)\mathcal{V}'(x_k) + (x - x_k)^2\mathcal{V}''(x_k)/2 + \dots \quad (5.2)$$

For a multidimensional function, the variable x is replaced by the vector \mathbf{x} and matrices are used for the various derivatives. Thus if the potential energy $\mathcal{V}(\mathbf{x})$ is a function of $3N$ Cartesian coordinates, the vector \mathbf{x} will have $3N$ components and \mathbf{x}_k corresponds to the current configuration of the system. $\mathcal{V}'(\mathbf{x}_k)$ is a $3N \times 1$ matrix (i.e. a vector), each element of which is the partial derivative of \mathcal{V} with respect to the appropriate coordinate, $\partial\mathcal{V}/\partial x_i$. We will also write the gradient at the point k as \mathbf{g}_k . Each element (i, j) of the matrix $\mathcal{V}''(\mathbf{x}_k)$ is the partial second derivative of the energy function with respect to the two coordinates x_i and x_j , $\partial^2\mathcal{V}/\partial x_i\partial x_j$. $\mathcal{V}''(\mathbf{x}_k)$ is thus of dimension $3N \times 3N$ and is known as the *Hessian* matrix or the *force constant* matrix. The Taylor series expansion can be written in the following form for the multidimensional case:

$$\mathcal{V}(\mathbf{x}) = \mathcal{V}(\mathbf{x}_k) + (\mathbf{x} - \mathbf{x}_k)\mathcal{V}'(\mathbf{x}_k) + (\mathbf{x} - \mathbf{x}_k)^T \cdot \mathcal{V}''(\mathbf{x}_k) \cdot (\mathbf{x} - \mathbf{x}_k)/2 + \dots \quad (5.3)$$

The energy functions used in molecular modelling are rarely quadratic and so the Taylor series expansion, Equation (5.3), can only be considered an approximation. There are two important consequences of this. The first consequence is that the performance of a given minimisation method will not be as good for a molecular mechanics or quantum mechanics energy surface as it is for a pure quadratic function. As we shall see, a second derivative method such as the Newton-Raphson algorithm can locate the minimum in a single step for a purely quadratic function, but several iterations are usually required for a typical molecular modelling energy function. The second consequence is that, far from the minimum, the harmonic approximation is a poor one and some of the less robust methods will fail, even though they may work very well close to a minimum, where the harmonic approximation is more valid. For this reason it is important to choose the minimisation

protocol with care, possibly using a robust (but perhaps inefficient) method at first, and then a less robust but more efficient method.

The derivative methods can be classified according to the highest-order derivative used. First-order methods use the first derivatives (i.e. the gradients) whereas second-order methods use both first and second derivatives. The simplex method can thus be considered a zeroth-order method as it does not use any derivatives.

5.4 First-order Minimisation Methods

Two first-order minimisation algorithms that are frequently used in molecular modelling are the method of *steepest descents* and the *conjugate gradient* method. These gradually change the coordinates of the atoms as they move the system closer and closer to the minimum point. The starting point for each iteration (k) is the molecular configuration obtained from the previous step, which is represented by the multidimensional vector \mathbf{x}_{k-1} . For the first iteration the starting point is the initial configuration of the system provided by the user, the vector \mathbf{x}_1 .

5.4.1 The Steepest Descents Method

The steepest descents method moves in the direction parallel to the net force, which in our geographical analogy corresponds to walking straight downhill. For $3N$ Cartesian coordinates this direction is most conveniently represented by a $3N$ -dimensional unit vector, \mathbf{s}_k . Thus:

$$\mathbf{s}_k = -\mathbf{g}_k / |\mathbf{g}_k| \quad (5.4)$$

Having defined the direction along which to move it is then necessary to decide how far to move along the gradient. Consider the two-dimensional energy surface of Figure 5.7. The gradient direction from the starting point is along the line indicated. If we imagine a cross-section through the surface along the line, the function will pass through a minimum and then increase, as shown in the figure. We can choose to locate the minimum point by performing a *line search* or we can take a step of arbitrary size along the direction of the force.

5.4.2 Line Search in One Dimension

The purpose of a line search is to locate the minimum along a specified direction (i.e. along a line through the multidimensional space). The first stage of the line search is to *bracket* the minimum. This entails finding three points along the line such that the energy of the middle point is lower than the energy of the two outer points. If three such points can be found, then at least one minimum must lie between the two outer points. An iterative procedure can then be used to decrease the distance between the three points, gradually restricting the minimum to an even smaller region. This is conceptually an easy process but it may require a considerable number of function evaluations, making it computationally expensive. An alternative is to fit a function such as a quadratic to the three points. Differentiation of the

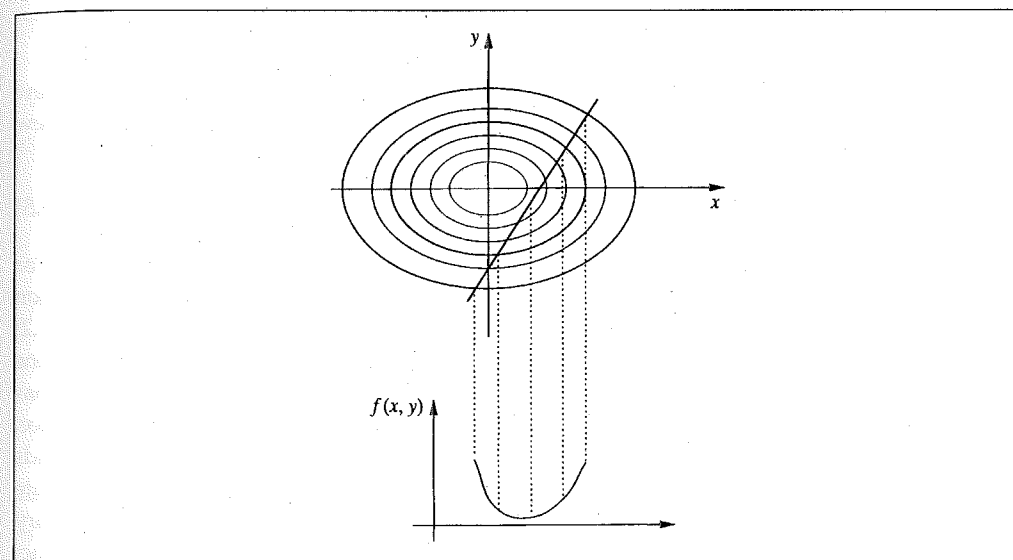


Fig. 5.7: A line search is used to locate the minimum in the function in the direction of the gradient.

fitted function enables an approximation to the minimum along the line to be identified analytically. A new function can then be fitted to give a better estimate, as shown in Figure 5.8. Higher-order polynomials may give a better fit to the bracketing points but these can give incorrect interpolations when used with functions that change sharply in the bracketed region.

The gradient at the minimum point obtained from the line search will be perpendicular to the previous direction. Thus, when the line search method is used to locate the minimum along the gradient then the next direction in the steepest descents algorithm will be orthogonal to the previous direction (i.e. $\mathbf{g}_k \cdot \mathbf{g}_{k-1} = 0$).

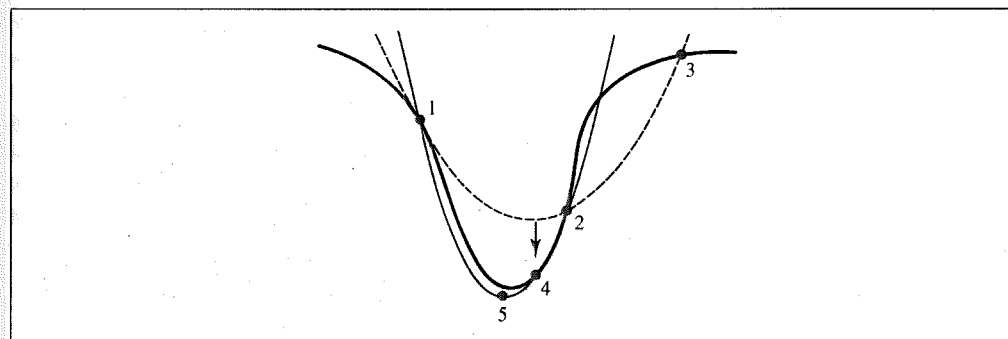


Fig. 5.8: The minimum in a line search may be found more effectively by fitting an analytical function such as a quadratic to the initial set of three points (1, 2 and 3). A better estimate of the minimum can then be found by fitting a new function to the points 1, 2 and 4 and finding its minimum. (Figure adapted from Press W H, B P Flannery, S A Teukolsky and W T Vetterling 1992. Numerical Recipes in Fortran. Cambridge, Cambridge University Press.)

5.4.3 Arbitrary Step Approach

As the line search may itself be computationally demanding we could obtain the new coordinates by taking a step of arbitrary length along the gradient unit vector \mathbf{s}_k . The new set of coordinates after step k would then be given by the equation:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \lambda_k \mathbf{s}_k \quad (5.5)$$

λ_k is the *step size*. In most applications of the steepest descents algorithm in molecular modelling the step size initially has a predetermined default value. If the first iteration leads to a reduction in energy, the step size is increased by a multiplicative factor (e.g. 1.2) for the second iteration. This process is repeated so long as each iteration reduces the energy. When a step produces an increase in energy, it is assumed that the algorithm has leapt across the valley which contains the minimum and up the slope on the opposite face. The step size is then reduced by a multiplicative factor (e.g. 0.5). The step size depends upon the nature of the energy surface; for a flat surface large step sizes would be appropriate but for a narrow, twisting gully a much smaller step would be more suitable. The arbitrary step method may require more steps to reach the minimum but it can often require fewer function evaluations (and thus less computer time) than the more rigorous line search approach.

The steepest descents method works as follows for our trial function, $f(x, y) = x^2 + 2y^2$. Differentiating the function gives $df = 2x dx + 4y dy$ and so the gradient at any point (x, y) equals $4y/2x$. The direction of the first move from the point $(9.0, 9.0)$ is $(-18.0, -36.0)$ and the equation of the line along which the search is performed is $y = 2x - 9$. The minimum of the function along this line can be obtained using Lagrange multipliers (see Section 1.10.5) and is at $(4.0, -1.0)$. The direction of the next move is the vector $(-8, 4)$ and the next line search is performed along the line $y = -0.5x + 1$. The minimum point along this line is $(2/3, 2/3)$ where the function has the value $4/3$. The third point found by the steepest descents method is at $(0.296, -0.074)$ where the function has the value 0.099. These moves are illustrated in Figure 5.9.

The direction of the gradient is determined by the largest interatomic forces and so steepest descents is a good method for relieving the highest-energy features in an initial configuration. The method is generally robust even when the starting point is far from a minimum, where the harmonic approximation to the energy surface is often a poor assumption. However, it suffers from the problem that many small steps will be performed when proceeding down a long narrow valley. The steepest descents method is forced to make a right-angled turn at each point, even though that might not be the best route to the minimum. The path oscillates and continually overcorrects itself, as illustrated in Figure 5.10; later steps reintroduce errors that were corrected by earlier moves.

5.4.4 Conjugate Gradients Minimisation

The conjugate gradients method produces a set of directions which does not show the oscillatory behaviour of the steepest descents method in narrow valleys. In the steepest descents method both the gradients and the direction of successive steps are orthogonal.

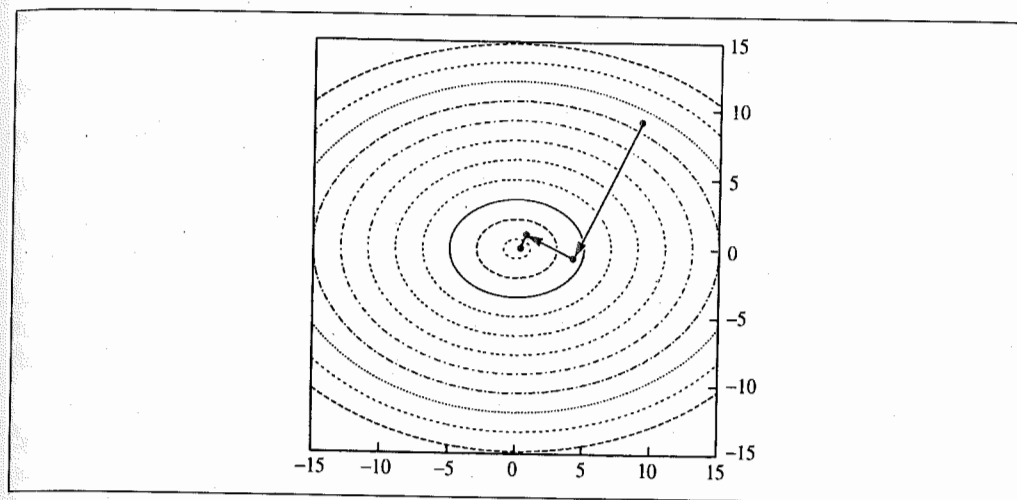


Fig. 5.9: Application of steepest descents to the function $x^2 + 2y^2$.

In conjugate gradients, the gradients at each point are orthogonal but the directions are *conjugate* (indeed, the method is more properly called the conjugate directions method). A set of conjugate directions has the property that for a quadratic function of M variables, the minimum will be reached in M steps. The conjugate gradients method moves in a direction \mathbf{v}_k from point \mathbf{x}_k where \mathbf{v}_k is computed from the gradient at the point and the previous direction vector \mathbf{v}_{k-1} :

$$\mathbf{v}_k = -\mathbf{g}_k + \gamma_k \mathbf{v}_{k-1} \quad (5.6)$$

γ_k is a scalar constant given by

$$\gamma_k = \frac{\mathbf{g}_k \cdot \mathbf{g}_k}{\mathbf{g}_{k-1} \cdot \mathbf{g}_{k-1}} \quad (5.7)$$

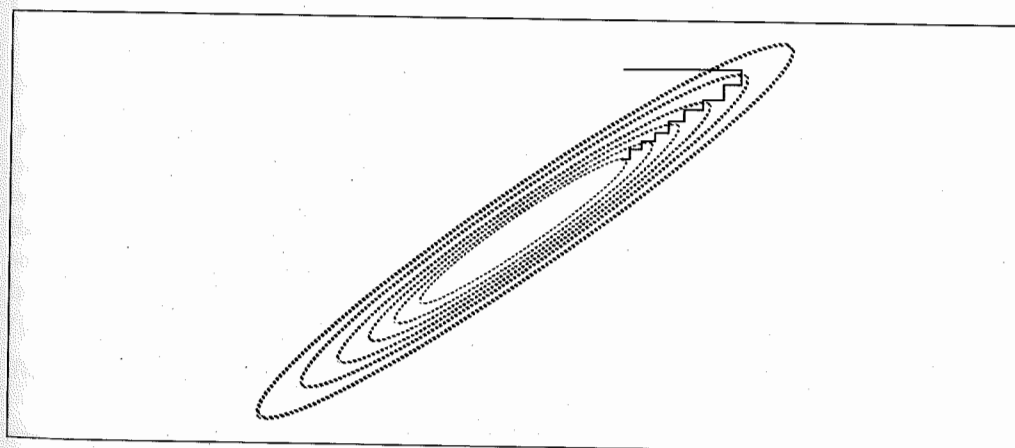


Fig. 5.10: The steepest descents method can give undesirable behaviour in a long narrow valley.

In the conjugate gradients method all of the directions and gradients satisfy the following relationships:

$$\mathbf{g}_i \cdot \mathbf{g}_j = 0 \quad (5.8)$$

$$\mathbf{v}_i \cdot \mathcal{V}_{ij}'' \cdot \mathbf{v}_j = 0 \quad (5.9)$$

$$\mathbf{g}_i \cdot \mathbf{v}_j = 0 \quad (5.10)$$

Clearly Equation (5.6) can only be used from the second step onwards and so the first step in the conjugate gradients method is the same as the steepest descents (i.e. in the direction of the gradient). The line search method should ideally be used to locate the one-dimensional minimum in each direction to ensure that each gradient is orthogonal to all previous gradients and that each direction is conjugate to all previous directions. However, an arbitrary step method is also possible.

The conjugate gradients method deals with our simple quadratic function $f(x, y) = x^2 + 2y^2$ as follows. From the initial point (9, 9) we move to the same point as in steepest descents, (4, -1). To find the direction of the next move, we first determine the negative gradient at the current point. This is the vector (-8, 4). This is then combined with the vector corresponding to minus the gradient at the initial point, (-18, -36) multiplied by γ :

$$\mathbf{v}_k = \begin{pmatrix} -8 \\ 4 \end{pmatrix} + \frac{(-8)^2 + (4)^2}{(-18)^2 + (-36)^2} \begin{pmatrix} -18 \\ -36 \end{pmatrix} = \begin{pmatrix} -80/9 \\ +20/9 \end{pmatrix} \quad (5.11)$$

To locate the second point we therefore need to perform a line search along the line with gradient -1/4 that passes through the point (4, -1). The minimum along this line is at the origin, at the true minimum of the function. The conjugate gradients method thus locates the exact minimum of the function exactly in just two moves, as illustrated in Figure 5.11.

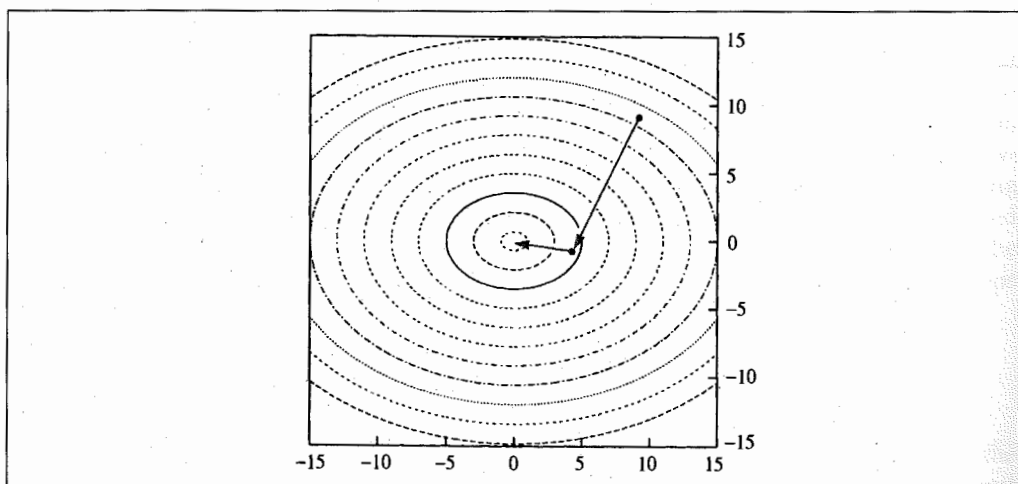


Fig. 5.11: Application of conjugate gradients method to the function $x^2 + 2y^2$.

Several variants of the conjugate gradients method have been proposed. The formulation given in Equation (5.7) is the original Fletcher-Reeves algorithm. Polak and Ribiere proposed an alternative form for the scalar constant γ_k :

$$\gamma_k = \frac{(\mathbf{g}_k - \mathbf{g}_{k-1}) \cdot \mathbf{g}_k}{\mathbf{g}_{k-1} \cdot \mathbf{g}_{k-1}} \quad (5.12)$$

For a purely quadratic function the Polak-Ribiere method is identical to the Fletcher-Reeves algorithm as all gradients will be orthogonal. However, most functions of interest, including those used in molecular modelling, are at best only approximately quadratic. Polak and Ribiere claimed that their method performed better than the original Fletcher-Reeves algorithm, at least for the functions that they examined.

5.5 Second Derivative Methods: The Newton-Raphson Method

Second-order methods use not only the first derivatives (i.e. the gradients) but also the second derivatives to locate a minimum. Second derivatives provide information about the curvature of the function. The *Newton-Raphson* method is the simplest second-order method. Recall our Taylor series expansion about the point x_k , Equation (5.2):

$$\mathcal{V}(x) = \mathcal{V}(x_k) + (x - x_k)\mathcal{V}'(x_k) + (x - x_k)^2\mathcal{V}''(x_k)/2 + \dots \quad (5.13)$$

The first derivative of $\mathcal{V}(x)$ is:

$$\mathcal{V}'(x) = x\mathcal{V}'(x_k) + (x - x_k)\mathcal{V}''(x_k) \quad (5.14)$$

If the function is purely quadratic, the second derivative is the same everywhere, and so $\mathcal{V}''(x) = \mathcal{V}''(x_k)$.

At the minimum ($x = x^*$) $\mathcal{V}'(x^*) = 0$ and so

$$x^* = x_k - \mathcal{V}'(x_k)/\mathcal{V}''(x_k) \quad (5.15)$$

For a multidimensional function: $\mathbf{x}^* = \mathbf{x}_k - \mathcal{V}'(\mathbf{x}_k)\mathcal{V}''^{-1}(\mathbf{x}_k)$.

$\mathcal{V}''^{-1}(\mathbf{x}_k)$ is the inverse Hessian matrix of second derivatives, which, in the Newton-Raphson method, must therefore be inverted. This can be computationally demanding for systems with many atoms and can also require a significant amount of storage. The Newton-Raphson method is thus more suited to small molecules (usually less than 100 atoms or so). For a purely quadratic function the Newton-Raphson method finds the minimum in one step from any point on the surface, as we will now show for our function $f(x, y) = x^2 + 2y^2$.

The Hessian matrix for this function is:

$$\mathbf{f}'' = \begin{pmatrix} 2 & 0 \\ 0 & 4 \end{pmatrix} \quad (5.16)$$

The inverse of this matrix is:

$$\mathbf{f}''^{-1} = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/4 \end{pmatrix} \quad (5.17)$$

The minimum is obtained using Equation (5.15):

$$\mathbf{x}^* = \begin{pmatrix} 9 \\ 9 \end{pmatrix} - \begin{pmatrix} 1/2 & 0 \\ 0 & 1/4 \end{pmatrix} \begin{pmatrix} 18 \\ 36 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (5.18)$$

In practice, of course, the surface is only quadratic to a first approximation and so a number of steps will be required, at each of which the Hessian matrix must be calculated and inverted. The Hessian matrix of second derivatives must be *positive definite* in a Newton-Raphson minimisation. A positive definite matrix is one for which all the eigenvalues are positive. When the Hessian matrix is not positive definite then the Newton-Raphson method moves to points (e.g. saddle points) where the energy increases. In addition, far from a minimum the harmonic approximation is not appropriate and the minimisation can become unstable. One solution to this problem is to use a more robust method to get near to the minimum (i.e. where the Hessian is positive definite) before applying the Newton-Raphson method.

5.5.1 Variants on the Newton-Raphson Method

There are a number of variations on the Newton-Raphson method, many of which aim to eliminate the need to calculate the full matrix of second derivatives. In addition, a family of methods called the quasi-Newton methods require only first derivatives and gradually construct the inverse Hessian matrix as the calculation proceeds. One simple way in which it may be possible to speed up the Newton-Raphson method is to use the same Hessian matrix for several successive steps of the Newton-Raphson algorithm with only the gradients being recalculated at each iteration.

A widely used algorithm is the *block-diagonal Newton-Raphson* method in which just one atom is moved at each iteration. Consequently all terms of the form $\partial^2 \mathcal{V} / \partial x_i \partial x_j$, where i and j refer to the Cartesian coordinates of atoms other than the atom being moved, will be zero. This only leaves those terms which involve the coordinates of the atom being moved and so reduces the problem to the trivial one of inverting a 3×3 matrix. However, the block-diagonal approach can be less efficient when the motions of some atoms are closely coupled, such as the concerted movements of connected atoms in a phenyl ring.

5.6 Quasi-Newton Methods

Calculation of the inverse Hessian matrix can be a potentially time-consuming operation that represents a significant drawback to the 'pure' second derivative methods such as Newton-Raphson. Moreover, one may not be able to calculate analytical second derivatives, which are preferable. The quasi-Newton methods (also known as variable metric methods) gradually build up the inverse Hessian matrix in successive iterations. That is, a sequence of

matrices \mathbf{H}_k is constructed that has the property

$$\lim_{k \rightarrow \infty} \mathbf{H}_k = \mathcal{V}''^{-1} \quad (5.19)$$

At each iteration k , the new positions \mathbf{x}_{k+1} are obtained from the current positions \mathbf{x}_k , the gradient \mathbf{g}_k and the current approximation to the inverse Hessian matrix \mathbf{H}_k :

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}_k \mathbf{g}_k \quad (5.20)$$

This formula is exact for a quadratic function, but for 'real' problems a line search may be desirable. This line search is performed along the vector $\mathbf{x}_{k+1} - \mathbf{x}_k$. It may not be necessary to locate the minimum in the direction of the line search very accurately, at the expense of a few more steps of the quasi-Newton algorithm. For quantum mechanics calculations the additional energy evaluations required by the line search may prove more expensive than using the more approximate approach. An effective compromise is to fit a function to the energy and gradient at the current point \mathbf{x}_k and at the point \mathbf{x}_{k+1} and determine the minimum in the fitted function.

Having moved to the new positions \mathbf{x}_{k+1} , \mathbf{H} is updated from its value at the previous step according to a formula depending upon the specific method being used. The methods of Davidon-Fletcher-Powell (DFP), Broyden-Fletcher-Goldfarb-Shanno (BFGS) and Murtaugh-Sargent (MS) are commonly encountered, but there are many others. These methods converge to the minimum, for a quadratic function of M variables, in M steps. The DFP formula is:

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\mathbf{x}_{k+1} - \mathbf{x}_k) \otimes (\mathbf{x}_{k+1} - \mathbf{x}_k)}{(\mathbf{x}_{k+1} - \mathbf{x}_k) \cdot (\mathbf{x}_{k+1} - \mathbf{x}_k)} - \frac{[\mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)] \otimes [\mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)]}{(\mathbf{g}_{k+1} - \mathbf{g}_k) \cdot \mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)} \quad (5.21)$$

The symbol \otimes when interposed between two vectors means that a matrix is to be formed. The ij th element of the matrix $\mathbf{u} \otimes \mathbf{v}$ is obtained by multiplying u_i by v_j .

The BFGS formula differs from the DFP equation by an additional term:

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\mathbf{x}_{k+1} - \mathbf{x}_k) \otimes (\mathbf{x}_{k+1} - \mathbf{x}_k)}{(\mathbf{x}_{k+1} - \mathbf{x}_k) \cdot (\mathbf{x}_{k+1} - \mathbf{x}_k)} - \frac{[\mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)] \otimes [\mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)]}{(\mathbf{g}_{k+1} - \mathbf{g}_k) \cdot \mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)} + [(\mathbf{g}_{k+1} - \mathbf{g}_k) \cdot \mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)] \mathbf{u} \otimes \mathbf{u} \quad (5.22)$$

where

$$\mathbf{u} = \frac{(\mathbf{x}_{k+1} - \mathbf{x}_k)}{(\mathbf{x}_{k+1} - \mathbf{x}_k) \cdot (\mathbf{x}_{k+1} - \mathbf{x}_k)} - \frac{[\mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)]}{(\mathbf{g}_{k+1} - \mathbf{g}_k) \cdot \mathbf{H}_k \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)} \quad (5.23)$$

The MS formula is:

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{[(\mathbf{x}_{k+1} - \mathbf{x}_k) - \mathbf{H}_k(\mathbf{g}_{k+1} - \mathbf{g}_k)] \otimes [(\mathbf{x}_{k+1} - \mathbf{x}_k) - \mathbf{H}_k(\mathbf{g}_{k+1} - \mathbf{g}_k)]}{[(\mathbf{x}_{k+1} - \mathbf{x}_k) - \mathbf{H}_k(\mathbf{g}_{k+1} - \mathbf{g}_k)] \cdot (\mathbf{g}_{k+1} - \mathbf{g}_k)} \quad (5.24)$$

All of these methods use just the new and current points to update the inverse Hessian. The default algorithm used in the Gaussian series of molecular orbital programs [Schlegel 1982] makes use of more of the previous points to construct the Hessian (and thence the inverse Hessian), giving better convergence properties. Another feature of this method is its use

of a quartic polynomial that is guaranteed to have just one local minimum in the line search. The DFP, BFGS and MS methods can also be used with numerical derivatives, but alternative approaches may be more effective under such circumstances.

The matrix H is often initialised to the unit matrix I . The performance of the quasi-Newton algorithms can be improved by using a better estimate of the inverse Hessian than just the unit matrix. The unit matrix gives no information about the bonding in the system, nor does it identify any coupling between the various degrees of freedom. For example, a molecular mechanics calculation can be used to provide an initial guess to H prior to a quantum mechanical calculation. Alternatively the matrix can be obtained from a quantum mechanical calculation at a lower level of theory (e.g. semi-empirical or with a smaller basis set).

5.7 Which Minimisation Method Should I Use?

The choice of minimisation algorithm is dictated by a number of factors, including the storage and computational requirements, the relative speeds with which the various parts of the calculation can be performed, the availability of analytical derivatives and the robustness of the method. Thus, any method that requires the Hessian matrix to be stored (let alone its inverse calculated) may present memory problems when applied to systems containing thousands of atoms. Calculations on systems of this size are invariably performed using molecular mechanics, and so the steepest descents and the conjugate gradients methods are very popular here. For molecular mechanics calculations on small molecules, the Newton-Raphson method may be used, although this algorithm can have problems with structures that are far from a minimum. For this reason it is usual to perform a few steps of minimisation using a more robust method such as the simplex or steepest descents before applying the Newton-Raphson algorithm. Analytical expressions for both first and second derivatives are available for most of the terms found in common force fields.

The performance of the steepest descents and conjugate gradients methods is contrasted in the following example. A model of the antibiotic netropsin (Figure 5.12) bound to DNA was constructed using an automated docking program. This initial model was then subjected to two stages of minimisation. In the first stage, the aim was to produce a structure that did not

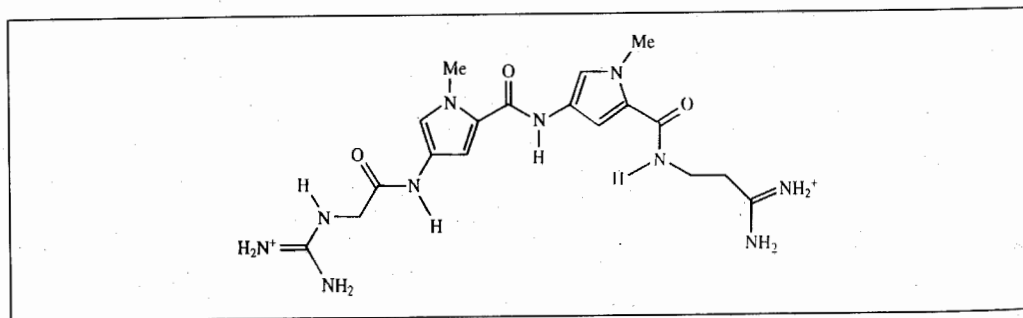


Fig. 5.12: The DNA inhibitor netropsin.

Method	Initial refinement (Av. gradient $< 1 \text{ kcal } \text{Å}^{-2}$)		Stringent minimisation (Av. gradient $< 0.1 \text{ kcal } \text{Å}^{-2}$)	
	CPU time (s)	Number of iterations	CPU time (s)	Number of iterations
Steepest descents	67	98	1405	1893
Conjugate gradients	149	213	257	367

Table 5.1 A comparison of the steepest descents and conjugate gradients methods for an initial refinement and a stringent minimisation.

have any significant high-energy interactions. The structure was then further minimised to give a structure much closer to the minimum. The results are shown in Table 5.1.

This study shows that the steepest descent method can actually be superior to conjugate gradients when the starting structure is some way from the minimum. However, conjugate gradients is much better once the initial strain has been removed.

Quantum mechanical calculations are restricted to systems with relatively small numbers of atoms, and so storing the Hessian matrix is not a problem. As the energy calculation is often the most time-consuming part of the calculation, it is desirable that the minimisation method chosen takes as few steps as possible to reach the minimum. For many levels of quantum mechanics theory analytical first derivatives are available. However, analytical second derivatives are only available for a few levels of theory and can be expensive to compute. The quasi-Newton methods are thus particularly popular for quantum mechanical calculations.

When using internal coordinates in a quantum mechanical minimisation it can be important to use an appropriate Z-matrix as input. For many systems the Z-matrix can often be written in many different ways as there are many combinations of internal coordinates. There should be no strong coupling between the coordinates. *Dummy atoms* can often help in the construction of an appropriate Z-matrix. A dummy atom is used solely to define the geometry and has no nuclear charge and no basis functions. A simple example of the use of dummy atoms is for a linear molecule such as HN_3 , where the angle of 180° would cause problems. The geometry of this molecule can be defined using a dummy atom as illustrated in Figure 5.13; the associated Z-matrix for this system would be:

1	N						
2	N	1	RN1N2				
3	X	1	1.0	2	90.0		
4	N	1	RN1N4	3	AN4N1X	2	180.0
5	H	4	RN4H	1	AHN4N1	3	180.0

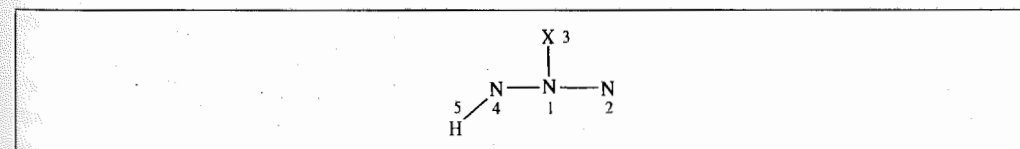


Fig. 5.13: Internal coordinates of HN_3 molecule defined using dummy atom X.

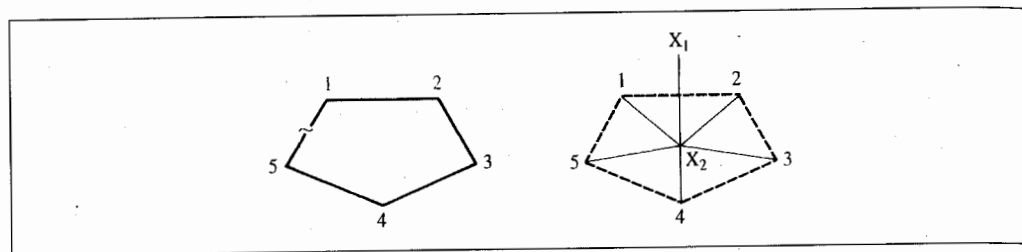


Fig. 5.14: The ring closure bond between atoms 1 and 5 would be strongly coupled to the other internal coordinates (left) unless dummy atoms are used to define the Z-matrix (right).

Strong coupling between coordinates can give long 'valleys' in the energy surface, which may also present problems. Care must be taken when defining the Z-matrix for cyclic systems in particular. The natural way to define a cyclic compound would be to number the atoms sequentially around the ring. However, this would then mean that the ring closure bond will be very strongly coupled to all of the other bonds, angles and torsion angles (Figure 5.14). A better definition uses a dummy atom placed at the centre of the ring (Figure 5.14). Some quantum mechanics programs are able to convert the input coordinates (be they Cartesian or internal) into the most efficient set for minimisation so removing from the user the problems of trying to decide what is an appropriate set of internal coordinates. For energy minimisations redundant internal coordinates have been shown to give significant improvements in efficiency compared with Cartesian coordinates or non-redundant internal coordinates, especially for flexible and polycyclic systems [Peng *et al.* 1996]. The redundant internal coordinates employed generally comprise the bond lengths, angles and torsion angles in the system. These methods obviously also require the means to interconvert between the internal coordinate representation and the Cartesian coordinates that are often used as input and desired as output. Of particular importance is the need to transform energy derivatives and the Hessian matrices (if appropriate).

5.7.1 Distinguishing Between Minima, Maxima and Saddle Points

A configuration at which all the first derivatives are zero need not necessarily be a minimum point; this condition holds at both maxima and saddle points as well. From simple calculus we know that the second derivative of a function of one variable, $f'(x)$ is positive at a minimum and negative at a maximum. It is necessary to calculate the eigenvalues of the Hessian matrix to distinguish between minima, maxima and saddle points. At a minimum point there will be six zero and $3N - 6$ positive eigenvalues if $3N$ Cartesian coordinates are used. The six zero eigenvalues correspond to the translational and rotational degrees of freedom of the molecule (thus these six zero eigenvalues are not obtained when internal coordinates are used). At a maximum point all eigenvalues are negative and at a saddle point one or more eigenvalues are negative. We will consider the uses of the eigenvalue and eigenvector information in Sections 5.8 and 5.9.

5.7.2 Convergence Criteria

In contrast to the simple analytical functions that we have used to illustrate the operation of the various minimisation methods, in 'real' molecular modelling applications it is rarely possible to identify the 'exact' location of minima and saddle points. We can only ever hope to find an approximation to the true minimum or saddle point. Unless instructed otherwise, most minimisation methods would keep going forever, moving ever closer to the minimum. It is therefore necessary to have some means to decide when the minimisation calculation is sufficiently close to the minimum and so can be terminated. Any calculation is of course limited by the precision with which numbers can be stored on the computer, but in most instances it is usual to stop well before this limit is reached. A simple strategy is to monitor the energy from one iteration to the next and to stop when the difference in energy between successive steps falls below a specified threshold. An alternative is to monitor the change in coordinates and to stop when the difference between successive configurations is sufficiently small. A third method is to calculate the root-mean-square gradient. This is obtained by adding the squares of the gradients of the energy with respect to the coordinates, dividing by the number of coordinates and taking the square root:

$$\text{RMS} = \sqrt{\frac{\mathbf{g}^T \mathbf{g}}{3N}} \quad (5.25)$$

It is also useful to monitor the maximum value of the gradient to ensure that the minimisation has properly relaxed all the degrees of freedom and has not left a large amount of strain in one or two coordinates.

5.8 Applications of Energy Minimisation

Energy minimisation is very widely used in molecular modelling and is an integral part of techniques such as conformational search procedures (Chapter 9). Energy minimisation is also used to prepare a system for other types of calculation. For example, energy minimisation may be used prior to a molecular dynamics or Monte Carlo simulation in order to relieve any unfavourable interactions in the initial configuration of the system. This is especially recommended for simulations of complex systems such as macromolecules or large molecular assemblies. In the following sections we will discuss some techniques that are specifically associated with energy minimisation methods.

5.8.1 Normal Mode Analysis

The molecular mechanics or quantum mechanics energy at an energy minimum corresponds to a hypothetical, motionless state at 0K. Experimental measurements are made on molecules at a finite temperature when the molecules undergo translational, rotational and vibration motion. To compare the theoretical and experimental results it is

necessary to make appropriate corrections to allow for these motions. These corrections are calculated using standard statistical mechanics formulae. The internal energy $U(T)$ at a temperature T is given by:

$$U(T) = U_{\text{trans}}(T) + U_{\text{rot}}(T) + U_{\text{vib}}(T) + U_{\text{vib}}(0) \quad (5.26)$$

If all translational and rotational modes are fully accessible in accordance with the equipartition theorem, then $U_{\text{trans}}(T)$ and $U_{\text{rot}}(T)$ are both equal to $\frac{3}{2}k_B T$ per molecule (except that $U_{\text{rot}}(T)$ equals $k_B T$ for a linear molecule); k_B is Boltzmann's constant. However, the vibrational energy levels are often only partially excited at room temperature. The vibrational contribution to the internal energy at a temperature T thus requires knowledge of the actual vibrational frequencies. The vibrational contribution equals the difference in the vibrational enthalpy at the temperature T and at 0K and is given by:

$$U_{\text{vib}}(T) = \sum_{i=1}^{N_{\text{nm}}} \left(\frac{h\nu_i}{2} + \frac{h\nu_i}{\exp[h\nu_i/k_B T] + 1} \right) \quad (5.27)$$

N_{nm} is the number of *normal vibrational modes* for the system. Even the zero-point energy ($U_{\text{vib}}(0)$, obtained by summing $\frac{1}{2}h\nu_i$ for each normal mode) can be quite substantial, amounting to about 100 kcal/mol for a six-carbon alkane. Other thermodynamic quantities such as entropies and free energies may also be calculated from the vibrational frequencies using the relevant statistical mechanics expressions.

Normal modes are useful because they correspond to collective motions of the atoms in a coupled system that can be individually excited. The three normal modes of water are schematically illustrated in Figure 5.15; a non-linear molecule with N atoms has $3N - 6$ normal modes. The frequencies of the normal modes together with the displacements of the individual atoms may be calculated from a molecular mechanics force field or from the wavefunction using the Hessian matrix of second derivatives (\mathcal{V}''). Of course, if we have used an appropriate minimisation algorithm then we already know the Hessian. The Hessian must first be converted to the equivalent force-constant matrix in *mass-weighted coordinates* (\mathbf{F}), as follows:

$$\mathbf{F} = \mathbf{M}^{-1/2} \mathcal{V}'' \mathbf{M}^{-1/2} \quad (5.28)$$

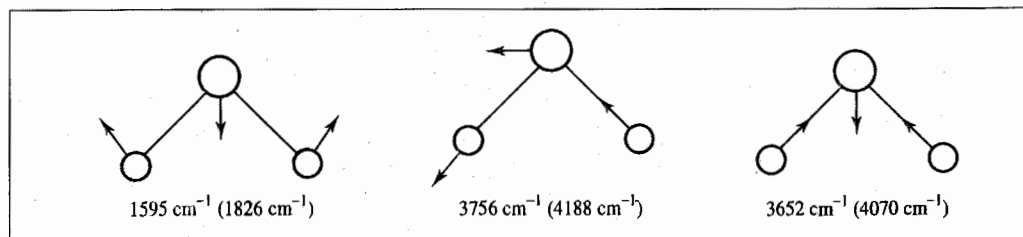


Fig. 5.15: Normal modes of water. Experimental and (calculated) frequencies are shown. Theoretical frequencies calculated using a 6-31G* basis set.

\mathbf{M} is a diagonal matrix of dimension $3N \times 3N$, containing the atomic masses. All elements of \mathbf{M} are zero except those on the diagonal; $M_{1,1} = m_1$, $M_{2,2} = m_1$, $M_{3,3} = m_1$, $M_{4,4} = m_2$, \dots , $M_{3N-2,3N-2} = m_N$, $M_{3N-1,3N-1} = m_N$, $M_{3N,3N} = m_N$. Each non-zero element of $\mathbf{M}^{-1/2}$ is thus the inverse square root of the mass of the appropriate atom. The masses of the atoms must be taken into account because a force of a given magnitude will have a different effect upon a larger mass than a smaller one. For example, the force constant for a bond to a deuterium atom is, to a good approximation, the same as to a proton, yet the different mass of the deuteron gives a different motion and a different zero-point energy. The use of mass-weighted coordinates takes care of these problems.

We next solve the secular equation $|\mathbf{F} - \mathbf{I}| = 0$ to obtain the eigenvalues and eigenvectors of the matrix \mathbf{F} . This step is usually performed using matrix diagonalisation, as outlined in Section 1.10.3. If the Hessian is defined in terms of Cartesian coordinates then six of these eigenvalues will be zero as they correspond to translational and rotational motion of the entire system. The frequency of each normal mode is then calculated from the eigenvalues using the relationship:

$$\nu_i = \frac{\sqrt{\lambda_i}}{2\pi} \quad (5.29)$$

As a simple example of a normal mode calculation consider the linear triatomic system in Figure 5.16. We shall just consider motion along the long axis of the molecule. The displacements of the atoms from their equilibrium positions along this axis are denoted by ξ_i . It is assumed that the displacements are small compared with the equilibrium values l_0 and the system obeys Hooke's law with bond force constants k . The potential energy is given by:

$$\mathcal{V} = \frac{1}{2}k(\xi_1 - \xi_2)^2 + \frac{1}{2}k(\xi_2 - \xi_3)^2 \quad (5.30)$$

We next calculate the first and then the second derivatives of the potential energy with respect to the three coordinates ξ_1 , ξ_2 and ξ_3 :

$$\frac{\partial \mathcal{V}}{\partial \xi_1} = k(\xi_1 - \xi_2); \quad \frac{\partial \mathcal{V}}{\partial \xi_2} = -k(\xi_1 - \xi_2) + k(\xi_2 - \xi_3); \quad \frac{\partial \mathcal{V}}{\partial \xi_3} = -k(\xi_2 - \xi_3) \quad (5.31)$$

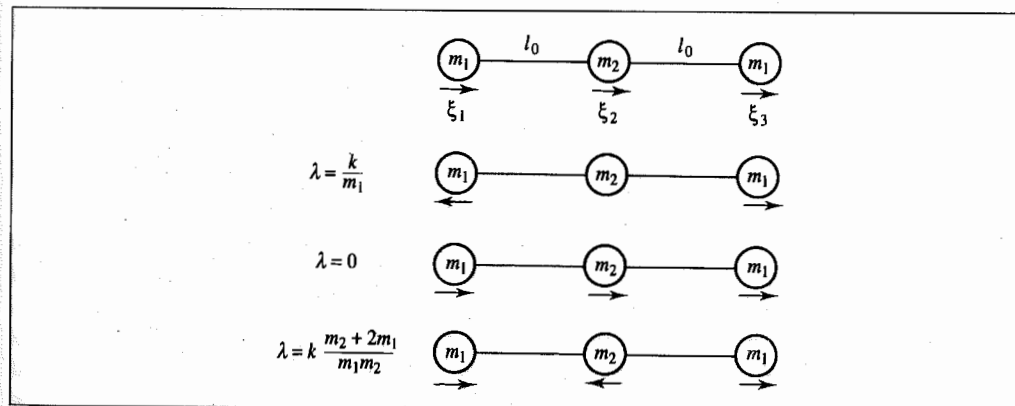


Fig. 5.16: Linear three-atom system with results of normal mode calculation.

The second derivatives are conveniently represented as a 3×3 matrix:

$$\begin{vmatrix} k & -k & 0 \\ -k & 2k & -k \\ 0 & -k & k \end{vmatrix} \quad (5.32)$$

The mass-weighted matrix is

$$\begin{vmatrix} m_1 & 0 & 0 \\ 0 & m_2 & 0 \\ 0 & 0 & m_3 \end{vmatrix} \quad (5.33)$$

The secular equation to be solved is thus:

$$\begin{vmatrix} \frac{k}{m_1} - \lambda & -\frac{k}{\sqrt{m_1}\sqrt{m_2}} & 0 \\ -\frac{k}{\sqrt{m_1}\sqrt{m_2}} & \frac{2k}{m_2} - \lambda & -\frac{k}{\sqrt{m_1}\sqrt{m_2}} \\ 0 & -\frac{k}{\sqrt{m_1}\sqrt{m_2}} & \frac{k}{m_1} - \lambda \end{vmatrix} = 0 \quad (5.34)$$

This determinant leads to a cubic in λ which has three roots (λ_k), each corresponding to a different mode of motion:

$$\lambda = \frac{k}{m_1}, \quad \lambda = 0, \quad \lambda = k \frac{m_2 + 2m_1}{m_1 m_2} \quad (5.35)$$

The corresponding frequencies can be obtained from Equation (5.29). The amplitudes (A) of each normal mode are given by the eigenvector solutions of the secular equation $\mathbf{FA} = \lambda \mathbf{A}$. If A_1 , A_2 and A_3 are the amplitudes of each atom then the amplitudes obtained for each eigenvalue are:

$$\lambda = \frac{k}{m_1}: \quad A_1 = -A_3; \quad A_2 = 0 \quad (5.36)$$

$$\lambda = 0: \quad A_1 = A_3; \quad A_2 = \sqrt{\frac{m_2}{m_1}} A_1 \quad (5.37)$$

$$\lambda = k \frac{m_2 + 2m_1}{m_1 m_2}: \quad A_1 = A_3; \quad A_2 = -2\sqrt{\frac{m_1}{m_2}} A_1 \quad (5.38)$$

These normal modes are schematically illustrated in Figure 5.16. They correspond to a symmetric stretch, a translation and an asymmetric stretch respectively.

We have already seen how the results of normal mode calculations can be used to calculate thermodynamic quantities. The frequencies themselves can also be compared with the results of spectroscopic experiments, information which can be used in the parametrisation of a force field. For example, the experimental frequencies for the normal modes of water are shown in Figure 5.15, together with the frequencies determined using a 6-31G* *ab initio* calculation. The calculated values clearly deviate from those obtained experimentally, but the ratio of the experimental and theoretical frequencies is

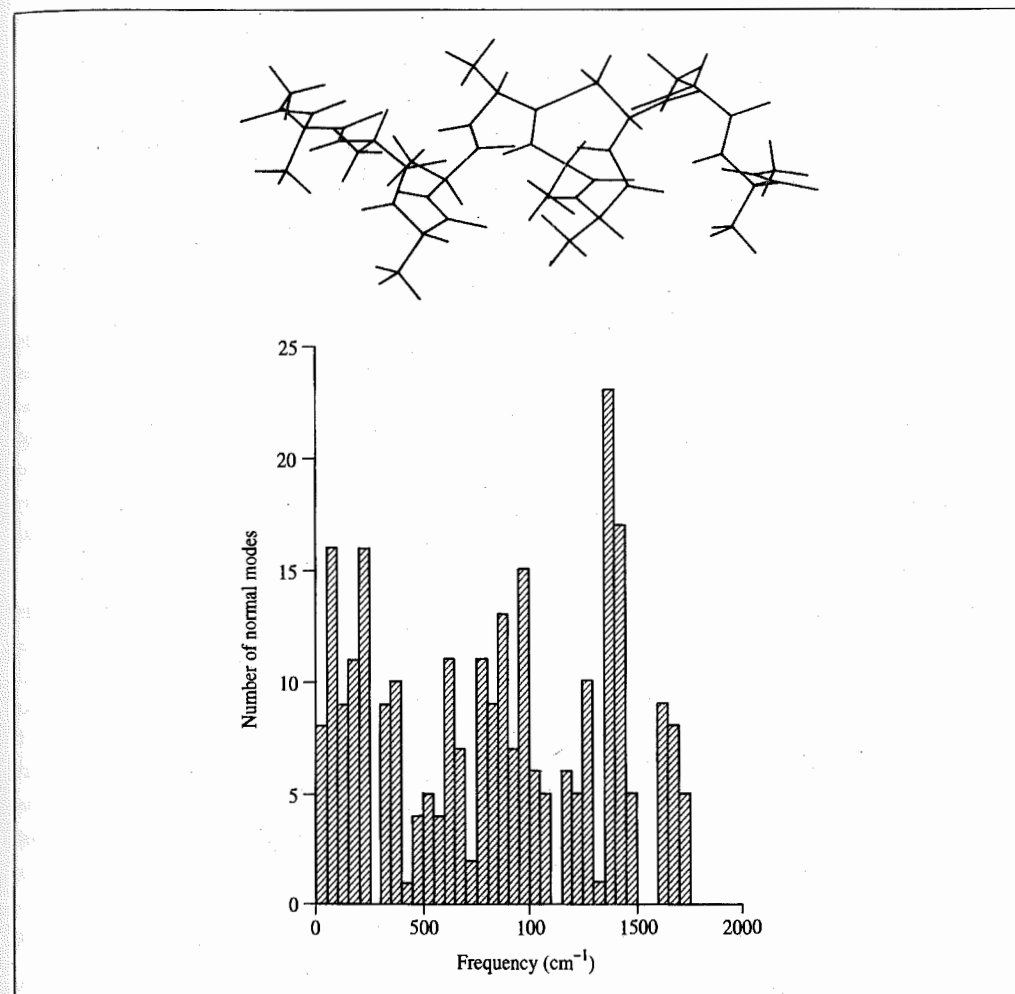


Fig. 5.17: Histogram of the normal modes calculated for a polyaniline polypeptide in an α -helical conformation. The height of each bar indicates the number of normal modes in each 50 cm^{-1} section.

remarkably consistent (at about 1.1). Such empirical scaling factors have been derived which enable frequencies obtained using a given level of theory to be converted to values for experiment or a higher level of theory [Pople *et al.* 1993]. The normal modes of much larger molecules can be calculated using molecular mechanics. For example, the vibrations of a helical polypeptide constructed from a sequence of ten alanine residues (112 atoms) are shown in Figure 5.17. In such cases it is usually the low-frequency vibrations that are of most interest as these correspond to the large-scale conformational motions of the molecule. The results of such analyses can be compared with molecular dynamics simulations from which vibrational contributions can also be extracted [Brooks and Karplus 1983].

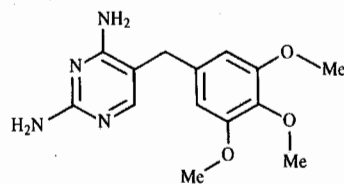


Fig. 5.18: Trimethoprim.

A normal mode calculation is based upon the assumption that the energy surface is quadratic in the vicinity of the energy minimum (the harmonic approximation). Deviations from the harmonic model can require corrections to calculated thermodynamic properties. One way to estimate anharmonic corrections is to calculate a force constant matrix using the atomic motions obtained from a molecular dynamics simulation; such simulations are not restricted to movements on a harmonic energy surface. The eigenvalues and eigenvectors are then calculated for this quasi-harmonic force-constant matrix in the normal way, giving a model which implicitly incorporates the anharmonic effects.

The harmonic approximation to the energy surface is found to be appropriate for well-defined energy minima such as the intramolecular degrees of freedom of small molecules and for some small intermolecular complexes. For larger systems such as liquids and large, 'floppy' molecules, the harmonic approximation breaks down. Such systems also have an extraordinarily large number of 'minima' on the energy surface. In such cases it is not possible to calculate accurately thermodynamic properties using energy minimisation and normal mode calculations. Rather, molecular dynamics or Monte Carlo simulations must be used to sample the energy surface from which properties can be derived, as we will discuss in Chapters 6–8.

5.8.2 The Study of Intermolecular Processes

One example of the use of minimisation methods and normal-mode analysis is the study by Hagler and co-workers of the binding of the antibacterial drug trimethoprim (Figure 5.18) to the enzyme dihydrofolate reductase (DHFR) [Dauber-Osguthorpe *et al.* 1988; Fisher *et al.* 1991]. DHFR catalyses the reduction of folic acid and dihydrofolic acid to tetrahydrofolic acid (Figure 5.19) and plays a vital metabolic role in the biosynthesis of nucleic acids in bacteria, protozoa, plants and animals. Trimethoprim exploits the structural differences between bacterial and vertebrate DHFR, binding much more strongly to the former, and is clinically used as an antibacterial agent. Inhibitors of human DHFR are used in cancer therapy. Hagler and colleagues applied energy minimisation to an isolated trimethoprim molecule, to the crystal structure of trimethoprim, to trimethoprim in the presence of water molecules, and to trimethoprim in intermolecular complexes with DHFR from both bacterial and vertebrate sources. An important observation was that the conformation of the trimethoprim, when bound to the enzyme, was significantly different from that obtained for the isolated molecule. This reinforces the view that the use of

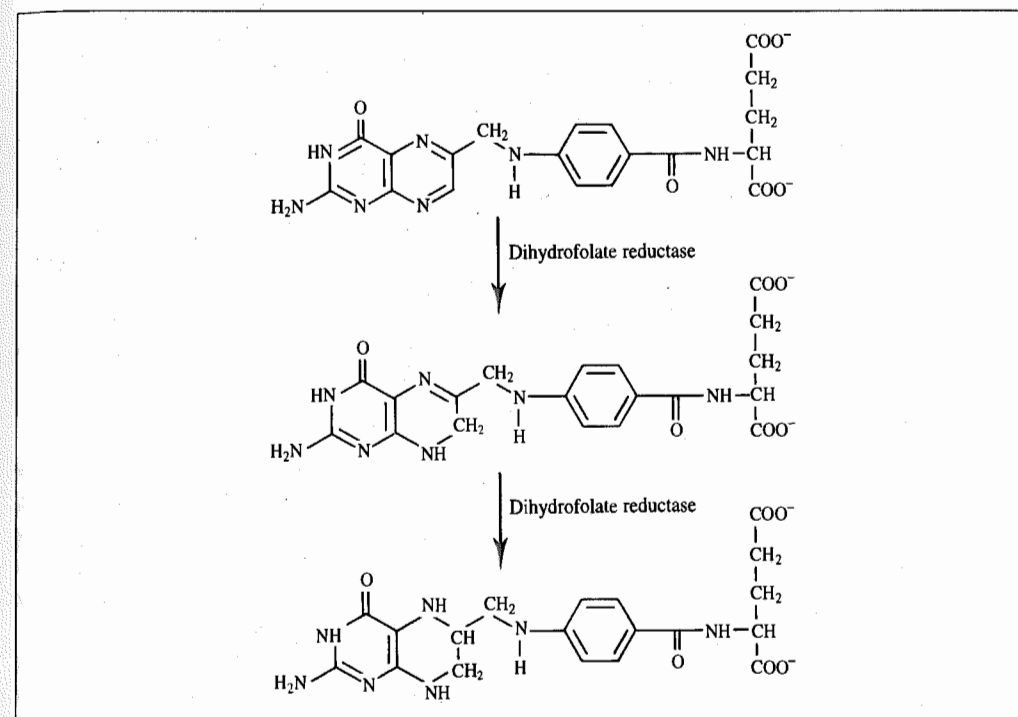


Fig. 5.19: DHFR catalyses the reduction of folic acid to tetrahydrofolic acid.

structures obtained from energy minimisation calculations on isolated molecules can lead to misleading conclusions. Intermolecular interactions with the receptor enable the ligand to adopt a conformation whose intramolecular energy is significantly higher than any of its minimum energy structures.

A normal mode analysis on the isolated and bound trimethoprim molecules enabled an estimate to be made of the entropic contribution to binding. Low-frequency modes for the isolated ligand were found to be shifted to higher frequencies for the ligand in the enzyme complex, reflecting a restriction of the motion of the ligand by the protein. This entropic contribution to the free energy of binding was predicted to be quite significant, indicating that conclusions based solely upon energies may be misleading.

5.9 Determination of Transition Structures and Reaction Pathways

Chemists are interested not only in the thermodynamics of a process (the relative stability of the various species) but also in its kinetics (the rate of conversion from one structure to another). Knowledge of the minimum points on an energy surface enables thermodynamic data to be interpreted, but for the kinetics it is necessary to investigate the nature of the

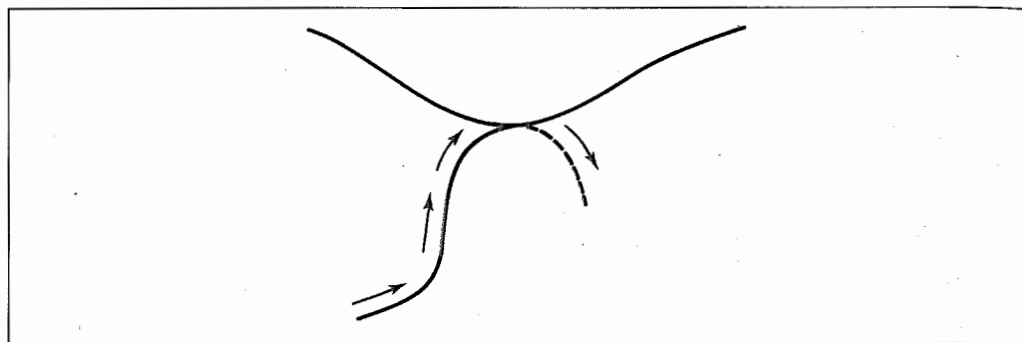


Fig. 5.20: The lowest-energy path from one minimum to another passes through a saddle point.

energy surface away from the minimum points. In particular, we would like to know how the system changes from one minimum to another. What changes in geometry are involved, and how does the energy vary during the transition? The minimum points on the energy surface may be the reactants and products of a chemical reaction, two conformations of a molecule, or two molecules that associate to form a non-covalently bound bimolecular complex. We shall use the term 'reaction pathway' to describe the path between two minima, but our use of the word 'reaction' does not necessarily mean that bond making and/or breaking is involved. Many methods have been proposed for finding transition structures and elucidating reaction pathways. We do not have space to cover all of the methods, and so we shall restrict our discussion to some of the more common approaches.

As a system moves from one minimum to another, the energy increases to a maximum at the transition structure and then falls. At a saddle point the first derivatives of the potential function with respect to the coordinates are all zero (just as they are at a minimum point). The number of negative eigenvalues in the Hessian matrix is used to distinguish different types of saddle point; an n th-order transition or saddle point has n negative eigenvalues. We are usually most interested in first-order saddle points, where the energy passes through a maximum for movement along the pathway that connects the two minima, but is a minimum for displacements in all other directions perpendicular to the path. This is shown schematically for a two-dimensional energy surface in Figure 5.20.

These negative eigenvalues of the Hessian matrix are often referred to as the 'imaginary' frequencies for motion of the system over the saddle point. We can illustrate this concept using the gas-phase S_N2 reaction between Cl^- and CH_3Cl . As the chloride ion approaches the methyl chloride along the line of the C–Cl bond the energy passes through an ion-dipole complex which is at an energy minimum. The energy then rises to a maximum at the pentagonal transition state. The energy profile is drawn in Figure 5.21. The geometries of the minimum and the pentagonal transition state, as determined by an *ab initio* HF/SCF calculation with the 6-31G⁺ basis set are shown in Figure 5.22. The lowest-frequency eigenvalues and a representation of the corresponding eigenvectors for the two geometries are also given in Figure 5.22. There are three frequencies in the ion-dipole minimum that are of particularly low energy; two of these correspond to degenerate 'wagging' motions of the

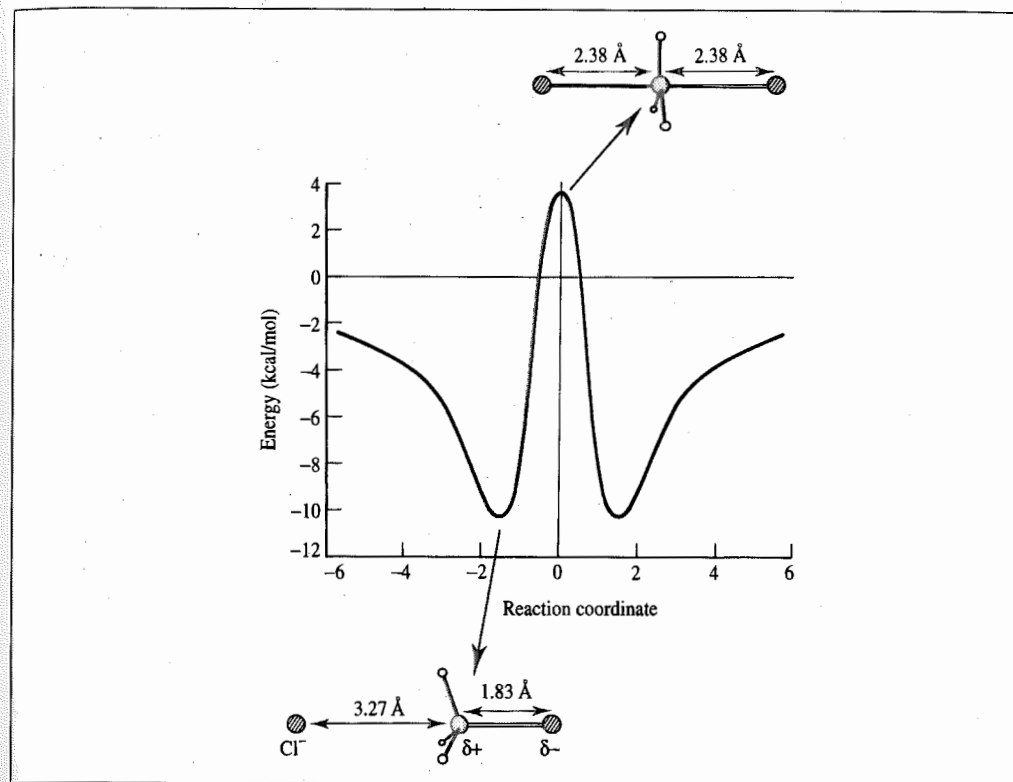


Fig. 5.21: The energy profile for the gas-phase $\text{Cl}^- + \text{MeCl}$ reaction. (Adapted in part from Chandrasekhar J, S F Smith and W L Jorgensen 1985. Theoretical Examination of the S_N2 Reaction Involving Chloride Ion and Methyl Chloride in the Gas Phase and Aqueous Solution. *Journal of the American Chemical Society* 107:154–163.)

system (at 71.3 cm^{-1}). The vibration at 101.0 cm^{-1} is the normal mode that corresponds to motion towards the transition state. At the saddle point there is a single negative eigenvalue (with an imaginary 'frequency' of -415.0 cm^{-1}) that corresponds to vibration along the Cl–C–Cl axis (i.e. motion along the reaction pathway). The other normal modes at the saddle point all have positive frequencies; the two lowest (at 204.2 cm^{-1}) correspond to wagging motions perpendicular to the Cl–C–Cl axis and the third is a symmetric stretch of the two chlorine atoms along the symmetry axis.

It is important to distinguish the transition *structure* from the transition *state*. The transition structure is the point of highest potential energy along the pathway. By contrast, the transition state is the geometry at the peak in the free energy profile. In many cases the geometry at the transition state is very similar to that of the transition structure. However, the transition state may be different as the free energy of activation includes contributions from sources other than just the potential energy. If the transition state geometry is temperature-dependent then entropic factors may be important. An example is the following radical

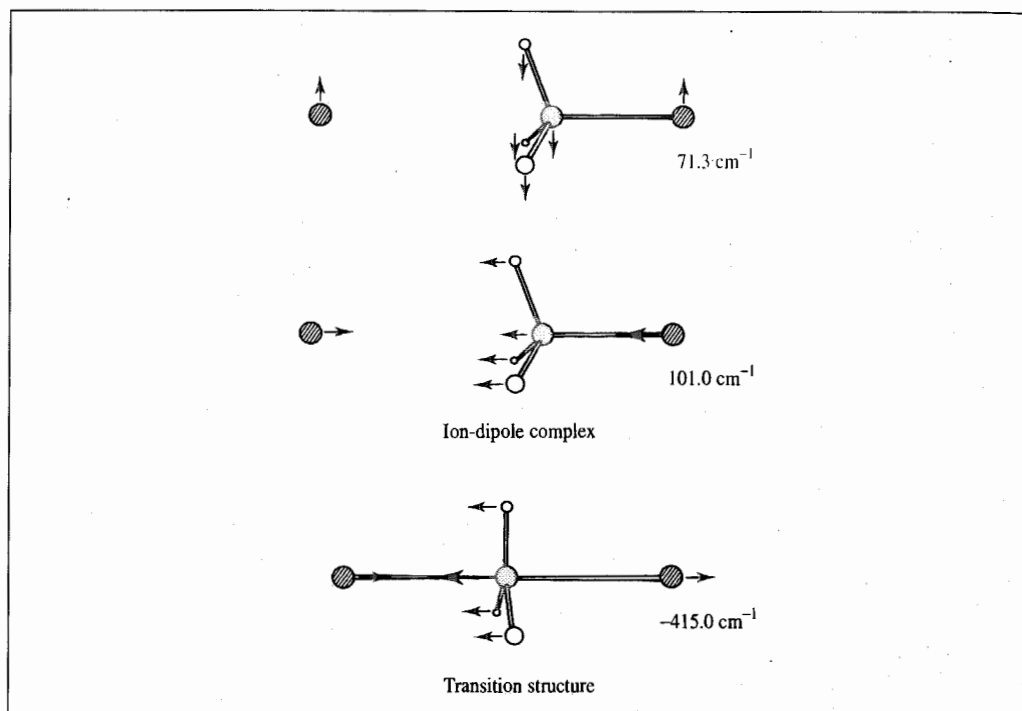
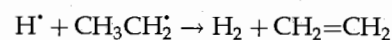


Fig. 5.22: Schematic representation of some of the lower frequencies in the ion-dipole complex for the $\text{Cl}^- + \text{MeCl}$ reaction and the imaginary frequency of the transition structure, calculated using a 6-31G* basis set.

reaction:



The calculated geometry of the transition structure resembles the ethyl radical (Figure 5.23) [Doubleday *et al.* 1985]. The entropy change for this reaction is negative and so, as the temperature is increased, the maximum in the free energy profile shifts more towards the products, in the direction of lower entropy.

Methods for finding transition structures and reaction pathways are often closely related. Thus, some methods for finding the reaction pathway start from the transition structure and move down towards a minimum. Such methods must be supplied with the transition structure geometry as the starting point. Conversely, some methods for locating transition structures do so by searching along the reaction pathway, or an approximation to it. Yet other methods require neither the transition structure nor the pathway, but can determine both simultaneously from the two minima. In general, it is more difficult to locate transition structures and determine reaction pathways than to find minimum points. It is therefore crucial to check that the Hessian matrix at any proposed saddle point has the required single negative eigenvalue. Methods for locating saddle points are usually most effective when given as input a geometry that is as close as possible to the transition structure. It

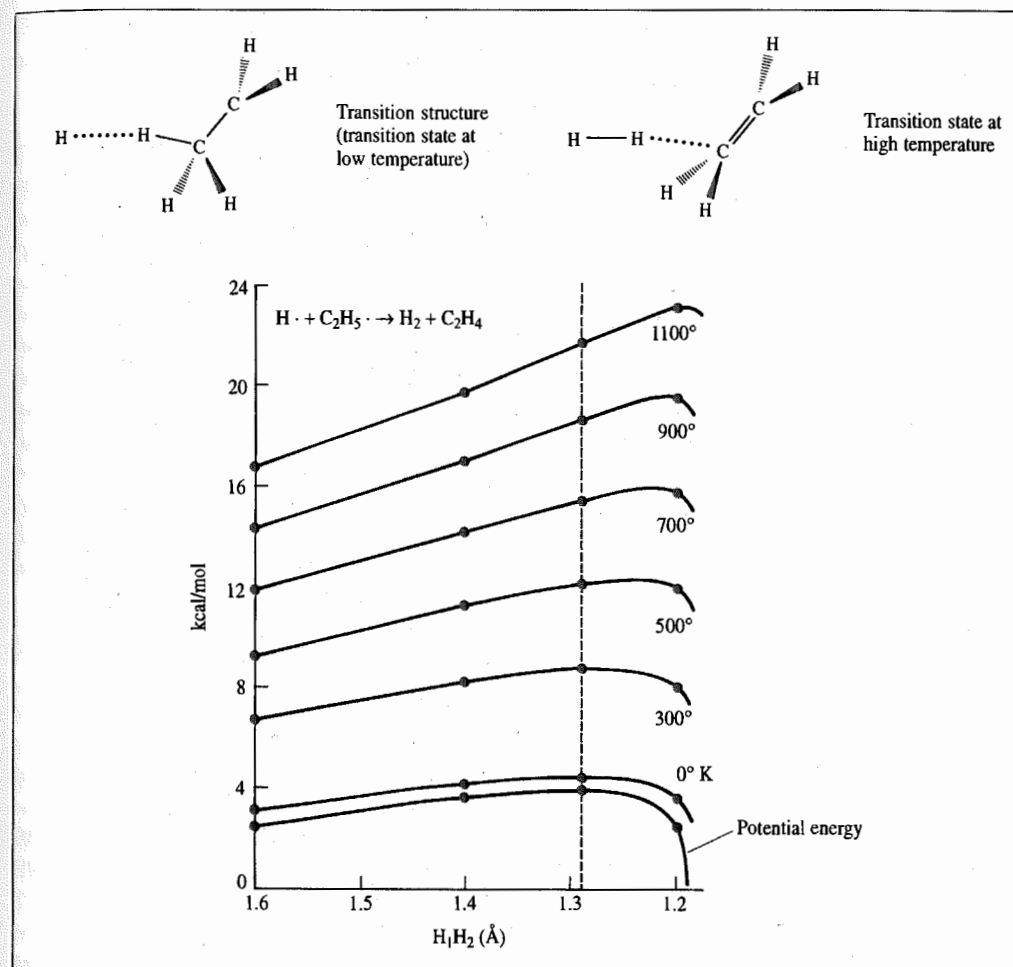


Fig. 5.23: The transition structure for the $\text{H} + \text{CH}_3\text{CH}_2^\bullet \rightarrow \text{H}_2 + \text{CH}_2=\text{CH}_2$ reaction. At low temperature the transition structure corresponds to the transition state (maximum of free energy). At high temperature the transition state moves closer to the products, as can be seen from the graph. (Redrawn from Doubleday C, J McIver, M Page and T Zielinski 1985. Temperature Dependence of the Transition-State Structure for the Disproportionation of Hydrogen Atom with Ethyl Radical. *Journal of the American Chemical Society* **107**:5800-5801.)

can also be helpful to examine the atomic displacements that correspond to the negative eigenvector, to ensure that it corresponds to the correct motion over the saddle point as for the $\text{Cl}^- + \text{CH}_3\text{Cl}$ reaction.

As one approaches the saddle point from a minimum, the Hessian matrix will change from having all positive eigenvalues to including one negative value. The *quadratic region* of a saddle point is that portion of the energy surface surrounding the point where the Hessian contains one negative eigenvalue. Similarly the quadratic region of a minimum is the

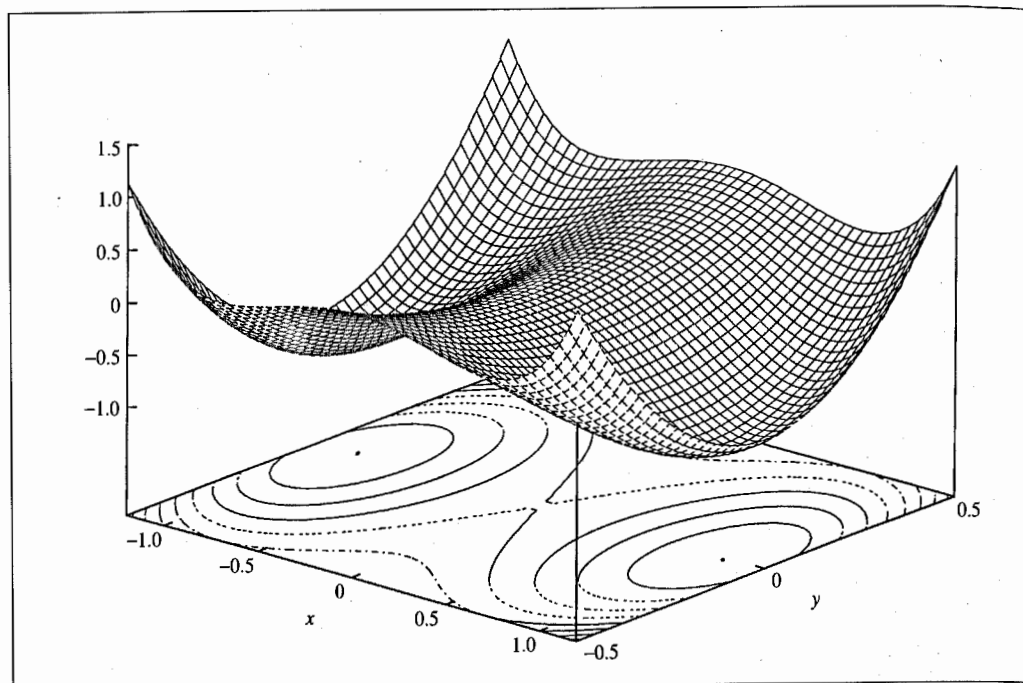


Fig. 5.24: The function $f(x, y) = x^4 + 4x^2y^2 - 2x^2 + 2y^2$ has a saddle point at $(0, 0)$ and minima at $(1, 0)$ and $(-1, 0)$.

region where all eigenvalues are positive and the Hessian is positive definite. Some algorithms for finding saddle points require a starting geometry within the quadratic region. We can illustrate the concept of a quadratic region by considering the function $f(x, y) = x^4 + 4x^2y^2 - 2x^2 + 2y^2$, which is drawn in Figure 5.24. This function has two minima at $(1, 0)$ and $(-1, 0)$ and one saddle point at $(0, 0)$. In this case it is possible to derive and characterise the stationary points analytically. The Hessian matrix of second derivatives for this function is:

$$\begin{pmatrix} 12x^2 + 8y^2 - 4 & 16xy \\ 16xy & 8x^2 + 4 \end{pmatrix} \quad (5.39)$$

At the point $(1, 0)$ the Hessian matrix is thus

$$\begin{pmatrix} 8 & 0 \\ 0 & 4 \end{pmatrix} \quad (5.40)$$

The eigenvalues of this matrix are obtained by setting the secular determinant to zero:

$$\begin{vmatrix} 8 - \lambda & 0 \\ 0 & 4 - \lambda \end{vmatrix} = 0 \quad (5.41)$$

The eigenvalues are $\lambda = 4$ and $\lambda = 8$. Thus both eigenvalues are positive and the point is a minimum. At the point $(0, 0)$ the Hessian matrix is

$$\begin{pmatrix} -4 & 0 \\ 0 & 4 \end{pmatrix} \quad (5.42)$$

with one negative and one positive eigenvalue (-4 and $+4$). The normalised eigenvectors corresponding to these eigenvalues are $(0, 1)$ for the eigenvalue $\lambda = 4$ and $(1, 0)$ for the eigenvalue $\lambda = -4$. These eigenvectors indicate the directions in which the gradient of the function changes sign. Thus along the line $x = 0$ the function passes through a minimum, as can be seen from Figure 5.24. By contrast, if one progresses from $(-1, 0)$ to $(1, 0)$ through the origin then the function passes through a maximum. As one progresses through a transition structure the eigenvector of the negative eigenvalue corresponds to the concerted motions of the atoms that give rise to motion through the saddle point. If we move along the x axis from the minimum at $(1, 0)$ to the saddle point at the origin, both eigenvalues will be positive so long as $12x^2 + 8y^2 - 4 > 0$. Thus, so long as x is larger than $1/\sqrt{3}$ the eigenvalues of the Hessian matrix will be positive. When x becomes smaller than $1/\sqrt{3}$ there will be one negative and one positive eigenvalue. In this case the quadratic region would correspond to all points where the absolute value of x was less than $1/\sqrt{3}$.

5.9.1 Methods to Locate Saddle Points

In some simple cases such as the chloride/methyl chloride reaction the geometry of the transition structure can be predicted by inspection. In other cases a *grid search* can be used to scan the energy surface in order to locate the approximate position of the transition state. In a grid search, the coordinates are systematically varied to generate a set of structures, for each of which the energy is calculated. It may then be possible to fit an analytical expression to these points, from which the saddle point can be predicted by standard calculus methods. The grid search method is widely used for constructing potential energy surfaces but is restricted to systems with a very small number of atoms or where only a limited number of degrees of freedom are being explored such as the $\text{H} + \text{H}_2 \rightarrow \text{H}_2 + \text{H}$ reaction. An advantage of the grid search is that it does provide information about the energy surface away from the pathway, which can be important if one wishes to investigate the dynamics of a reaction and the interconversion of energy between different modes. The grid search method is not the method of choice for all but the smallest systems due to the number of energy evaluations that are required. In any case it does not directly provide the transition structure.

The conversion of one minimum-energy structure into another may sometimes occur primarily along just one or two coordinates. In such cases, an approximation to the reaction pathway can be obtained by gradually changing the coordinate(s), allowing the system to relax at each stage using minimisation while keeping the chosen coordinate(s) fixed. The point of highest energy on the path is an approximation to the saddle point and the structures generated during the course of the calculation can be considered to represent a sequence of points on the interconversion pathway. When such coordinate driving methods are applied to conformational changes that occur primarily via rotation about bonds, the

the Hessian matrix of second derivatives is available then the appropriate direction to take is uphill along the eigenvector of the smallest eigenvalue when all eigenvalues are positive and downhill along the eigenvector corresponding to the negative eigenvalue when within the quadratic region of the saddle point [Baker 1986].

As we have stated frequently, at a saddle point the gradient is zero (as it is for a minimum). It might therefore be imagined that a minimisation algorithm (or some variant) could be used to locate saddle points. Some minimisation algorithms can occasionally incorrectly converge to a saddle point, especially if the starting structure is close to the transition structure. A simple example is the Newton-Raphson method, which will converge to a transition structure when giving a starting position that is within the quadratic region. Other minimisation algorithms can also be modified so that they consistently locate saddle points when provided with an initial structure within the quadratic region [Schlegel 1982].

5.9.2 Reaction Path Following

The traditional way to elucidate the reaction path is to move downhill from a saddle point to the two associated minima. There may be many different paths that could be followed from the saddle point to the associated minima. The *intrinsic reaction coordinate* (IRC) is the path that would be followed by a particle moving along the steepest descents path with an infinitely small step from the transition structure down to each minimum when the system is described using mass-weighted coordinates (as in a normal mode calculation) [Fukui 1981]. The initial directions towards each minimum can be obtained directly from the eigenvector that corresponds to the imaginary frequency at the transition structure. A simple steepest descents algorithm with a reasonable step size will usually give a path that oscillates about the true minimum energy path, as illustrated in Figure 5.28. This is perfectly acceptable in a minimisation, where the objective is to locate the minimum as efficiently as possible and where we are not interested in the intermediate structures. To determine the true reaction pathway (or a better approximation to it) it is necessary to 'correct' the path taken by the steepest descents algorithm. These corrective methods are especially useful when the path is curved.

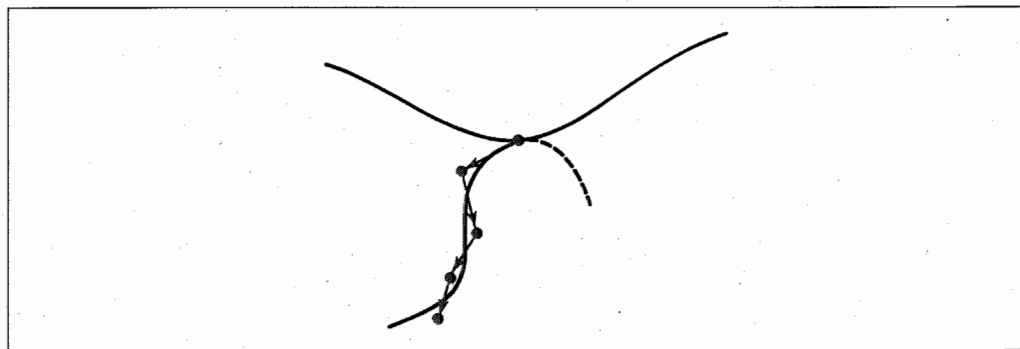


Fig. 5.28: A steepest descents minimisation algorithm produces a path that oscillates about the true reaction pathway from the transition structure to a minimum.

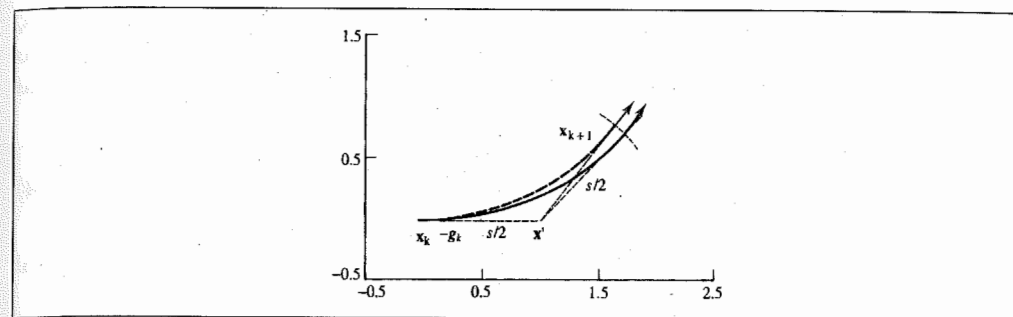


Fig. 5.29: Method for correcting the path followed by a steepest descents algorithm to generate the intrinsic reaction coordinate. The solid line shows the real path and the dotted line shows the algorithmic approximation to it. (Figure redrawn from Gonzalez C and H B Schlegel 1988. An Improved Algorithm for Reaction Path Following. *Journal of Chemical Physics* 90:2154-2161.)

Many different algorithms have been suggested for determining reaction paths. The real challenge is to find an approach that is sufficiently general to work well in many (if not all) situations and with relatively little computational expense. One widely used method was devised by Gonzalez and Schlegel [Gonzalez and Schlegel 1988] and is illustrated in Figure 5.29. First it calculates the gradient at the current point, x_k . A step of length $s/2$ is taken along the direction of this gradient to give a new point (x'). The next point on the reaction path is obtained by minimising the energy subject to the constraint that the distance between x' and the new point on the reaction path (x_{k+1}) is $s/2$. The reaction path is then approximated by a circle that passes through both x_k and x_{k+1} and whose tangents at those two points are in the directions of the gradients. A refined version of this path-following algorithm has been incorporated into an efficient combined procedure which can determine reaction paths, minima and transition state geometries [Ayala and Schlegel 1997] without the need for second derivatives to be calculated.

5.9.3 Transition Structures and Reaction Pathways for Large Systems

Most of the algorithms we have discussed so far, with the possible exception of adiabatic mapping, were originally designed to be used with quantum mechanics where relatively small numbers of atoms are involved. It is often difficult to apply these methods to the study of conformational transitions. There are several reasons for this, but one important feature is that it is assumed that there is only one saddle point between the initial and final states. There may be a number of transition structures along the pathway between two conformations of a complex molecule. Here we will discuss two related methods that were originally designed to tackle this problem using molecular mechanics.

In the self-penalty walk (SPW) method of Czerminski and Elber [Czerminski and Elber 1990; Nowak *et al.* 1991] a 'polymer' is constructed that consists of a series of $M + 2$ 'monomers'. Each monomer is a complete copy of the actual system and so there are $(M + 2)N$ atoms present in the calculation. The two ends of the polymer correspond to the two minima between which we are trying to elucidate the pathway (the 'reactant' and the 'product').

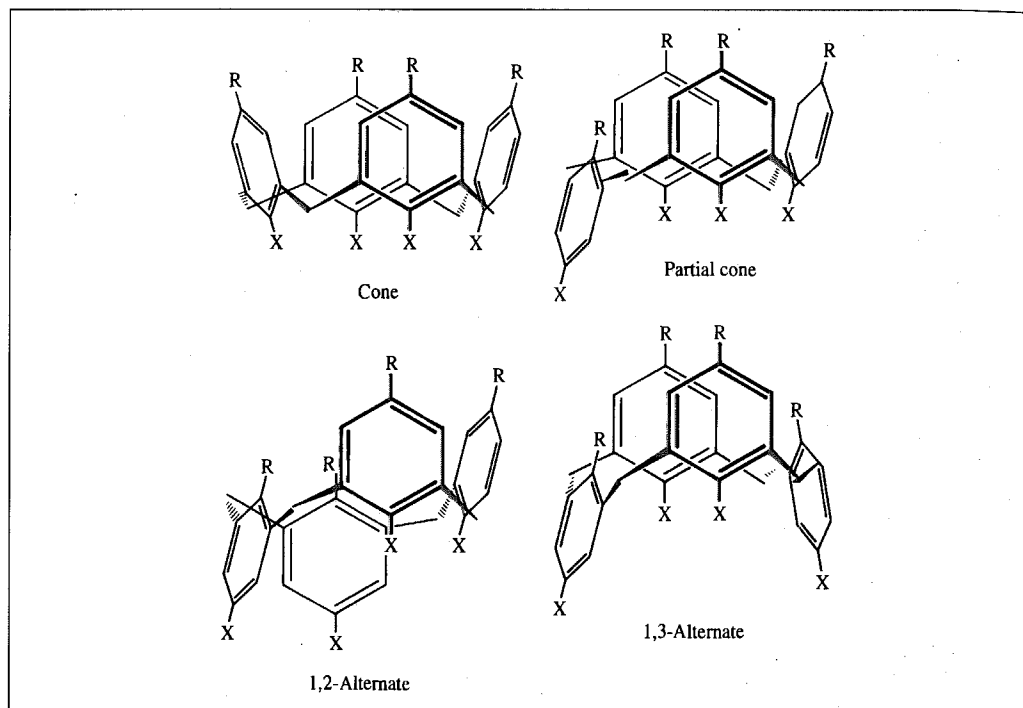


Fig. 5.31: Possible conformations of the calix[4]arene systems. (Figure adapted from Fischer S, P D J Groothuis, L C Groenen, W P van Hoorn, F C J M van Geggel, D N Reinhoudt and M Karplus 1995. Pathways for Conformational Interconversion of Calix[4]arenes. *Journal of the American Chemical Society* 117:1611–1620.)

involved, giving the energy diagram in Figure 5.32. The predicted activation barrier of 14.5 kcal/mol for the cone \rightarrow inverted cone transition was in very good agreement with the experimentally determined value of 14.2 kcal/mol. Much of the barrier (9.1 kcal/mol) was due to the need to break two hydrogen bonds; the remainder was due to the need to deform some bond angles such as those of the bridging methylene carbons.

5.9.4 The Transition Structures of Pericyclic Reactions

One of the most celebrated examples of the use of quantum mechanical methods in understanding chemical reactivity is the work of Woodward and Hoffmann [Woodward and Hoffmann 1969] who were able to explain the experimentally observed nature of certain types of concerted reaction. The reactions which they studied include cycloadditions, sigmatropic rearrangements, cheletropic reactions, electrocyclic reactions and the ene reaction (Figure 5.33) and are collectively known as pericyclic reactions. The products obtained from such reactions can be understood in terms of simple mechanistic arguments, but such arguments cannot explain some aspects. In particular, the reactions are often highly stereospecific with the reaction rates and the stereoselectivity changes dramatically with

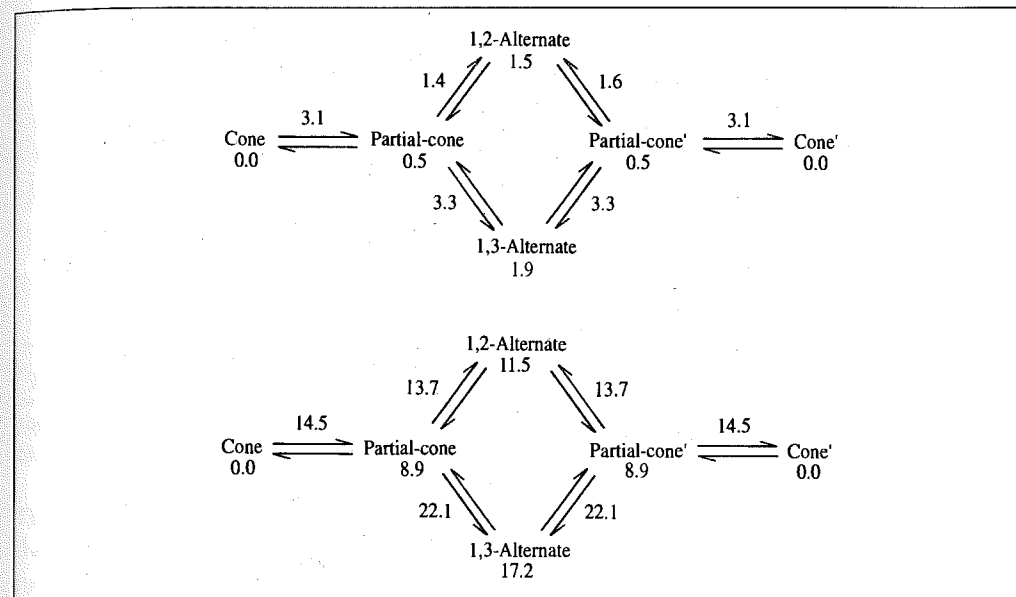


Fig. 5.32: Interconversion between various conformations of calix[4]arenes. X = H, R = H (top); X = OH, R = H (bottom). Energies in kcal/mol.

the reaction conditions. Woodward and Hoffmann successfully employed molecular orbital theory to rationalise the existing data and their theory has also been very successful in predicting the outcome of similar reactions. The basic principle applied by Woodward and Hoffmann was that of the conservation of orbital symmetry and as a consequence of their work a series of rules (often called the Woodward–Hoffmann rules) were developed. The Woodward–Hoffmann rules apply only to concerted reactions and are based upon the principle that maximum bonding is maintained throughout the course of a reaction. Fukui also discovered the importance of orbital symmetry and suggested that the majority of chemical reactions should take place at the position of, and in the direction of, maximum overlap between the highest occupied molecular orbital (HOMO) of one species and the lowest unoccupied molecular orbital (LUMO) of the other component [Fukui 1971]. These orbitals are collectively known as the *frontier orbitals*.

The HOMO–LUMO interaction depends on various factors, including the geometry of approach (which affects the amount of overlap), the phase relationship of the orbitals and their energy separation. For example, the HOMO and LUMO of ethene are illustrated pictorially in Figure 5.34. The most obvious mode of interaction between the two molecules involves suprafacial attack shown in Figure 5.34 to give cyclobutane. However, the symmetries of the overlapping orbitals must have the same phase for a favourable interaction to occur and this is not possible for ethene unless an energetically unfavourable antarafacial approach is adopted. By contrast, the interaction between ethene and the butadiene does occur in a suprafacial sense with both HOMO/LUMO pairs of orbitals having the appropriate phase relationship (Figure 5.34).

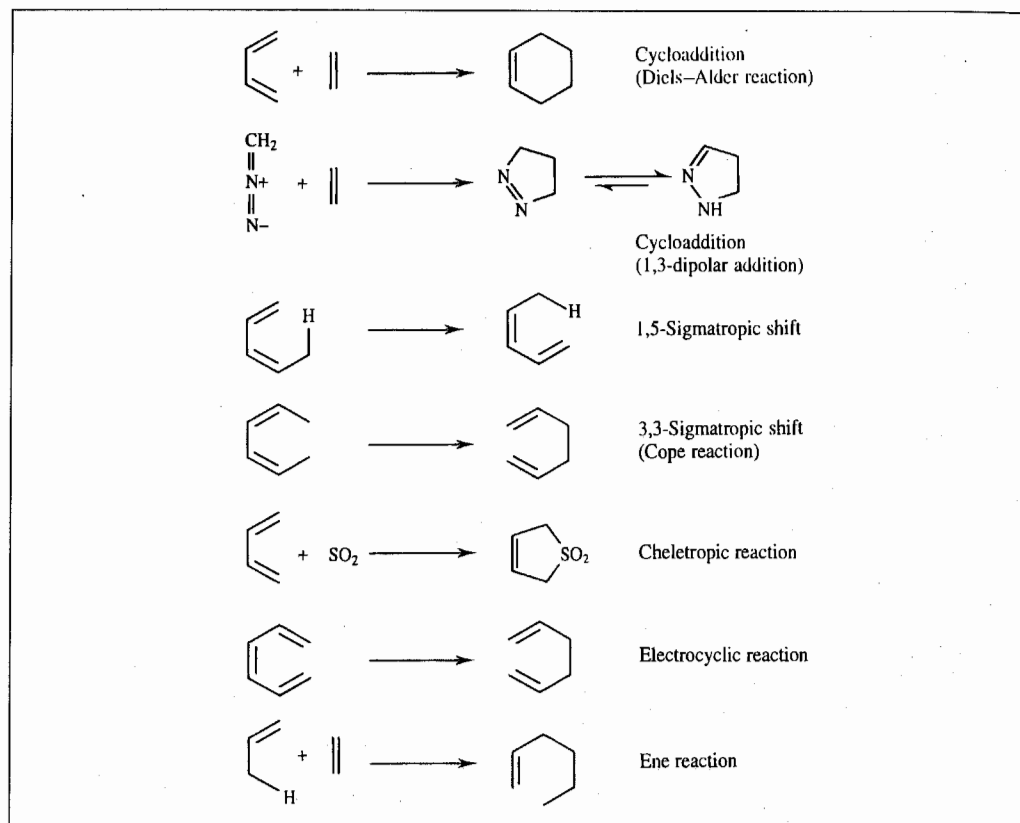


Fig. 5.33: Typical pericyclic reactions.

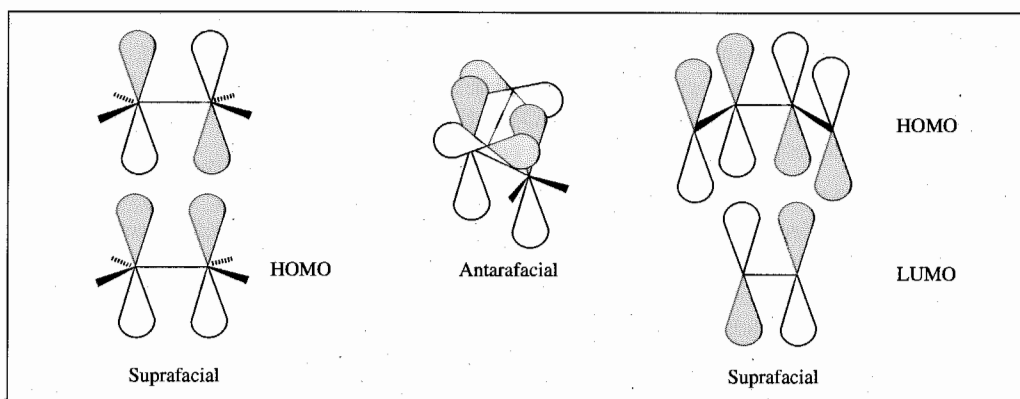


Fig. 5.34: Suprafacial attack of one ethene molecule on another (left) is not permitted by the Woodward-Hoffmann rules and the alternative antarafacial mode of attack is sterically unfavourable. Suprafacial attack is however permitted for the Diels-Alder reaction between butadiene and ethene (right).

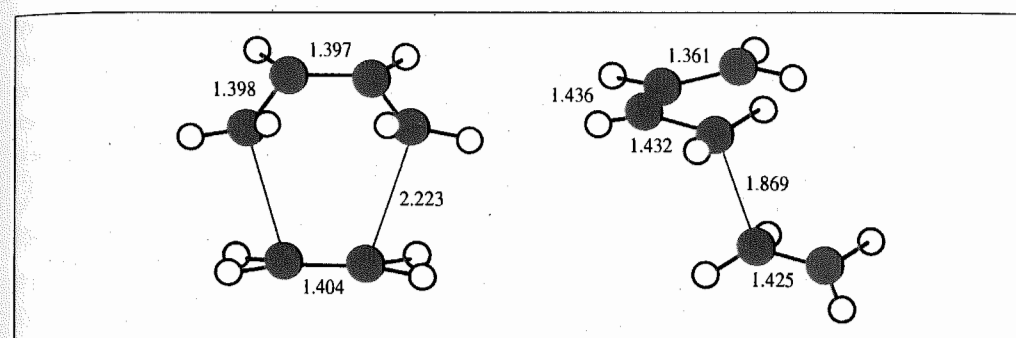


Fig. 5.35: Geometry predicted by CASSCF *ab initio* calculations of the two possible transition structure geometries for the Diels-Alder reaction between ethene and butadiene. (Figure adapted from Houk K N, J González and Y Li 1995. *Pericyclic Reaction Transition States: Passions and Punctilios 1935-1995*. Accounts of Chemical Research 28:81-90.)

The Woodward-Hoffmann rules state what the outcome of a pericyclic reaction will be, but they do not define the mechanism by which the reaction occurs. Many theoretical techniques have been applied to the study of these problems over the years [Houk *et al.* 1992] and a passionate debate has ensued on the nature of the transition structures involved in these reactions. The debate has been fuelled by the fact that different theoretical treatments (especially semi-empirical methods) give different results. For example, at one extreme the Diels-Alder reaction between butadiene and ethene would proceed via a two-step mechanism involving a biradial transition structure. At the other extreme the reaction would involve a symmetrical transition state formed in a concerted, synchronous reaction. *Ab initio* calculations at various levels of theory suggest the concerted transition structure. The geometry obtained for the prototypical Diels-Alder reaction between butadiene/ethene using a CASSCF calculation and a 6-31G* basis set is shown in Figure 5.35 [Houk *et al.* 1995]. The alternative biradial structure is also shown in Figure 5.35; this is predicted to be 6 kcal/mol higher in energy than the symmetrical transition structure.

5.10 Solid-state Systems: Lattice Statics and Lattice Dynamics

Energy minimisation and normal mode analysis have an important role to play in the study of the solid state. Algorithms similar to those discussed above are employed but an extra feature of such systems, at least when they form a perfect lattice, is that it is possible to exploit the space group symmetry of the lattice to speed up the calculations. It is also important to properly take the interactions with atoms in neighbouring cells into account.

The most straightforward type of lattice minimisation is performed at constant volume, where the dimensions of the basic unit cell do not change. A more advanced type of calculation is one performed at constant pressure, in which case there are forces on both the atoms and the unit cell as a whole. The lattice vectors are considered as additional variables along with the atomic coordinates. The laws of elasticity describe the behaviour of a material when

subjected to a *stress* (defined as the force per unit area). One obvious source of stress is any external pressure, but stress may also arise from other sources, especially from interatomic forces within the cell, which give rise to 'internal stress'. The concept of *strain* is also key to this subject; the strain is the fractional change in the dimension (for example, the change per unit length when a steel rod is stretched). In the general case we consider a situation where a point \mathbf{r} in the unstrained material moves to a new point \mathbf{r}' under the effect of some strain:

$$\mathbf{u} = \mathbf{r}' - \mathbf{r} \quad (5.48)$$

If we apply the strain uniformly in one dimension (e.g. the x axis) then the x coordinate of a point that was initially at x will change by an amount proportional to x . This is written:

$$u_x = \varepsilon_{xx}x \quad (5.49)$$

In the general case the constant of proportionality is written as the first derivative:

$$\varepsilon_{xx} = \partial u_x / \partial x \quad (5.50)$$

Deformation in the y and z directions is described in an analogous manner. In order to cater for shear-type strains additional elements are defined:

$$\varepsilon_{xy} = \varepsilon_{yx} = \frac{1}{2}(\partial u_y / \partial x + \partial u_x / \partial y) \quad (5.51)$$

These values ε give rise to a *strain tensor* (see Section 4.9.1 for more discussion on tensors), which is symmetric and is often written in the following form:

$$\boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 & \frac{1}{2}\varepsilon_6 & \frac{1}{2}\varepsilon_5 \\ \frac{1}{2}\varepsilon_6 & \varepsilon_2 & \frac{1}{2}\varepsilon_4 \\ \frac{1}{2}\varepsilon_5 & \frac{1}{2}\varepsilon_4 & \varepsilon_3 \end{bmatrix} \quad \begin{matrix} \varepsilon_1 \equiv \varepsilon_{xx}, & \varepsilon_2 \equiv \varepsilon_{yy}, & \varepsilon_3 \equiv \varepsilon_{zz} \\ \varepsilon_4 \equiv \varepsilon_{yz}, & \varepsilon_5 \equiv \varepsilon_{xz}, & \varepsilon_6 \equiv \varepsilon_{xy} \end{matrix} \quad (5.52)$$

There are thus six different numbers present in the strain tensor. The symmetric form of the strain tensor prevents rotation of the unit cell with respect to the Cartesian axis system. It is possible to use this matrix to relate how a vector \mathbf{r} in the unstrained matrix is related to one \mathbf{r}' in the strained structure as follows:

$$\mathbf{r}' = (\mathbf{I} + \boldsymbol{\varepsilon})\mathbf{r} \quad (5.53)$$

\mathbf{I} is the identity matrix. The six first derivatives of the energy with respect to the strain components ε_i measure the forces acting on the unit cell. When combined with the atomic coordinates we get a matrix with $3N + 6$ dimensions. At a minimum not only should there be no force on any of the atoms but the forces on the unit cell should also be zero. Application of a standard iterative minimisation procedure such as the Davidon-Fletcher-Powell method will optimise all these degrees of freedom to give a strain-free final structure. In such procedures a reasonably accurate estimate of the initial inverse Hessian matrix is usually required to ensure that the changes in the atomic positions and in the cell dimensions are matched.

Two common properties which can be calculated from the minimum-energy structure are the elastic and dielectric constants. The elastic constant matrix is used to relate the strains of a material to the internal forces, or stresses. It is defined as the second derivative of the energy with respect to the strain, normalised by the cell volume. The inverse of the elastic

constant matrix gives the constant of proportionality between the stress and the strain. The elastic constant matrix has dimensions 6×6 and is given by the following expression:

$$\mathbf{C} = \frac{1}{V} [\mathcal{V}_{\varepsilon\varepsilon}'' - (\mathcal{V}_{\varepsilon\tau}'' \cdot \mathcal{V}_{\tau\tau}''^{-1} \cdot \mathcal{V}_{\tau\varepsilon}'')] \quad (5.54)$$

In this equation $\mathcal{V}_{\varepsilon\varepsilon}''$ is the 6×6 matrix of second derivatives (elements $\partial^2 \mathcal{V} / \partial \varepsilon_{ij}^2$), $\mathcal{V}_{\varepsilon\tau}''$ and $\mathcal{V}_{\tau\varepsilon}''$ are the corresponding $3N \times 6$ and $6 \times 3N$ mixed coordinate/strain matrices, $\mathcal{V}_{\tau\tau}''$ is the $3N \times 3N$ second-derivative coordinate matrix and V is the unit cell volume. It is the second term in Equation (5.54) that accounts for internal atomic relaxations as the cell distorts.

The strains on the lattice are equal to the stress divided by the elastic constant matrix:

$$\boldsymbol{\varepsilon} = (P_{\text{static}} + P_{\text{applied}}) \cdot \mathbf{C}^{-1} \quad (5.55)$$

Here we have expressed the stress as the sum of the (external) applied pressure P_{applied} together with a static pressure P_{static} , which arises from the internal forces acting on the unit cell.

The dielectric constant is concerned with the electrical properties of a material. The dielectric constant for a solid is a 3×3 matrix with different components according to the Cartesian axes. These elements are given by:

$$D_{ij} = \delta_{ij} + \frac{4\pi}{V} \mathbf{q}^T \cdot \mathcal{V}_{\tau\tau}''^{-1} \cdot \mathbf{q} \quad (5.56)$$

In Equation (5.56) i and j are one of x , y or z ; δ_{ij} is the delta function (i.e. equal to one when $i = j$ and zero otherwise) and \mathbf{q} is a vector containing the charges of each species. It is well known that the effect of a dielectric changes in an oscillating electric field (at high enough frequencies the permanent dipoles in the material are unable to keep up with the rapidly changing field). Thus one usually calculates two sets of dielectric constant matrices, corresponding to the low- and high-frequency regimes. If polarisation is included via a shell model (see Section 4.22.2) then both the cores and the shells are used to determine the low-frequency dielectric matrix; at high frequency only the shells are considered.

Comparison of the relative energies following a minimisation calculation can enable predictions to be made of the likely structure for a given material. In the same way that an organic molecule may be able to exist in more than one three-dimensional structure (or conformation - see Chapter 9) so a solid may (in principle) be able to adopt more than one three-dimensional arrangement of its atoms whilst still maintaining a periodic lattice structure. Silica, SiO_2 , has been the subject of considerable attention using these methods. The lowest-energy form is $\alpha\text{-SiO}_2$, or quartz. However, it can also form more open structures. A number of such microporous structures are in principle available, three being silicalite, mordenite and faujasite. In one study the energies of these structures relative to the quartz structure were found to be approximately 2.6, 4.9 and 5.1 kcal/mol, respectively [Ooms *et al.* 1988]. Indeed, the silicalite structure is the only one which can be prepared as the pure silicon oxide; the other forms usually require a high aluminium content and are more traditional zeolites. In an extension of this work two slightly different forms of the silicalite structure were simulated. The normal form at room temperature has orthorhombic

symmetry but at low temperatures this changes to monoclinic. These two forms are very closely related, differing only by the distortion of a key angle by 0.64° . Nevertheless, an energy-minimisation calculation starting from the orthorhombic structure did indeed change to the monoclinic, in agreement with the experimental data [Bell *et al.* 1990]. The orthorhombic \rightarrow monoclinic transition could only be observed using a force field which included polarisation effects (i.e. the shell model). Lattice minimisation methods can sometimes be very useful in helping to solve the structure of materials, a noteworthy example being the determination of the structure of a zeolite NU-87 [Shannon *et al.* 1991]. This synthetic material is of particular interest as a catalyst as it contains a multidimensional channel system. Multidimensional systems permit more complex catalytic reactions to occur and are also less prone to deactivation than one-dimensional systems. In this case, there are rings containing ten and twelve oxygen atoms (Figure 5.36 (colour plate section)). NU-87 also has a high silica content, which confers improved stability to heat. A number of experimental techniques were used to try to determine the structure, including electron diffraction and powder synchrotron X-ray diffraction, as a result of which an approximate structure was deduced but there remained some features in the powder diffraction spectrum that could not be accounted for. These were initially believed to be due to impurities but after energy-minimisation studies some subtle changes in the structure occurred to give a related structure that was a better match to the experimental data. A key feature of this particular minimisation was that the structure was not forced to adopt any specific symmetry but rather each atom was able to move independently of the others.

The calculation of vibrational frequencies (called *phonons*) is important to the study of the solid state. Indeed, the calculation of and study of phonons is often given a special name, *lattice dynamics*. To calculate the vibrational frequencies for a solid one follows a very similar approach to that described earlier for molecules, with the exception that when a shell model is being used* then their effect must be incorporated into the mass-weighted matrix of second derivatives (though not directly as they have no mass):

$$\mathcal{V}'' = \mathcal{V}''_{\text{core-core}} - \mathcal{V}''_{\text{core-shell}} \cdot \mathcal{V}''_{\text{shell-shell}}^{-1} \cdot \mathcal{V}''_{\text{core-shell}} \quad (5.57)$$

Of additional importance is that the vibrational modes are dependent upon the reciprocal lattice vector \mathbf{k} . As with calculations of the electronic structure of periodic lattices these calculations are usually performed by selecting a suitable set of points from within the Brillouin zone. For periodic solids it is necessary to take this periodicity into account; the effect on the second-derivative matrix is that each element ij needs to be multiplied by the phase factor $\exp(i\mathbf{k} \cdot \mathbf{r}_{ij})$. A *phonon dispersion curve* indicates how the phonon frequencies vary over the Brillouin zone, an example being shown in Figure 5.37. The phonon density of states is the variation in the number of frequencies as a function of frequency. A purely transverse vibration is one where the displacement of the atoms is perpendicular to the direction of motion of the wave; in a purely longitudinal vibration the atomic displacements are in the direction of the wave motion. Such motions can be observed in simple systems (e.g. those that contain just one or two atoms per unit cell) but for general three-dimensional lattices most of the vibrations are a mixture of transverse and longitudinal motions, the exceptions

* The use of a shell model is generally recommended otherwise the resulting frequencies are too high.

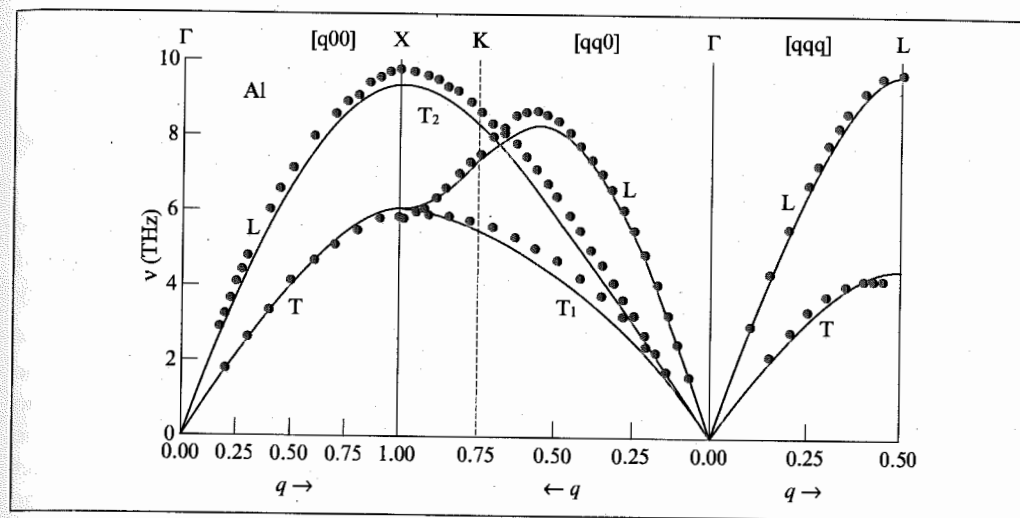


Fig. 5.37: Comparison of the calculated phonon dispersion curve for Al with the experimental values measured using neutron diffraction. (Figure redrawn from Michin Y, D Farkas, M J Mehl and D A Papaconstantopoulos 1999. *Interatomic Potentials for Monatomic Metals from Experimental Data and ab initio Calculations*. Physical Review B59:3393–3407.)

being those along directions of high symmetry. The phonons are additionally classified as acoustic or optical; the former are typically of longer wavelength (lower-frequency oscillations) where the atoms move as a unit. The name arises from the fact that these are often measured as sound waves. At the point $\mathbf{k} = 0$ (the gamma point) the first three vibrational frequencies correspond to translation of the entire lattice. The optical phonons are typically higher in frequency. Various experimental techniques can be used to investigate lattice vibrations and to determine the phonon dispersion curves, the most powerful of which is inelastic scattering using thermal neutrons. These often allow the entire range of \mathbf{k} to be sampled, in contrast to some of the alternative types of radiation.

Once the phonon frequencies are known it becomes possible to determine various thermodynamic quantities using statistical mechanics (see Appendix 6.1). Here again some slight modifications are required to the standard formulae. These modifications are usually a consequence of the need to sum over the points sampled in the Brillouin zone. For example, the zero-point energy is:

$$U_{\text{vib}}(0) = \sum_{q=1}^p w_q \sum_{i=1}^{N_{\text{nm}}} \frac{h\nu_i}{2} \quad (5.58)$$

In Equation (5.58) the outer summation is over the p points q which are used to sample the Brillouin zone, w_q is the fractional weight associated with each point (related to the volume of Brillouin zone space surrounding q) and ν_i are the phonon frequencies. In addition to the internal energy due to the vibrational modes it is also possible to calculate the vibrational entropy, and hence the free energy. The Helmholtz free energy at a temperature

T is thus given in the quasi-harmonic approximation by the sum of static and vibrational contributions:

$$A = \mathcal{V} + \sum_{q=1}^p w_q \sum_{i=1}^{N_{\text{nm}}} \left(\frac{h\nu_i}{2} + k_{\text{B}}T \ln \left[1 - \exp \left(-\frac{h\nu_i}{k_{\text{B}}T} \right) \right] \right) \quad (5.59)$$

Here, \mathcal{V} is the internal energy calculated from the potential energy model. The heat capacity at constant volume is another useful thermodynamic quantity that can be determined directly from the frequencies as it equals the derivative of the vibrational internal energy with respect to temperature.

An extension of these ideas involves minimisation of the free energy as a function of the coordinates and the temperature. The function to be minimised is sometimes referred to as an *availability*, given by $G^* = A + P_{\text{ext}}V$ where P_{ext} is the external pressure and V is the volume. Such a free energy minimisation requires derivatives of the free energy with respect to the coordinates. Early implementations used approximations such as separating the changes in the external coordinates (i.e. the dimensions of the unit cell) from the internal coordinates (i.e. the locations of the ions within the unit cell). In addition, true free energy derivatives might be calculated for the external coordinates only (due to the computational cost) with the internal coordinates being changed using the static potential energy. It is now possible to calculate a full set of analytical free energy first derivatives and hence to perform a full minimisation of the free energy with respect to external and internal coordinates simultaneously [Taylor *et al.* 1998].

Free energy minimisation provides information that is in many ways complementary to molecular dynamics simulations [Allan *et al.* 2000]. The former is particularly useful for investigating materials at lower temperatures where the harmonic assumption is valid; moreover it includes zero-point energy and quantisation effects, which are ignored by molecular dynamics. In addition, free energy minimisation provides free energies directly, rather than free energy differences, and is computationally significantly less expensive. Conversely anharmonic effects become important at higher temperatures, making molecular dynamics and Monte Carlo more suitable. One property that can be calculated via free energy minimisation is the thermal expansivity. This involves a series of free energy minimisations at different temperatures. It is also possible to calculate the free energy of disordered solids, and thus enthalpies and entropies of mixing.

Further Reading

- Catlow C R A 1998. Solids: Computer Modelling. In Schleyer, P v R, N L Allinger, T Clark, J Gasteiger, P A Kollman, H F Schaefer III and P R Schreiner (Editors) *The Encyclopedia of Computational Chemistry*, John Wiley & Sons, Chichester.
- Gill P E and W Murray 1981. *Practical Optimization*. London, Academic Press.
- McKee M L and M Page 1993. Computing Reaction Pathways on Molecular Potential Energy Surfaces. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 4. New York, VCH Publishers, pp. 35–65.

- Press W H, B P Flannery, S A Teukolsky and W T Vetterling 1992. *Numerical Recipes in Fortran*. Cambridge, Cambridge University Press.
- Schlegel H B 1987. Optimization of Equilibrium Geometries and Transition Structures. In Lawley K P (Editor) *Ab Initio Methods in Quantum Chemistry – I*. New York, John Wiley & Sons, pp. 249–286.
- Schlegel H B 1989. Some Practical Suggestions for Optimizing Geometries and Locating Transition States. In Bertrán J and I G Csizmadia (Editors). *New Theoretical Concepts for Understanding Organic Reactions*. Dordrecht, Kluwer, pp. 33–53.
- Schlick T 1992. Optimization Methods in Computational Chemistry. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 3. New York, VCH Publishers, pp. 1–71.
- Stassis C 19. Lattice Dynamics. In Sköld and D L Price (Editors) *Methods of Experimental Physics* Volume 23: *Neutron Scattering Part A*. Orlando, Academic Press, pp. 369–440.
- Watson G W, P Tschaufeser, A Wall, R A Jackson and S C Parker 1997. Lattice Energy and Free Energy Minimisation Techniques. *Computer Modelling in Inorganic Crystallography*. San Diego, Academic Press, pp. 55–81.
- Williams I H 1993. Interplay of Theory and Experiment in the Determination of Transition-State Structures. *Chemical Society Reviews* 1:277–283.

References

- Allan N L, G D Barrera, J A Purton, C E Sims and M B Taylor 2000. Ionic Solids at High Temperatures and Pressures: *Ab initio*, Lattice Dynamics and Monte Carlo Studies. *Physical Chemistry Chemical Physics* 2:1099–1111.
- Ayala P Y and H B Schlegel 1997. A Combined Method for Determining Reaction Paths, Minima and Transition State Geometries. *Journal of Chemical Physics* 107:375–384.
- Baker J 1986. An Algorithm for the Location of Transition States. *Journal of Computational Chemistry* 7:385–395.
- Bell R G, R A Jackson and C R A Catlow 1990. Computer Simulation of the Monoclinic Distortion in Silicalite. *Journal of the Chemical Society Chemical Communications* 10:782–783.
- Brooks B and M Karplus 1983. Harmonic Dynamics of Proteins: Normal Modes and Fluctuations in Bovine Pancreatic Trypsin Inhibitor. *Proceedings of the National Academy of Sciences USA* 80:6571–6575.
- Czerminski R and R Elber 1990. Self-Avoiding Walk Between 2 Fixed-Points as a Tool to Calculate Reaction Paths in Large Molecular Systems. *International Journal of Quantum Chemistry* 52:167–186.
- Dauber-Osguthorpe P, V A Roberts, D J Osguthorpe, J Wolff, M Genest and A T Hagler 1988. Structure and Energetics of Ligand Binding to Proteins: *Escherichia coli* Dihydrofolate Reductase-Tri-methoprim, A Drug-Receptor System. *Proteins: Structure, Function and Genetics* 4:31–47.
- Doubleday C, J McIver, M Page and T Zielinski 1985. Temperature Dependence of the Transition-State Structure for the Disproportionation of Hydrogen Atom with Ethyl Radical. *Journal of the American Chemical Society* 107:5800–5801.
- Elber R and M Karplus 1987. A Method for Determining Reaction Paths in Large Molecules: Application to Myoglobin. *Chemical Physics Letters* 139:375–380.
- Fischer S, P D J Groothuis, L C Groenen, W P van Hoorn, F C J M van Geggel, D N Reinhoudt and M Karplus 1995. Pathways for Conformational Interconversion of Calix[4]arenes. *Journal of the American Chemical Society* 117:1611–1620.
- Fischer S and M Karplus 1992. Conjugate Peak Refinement: An Algorithm for Finding Reaction Paths and Accurate Transition States in Systems with Many Degrees of Freedom. *Chemical Physics Letters* 194:252–261.

- Fisher C L, V A Roberts and A T Hagler 1991. Influence of Environment on the Antifolate Drug Trimethoprim: Energy Minimization Studies. *Biochemistry* **30**:3518–3526.
- Fukui K 1971. Recognition of Stereochemical Paths by Orbital Interaction. *Accounts of Chemical Research* **4**:57–64.
- Fukui K 1981. The Path of Chemical Reactions – The IRC Approach. *Accounts of Chemical Research* **14**:368–375.
- Gelin B R and M Karplus 1975. Sidechain Torsional Potential and Motion of Amino Acids in Proteins: Bovine Pancreatic Trypsin Inhibitor. *Proceedings of the National Academy of Sciences USA* **72**:2002–2006.
- González C and H B Schlegel 1988. An Improved Algorithm for Reaction Path Following. *Journal of Chemical Physics* **90**:2154–2161.
- Houk K N, J González and Y Li 1995. Pericyclic Reaction Transition States: Passions and Punctilios 1935–1995. *Accounts of Chemical Research* **28**:81–90.
- Houk K N, Y Li and J D Evanseck 1992. Transition Structures of Hydrocarbon Pericyclic Reactions. *Angewandte Chemie International Edition in English* **31**:682–708.
- Nowak W, R Czerminski and R Elber 1991. Reaction Path Study of Ligand Diffusion in Proteins: Application of the Self Penalty Walk (SPW) Method to Calculate Reaction Coordinates for the Motion of CO through Leghemoglobin. *Journal of the American Chemical Society* **113**:5627–5737.
- Ooms G, R A van Santen, C J J Den Ouden, R A Jackson and C R A Catlow 1988. Relative Stabilities of Zeolitic Aluminosilicates. *Journal of Physical Chemistry* **92**:4462–4465.
- Peng C, P Y Ayala, H B Schlegel and M J Frisch 1996. Using Redundant Internal Coordinates to Optimise Equilibrium Geometries and Transition States. *Journal of Computational Chemistry* **17**:49–56.
- Pople J A, A P Scott, M W Wong and L Radom 1993. Scaling Factors for Obtaining Fundamental Vibrational Frequencies and Zero-Point Energies from HF/6-31G* and MP2/6-31G* Harmonic Frequencies. *Israel Journal of Chemistry* **33**:345–350.
- Schlegel H B 1982. Optimisation of Equilibrium Geometries and Transition Structures. *Journal of Computational Chemistry* **3**:214–218.
- Shannon M D, J L Casci, P A Cox and S J Andrews 1991. Structure of the Two-dimensional Medium-pore High-silica Zeolite NU-87. *Nature* **353**:417–420.
- Taylor M B, G D Barrera, N L Allan, T H K Barron and W C Mackrodt 1998. Shell: A Code for Lattice Dynamics and Structure Optimisation of Ionic Crystals. *Computer Physics Communications* **109**: 135–143.
- Woodward R B and R Hoffmann 1969. The Conservation of Orbital Symmetry. *Angewandte Chemie International Edition in English* **8**:781–853.

CHAPTER SIX

Computer Simulation Methods

6.1 Introduction

Energy minimisation generates individual minimum energy configurations of a system. In some cases the information provided by energy minimisation can be sufficient to predict accurately the properties of a system. If all minimum configurations on an energy surface can be identified then statistical mechanical formulae can be used to derive a partition function from which thermodynamic properties can be calculated. However, this is possible only for relatively small molecules or small molecular assemblies in the gas phase. The molecular modeller more often wants to understand and to predict the properties of liquids, solutions and solids, to study complex processes such as the adsorption of molecules onto surfaces and into solids and to investigate the behaviour of macromolecules which have many closely separated minima. In such systems the experimental measurements are made on macroscopic samples that contain extremely large numbers of atoms or molecules, with an enormous number of minima on their energy surfaces. A full quantification of the energy surfaces of such systems is not possible, nor is it ever likely to be. Computer simulation methods enable us to study such systems and predict their properties through the use of techniques that consider small replications of the macroscopic system with manageable numbers of atoms or molecules. A simulation generates representative configurations of these small replications in such a way that accurate values of structural and thermodynamic properties can be obtained with a feasible amount of computation. Simulation techniques also enable the time-dependent behaviour of atomic and molecular systems to be determined, providing a detailed picture of the way in which a system changes from one conformation or configuration to another. Simulation techniques are also widely used in some experimental procedures, such as the determination of protein structures from X-ray crystallography.

In this chapter we shall discuss some of the general principles involved in the two most common simulation techniques used in molecular modelling: the molecular dynamics and the Monte Carlo methods. We shall also discuss several concepts that are common to both of these methods. A more detailed discussion of the two simulation methods can be found in Chapters 7 and 8.

6.1.1 Time Averages, Ensemble Averages and Some Historical Background

Suppose we wish to determine experimentally the value of a property of a system such as the pressure or the heat capacity. In general, such properties will depend upon the positions and

momenta of the N particles that comprise the system. The instantaneous value of the property A can thus be written as $A(\mathbf{p}^N(t), \mathbf{r}^N(t))$, where $\mathbf{p}^N(t)$ and $\mathbf{r}^N(t)$ represent the N momenta and positions respectively at time t (i.e. $A(\mathbf{p}^N(t), \mathbf{r}^N(t)) \equiv A(p_{1x}, p_{1y}, p_{1z}, p_{2x}, \dots, x_1, y_1, z_1, x_2, \dots, t)$ where p_{1x} is the momentum of particle 1 in the x direction and x_1 is its x coordinate). Over time, the instantaneous value of the property A fluctuates as a result of interactions between the particles. The value that we measure experimentally is an average of A over the time of the measurement and is therefore known as a *time average*. As the time over which the measurement is made increases to infinity, so the value of the following integral approaches the 'true' average value of the property:

$$A_{\text{ave}} = \lim_{\tau \rightarrow \infty} \frac{1}{\tau} \int_{t=0}^{\tau} A(\mathbf{p}^N(t), \mathbf{r}^N(t)) dt \quad (6.1)$$

To calculate average values of the properties of the system, it would therefore appear to be necessary to simulate the dynamic behaviour of the system (i.e. to determine values of $A(\mathbf{p}^N(t), \mathbf{r}^N(t))$, based upon a model of the intra- and intermolecular interactions present). In principle, this is relatively straightforward to do. For any arrangement of the atoms in the system, the force acting on each atom due to interactions with other atoms can be calculated by differentiating the energy function. From the force on each atom it is possible to determine its acceleration via Newton's second law. Integration of the equations of motion should then yield a trajectory that describes how the positions, velocities and accelerations of the particles vary with time, and from which the average values of properties can be determined using the numerical equivalent of Equation (6.1). The difficulty is that for 'macroscopic' numbers of atoms or molecules (of the order of 10^{23}) it is not even feasible to determine an initial configuration of the system, let alone integrate the equations of motion and calculate a trajectory. Recognising this problem, Boltzmann and Gibbs developed statistical mechanics, in which a single system evolving in time is replaced by a large number of replications of the system that are considered simultaneously. The time average is then replaced by an *ensemble average*:

$$\langle A \rangle = \iint d\mathbf{p}^N d\mathbf{r}^N A(\mathbf{p}^N, \mathbf{r}^N) \rho(\mathbf{r}^N, \mathbf{p}^N) \quad (6.2)$$

The angle brackets $\langle \rangle$ indicate an ensemble average, or *expectation value*; that is, the average value of the property A over all replications of the ensemble generated by the simulation. Equation (6.2) is written as a double integral for convenience but in fact there should be $6N$ integral signs on the integral for the $6N$ positions and momenta of all the particles. $\rho(\mathbf{p}^N, \mathbf{r}^N)$ is the *probability density* of the ensemble; that is, the probability of finding a configuration with momenta \mathbf{p}^N and positions \mathbf{r}^N . The ensemble average of the property A is then determined by integrating over all possible configurations of the system. In accordance with the *ergodic hypothesis*, which is one of the fundamental axioms of statistical mechanics, the ensemble average is equal to the time average. Under conditions of constant number of particles, volume and temperature, the probability density is the familiar Boltzmann distribution:

$$\rho(\mathbf{p}^N, \mathbf{r}^N) = \exp(-E(\mathbf{p}^N, \mathbf{r}^N)/k_B T) / Q \quad (6.3)$$

In Equation (6.3), $E(\mathbf{p}^N, \mathbf{r}^N)$ is the energy, Q is the partition function, k_B is Boltzmann's constant and T is the temperature. The partition function is more generally written in terms of the Hamiltonian, \mathcal{H} ; for a system of N identical particles the partition function

for the canonical ensemble is as follows:

$$Q_{NVT} = \frac{1}{N!} \frac{1}{h^{3N}} \iint d\mathbf{p}^N d\mathbf{r}^N \exp \left[-\frac{\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{k_B T} \right] \quad (6.4)$$

The canonical ensemble is the name given to an ensemble for constant temperature, number of particles and volume. For our purposes \mathcal{H} can be considered the same as the total energy, $E(\mathbf{p}^N, \mathbf{r}^N)$, which equals the sum of the kinetic energy ($\mathcal{K}(\mathbf{p}^N)$) of the system, which depends upon the momenta of the particles, and the potential energy ($\mathcal{V}(\mathbf{r}^N)$), which depends upon the positions. The factor $N!$ arises from the indistinguishability of the particles and the factor $1/h^{3N}$ is required to ensure that the partition function is equal to the quantum mechanical result for a particle in a box. A short discussion of some of the key results of statistical mechanics is provided in Appendix 6.1 and further details can be found in standard textbooks.

The first computer simulations of fluids were performed in 1952 by Metropolis, Rosenbluth, Rosenbluth, Teller and Teller, who developed a scheme for sampling from the Boltzmann distribution to give ensemble averages. This gave rise to the Monte Carlo simulation method. Not long afterwards (in 1957) Alder recognised that it was, in fact, possible to integrate the equations of motion for a relatively small number of particles, and to mimic the behaviour of a real system using periodic boundary conditions. This led to the first molecular dynamics simulations of molecular systems.

6.1.2 A Brief Description of the Molecular Dynamics Method

Molecular dynamics calculates the 'real' dynamics of the system, from which time averages of properties can be calculated. Sets of atomic positions are derived in sequence by applying Newton's equations of motion. Molecular dynamics is a *deterministic* method, by which we mean that the state of the system at any future time can be predicted from its current state. The first molecular dynamics simulations were performed using very simple potentials such as the hard-sphere potential. The behaviour of the particles in this potential is similar to that of billiard or snooker balls: the particles move in straight lines at constant velocity between collisions. The collisions are perfectly elastic and occur when the separation between a pair of spheres equals the sum of their radii. After a collision, the new velocities of the colliding spheres are calculated using the principle of conservation of linear momentum. The hard-sphere model has provided many useful results but is obviously not ideal for simulating atomic or molecular systems. In potentials such as the Lennard-Jones potential the force between two atoms or molecules changes continuously with their separation. By contrast, in the hard-sphere model there is no force between particles until they collide. The continuous nature of the more realistic potentials requires the equations of motion to be integrated by breaking the calculation into a series of very short time steps (typically between 1 femtosecond and 10 femtoseconds; 10^{-15} s to 10^{-14} s). At each step, the forces on the atoms are computed and combined with the current positions and velocities to generate new positions and velocities a short time ahead. The force acting on each atom is assumed to be constant during the time interval. The atoms are then moved to the new positions, an updated set of forces is computed, and so on. In this way a molecular dynamics simulation generates a

trajectory that describes how the dynamic variables change with time. Molecular dynamics simulations are typically run for tens or hundreds of picoseconds (a 100 ps simulation using a 1 fs time step requires 100 000 steps). Thermodynamic averages are obtained from molecular dynamics as time averages using numerical integration of Equation (6.2):

$$\langle A \rangle = \frac{1}{M} \sum_{i=1}^M A(\mathbf{p}^N, \mathbf{r}^N) \quad (6.5)$$

M is the number of time steps. Molecular dynamics is also extensively used to investigate the conformational properties of flexible molecules as will be discussed in Chapters 7 and 9.

6.1.3 The Basic Elements of the Monte Carlo Method

In a molecular dynamics simulation the successive configurations of the system are connected in time. This is not the case in a Monte Carlo simulation, where each configuration depends only upon its predecessor and not upon any other of the configurations previously visited. The Monte Carlo method generates configurations randomly and uses a special set of criteria to decide whether or not to accept each new configuration. These criteria ensure that the probability of obtaining a given configuration is equal to its Boltzmann factor, $\exp\{-\mathcal{V}(\mathbf{r}^N)/k_B T\}$, where $\mathcal{V}(\mathbf{r}^N)$ is calculated using the potential energy function. States with a low energy are thus generated with a higher probability than configurations with a higher energy. For each configuration that is accepted the values of the desired properties are calculated and at the end of the calculation the average of these properties is obtained by simply averaging over the number of values calculated, M :

$$\langle A \rangle = \frac{1}{M} \sum_{i=1}^M A(\mathbf{r}^N) \quad (6.6)$$

Most Monte Carlo simulations of molecular systems are more properly referred to as Metropolis Monte Carlo calculations after Metropolis and his colleagues, who reported the first such calculation. The distinction can be important because there are other ways in which an ensemble of configurations can be generated. As we shall see in Chapter 7, the Metropolis scheme is only one of a number of possibilities, though it is by far the most popular.

In a Monte Carlo simulation each new configuration of the system may be generated by randomly moving a single atom or molecule. In some cases new configurations may also be obtained by moving several atoms or molecules or by rotating about one or more bonds. The energy of the new configuration is then calculated using the potential energy function. If the energy of the new configuration is lower than the energy of its predecessor then the new configuration is accepted. If the energy of the new configuration is higher than the energy of its predecessor then the Boltzmann factor of the energy difference is calculated: $\exp[-(\mathcal{V}_{\text{new}}(\mathbf{r}^N) - \mathcal{V}_{\text{old}}(\mathbf{r}^N))/k_B T]$. A random number is then generated between 0 and 1 and compared with this Boltzmann factor. If the random number is higher than the Boltzmann factor then the move is rejected and the original configuration is retained for the next iteration; if the random number is lower then the move is accepted and the new

configuration becomes the next state. This procedure has the effect of permitting moves to states of higher energy. The smaller the uphill move (i.e. the smaller the value of $\mathcal{V}_{\text{new}}(\mathbf{r}^N) - \mathcal{V}_{\text{old}}(\mathbf{r}^N)$) the greater is the probability that the move will be accepted.

6.1.4 Differences Between the Molecular Dynamics and Monte Carlo Methods

The molecular dynamics and Monte Carlo simulation methods differ in a variety of ways. The most obvious difference is that molecular dynamics provides information about the time dependence of the properties of the system whereas there is no temporal relationship between successive Monte Carlo configurations. In a Monte Carlo simulation the outcome of each trial move depends only upon its immediate predecessor, whereas in molecular dynamics it is possible to predict the configuration of the system at any time in the future – or indeed at any time in the past. Molecular dynamics has a kinetic energy contribution to the total energy whereas in a Monte Carlo simulation the total energy is determined directly from the potential energy function. The two simulation methods also sample from different ensembles. Molecular dynamics is traditionally performed under conditions of constant number of particles (N), volume (V) and energy (E) (the microcanonical or constant NVE ensemble) whereas a traditional Monte Carlo simulation samples from the canonical ensemble (constant N , V and temperature, T). Both the molecular dynamics and Monte Carlo techniques can be modified to sample from other ensembles; for example, molecular dynamics can be adapted to simulate from the canonical ensemble. Two other ensembles are common:

- isothermal–isobaric: fixed N , T , P (pressure)
- grand canonical: fixed μ (chemical potential), V , T

In the canonical, microcanonical and isothermal–isobaric ensembles the number of particles is constant but in a grand canonical simulation the composition can change (i.e. the number of particles can increase or decrease). The equilibrium states of each of these ensembles are characterised as follows:

- canonical ensemble: minimum Helmholtz free energy (A)
- microcanonical ensemble: maximum entropy (S)
- isothermal–isobaric ensemble: minimum Gibbs function (G)
- grand canonical ensemble: maximum pressure \times volume (PV)

6.2 Calculation of Simple Thermodynamic Properties

A wide variety of thermodynamic properties can be calculated from computer simulations; a comparison of experimental and calculated values for such properties is an important way in which the accuracy of the simulation and the underlying energy model can be quantified. Simulation methods also enable predictions to be made of the thermodynamic properties of systems for which there is no experimental data, or for which experimental data is difficult or impossible to obtain. Simulations can also provide structural information about the

conformational changes in molecules and the distributions of molecules in a system. The emphasis in our discussion will be on those properties that are routinely calculated in computer simulations and on the way in which they are obtained. It is important to recognise that the results we derive are for the canonical ensemble. Sometimes the equivalent expressions in other ensembles are provided. The result obtained from one ensemble may also be transformed to another ensemble, though this is strictly only possible in the limit of an infinitely large system. The expressions follow from standard statistical mechanical relationships, which are given in standard texts and summarised in Appendix 6.1.

6.2.1 Energy

The internal energy is easily obtained from a simulation as the ensemble average of the energies of the states that are examined during the course of the simulation:

$$U = \langle E \rangle = \frac{1}{M} \sum_{i=1}^M E_i \quad (6.7)$$

6.2.2 Heat Capacity

At a phase transition the heat capacity will often show a characteristic dependence upon the temperature (a first-order phase transition is characterised by an infinite heat capacity at the transition but in a second-order phase transition the heat capacity changes discontinuously). Monitoring the heat capacity as a function of temperature may therefore enable phase transitions to be detected. Calculations of the heat capacity can also be compared with experimental results and so be used to check the energy model or the simulation protocol.

The heat capacity is formally defined as the partial derivative of the internal energy with respect to temperature:

$$C_V = \left(\frac{\partial U}{\partial T} \right)_V \quad (6.8)$$

The heat capacity can therefore be calculated by performing a series of simulations at different temperatures, and then differentiating the energy with respect to the temperature. The differentiation can be done numerically or by fitting a polynomial to the data and then analytically differentiating the fitted function. The heat capacity may also be calculated from a single simulation by considering the instantaneous fluctuations in the energy as follows:

$$C_V = \{ \langle E^2 \rangle - \langle E \rangle^2 \} / k_B T^2 \quad (6.9)$$

An alternative way to write this expression uses the relationship:

$$\langle (E - \langle E \rangle)^2 \rangle = \langle E^2 \rangle - \langle E \rangle^2 \quad (6.10)$$

giving:

$$C_V = \langle (E - \langle E \rangle)^2 \rangle / k_B T^2 \quad (6.11)$$

A derivation of this result is provided in Appendix 6.2.

The heat capacity can therefore be obtained by keeping a running count of E^2 and E during the simulation, from which their expectation values $\langle E^2 \rangle$ and $\langle E \rangle$ can be calculated at the end of the calculation. Alternatively, if the energies are stored during the simulation then the value of $\langle (E - \langle E \rangle)^2 \rangle$ can be calculated once the simulation has finished. This second approach may be more accurate due to round-off errors; $\langle E^2 \rangle$ and $\langle E \rangle^2$ are usually both large numbers and so there may be a large uncertainty in their difference.

6.2.3 Pressure

The pressure is usually calculated in a computer simulation via the virial theorem of Clausius. The *virial* is defined as the expectation value of the sum of the products of the coordinates of the particles and the forces acting on them. This is usually written $W = \sum x_i \dot{p}_x$, where x_i is a coordinate (e.g. the x or y coordinate of an atom) and \dot{p}_x is the first derivative of the momentum along that coordinate (\dot{p}_x is the force, by Newton's second law). The virial theorem states that the virial is equal to $-3Nk_B T$.

In an ideal gas, the only forces are those due to interactions between the gas and the container and it can be shown that the virial in this case equals $-3PV$. This result can also be obtained directly from $PV = Nk_B T$.

Forces between the particles in a real gas or liquid affect the virial, and thence the pressure. The total virial for a real system equals the sum of an ideal gas part ($-3PV$) and a contribution due to interactions between the particles. The result obtained is:

$$W = -3PV + \sum_{i=1}^N \sum_{j=i+1}^N r_{ij} \frac{d_v(r_{ij})}{dr_{ij}} = -3Nk_B T \quad (6.12)$$

The real gas part is derived in Appendix 6.3. If $d_v(r_{ij})/dr_{ij}$ is written as f_{ij} , the force acting between atoms i and j , then we have the following expression for the pressure:

$$P = \frac{1}{V} \left[Nk_B T - \frac{1}{3} \sum_{i=1}^N \sum_{j=i+1}^N r_{ij} f_{ij} \right] \quad (6.13)$$

The forces are calculated as part of a molecular dynamics simulation, and so little additional effort is required to calculate the virial and thus the pressure. The forces are not routinely calculated during a Monte Carlo simulation, and so some additional effort is required to determine the pressure by this route. When calculating the pressure it is also important to check that the components of the pressure in all three directions are equal.

6.2.4 Temperature

In a canonical ensemble the total temperature is constant. In the microcanonical ensemble, however, the temperature will fluctuate. The temperature is directly related to the kinetic energy of the system as follows:

$$\mathcal{K} = \sum_{i=1}^N \frac{|p_i|^2}{2m_i} = \frac{k_B T}{2} (3N - N_c) \quad (6.14)$$

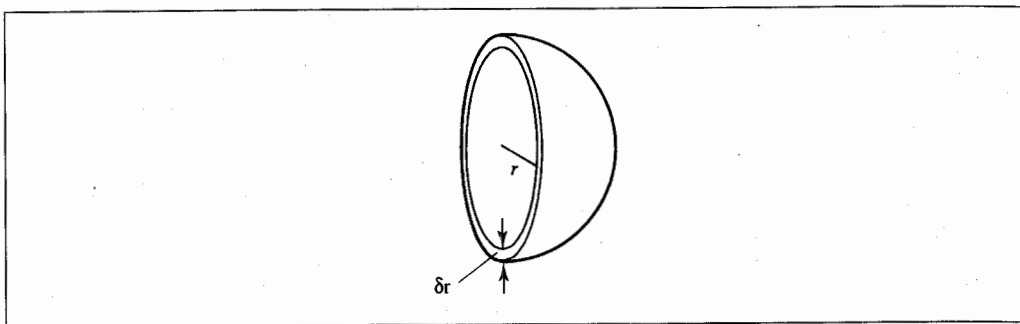


Fig. 6.1: Radial distribution functions use a spherical shell of thickness δr .

In this equation, \mathbf{p}_i is the total momentum of particle i and m_i is its mass. According to the theorem of the equipartition of energy each degree of freedom contributes $k_B T/2$. If there are N particles, each with three degrees of freedom, then the kinetic energy should equal $3Nk_B T/2$. N_c in Equation (6.14) is the number of constraints on the system. In a molecular dynamics simulation the total linear momentum of the system is often constrained to a value of zero, which has the effect of removing three degrees of freedom from the system and so N_c would be equal to 3. Other types of constraint are also possible as we shall discuss in Section 7.5.

6.2.5 Radial Distribution Functions

Radial distribution functions are a useful way to describe the structure of a system, particularly of liquids. Consider a spherical shell of thickness δr at a distance r from a chosen atom (Figure 6.1). The volume of the shell is given by:

$$\begin{aligned} V &= \frac{4}{3}\pi(r + \delta r)^3 - \frac{4}{3}\pi r^3 \\ &= 4\pi r^2 \delta r + 4\pi r \delta r^2 + \frac{4}{3}\pi \delta r^3 \approx 4\pi r^2 \delta r \end{aligned} \quad (6.15)$$

If the number of particles per unit volume is ρ , then the total number in the shell is $4\pi\rho r^2 \delta r$ and so the number of atoms in the volume element varies as r^2 .

The pair distribution function, $g(r)$, gives the probability of finding an atom (or molecule, if simulating a molecular fluid) a distance r from another atom (or molecule) compared to the ideal gas distribution. $g(r)$ is thus dimensionless. Higher radial distribution functions (e.g. the triplet radial distribution function) can also be defined but are rarely calculated and so references to the 'radial distribution function' are usually taken to mean the pairwise version. In a crystal, the radial distribution function has an infinite number of sharp peaks whose separations and heights are characteristic of the lattice structure.

The radial distribution function of a liquid is intermediate between the solid and the gas, with a small number of peaks at short distances, superimposed on a steady decay to a constant value at longer distances. The radial distribution function calculated from a molecular dynamics simulation of liquid argon (shown in Figure 6.2) is typical. For short distances (less

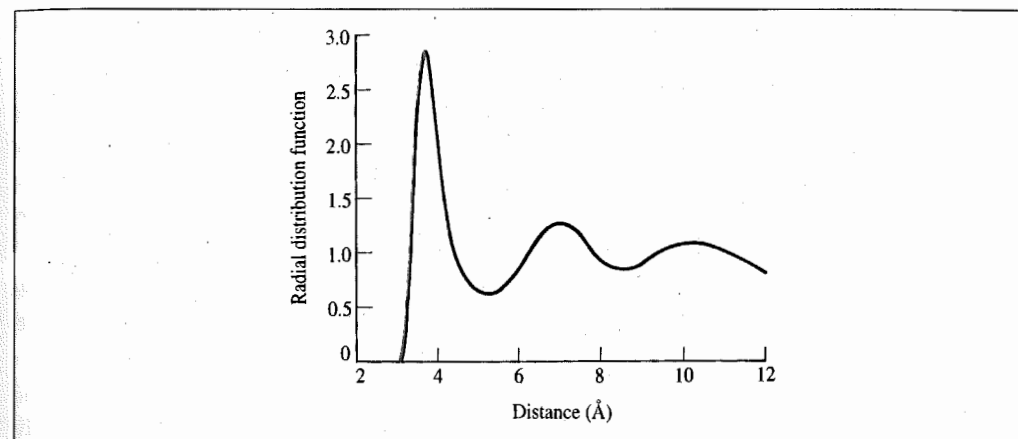


Fig. 6.2: Radial distribution function determined from a 100 ps molecular dynamics simulation of liquid argon at a temperature of 100 K and a density of 1.396 g cm^{-3} .

than the atomic diameter) $g(r)$ is zero. This is due to the strong repulsive forces. The first (and largest) peak occurs at $r \approx 3.7 \text{\AA}$, with $g(r)$ having a value of about 3. This means that it is three times more likely that two molecules would have this separation than in the ideal gas. The radial distribution function then falls and passes through a minimum value around $r \approx 5.4 \text{\AA}$. The chances of finding two atoms with this separation are less than for the ideal gas. At long distances, $g(r)$ tends to the ideal gas value, indicating that there is no long-range order.

To calculate the pair distribution function from a simulation, the neighbours around each atom or molecule are sorted into distance 'bins', or histograms. The number of neighbours in each bin is then averaged over the entire simulation. For example, a count is made of the number of neighbours between (say) 2.5\AA and 2.75\AA , 2.75\AA and 3.0\AA and so on for every atom or molecule in the simulation. This count can be performed during the simulation itself or by analysing the configurations that are generated.

Radial distribution functions can be measured experimentally using X-ray diffraction. The regular arrangement of the atoms in a crystal gives the characteristic X-ray diffraction pattern with bright, sharp spots. For liquids, the diffraction pattern has regions of high and low intensity but no sharp spots. The X-ray diffraction pattern can be analysed to calculate an experimental distribution function, which can then be compared with that obtained from the simulation.

Thermodynamic properties can be calculated using the radial distribution function, if pairwise additivity of the forces is assumed. These properties are usually given as an ideal gas part plus a real gas part. For example, to calculate the energy of a real gas, we consider the spherical shell of volume $4\pi r^2 \delta r$ that contains $4\pi r^2 \rho g(r) \delta r$ particles. If the pair potential at a distance r has a value $v(r)$ then the energy of interaction between the particles in the shell and the central particle is $4\pi r^2 \rho g(r) v(r) \delta r$. The total potential energy of the real gas is obtained by integrating this between 0 and ∞ and multiplying the result

by $N/2$ (the factor $1/2$ ensures that we only count each interaction once). The total energy is then given by:

$$E = \frac{3}{2}Nk_B T + 2\pi N\rho \int_0^\infty r^2 v(r)g(r) dr \quad (6.16)$$

In a similar way the following expression for the pressure can be derived:

$$PV = Nk_B T - \frac{2\pi N\rho}{3k_B T} \int_0^\infty r^2 r \frac{dv(r)}{dr} g(r) dr \quad (6.17)$$

It is usually more accurate to calculate such properties directly, partly because the radial distribution function is not obtained as a continuous function but is derived by dividing the space into small but discrete bins.

For molecules, the orientation must be taken into account if the true nature of the distribution is to be determined. The radial distribution function for molecules is usually measured between two fixed points, such as between the centres of mass. This may then be supplemented by an orientational distribution function. For linear molecules, the orientational distribution function may be calculated as the angle between the axes of the molecule, with values ranging from -180° to $+180^\circ$. For more complex molecules it is usual to calculate a number of site-site distribution functions. For example, for a three-site model of water, three functions can be defined ($g(\text{O}-\text{O})$, $g(\text{O}-\text{H})$ and $g(\text{H}-\text{H})$). An advantage of the site-site models is that they can be directly related to information obtained from the X-ray scattering experiments. The O-O, O-H and H-H radial distribution functions have been particularly useful for refining the various potential models for simulating liquid water.

6.3 Phase Space

An important concept in computer simulation is that of the *phase space*. For a system containing N atoms, $6N$ values are required to define the state of the system (three coordinates per atom and three components of the momentum). Each combination of $3N$ positions and $3N$ momenta (usually denoted by Γ_N) defines a point in the $6N$ -dimensional phase space; an ensemble can thus be considered to be a collection of points in phase space. The way in which the system moves through phase space is governed by Hamiltonian's equations:

$$\frac{d\mathbf{r}_i}{dt} = \frac{\partial \mathcal{H}}{\partial \mathbf{p}_i} \quad (6.18)$$

$$\frac{d\mathbf{p}_i}{dt} = -\frac{\partial \mathcal{H}}{\partial \mathbf{r}_i} \quad (6.19)$$

where i varies from 1 to N . Molecular dynamics generates a sequence of points in phase space that are connected in time. These points correspond to the successive configurations of the system generated by the simulation. A molecular dynamics simulation performed in the microcanonical (constant NVE) ensemble will sample phase space along a contour of constant energy. There is no momentum component in a Monte Carlo simulation and such simulations sample from the $3N$ -dimensional space corresponding to the positions of

the atoms. It might seem odd that thermodynamic properties can be obtained from Monte Carlo simulations, given that there is no momentum contribution and so $3N$ degrees of freedom are not explored. In fact, all of the deviations from ideal gas behaviour are a consequence of interactions between the atoms and are encapsulated in the potential function, $v(\mathbf{r}^N)$, which only depends upon the positions of the atoms. A Monte Carlo simulation does sample from the positional degrees of freedom and so can be used to provide the deviations of thermodynamic properties from ideal gas behaviour, which is what we want to calculate. We shall return to this point in Chapter 8.

If it were possible to visit all the points in phase space then the partition function could be calculated by summing the values of $\exp(-E/k_B T)$. The phase-space trajectory in such a case would be termed *ergodic* and the results would be independent of the initial configuration. For the systems that are typical of those studied using simulation methods the phase space is immense, and an ergodic trajectory is not achievable (indeed, even for relatively small systems with only a few tens of atoms the time that would be required to cycle round all of the points in phase space is longer than the age of the universe). A simulation can thus only ever provide an estimate of the 'true' energies and other thermodynamic properties and so a sequence of simulations using different starting conditions would be expected to give similar, but different, results.

The thermodynamic properties that we have considered so far, such as the internal energy, the pressure and the heat capacity are collectively known as the mechanical properties and can be routinely obtained from a Monte Carlo or molecular dynamics simulation. Other thermodynamic properties are difficult to determine accurately without resorting to special techniques. These are the so-called entropic or thermal properties: the free energy, the chemical potential and the entropy itself. The difference between the mechanical and thermal properties is that the mechanical properties are related to the derivative of the partition function whereas the thermal properties are directly related to the partition function itself. To illustrate the difference between these two classes of properties, let us consider the internal energy, U , and the Helmholtz free energy, A . These are related to the partition function by:

$$U = k_B T^2 \frac{\partial Q}{\partial T} \quad (6.20)$$

$$A = -k_B T \ln Q \quad (6.21)$$

Q is given by Equation (6.4) for a system of identical particles. We shall ignore any normalisation constants in our treatment here to enable us to concentrate on the basics, and so it does not matter whether the system consists of identical or distinguishable particles. We also replace the Hamiltonian by the energy, E . The internal energy is obtained via Equation (6.20):

$$\begin{aligned} U &= k_B T^2 \frac{1}{Q} \iint d\mathbf{p}^N d\mathbf{r}^N \frac{E(\mathbf{p}^N, \mathbf{r}^N)}{k_B T^2} \exp(-E(\mathbf{p}^N, \mathbf{r}^N)/k_B T) \\ &= \iint d\mathbf{p}^N d\mathbf{r}^N E(\mathbf{p}^N, \mathbf{r}^N) \frac{\exp(-E(\mathbf{p}^N, \mathbf{r}^N)/k_B T)}{Q} \end{aligned} \quad (6.22)$$

Now consider the probability of the state with energy $E(\mathbf{p}^N, \mathbf{r}^N)$:

$$\frac{\exp(-E(\mathbf{p}^N, \mathbf{r}^N)/k_B T)}{Q} \quad (6.23)$$

This probability is written $\rho(\mathbf{p}^N, \mathbf{r}^N)$; the internal energy is thus given by

$$U = \iint d\mathbf{p}^N d\mathbf{r}^N E(\mathbf{p}^N, \mathbf{r}^N) \rho(\mathbf{p}^N, \mathbf{r}^N) \quad (6.24)$$

The crucial point about Equation (6.24) is that high values of $E(\mathbf{p}^N, \mathbf{r}^N)$ have a very low probability and make an insignificant contribution to the integral. The Monte Carlo and molecular dynamics methods preferentially generate states of low energy, which are the states that make a significant contribution to the integral in Equation (6.24). These methods sample from phase space in a way that is representative of the equilibrium state and are able to generate accurate estimates of properties such as the internal energy, heat capacity, and so on.

Let us now consider the problem of calculating the Helmholtz free energy of a molecular liquid. Our aim is to express the free energy in the same functional form as the internal energy, that is as an integral which incorporates the probability of a given state. First, we substitute for the partition function in Equation (6.21):

$$A = -k_B T \ln Q = k_B T \ln \left(\frac{N! h^{3N}}{\iint d\mathbf{p}^N d\mathbf{r}^N \exp(-E(\mathbf{p}^N, \mathbf{r}^N)/k_B T)} \right) \quad (6.25)$$

Next we recognise that the following integral is equal to 1:

$$1 = \frac{1}{(8\pi^2 V)^N} \iint d\mathbf{p}^N d\mathbf{r}^N \exp\left(-\frac{E(\mathbf{p}^N, \mathbf{r}^N)}{k_B T}\right) \exp\left(\frac{E(\mathbf{p}^N, \mathbf{r}^N)}{k_B T}\right) \quad (6.26)$$

Inserting this into the expression for the free energy and ignoring the constants (which act to change the zero point from which the free energy is calculated) gives:

$$A = k_B T \ln \left(\frac{\iint d\mathbf{p}^N d\mathbf{r}^N \exp\left(-\frac{E(\mathbf{p}^N, \mathbf{r}^N)}{k_B T}\right) \exp\left(+\frac{E(\mathbf{p}^N, \mathbf{r}^N)}{k_B T}\right)}{\iint d\mathbf{p}^N d\mathbf{r}^N \exp(-E(\mathbf{p}^N, \mathbf{r}^N)/k_B T)} \right) \quad (6.27)$$

We can now substitute for the probability density, $\rho(\mathbf{p}^N, \mathbf{r}^N)$ in this equation, leading to the final result (in which we have again ignored the normalisation factors):

$$A = k_B T \ln \left(\iint d\mathbf{p}^N d\mathbf{r}^N \exp\left(+\frac{E(\mathbf{p}^N, \mathbf{r}^N)}{k_B T}\right) \rho(\mathbf{p}^N, \mathbf{r}^N) \right) \quad (6.28)$$

The important feature of this result is that the configurations with a high energy make a significant contribution to the integral due to the presence of the exponential term $\exp(+E(\mathbf{p}^N, \mathbf{r}^N)/k_B T)$. A Monte Carlo or molecular dynamics simulation preferentially samples the *lower-energy* regions of phase space. An ergodic trajectory would, of course, visit all of these high-energy regions, but in practice these will never be adequately sampled

by a real simulation. The results for the free energy and other entropic properties will as a consequence be poorly converged and inaccurate.

To reiterate a point that we made earlier, these problems of accurately calculating the free energy and entropy do not arise for isolated molecules that have a small number of well-characterised minima which can all be enumerated. The partition function for such systems can be obtained by standard statistical mechanical methods involving a summation over the minimum energy states, taking care to include contributions from internal vibrational motion.

6.4 Practical Aspects of Computer Simulation

6.4.1 Setting Up and Running a Simulation

There are significant differences between the molecular dynamics and Monte Carlo simulation methods, but the same general strategies are used to set up and run either type of simulation. The first task is to decide which energy model is to be used to describe the interactions within the system. Simulations are usually performed with relatively large numbers of atoms over many iterations or time steps. The intra- and intermolecular interactions are therefore almost always described using an empirical (i.e. molecular mechanics) energy model. Faster computers and new theoretical techniques do now enable simulations to be performed using models based only on quantum mechanics or mixed models based on molecular mechanics/quantum mechanics as discussed in Section 11.13. Having chosen an energy model, the simulation itself can be broken into four distinct stages. First, an initial configuration for the system must be established. An *equilibration phase* is then performed, during which the system evolves from the initial configuration. Thermodynamic and structural properties are monitored during the equilibration until stability is achieved. Several distinct steps may be required during the equilibration, particularly for inhomogeneous systems. At the end of the equilibration the *production phase* commences. It is during the production phase that simple properties of the system are calculated. At regular intervals the configuration of the system (i.e. the atomic coordinates) is output to a disk file. Finally, the simulation is analysed; properties not calculated during the simulation are determined and the configurations are examined, not only to discover how the structure of the system changed but also to check for any unusual behaviour that might indicate a problem with the simulation.

6.4.2 Choosing the Initial Configuration

Before a simulation can be performed it is obviously necessary to select an initial configuration of the system. This should be done with some care, as the initial arrangement can often determine the success or failure of a simulation. For simulations of systems at equilibrium (the most common sort) it is wise to choose an initial configuration that is close to the state which it is desired to simulate. For example, it would be unwise to initiate a simulation of a face-centred cubic solid from a body-centred cubic starting point. It is also good practice to ensure that the initial configuration does not contain any high-energy

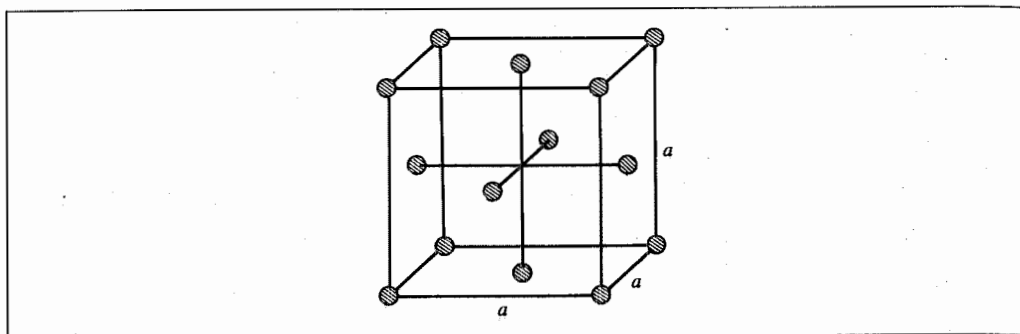


Fig. 6.3: The face-centred cubic cell.

interactions as these may cause instabilities in the simulation. Such 'hot spots' can often be eradicated by performing energy minimisation prior to the simulation itself.

To simulate homogeneous liquids which contain large numbers of the same molecule, a standard lattice structure is often chosen as the starting configuration. If an experimentally determined arrangement is available (e.g. an X-ray structure) then this could be used, provided that it was appropriate to the simulation being performed. When no experimental structure is available the initial configuration can be chosen from one of the common crystallographic lattices (simply placing molecules at random can often give rise to high-energy overlaps and instabilities). The most common lattice is the face-centred cubic lattice (fcc), shown in Figure 6.3. This structure contains $4M^3$ points ($M = 2, 3, 4, \dots$). For this reason, simulations are often performed using 108, 256, 525, 784, \dots , atoms or molecules. The lattice size is chosen so that the density is appropriate to that of the system under study. For simulations of molecules it is also necessary to assign an orientation to each molecule. For small linear molecules, the solid structure of CO_2 is often chosen as the initial configuration. This is a face-centred cubic lattice with the molecules oriented in a regular fashion along the four diagonals of the unit cell. Alternatively, the orientations may be chosen completely at random or by making small random changes from the orientation in a regular lattice. At high densities non-physical overlaps may result, particularly if the molecules are large; in such cases it is more important to use an initial configuration that is close to the expected equilibrium distribution. For example, simulations of rod-shaped molecules such as liquid crystals are usually initiated from a configuration in which the molecules are all aligned approximately in the same direction.

For simulations of inhomogeneous systems comprising a solute molecule or intermolecular complex immersed in a solvent, the starting conformation of the solute may be obtained from an experimental technique such as X-ray crystallography or NMR, or may be generated by theoretical modelling. The coordinates of some solvent molecules may be known if the structure is obtained from X-ray crystallography, but it is usually necessary to add other solvent molecules to give the appropriate solvent density. A typical approach is to use the coordinates obtained from a previous simulation of the pure solvent. The solute is immersed in the solvent 'bath' and any solvent molecules that are too close to the solute are then discarded before the calculation proceeds.

6.5 Boundaries

The correct treatment of boundaries and boundary effects is crucial to simulation methods because it enables 'macroscopic' properties to be calculated from simulations using relatively small numbers of particles. The importance of boundary effects can be illustrated by considering the following simple example. Suppose we have a cube of volume 1 litre which is filled with water at room temperature. The cube contains approximately 3.3×10^{25} molecules. Interactions with the walls can extend up to 10 molecular diameters into the fluid. The diameter of the water molecule is approximately 2.8 \AA and so the number of water molecules that are interacting with the boundary is about 2×10^{19} . So only about one in 1.5 million water molecules is influenced by interactions with the walls of the container. The number of particles in a Monte Carlo or molecular dynamics simulation is far fewer than 10^{25} – 10^{26} and is frequently less than 1000. In a system of 1000 water molecules most, if not all of them, would be within the influence of the walls of the boundary. Clearly, a simulation of 1000 water molecules in a vessel would not be an appropriate way to derive 'bulk' properties. The alternative is to dispense with the container altogether. Now, approximately three-quarters of the molecules would be at the surface of the sample rather than being in the bulk. Such a situation would be relevant to studies of liquid drops, but not to studies of bulk phenomena.

6.5.1 Periodic Boundary Conditions

Periodic boundary conditions enable a simulation to be performed using a relatively small number of particles, in such a way that the particles experience forces as if they were in bulk fluid. Imagine a cubic box of particles which is replicated in all directions to give a periodic array. A two-dimensional box is shown in Figure 6.4. In the two-dimensional example each box is surrounded by eight neighbours; in three dimensions each box would have 26 nearest neighbours. The coordinates of the particles in the image boxes can be computed simply by adding or subtracting integral multiples of the box sides. Should a particle leave the box during the simulation then it is replaced by an image particle that enters from the opposite side, as illustrated in Figure 6.4. The number of particles within the central box thus remains constant.

The cubic cell is the simplest periodic system to visualise and to program. However, a cell of a different shape might be more appropriate for a given simulation. This may be particularly important for simulations of systems which comprise a single molecule or intermolecular complex surrounded by solvent molecules. In such systems it is usually the behaviour of the central solute molecule that is of most interest and so it is desirable that as little of the computer time as possible is spent simulating the solvent far from the solute. In principle, any cell shape can be used provided it fills all of space by translation operations of the central box in three dimensions. Five shapes satisfy this condition: the cube (and its close relation, the parallelepiped), the hexagonal prism, the truncated octahedron, the rhombic dodecahedron and the 'elongated' dodecahedron (Figure 6.5) [Adams 1983]. It is often sensible to choose a periodic cell that reflects the underlying geometry of the system. For example, a rectangular cell is not the ideal choice to simulate an approximately spherical molecule.

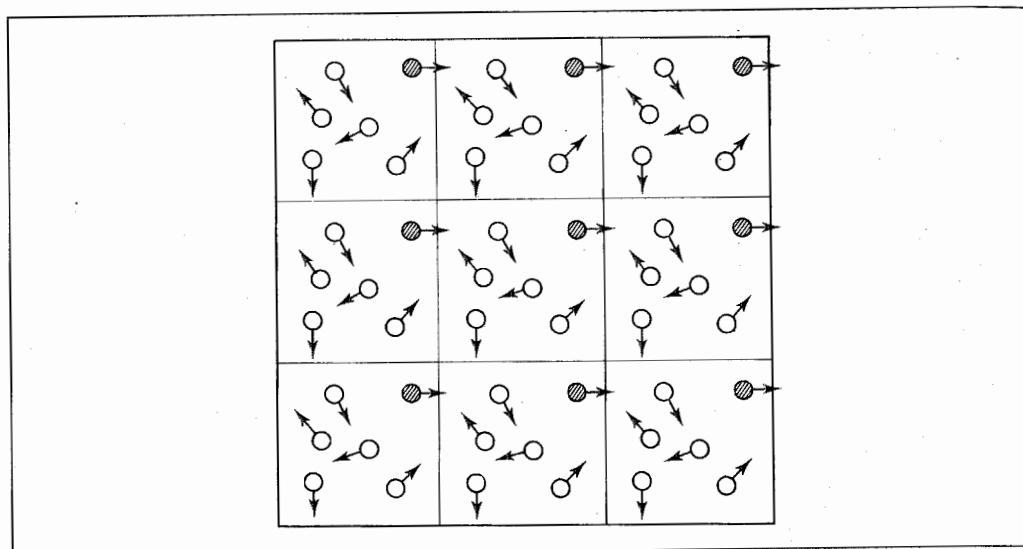


Fig. 6.4: Periodic boundary conditions in two dimensions.

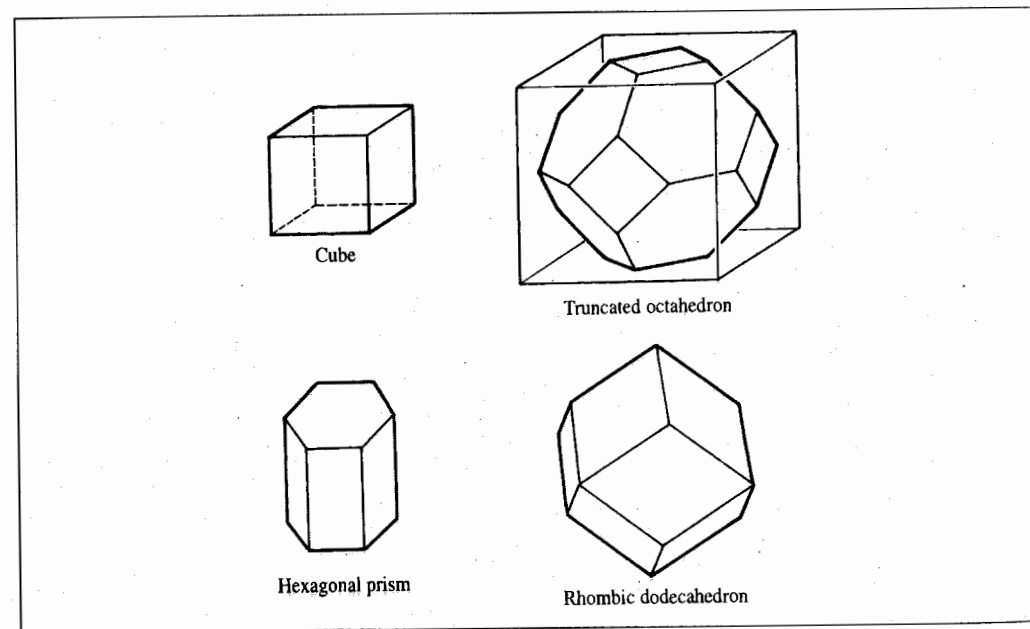


Fig. 6.5: Periodic cells used in computer simulations: the cube, truncated octahedron, hexagonal prism and rhombic dodecahedron.

The truncated octahedron and the rhombic dodecahedron provide periodic cells that are approximately spherical and so may be more appropriate for simulations of spherical molecules. The distance between adjacent cells in the truncated octahedron or the rhombic dodecahedron is larger than the conventional cube for a system with a given number of particles and so a simulation using one of the spherical cells will require fewer particles than a comparable simulation using a cubic cell. Of the two approximately spherical cells, the truncated octahedron is often preferred as it is somewhat easier to program. The hexagonal prism can be used to simulate molecules with a cylindrical shape such as DNA.

Of the five possible shapes, the cube/parallelepiped and the truncated octahedron have been most widely used, with some simulations in the hexagonal prism. The formulae used to translate a particle back into the central simulation box for these three shapes are given in Appendix 6.4. It may be preferable to use one of the more common periodic cells even if there are aesthetic reasons for using an alternative. This is because the expressions for calculating the images may be difficult and inefficient to implement, even though the simulation would use fewer atoms.

For some simulations it is inappropriate to use standard periodic boundary conditions in all directions. For example, when studying the adsorption of molecules onto a surface, it is clearly inappropriate to use the usual periodic boundary conditions for motion perpendicular to the surface. Rather, the surface is modelled as a true boundary, for example by explicitly including the atoms in the surface. The opposite side of the box must still be treated; when a molecule strays out of the top side of the box it is reflected back into the simulation cell, as indicated in Figure 6.6. Usual periodic boundary conditions apply to motion parallel to the surface.

Periodic boundaries are widely used in computer simulations, but they do have some drawbacks. A clear limitation of the periodic cell is that it is not possible to achieve fluctuations that have a wavelength greater than the length of the cell. This can cause problems in certain situations, such as near the liquid-gas critical point. The range of the interactions present in the system is also important; if the cell size is large compared with the range over which the interactions act then there should be no problems. For example, for the relatively short-range Lennard-Jones potential the cell should have a side greater than approximately 6σ , which corresponds to about 20 Å for argon. For longer-range electrostatic interactions the situation is more difficult and it is often necessary to accept that some long-range order will be imposed upon the system. The effects of imposing a periodic boundary can be evaluated empirically by comparing the results of simulations performed using a variety of cell shapes and sizes.

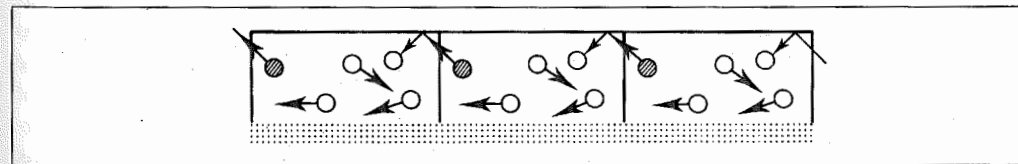


Fig. 6.6: Periodic boundary conditions for surface simulations. (Figure adapted from Allen M P and D J Tildesley 1987. *Computer Simulation of Liquids*. Oxford, Oxford University Press.)

6.5.2 Non-periodic Boundary Methods

Periodic boundary conditions are not always used in computer simulations. Some systems, such as liquid droplets or van der Waals clusters, inherently contain a boundary. Periodic boundary conditions may also cause difficulties when simulating inhomogeneous systems or systems that are not at equilibrium. In other cases the use of periodic boundary conditions would require a prohibitive number of atoms to be included in the simulation. This particularly arises in the study of the structural and conformational behaviour of macromolecules such as proteins and protein-ligand complexes. The first simulations of such systems ignored all solvent molecules due to the limited computational resources then available. This corresponds to the unrealistic situation of simulating an isolated protein *in vacuo* and then comparing the results with experimental data obtained in solution. Vacuum calculations can lead to significant problems. A vacuum boundary tends to minimise the surface area and so may distort the shape of the system if it is non-spherical. Small molecules may adopt more compact conformations when simulated *in vacuo* due to favourable intramolecular electrostatic and van der Waals interactions, which would be dampened in the presence of a solvent.

As computer power has increased it has become possible to incorporate explicitly some solvent molecules and thereby simulate a more realistic system. The simplest way to do this is to surround the molecule with a 'skin' of solvent molecules. If the skin is sufficiently deep then the system is equivalent to a solute molecule inside a 'drop' of solvent. The number of solvent molecules in such cases is usually significantly fewer than would be required in the analogous periodic boundary simulation, where the solute molecule is positioned at the centre of the cell and the empty space is filled with solvent. Boundary effects should be transferred from the molecule-vacuum interface to the solvent-vacuum interface and so might be expected to result in a more realistic treatment of the solute. To illustrate these three situations, we can consider dihydrofolate reductase, which is a small enzyme that contains approximately 2500 atoms. If this enzyme is surrounded by water molecules in a cubic periodic system such that the surface of the protein is at least 10 Å from any side of the box, then the number of atoms rises to almost 20000. If a shell 10 Å thick is used then the number of atoms falls to 14700, and with a 5 Å shell the system contains 8900 atoms.

Sometimes we are only interested in a specific part of the solute, such as the active site of an enzyme. It has been common practice in such cases to divide the system into two regions (Figure 6.7). One region, often called the *reaction zone*, contains all atoms or groups within a given radius R_1 of the site of interest. The atoms in the reaction zone are subjected to the full simulation method. The second region (the *reservoir region*) contains all atoms outside the reaction zone but within a distance R_2 of the active site. The atoms in the reservoir region may be kept fixed in their initial positions, or may be restrained so that they stay within the shell defined by R_1 and R_2 . Alternatively, they may be restrained to their initial positions using a harmonic potential. Any atoms further away from the active site than R_2 are discarded or may be kept fixed in their initial positions. It is important to be aware that restraining or fixing atoms in this way may prevent natural changes occurring and so lead to artificial behaviour. A variety of schemes for performing simulations using such *stochastic boundary conditions* have been proposed. However, such methods can be rather complicated to implement and if not used properly can give anomalous results.

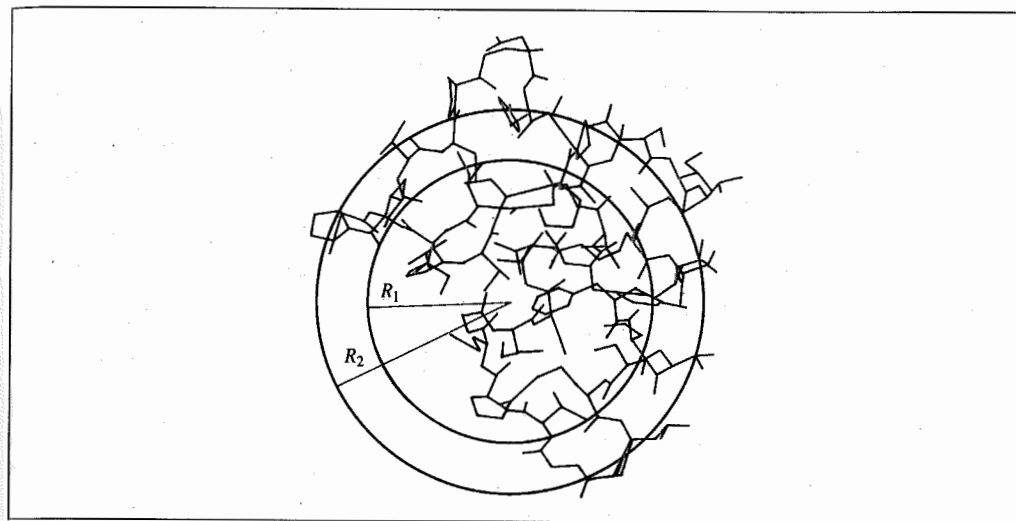


Fig. 6.7: Division into reaction zone and reservoir regions in a simulation using stochastic boundary conditions.

If at all possible, a periodic boundary is the 'safest' way to ensure that boundary effects are minimised, but sometimes an alternative may be the only practical course.

6.6 Monitoring the Equilibration

The purpose of the equilibration phase is to enable the system to evolve from the starting configuration to reach equilibrium. Equilibration should continue until the values of a set of monitored properties become stable. The properties to be monitored usually include thermodynamic quantities such as the energy, temperature and pressure and also structural properties. Many simulations of the liquid state involve a starting configuration that corresponds to a solid lattice. It is therefore important to establish that the lattice has 'melted' before the production phase begins. *Order parameters* can be used to determine that the liquid state has been reached. An order parameter is a measure of the degree of order (or, equivalently, disorder) in the system. During a simulation of a crystal lattice the atoms would be expected to remain in approximately the same positions throughout and thereby maintain a high degree of order. In a liquid, however, we would expect considerable mobility of the species present, giving rise to translational disorder. One way to measure translational order in a system initially in a face-centred cubic lattice was suggested by Verlet, whose order parameter λ is:

$$\lambda = \frac{1}{3}[\lambda_x + \lambda_y + \lambda_z] \quad (6.29)$$

$$\lambda_x = \frac{1}{N} \sum_{i=1}^N \cos\left(\frac{4\pi x_i}{a}\right) \quad (6.30)$$

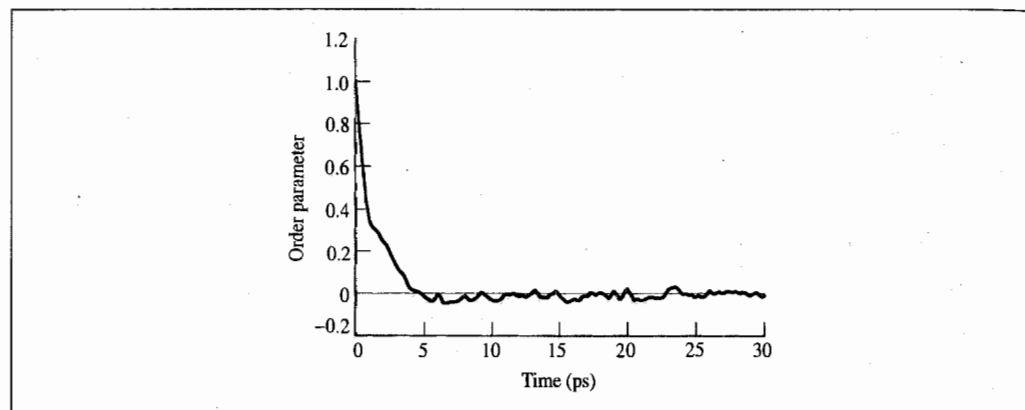


Fig. 6.8: Variation in Verlet order parameter during the equilibration phase of a molecular dynamics simulation of argon.

where a is the length of one edge of the unit cell. Initially, all of the coordinates x_i , y_i and z_i are multiples of $a/2$ and so the order parameter has a value of 1. As the simulation proceeds the order parameter should gradually decrease to a value of zero, indicating that the atoms are distributed randomly. When equilibrium has been reached the fluctuations in the order parameter should be proportional to $1/\sqrt{N}$, where N is the size of the system. A typical result is shown in Figure 6.8 for an argon simulation.

For molecules, it is also necessary to consider their orientations, which can be monitored using a rotational order parameter. For some systems, such as carbon monoxide or water, complete disorder would be expected in the liquid state at equilibrium. However, if we were simulating a dense fluid of rod-shaped molecules which form a liquid crystalline phase then we might expect that, on average, the molecules would tend to line up in a common direction. The Viellard-Baron rotational order parameter for linear molecules is calculated using the following formula:

$$P_1 = \frac{1}{N} \sum_{i=1}^N \cos \gamma_i \quad (6.31)$$

where γ_i is the angle between the current and original direction of the molecular axis of molecule i . A value of 1 indicates that the molecules are perfectly aligned. Rotational disorder is indicated by a value of zero. The fluctuations about the average value should again be proportional to $1/\sqrt{N}$. For non-linear molecules, a number of rotational order parameters can be defined and each monitored.

The mean squared displacement also provides a means to establish whether a solid lattice has melted. The mean squared displacement is given by:

$$\Delta r^2(t) = \frac{1}{N} \sum_{i=1}^N [\mathbf{r}_i(t) - \mathbf{r}_i(0)]^2 \quad (6.32)$$

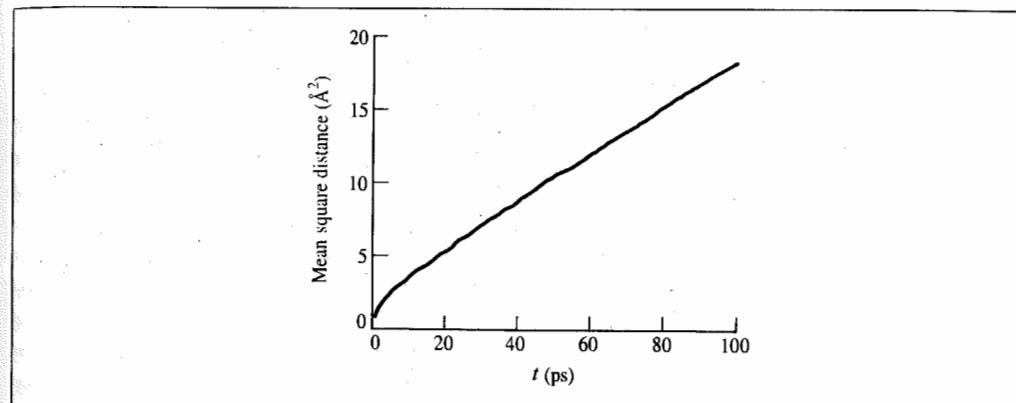


Fig. 6.9: Variation in mean squared displacement during the initial steps of a molecular dynamics simulation of argon.

For a fluid, with no underlying regular structure, the mean squared displacement gradually increases with time (Figure 6.9). For a solid, however, the mean squared displacement typically oscillates about a mean value. However, if there is diffusion within a solid then this can be detected from the mean squared displacement and may be restricted to fewer than three dimensions. For example, Figure 6.10 shows the mean squared displacement calculated for Li^+ ions in Li_3N at 400 K [Wolf *et al.* 1984]. This material contains layers of Li_2N ; mobility of the Li^+ ions is much greater within these planes than perpendicular to them.

The radial distribution function can also be used to monitor the progress of the equilibration. This function is particularly useful for detecting the presence of two phases. Such a situation is characterised by a larger than expected first peak and by the fact that $g(r)$ does not decay towards a value of 1 at long distances. If two-phase behaviour is inappropriate then the simulation should probably be terminated and examined. If, however, a two-phase system is desired, then a long equilibration phase is usually required.

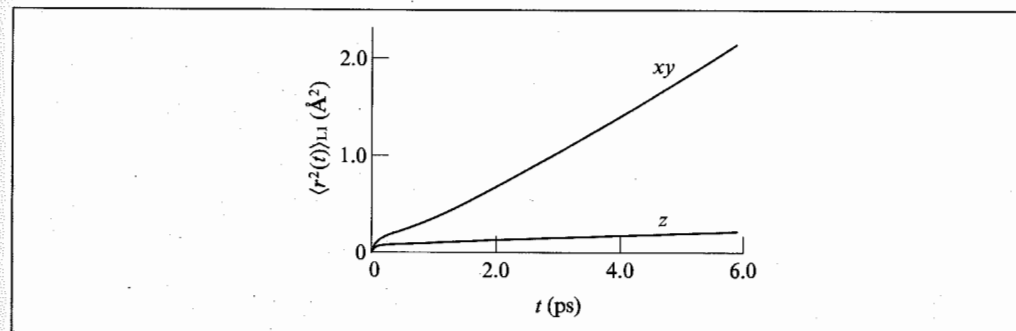


Fig. 6.10: Mean squared displacement for Li^+ ions in Li_3N for motion parallel (xy) and perpendicular (z) to the Li_2N layers [Wolf *et al.* 1984].

6.7 Truncating the Potential and the Minimum Image Convention

The most time-consuming part of a Monte Carlo or molecular dynamics simulation (or, indeed, of an energy minimisation) is the calculation of the non-bonded energies and/or forces. The numbers of bond-stretching, angle-bending and torsional terms in a force field model are all proportional to the number of atoms but the number of non-bonded terms that need to be evaluated increases as the square of the number of atoms (for a pairwise model) and is thus of order N^2 . In principle, the non-bonded interactions are calculated between every pair of atoms in the system. However, for many interaction models this is not justified. The Lennard-Jones potential falls off very rapidly with distance: at 2.5σ the Lennard-Jones potential has just 1% of its value at σ . This reflects the r^{-6} distance dependence of the dispersion interaction. The most popular way to deal with the non-bonded interactions is to use a *non-bonded cutoff* and to apply the *minimum image convention*. In the minimum image convention, each atom 'sees' at most just one image of every other atom in the system (which is repeated infinitely via the periodic boundary method). The energy and/or force is calculated with the closest atom or image, as illustrated in Figure 6.11. When a cutoff is employed, the interactions between all pairs of atoms that are further apart than the cutoff value are set to zero, taking into account the closest image. When periodic boundary conditions are being used, the cutoff should not be so large that a particle sees its own image, or indeed the same molecule twice. This has the effect of limiting the cutoff to no more than half the length of the cell when simulating atomic fluids in a cubic cell. For rectangular cells the cutoff should be no greater than half the length of the shortest side. For simulations of molecules the upper limit on the cutoff may also be affected by the size of the molecules, as we shall see below in Section 6.7.2. In simulations where the Lennard-Jones potential is the only non-bonded interaction, a cutoff of 2.5σ gives rise to a relatively small error. However, when long-range electrostatic

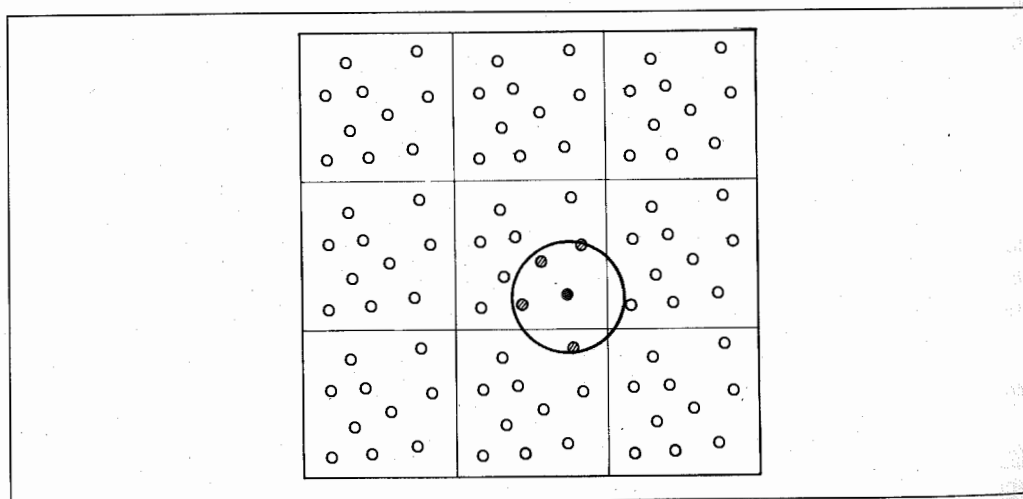


Fig. 6.11: The spherical cutoff and the minimum image convention.

interactions are involved, the cutoff should be much greater and indeed there is evidence to suggest that using any cutoff leads to errors. A value of at least 10 \AA is generally recommended but even this may be insufficient. More comprehensive methods have been devised for dealing with the electrostatic interactions, which are considered in Section 6.8.

6.7.1 Non-bonded Neighbour Lists

By itself, the use of a cutoff may not dramatically reduce the time taken to compute the number of non-bonded interactions. This is because we would still have to calculate the distance between every pair of atoms in the system simply to decide whether they are close enough to calculate their interaction energy. Calculating all the $N(N-1)$ distances takes almost as much time as calculating the energy itself.

In simulations of fluids, an atom's neighbours (i.e. those atoms that are within the cutoff distance) do not change significantly over 10 or 20 molecular dynamics time steps or Monte Carlo iterations. If we 'knew' which atoms to include in the non-bonded calculation (for example, by storing them in an array), then it would be possible to identify directly each atom's neighbours without having to calculate the distances to all the other atoms in the system. The *non-bonded neighbour list* (first proposed by Verlet) is just such a device. The Verlet neighbour list [Verlet 1967] stores all atoms within the cutoff distance, together with all atoms that are slightly further away than the cutoff distance. This is most efficiently done using a large neighbour list array, L , and a pointer array, P . The pointer array indicates where in the neighbour list the first neighbour for that atom is located. The last neighbour of atom i is stored in element $P[i+1] - 1$ of the neighbour list as shown in Figure 6.12. Thus the neighbours of atom i are stored in elements $L[P[i]]$ through $L[P[i+1] - 1]$ of the array L . The neighbour list is updated at regular intervals throughout the simulation. Between updates, the neighbour and pointer lists are used to directly identify the nearest neighbours of each atom i . The distance used to calculate each atom's neighbours should be larger than the actual non-bonded cutoff distance so that no atom, initially outside the neighbour cutoff, approaches closer than the non-bonded cutoff distance before the neighbour list is updated again.

It is important to update the neighbour list at the correct frequency. If the update frequency is too high the procedure is inefficient; too low and the energies and forces may be calculated incorrectly due to atoms moving within the non-bonded cutoff region. An update frequency between 10 and 20 steps is common. An algorithm that can automatically update the neighbour list and so circumvent these problems is as follows [Thompson 1983]. An array element is set to zero for each atom whenever the neighbour list is updated. This array is used to store the displacement of each atom or molecule in subsequent steps. When the sum of the maximum displacements of any two atoms exceeds the difference between the non-bonded cutoff distance and the neighbour list distance, then it is time to update the neighbour list again.

There are no fixed rules that determine how much larger the neighbour list cutoff should be than the non-bonded cutoff. Clearly there will be a trade-off between the size of the cutoff and the frequency at which the neighbour list must be updated: the larger the difference, the less frequently will the neighbour list have to be updated. There may also be storage implications if the list is too large.

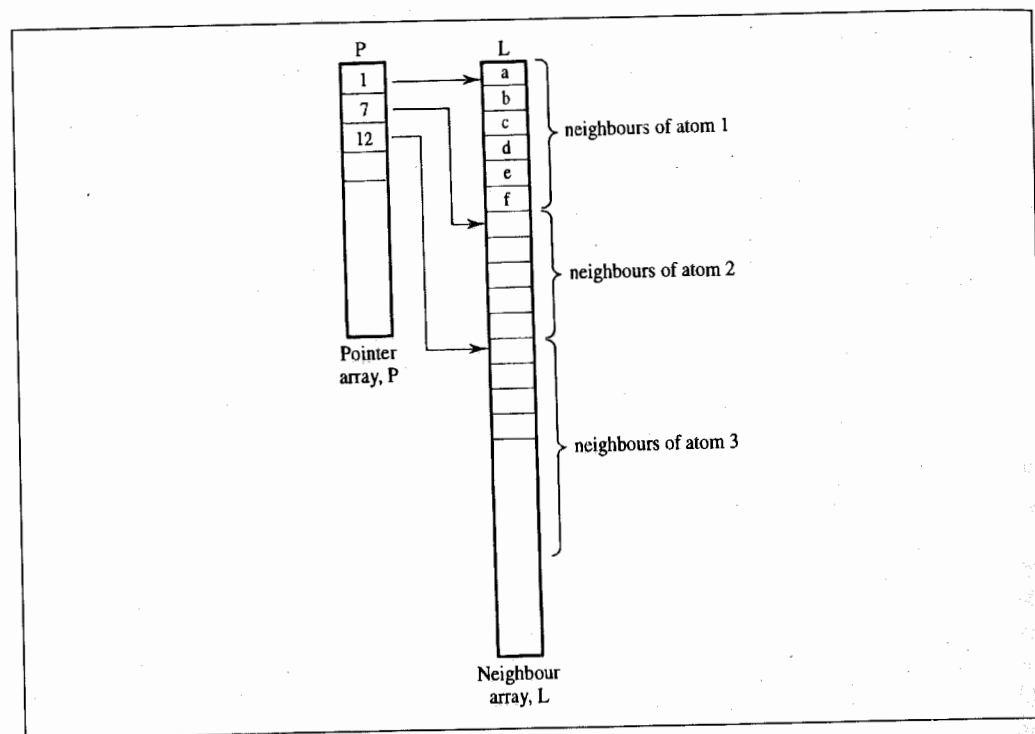


Fig. 6.12: Pointer and neighbour arrays can be used to implement the Verlet neighbour list.

When the number of molecules in the simulation is very large, it can require a significant computational effort just to update the neighbour list. This is because the standard way to update the neighbour list requires the distance between all pairs of atoms in the system to be calculated. When the size of the system is much larger than the cutoff distance, a *cell index method* can be used to make the updating procedure more efficient. In the cell index method, the simulation box is divided into a number of cells. The length of each cell is longer than the non-bonded cutoff distance. All of the neighbours of an atom will then be found either in the cell containing the atom or in one of the surrounding cells. If the entire system is divided into M^3 cells, there will be an average of N/M^3 molecules in each cell. To determine the neighbours of a given atom or molecule, it is then necessary to consider just $27N/M^3$ atoms rather than N . The cell index method requires a mechanism for identifying the atoms or molecules in each cell. Two arrays can be used to do this: a linked list array L and a pointer array P . The pointer array indicates the location of one of the atoms or molecules in a given cell. Thus, $P[1]$ would indicate the number of the 'first' atom or molecule in cell 1 and $P[2]$ is the number of the 'first' atom in cell 2. Each element of the linked list array then gives the number of the 'next' atom or molecule in the cell. Thus the value stored in $L[1]$ is the number of the second atom in the first cell. Suppose $P[1]$ is atom 10. Then the value stored in $L[10]$ is the second atom in the cell. If this second atom is number 15, then $L[15]$ contains the third atom in the cell. The last molecule in the sequence is identified by the fact that its array element is zero. The cell index method clearly

requires a mechanism for updating the pointer and linked list arrays when atoms or molecules move from one cell to another, which can add to the complexity.

When simulating species with a significant electrostatic contribution, it may be desirable to use different cutoffs for the electrostatic and van der Waals interactions. This is because the electrostatic interaction has a much longer range. Using a longer cutoff for the electrostatic interactions will, of course, significantly increase the number of pairs that must be calculated. A compromise is to use a *twin-range method*, in which two cutoffs are specified. All interactions below the lower cutoff are calculated as normal at each step. Interactions due to atoms between the lower and upper cutoffs are evaluated only when the neighbour list is updated and are kept constant between these updates. The rationale here is that the contribution of the atoms that are further away will not vary much between updates.

The use of a cutoff is amply justified in many cases, if only on the grounds of expediency, but its use will always lead to some fraction of the potential energy being ignored. This lost energy can be easily captured at the end of the simulation if it is assumed that the radial distribution function takes the value of 1 at distances greater than the cutoff. The calculation is analogous to that used to determine the total energy from the radial distribution function, Equation (6.16), but the integration is now performed between the cutoff distance r_c and infinity and $g(r)$ is now taken to be 1 in this range. For N particles, the correction is:

$$E_{\text{correction}} = 2\pi\rho N \int_{r_c}^{\infty} r^2 v(r) dr \quad (6.33)$$

For the Lennard-Jones potential the long-range contribution can be determined analytically:

$$E_{\text{correction}} = 8\pi\rho N \epsilon \left[\frac{\sigma^{12}}{9r^9} - \frac{\sigma^6}{3r^3} \right] \quad (6.34)$$

6.7.2 Group-based Cutoffs

When simulating large molecular systems, it is often advantageous to use a group-based cutoff (sometimes called a residue-based cutoff). Here, the large molecules are divided into 'groups', each of which contains a relatively small number of connected atoms. If the calculation involves small solvent molecules then each solvent molecule is also conveniently regarded as a single unconnected group. Why are groups useful? Let us consider the electrostatic interaction between two water molecules. In the popular TIP3P model there is a charge of $-0.834e$ on the oxygen and $0.417e$ on each hydrogen. The electrostatic interaction between two water molecules is calculated as the sum of nine distinct site-site interactions. If we start from the minimum energy arrangement for the water dimer shown in Figure 6.13 and

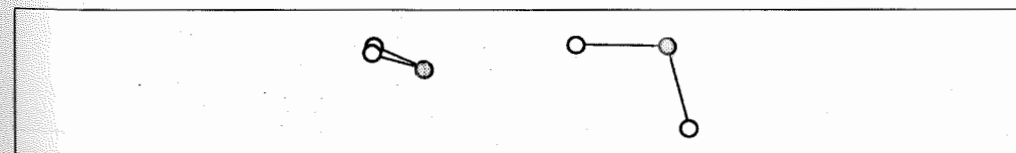


Fig. 6.13: Minimum energy structure for water dimer with TIP3P model.

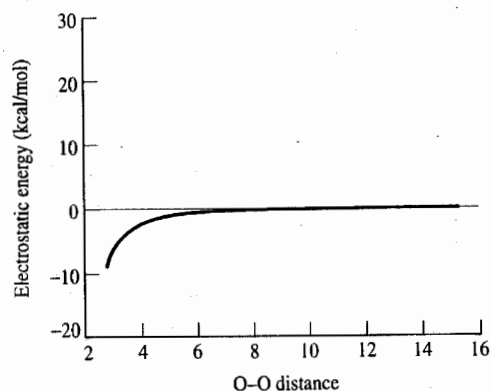


Fig. 6.14: The variation in the electrostatic interaction energy of the water dimer as a function of the O-O distance without a cutoff.

gradually move one water molecule relative to the other as indicated then the electrostatic energy varies as shown in the graph in Figure 6.14.

Although the overall interaction energy is relatively small beyond 6 Å or so, each of these energies is the sum of several rather large terms; for example, at an O-O separation of 8 Å, the overall interaction energy is about -0.27 kcal/mol but this comprises an oxygen-oxygen interaction of approximately 29 kcal/mol, oxygen-hydrogen interactions of -59.4 kcal/mol and hydrogen-hydrogen interactions of 29.2 kcal/mol. Suppose that a simple atom-based non-bonded cutoff is applied to the water dimer. The interaction energy then fluctuates violently near the cutoff distance, as shown in Figure 6.15 for a cutoff of 8 Å. This is because only some of the pairwise interactions are included in this case. Clearly such a model would almost certainly lead to serious problems for any

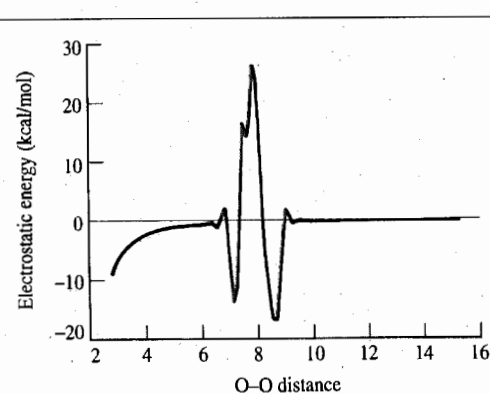


Fig. 6.15: The variation in interaction of the water dimer as a function of the O-O distance with an 8 Å atom-based cutoff.

simulation. This problem can be avoided by collecting all of the atoms from each water molecule into a single group, and by calculating the interactions on a group-group basis, even though some of the atom pairs may have a separation larger than the cutoff.

How should a molecule be divided into groups? In some cases there may appear to be a chemically obvious way to define the groups, especially when the molecule is a polymer that is constructed from distinct chemical residues. A particularly desirable feature is that each group should, if possible, be of zero charge. The reason for this can be understood if we recall how the different electrostatic interactions vary with distance from Section 4.9.1:

$$\begin{aligned} \text{charge-charge} &\sim 1/r \\ \text{charge-dipole} &\sim 1/r^2 \\ \text{dipole-dipole} &\sim 1/r^3 \\ \text{dipole-quadrupole} &\sim 1/r^4 \\ \text{charge-induced dipole} &\sim 1/r^4 \\ \text{dipole-induced dipole} &\sim 1/r^6 \end{aligned}$$

If the groups are electrically neutral, then the leading term in the electrostatic interaction between a pair of groups is the dipole-dipole interaction, which is dependent upon $1/r^3$. By comparison, the charge-charge terms vary as $1/r$. Of course, it is not always possible to arrange atoms in neutral groups as occurs when some of the species are charged.

A further question with the group-based scheme is: how do we decide whether a particular group-group interaction needs to be considered? In other words, how are cutoffs included in the group scheme? One strategy is to include a particular group-group interaction if any pair of atoms in the two groups is closer than the cutoff distance. Alternatively, a 'marker' atom may be nominated within the group; when the marker atoms come closer than the cutoff then the appropriate group-group interaction is included. When using marker atoms, it is important that the groups are not too large; thus the groups used by some simulation programs are much smaller than the 'chemically obvious' groupings. For example, the most obvious choice for proteins and peptides is to define each entire amino acid residue as a single group. However, this is not necessarily the most appropriate strategy. Consider the situation in which two arginine residues are spatially close together (Figure 6.16). Arginine has a long side chain that is comparable in length to the non-bonded cutoff distances often employed. Suppose the alpha-carbon atom (marked C_α in Figure 6.16)

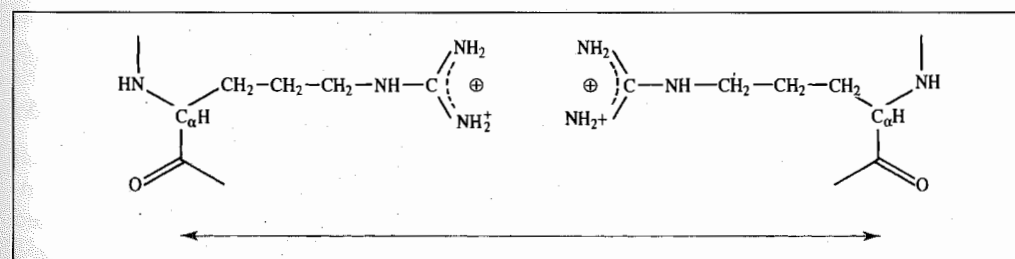


Fig. 6.16: The use of a marker atom on the alpha-carbon in an arginine residue may lead to a significant electrostatic interaction being neglected because the distance between the marker atoms exceeds the cutoff.

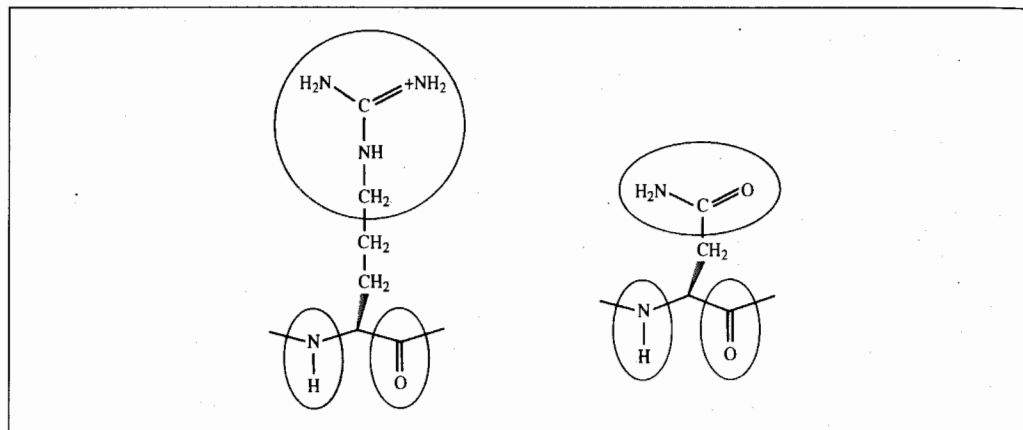


Fig. 6.17: The charge groups used in the GROMOS simulation program for simulating proteins [van Gunsteren and Berendsen 1986], illustrated using the amino acids arginine and asparagine. The CH_2 groups have zero charge.

in each arginine residue is chosen as the marker atom. If the distance between the alpha-carbons of two arginine residues is greater than the cutoff, no interactions between any atoms in the arginine residues would be calculated, despite the fact that the positively charged ends of the residues could be very close, as shown in Figure 6.16. Were the alpha-carbons to approach closer than the cutoff, there would then be a dramatic increase in the energy due to the unfavourable interaction between the two side chains, inevitably leading to an unstable simulation. It may therefore be appropriate to define 'charge groups' that contain smaller numbers of atoms than are in the chemically obvious scheme. For example, the groups that are used by the GROMOS simulation program for the amino acids arginine and asparagine are shown in Figure 6.17.

6.7.3 Problems with Cutoffs and How to Avoid Them

A cutoff introduces a discontinuity in both the potential energy and the force near the cutoff value. This creates problems, especially in molecular dynamics simulations where energy conservation is required. There are several ways that the effects of this discontinuity can be counteracted. One approach is to use a shifted potential, in which a constant term is subtracted from the potential at all values (Figure 6.18):

$$v'(r) = v(r) - v_c \quad r \leq r_c \quad (6.35)$$

$$v'(r) = 0 \quad r > r_c \quad (6.36)$$

where r_c is the cutoff distance and v_c is equal to the value of the potential at the cutoff distance. As the additional term is constant, it disappears when the potential is differentiated and so does not affect the force calculation in molecular dynamics. Use of the shifted potential does improve energy conservation, though as the number of atom pairs separated by a distance smaller than the cutoff varies, so the contribution of the shifted potential to the total energy will change. An additional problem is that there is a discontinuity in the force

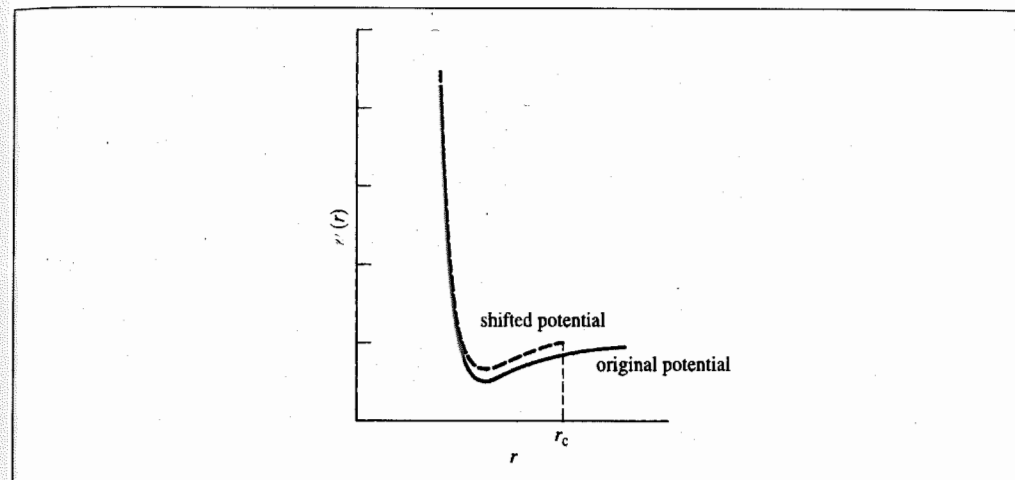


Fig. 6.18: A shifted Lennard-Jones potential.

with the shifted potential; at the cutoff distance, the force will have a finite value which drops suddenly to zero just beyond the cutoff. This can also give instabilities in a simulation. To avoid this, a linear term can be added to the potential, making the derivative zero at the cutoff:

$$v'(r) = v(r) - v_c - \left(\frac{dv(r)}{dr} \right)_{r=r_c} (r - r_c) \quad r \leq r_c \quad (6.37)$$

$$v'(r) = 0 \quad r > r_c \quad (6.38)$$

The shift makes the potential deviate from the 'true' potential, and so any calculated thermodynamic properties will be changed. The 'true' values can be retrieved but it is difficult to do so, and the shifted potential is thus rarely used in 'real' simulations. Moreover, while it is relatively straightforward to implement for a homogeneous system under the influence of a simple potential such as the Lennard-Jones potential, it is not easy for inhomogeneous systems containing many different types of atom.

An alternative way to eliminate discontinuities in the energy and force equations is to use a *switching function*. A switching function is a polynomial function of the distance by which the potential energy function is multiplied. Thus the switched potential $v'(r)$ is related to the true potential $v(r)$ by $v'(r) = v(r)S(r)$. Some switching functions are applied to the entire range of the potential up to the cutoff point. One such function is:

$$v'(r) = v(r) \left[1 - 2 \left(\frac{r}{r_c} \right)^2 + \left(\frac{r}{r_c} \right)^4 \right] \quad (6.39)$$

The switching function has a value of 1 at $r = 0$ and a value of 0 at $r = r_c$, the cutoff distance. Between these two values it varies as shown in Figure 6.19, which also shows how the potential function is affected.

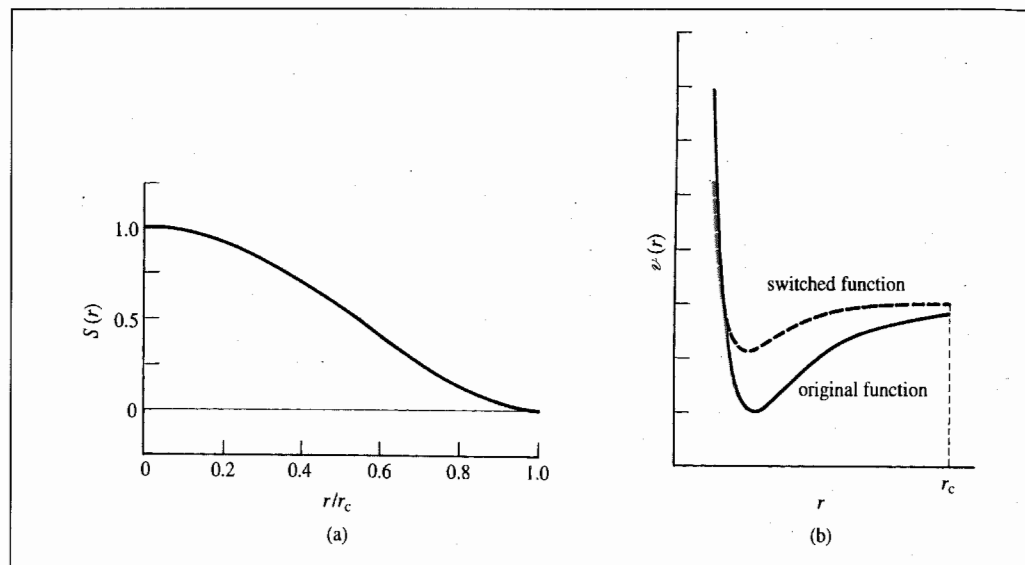


Fig. 6.19: (a) The effect of a switching function that applies over the entire range and (b) its effect on the Lennard-Jones potential.

A switching function applied to the potential function over the entire range does have drawbacks; for example, equilibrium structures are affected (the minimum energy separation for the argon dimer decreases slightly). A more acceptable alternative is to gradually taper the potential between two cutoff values. The potential takes its usual value until the lower cutoff distance. Between the lower (r_1) and upper cutoff distance (r_u) the potential is multiplied by the switching function, which takes the value 1 at the lower cutoff distance and 0 at the upper cutoff distance. The lower cutoff distance is typically relative close to the upper cutoff distance (for example, r_1 might be 9 Å and r_u 10 Å). A simple switching function has the following linear form:

$$S = 1.0 \quad r_{ij} < r_1 \quad (6.40)$$

$$S = (r_u - r_{ij}) / (r_u - r_1) \quad r_1 \leq r_{ij} \leq r_u \quad (6.41)$$

$$S = 0.0 \quad r_u < r_{ij} \quad (6.42)$$

This suffers from discontinuities in both the energy and the force at the two cutoff values. A more acceptable switching function smoothly changes from a value of 1 to a value of 0 (Figure 6.20) between r_1 and r_u and satisfies the following requirements:

$$S_{r=r_1} = 1; \quad \left(\frac{dS}{dr}\right)_{r=r_1} = 0; \quad \left(\frac{d^2S}{dr^2}\right)_{r=r_1} = 0 \quad (6.43)$$

$$S_{r=r_u} = 0; \quad \left(\frac{dS}{dr}\right)_{r=r_u} = 0; \quad \left(\frac{d^2S}{dr^2}\right)_{r=r_u} = 0 \quad (6.44)$$

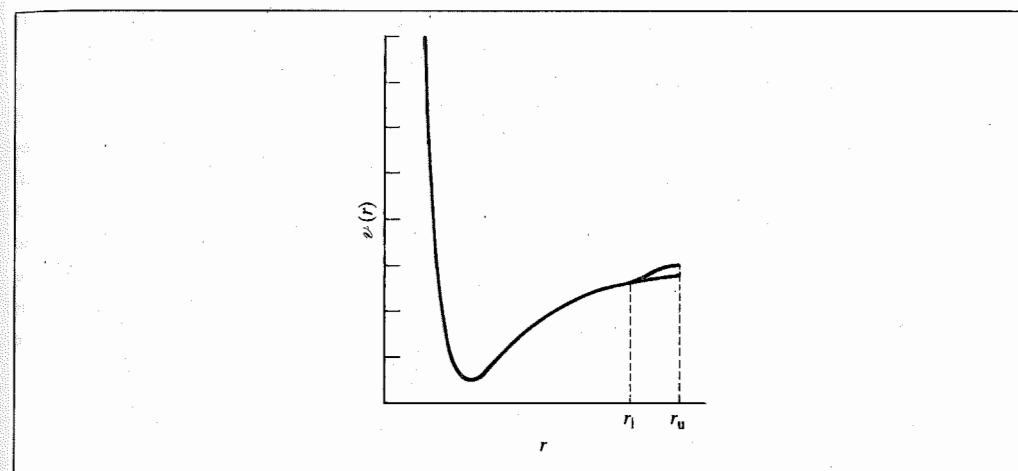


Fig. 6.20: A switching function that applies over a narrow range near the cutoff and its effect on the Lennard-Jones potential.

By ensuring that the first derivative is zero at the endpoints the force also approaches zero smoothly. A continuous second derivative is required to ensure that the integration algorithm works properly. If the switch function is assumed to take the following form:

$$S(r) = c_0 + c_1 \left[\frac{r-r_1}{r_u-r_1}\right] + c_2 \left[\frac{r-r_1}{r_u-r_1}\right]^2 + c_3 \left[\frac{r-r_1}{r_u-r_1}\right]^3 + c_4 \left[\frac{r-r_1}{r_u-r_1}\right]^4 + c_5 \left[\frac{r-r_1}{r_u-r_1}\right]^5 \quad (6.45)$$

then the following values of the coefficients $c_0 \dots c_5$ satisfy the six requirements in equations (6.43) and (6.44):

$$c_0 = 1; \quad c_1 = 0; \quad c_2 = 0; \quad c_3 = -10; \quad c_4 = 15; \quad c_5 = -6 \quad (6.46)$$

When using a switching function in a molecular simulation with a residue-based cutoff it is important that the function has the same value for all pairs of atoms in the two interacting groups. Otherwise, severe fluctuations in the energy can arise when the separation is within the cutoff region. These two contrasting situations can be formally expressed as follows:

$$\text{atom based: } \mathcal{V}_{AB} = \sum_{i=1}^{N_A} \sum_{j=1}^{N_B} S_{ij}(r_{ij}) v_{ij}(r_{ij}) \quad (6.47)$$

$$\text{residue or molecule based: } \mathcal{V}_{AB} = S_{AB}(|\mathbf{r}_A - \mathbf{r}_B|) \sum_{i=1}^{N_A} \sum_{j=1}^{N_B} v_{ij}(r_{ij}) \quad (6.48)$$

N_A and N_B are the numbers of atoms in the two groups A and B and S is the switching function. With the group-based switching function, it is necessary to define the 'distance' between the two groups (i.e. the two points \mathbf{r}_A and \mathbf{r}_B). There is no definitive way to do this. As with cutoffs, a special marker atom can be nominated within each residue, or the centre of mass, centre of geometry or centre of charge may be used.

Group-based switching functions have several advantages. Better energy conservation can be achieved, and there are advantages when performing energy minimisation, since the potential is defined analytically at all points. However, it is important to beware of possible problems with group-based switching functions when the groups are large. We have already seen how this can arise when an ordinary group-based cutoff is used. Let us re-examine our arginine problem (Figure 6.16) when a switching function is employed. When the two marker atoms have a separation only slightly less than the upper switch cutoff, the switching function will be close to zero, and so there will not be the dramatic increase in energy that is observed with the simple cutoff. Nevertheless, although the switching function does help to prevent the simulation from 'blowing up', the representation of the energy and the forces in the system is still unsatisfactory. The only real alternative is to make the groups smaller or dispense with cutoffs altogether.

6.8 Long-range Forces

Those interactions that decay no faster than r^{-n} , where n is the dimensionality of the system, can be a problem as their range is often greater than half the box length. The charge-charge interaction, which decays as r^{-1} , is particularly problematic in molecular simulations. There is much evidence that it is important to properly model these long-range forces, which are particularly acute when simulating charged species such as molten salts (when it is not possible to construct neutral groups). A proper treatment of long-range forces can also be important when calculating certain properties, such as the dielectric constant. One way to tackle the errors introduced by an inadequate treatment of long-range forces would be to use a much larger simulation cell, but this is usually impractical. Nevertheless, increasing computer power does mean that more rigorous ways of dealing with long-range forces can be considered, even in simulations of large systems. A variety of methods have been developed to handle long-range forces. The methods that we will discuss in detail are the Ewald summation, the reaction field method and the cell multiple method.

6.8.1 The Ewald Summation Method

The Ewald sum was first devised by Ewald [Ewald 1921] to study the energetics of ionic crystals. In this method, a particle interacts with all the other particles in the simulation box and with all of their images in an infinite array of periodic cells. Figure 6.21 illustrates how the array of simulation cells is constructed; in the limit, the cell array is considered to have a spherical shape. The position of each image box (assumed for simplicity to be a cube of side L containing N charges) can be related to the central box by specifying a vector, each of whose components is an integral multiple of the length of the box, $(\pm iL, \pm jL, \pm kL)$; $i, j, k = 0, 1, 2, 3$, etc. The charge-charge contribution to the potential energy due to all pairs of charges in the central simulation box can be written:

$$\mathcal{V} = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} \quad (6.49)$$

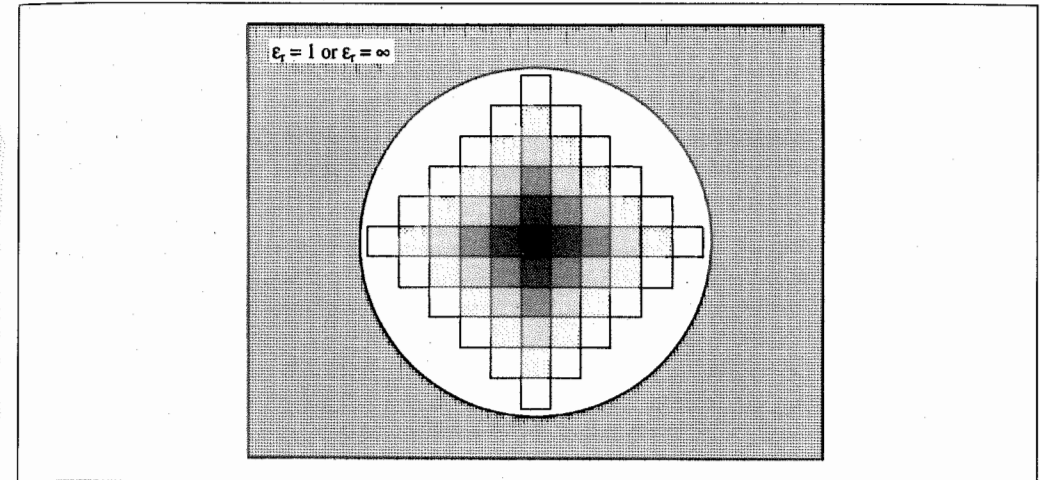


Fig. 6.21: The construction of a system of periodic cells in the Ewald method. (Figure adapted from Allen M P and D J Tildesley 1987. *Computer Simulation of Liquids*. Oxford, Oxford University Press.)

where r_{ij} is the minimum distance between the charges i and j . There are six boxes at a distance L from the central box with coordinates $(\mathbf{r}_{\text{box}})$ given by $(0, 0, L)$, $(0, 0, -L)$, $(0, L, 0)$, $(0, -L, 0)$, $(L, 0, 0)$ and $(-L, 0, 0)$ (only four of these are shown in the two-dimensional picture in Figure 6.21). The contribution of the charge-charge interaction between the charges in the central box and all images of all particles in these six surrounding boxes is given by:

$$\mathcal{V} = \frac{1}{2} \sum_{\mathbf{n} \neq 0} \sum_{i=1}^N \sum_{j=1}^N \frac{q_i q_j}{4\pi\epsilon_0 |\mathbf{r}_{ij} + \mathbf{n}|} \quad (6.50)$$

In general, for a box which is positioned at a cubic lattice point $\mathbf{n} (= (n_x L, n_y L, n_z L)$ with n_x, n_y, n_z being integers):

$$\mathcal{V} = \frac{1}{2} \sum_{\mathbf{n}} \sum_{i=1}^N \sum_{j=1}^N \frac{q_i q_j}{4\pi\epsilon_0 |\mathbf{r}_{ij} + \mathbf{n}|} \quad (6.51)$$

$|\mathbf{n}|$ thus takes the values $1, \sqrt{2}, \dots$. This expression is often written in such a way to incorporate the interactions between pairs of charges in the central box (for which $|\mathbf{n}| = 0$):

$$\mathcal{V} = \frac{1}{2} \sum'_{|\mathbf{n}|=0} \sum_{i=1}^N \sum_{j=1}^N \frac{q_i q_j}{4\pi\epsilon_0 |\mathbf{r}_{ij} + \mathbf{n}|} \quad (6.52)$$

The prime on the first summation indicates that the series does not include the interaction $i = j$ for $\mathbf{n} = 0$.

There is thus a contribution to the total energy from the interactions in the central box together with the interactions between the central box and all image boxes. There is also a contribution from the interaction between the spherical array of boxes and the surrounding

medium. The problem is that the summation in Equation (6.52) converges extremely slowly and in fact is *conditionally convergent*. A conditionally convergent series contains a mixture of positive and negative terms such that the positive terms alone form a divergent series (i.e. a series which does not have a finite sum) as do the negative terms when taken alone. The sum of a conditionally convergent series depends on the order in which its terms are considered. An additional problem with the Coulomb interaction is that it can vary rapidly at small distances.

The trick when calculating the Ewald sum is to convert the summation into two series, each of which converges much more rapidly. The mathematical foundation for this is the following identity:

$$\frac{1}{r} = \frac{f(r)}{r} + \frac{1-f(r)}{r} \quad (6.53)$$

The aim is thus to choose an appropriate function $f(r)$ which will deal with the rapid variation of $1/r$ at small r and the slow decay at long r . In the Ewald method each charge is considered to be surrounded by a neutralising charge distribution of equal magnitude but of opposite sign, as shown in Figure 6.22. A Gaussian charge distribution of the following functional form is commonly used:

$$\rho_i(\mathbf{r}) = \frac{q_i \alpha^3}{\pi^{3/2}} \exp(-\alpha^2 r^2) \quad (6.54)$$

The sum over point charges is now converted to a sum of the interactions between the charges *plus* the neutralising distributions. This dual summation (the 'real space'

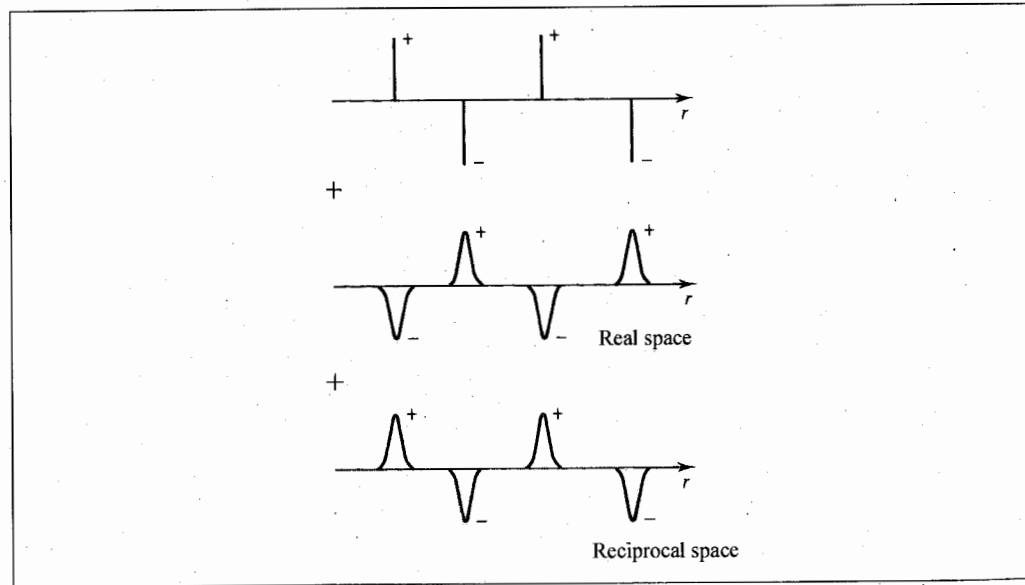


Fig. 6.22: In the Ewald summation method the initial set of charges are surrounded by a Gaussian distribution (calculated in real space) to which a cancelling charge distribution must be added (calculated in reciprocal space).

summation) is given by:

$$\mathcal{V} = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \sum_{|\mathbf{n}|=0} \frac{q_i q_j}{4\pi\epsilon_0} \frac{\text{erfc}(\alpha|\mathbf{r}_{ij} + \mathbf{n}|)}{|\mathbf{r}_{ij} + \mathbf{n}|} \quad (6.55)$$

erfc is the complementary error function, which is:

$$\text{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^{\infty} \exp(-t^2) dt \quad (6.56)$$

The Ewald method thus uses $\text{erfc}(r)$ for the function $f(r)$ in Equation (6.53). The crucial point is that this new summation involving the error function converges very rapidly and beyond some cutoff distance its value can be considered negligible. The rate of convergence depends upon the width of the cancelling Gaussian distributions; the wider the Gaussian, the faster the series converges. Specifically, α should be chosen so that the only terms in the series (6.55) are those for which $|\mathbf{n}| = 0$ (i.e. only pairwise interactions involving charges in the central box, or if a cutoff is used α is chosen so that only interactions with other charges within the cutoff are included). A second charge distribution is now added to the system which exactly counteracts the first neutralising distribution (Figure 6.22). The contribution from this second charge distribution is:

$$\mathcal{V} = \frac{1}{2} \sum_{\mathbf{k} \neq 0} \sum_{i=1}^N \sum_{j=1}^N \frac{1}{\pi L^3} \frac{q_i q_j}{4\pi\epsilon_0} \frac{4\pi^2}{k^2} \exp\left(-\frac{k^2}{4\alpha^2}\right) \cos(\mathbf{k} \cdot \mathbf{r}_{ij}) \quad (6.57)$$

This summation is performed in *reciprocal space*, the details of which need not concern us here. The vectors \mathbf{k} are reciprocal vectors and are given by $\mathbf{k} = 2\pi\mathbf{n}/L$. This reciprocal sum also converges much more rapidly than the original point-charge sum. However, the number of terms that must be included increases with the width of the Gaussians. There is thus a clear need to balance the real-space and reciprocal-space summations; the former converges more rapidly for large α , whereas the latter converges more rapidly for small α . A value for α of $5/L$ and 100-200 reciprocal vectors \mathbf{k} have been suggested as providing acceptable results. This reciprocal space summation corresponds to the second term $([1-f(r)]/r)$ in Equation (6.53); the requirement for this term is that it is a slowly varying function for all r . As such, its Fourier transform (which is what the summation is) can be represented by a small number of reciprocal vectors. The sum of Gaussian functions in real space includes the interaction of each Gaussian with itself. A third self-term must therefore be subtracted:

$$\mathcal{V} = -\frac{\alpha}{\sqrt{\pi}} \sum_{k=1}^N \frac{q_k^2}{4\pi\epsilon_0} \quad (6.58)$$

A fourth correction term may also be required, depending upon the medium that surrounds the sphere of simulation boxes. If the surrounding medium has an infinite relative permittivity (e.g. if it is a conductor) then no correction term is required. However, if the surrounding medium is a vacuum (with a relative permittivity of 1) then the following energy must be added:

$$\mathcal{V}_{\text{correction}} = \frac{2\pi}{3L^3} \left| \sum_{i=1}^N \frac{q_i}{4\pi\epsilon_0} \mathbf{r}_i \right|^2 \quad (6.59)$$

The final expression is thus:

$$\psi = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \left\{ \begin{aligned} & \sum_{|\mathbf{n}|=0}^{\infty} \frac{q_i q_j}{4\pi\epsilon_0} \frac{\operatorname{erfc}(\alpha|\mathbf{r}_{ij} + \mathbf{n}|)}{|\mathbf{r}_{ij} + \mathbf{n}|} \\ & + \sum_{k \neq 0} \frac{1}{\pi L^3} \frac{q_i q_j}{4\pi\epsilon_0} \frac{4\pi^2}{k^2} \exp\left(-\frac{k^2}{4\alpha^2}\right) \cos(\mathbf{k} \cdot \mathbf{r}_{ij}) \\ & - \frac{\alpha}{\sqrt{\pi}} \sum_{k=1}^N \frac{q_k^2}{4\pi\epsilon_0} + \frac{2\pi}{3L^3} \left| \sum_{k=1}^N \frac{q_k}{4\pi\epsilon_0} r_k \right|^2 \end{aligned} \right. \quad (6.60)$$

The Ewald sum is the most 'correct' way yet devised to accurately include all the effects of long-range forces in a computer simulation. It has been extensively used in simulations involving highly charged systems (such as ionic melts and in studies of processes in and on solids) and is increasingly being applied to other systems where electrostatic effects are important, such as lipid bilayers, proteins and DNA. Nevertheless, the Ewald method is not without problems and it tends to reinforce artefacts that arise from imposing periodic boundary conditions. For example, the method artificially results in each charge-charge interaction being minimised at a separation of half the box length. Instantaneous fluctuations in the simulation cell tend to be replicated throughout the infinite system rather than being damped out.

The Ewald summation is computationally quite expensive to implement. Under conditions of constant α (which will give the same density of reciprocal vectors, \mathbf{k}) then it scales as the square of the number of particles in the central simulation cell. If α is allowed to vary then the algorithm can be made to scale as $N^{3/2}$ though the consequent value of α might make the range of the Coulomb potential incompatible with the range of the van der Waals interactions. Several methods have been proposed to speed up the computationally demanding reciprocal space part of the calculation, such as the use of polynomial approximations but these do not solve the unfavourable N^2 scaling. The most promising way to tackle this difficulty is to modify the problem so that the fast Fourier transform (FFT) can be used to compute the reciprocal space summation. The fast Fourier algorithm scales as $N \ln N$, which gives considerable advantages over the N^2 alternative. If, in addition, a sufficiently large value of α is chosen such that the interatomic interaction is negligible for r_{ij} greater than a cutoff (e.g. 9 Å) then the real-space summation is reduced to order N and the order of the entire algorithm becomes $N \ln N$.

As outlined in Section 1.10.8, the FFT method requires that the data are not continuous but are discrete values. In order to employ the fast Fourier transform in the Ewald summation the point charges with their continuous coordinates must be replaced by a grid-based charge distribution. Each of the atomic point charges must thus be distributed among the surrounding grid points in some fashion so as to reproduce the potential of the charge at the original location. As usual an element of compromise is required; the more surrounding points that are used the more accurately the potential of the charge at the original location can be approximated but the greater the computational cost per particle. A popular approach is the particle-mesh method of Hockney and Eastwood [Hockney and Eastwood 1988], which uses the nearest 27 points in three dimensions. From this gridded charge density it is possible

to calculate (through use of the FFT algorithm) the potential due to the Gaussian distributions at the grid points, which by interpolation gives rise to the desired potential at (and thus the forces on) each of the particles. A number of variants on this general theme have been described, all of which use the fast Fourier transform algorithm but which differ in other aspects of their implementation. These include the particle-mesh Ewald method [Darden *et al.* 1993] and the particle-particle-particle-mesh approach [Hockney and Eastwood 1988; Luty *et al.* 1994, 1995]. Deserno and Holm presented a unification of these methods and also demonstrated that although very similar in spirit they could have very different accuracies [Deserno and Holm 1998a, b]. The particle-particle-particle-mesh approach was generally preferred as it was believed to be more flexible.

The Ewald method has been widely used to study highly polar or charged systems. Its use is considered routine for many types of solid-state materials. It is increasingly used for calculations on much larger molecular systems, such as proteins and DNA, due both to the increases in computer performance and to the new methodological advances we have just discussed [Darden *et al.* 1999]. For example, an early application of the particle-mesh Ewald method was the molecular dynamics simulation of a crystal of the protein bovine pancreatic trypsin inhibitor [York *et al.* 1994]. The full crystal environment was reproduced, with four protein molecules in the unit cell, together with associated water molecules and chloride counterions. Over the course of the 1 ns simulation the deviation of the simulated structures from the initial crystallographic structure was monitored. Once equilibrium was achieved this deviation (measured as the root-mean-square positional deviation) settled down to a value of 0.63 Å for all non-hydrogen atoms and 0.52 Å for the backbone atoms alone. By contrast, an equivalent simulation run with a 9 Å residue-based cutoff showed a deviation of more than 1.8 Å. In addition, the atomic fluctuations calculated from the Ewald simulation were in close agreement with those derived from the crystallographic temperature factors, unlike the non-Ewald simulation, which was significantly overestimated due to the use of the electrostatic cutoff. The highly charged nature of DNA makes it particularly important to deal properly with the electrostatic interactions and simulations using the particle-mesh Ewald approach are often much more stable, with the trajectories remaining much closer to the experimental structures [Cheatham *et al.* 1995].

6.8.2 The Reaction Field and Image Charge Methods

In the reaction field method, a sphere is constructed around the molecule with a radius equal to the cutoff distance. The interaction with molecules that are within the sphere is calculated explicitly. To this is added the energy of interaction with the medium beyond the sphere, which is modelled as a homogeneous medium of dielectric constant ϵ_s (Figure 6.23). The electrostatic field due to the surrounding dielectric is given by:

$$\mathbf{E}_i = \frac{2(\epsilon_s - 1)}{\epsilon_s + 1} \left(\frac{1}{r_c^3} \right) \sum_{j: r_{ij} \leq r_c} \boldsymbol{\mu}_j \quad (6.61)$$

where $\boldsymbol{\mu}_j$ are the dipoles of the neighbouring molecules that are within the cutoff distance (r_c) of the molecule i . The interaction between the molecule i and the reaction field equals $\mathbf{E}_i \cdot \boldsymbol{\mu}_i$.

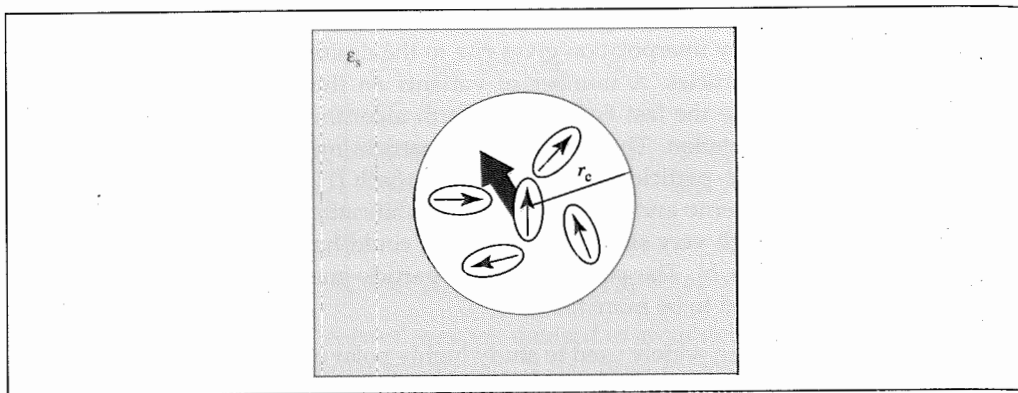


Fig. 6.23: The reaction field method. The shaded arrow represents the sum of the dipoles of the other molecules within the cutoff sphere.

which is added to the short-range molecule-molecule interaction. Problems with the reaction field method may arise from discontinuities in the energy and/or force when the number of molecules j within the cavity of the molecule i changes. These problems can be avoided by employing a switching function for molecules that are near the reaction field boundary.

Similar approaches employ a single boundary for the entire system. This boundary may be spherical or may have a more complicated shape that better approximates the true molecular surface of the molecule. In the *image charge method*, a spherical boundary is employed and the reaction field due to a charge inside the boundary is considered to arise from a so-called image charge situated in the continuous dielectric beyond the sphere (Figure 6.24) [Friedman 1975]. If the position of the charge is \mathbf{r}_i , then the image charge is located at $(R/r_i)^2 \mathbf{r}_i$ (where R is the radius of the bounding sphere) and has magnitude:

$$q_{\text{im}} = -\frac{(\epsilon_s - \epsilon_r)}{(\epsilon_s + \epsilon_r)} \frac{q_i R}{r_i} \quad (6.62)$$

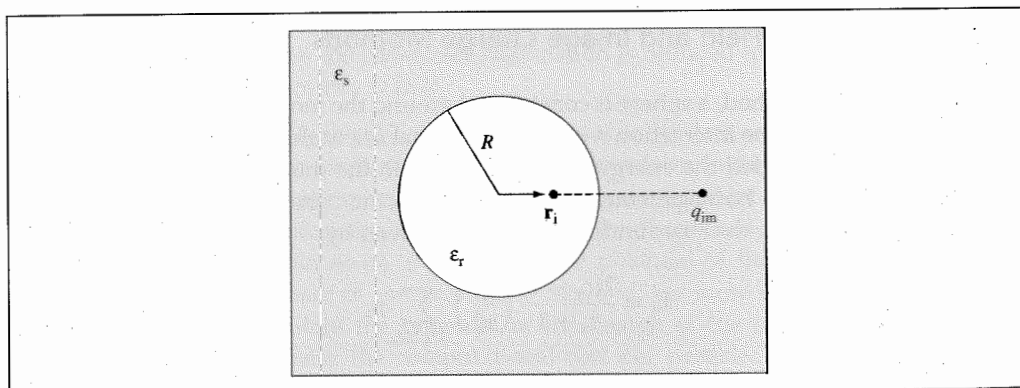


Fig. 6.24: The image charge method.

where ϵ_r and ϵ_s are the dielectric constants inside and outside the boundary, respectively. This expression holds if the dielectric constant beyond the boundary is much greater than that inside ($\epsilon_s \gg \epsilon_r$). A drawback with this method is that as a charge approaches the boundary, so too does its image and the method breaks down.

The reaction field and image charge methods have the advantages of being conceptually simple, relatively easy to implement and computationally efficient. However, they do rely upon the assumption that molecules beyond the cutoff can be modelled as a continuous dielectric. This is not necessarily the case but is often a reasonable assumption for homogeneous fluids. A value for the dielectric of the surrounding continuum must also be specified. This can be taken from experimental data, but the dielectric constant may be the property that one is trying to calculate! In practice, it is often necessary only to ensure that $\epsilon_r \leq \epsilon_s \leq \infty$. There are several ways in which the dielectric constant can be calculated from a computer simulation. A common approach is via the average of the square of the total dipole moment of the system, $\langle \mathbf{M}^2 \rangle$. With a reaction field boundary the dielectric constant ϵ_r is given by:

$$\frac{4\pi \langle \mathbf{M}^2 \rangle}{9 V k_B T} = \frac{(\epsilon_r - 1)(2\epsilon_s + 1)}{3(2\epsilon_s + \epsilon_r)} \quad (6.63)$$

where V is the volume of the simulation system. Even though the value of $\langle \mathbf{M}^2 \rangle$ can vary quite considerably with the reaction field dielectric ϵ_s , almost identical values of ϵ_r are obtained. An alternative approach is to determine the polarisation response of the liquid to an electric field \mathbf{E}_0 . If the average dipole moment per unit volume along the direction of the applied field is $\langle \mathbf{P} \rangle$ then the dielectric constant is given by:

$$\frac{4\pi \langle \mathbf{P} \rangle}{3 \mathbf{E}_0} = \frac{(\epsilon_r - 1)(2\epsilon_s + 1)}{3(2\epsilon_s + \epsilon_r)} \quad (6.64)$$

This perturbation method is claimed to be more efficient than the fluctuating dipole method, at least for certain water models [Alper and Levy 1989], but it is important to ensure that the polarisation $\langle \mathbf{P} \rangle$ is linear in the electric field strength to avoid problems with dielectric saturation.

6.8.3 The Cell Multipole Method for Non-bonded Interactions

The cell multipole method (also called the fast multipole method) is an algorithm that enables *all* $N(N-1)$ pairwise non-bonded interactions to be enumerated in a time that scales linearly with N , rather than N^2 , as in the standard Ewald approach [Greengard and Roklin 1987; Ding *et al.* 1992a, b; Greengard 1994]. The cell multipole method can be used to evaluate interactions that can be expressed in the following general form:

$$\sum_i \sum_{j>i} \frac{q_i q_j}{|\mathbf{r}_i - \mathbf{r}_j|^p} \quad (6.65)$$

Both the Coulomb and Lennard-Jones potentials can be considered examples of this type. In the cell multipole method the simulation space is divided into uniform cubic

cells. The multipole moments (charge, dipole, quadrupole) of each cell are then calculated by summing over the atoms contained within the cell. The interaction between all of the atoms in the cell and another atom outside the cell (or indeed another cell) can then be calculated using an appropriate multipole expansion (see Section 4.9.1).

This multipole expansion is only valid if the separation between the interacting particles (be they atoms, molecules or cells) is larger than the sum of the radii of convergence of the multipoles. In the cell multipole method, the multipole expansion is used for interactions that are more than one cell distance away. For interactions that are within one cell distance the usual atomic pairwise interaction method is employed.

Consider an atom in a cell, C_0 . The interactions with atoms in nearby cells are calculated using the usual pairwise formulae. There are 27 such cells (i.e. the cell in which the atom is positioned and the surrounding 26 cells). The interaction between the atom and all of the atoms in each of the faraway cells is then calculated using the multipole expansion. The potential due to a faraway cell will be approximately constant for all atoms in the cell of current interest, C_0 (the cells are usually small, containing on average four atoms). Thus the potential due to each faraway cell can be represented as a Taylor series expansion about the centre of C_0 . If there are M cells in total then there are $M - 27$ faraway cells; then the calculation of these cell-cell interactions for the entire system will be of order $M(M - 27)$. As the number of cells is approximately equal to the number of atoms, this still leaves us with a quadratic dependency upon the number of atoms present (though it does now vary as about $N^2/16$, if there is an average of four atoms per cell).

The algorithm can be converted to one which shows linear dependency by recognising that in the method we have just described, the interactions due to very faraway cells are calculated with the same accuracy as interactions with cells that are much closer. This level of accuracy can be considered unnecessary as any error is largely due to the closer cells. The small cells are thus grouped into larger cells, with the cell size increasing with the distance from the cell of interest, C_0 . The accuracy of the calculation remains approximately constant if the ratio of the cell size to the distance remains constant. This grouping scheme is illustrated in Figure 6.25. The multipoles for each of the larger cells are calculated by translating and adding the moments of its constituent smaller cells. The use of multipole expansions and Taylor series approximations does mean that there will be a degree of truncation error, though this can be reduced by simply including more terms in the multipole expansion. The cell multipole algorithm requires an amount of bookkeeping to keep track of the hierarchy of the cells, which means that up to a certain size of problem the exact N^2 algorithm is faster. The algorithm then suddenly switches to a linear dependence. There is some debate over the break-even point at which the cell multipole method is equally as fast as an N^2 method, with estimates ranging from 300 particles to 100 000. Another complication to the debate is the introduction of the fast Fourier transform Ewald methods with their $N \ln N$ scaling. Nevertheless, for calculations on systems with thousands, if not millions, of atoms, the cell multipole methods appear promising, especially as enhanced versions are developed [Petersen *et al.* 1994].

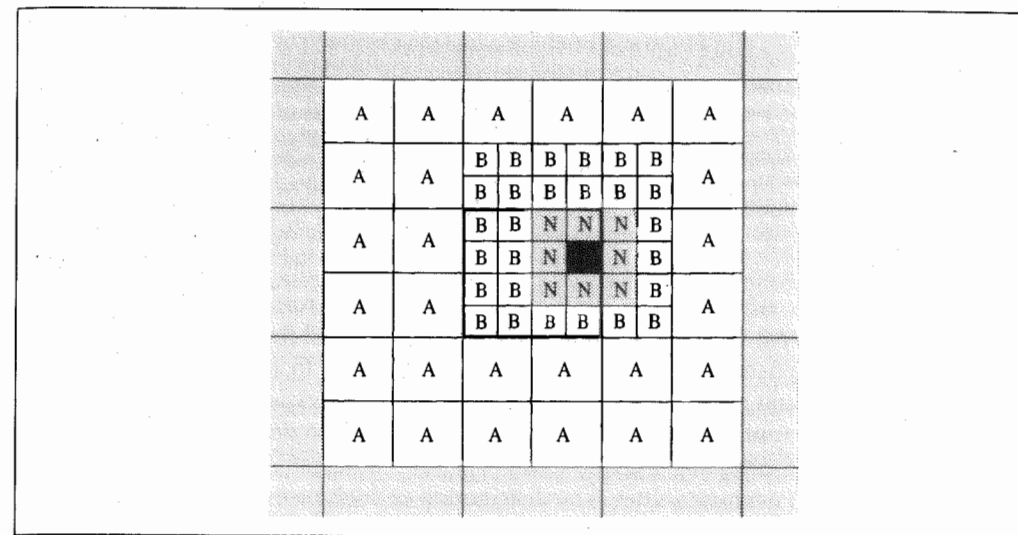


Fig. 6.25: The hierarchy of cells used in the cell multipole method. For an atom in the black cell, the interactions with atoms in the 26 nearby cells (N) are calculated explicitly. Interactions with the atoms in cells labelled A and B are calculated using a Taylor series multipole expansion. (Figure adapted from Ding H-Q, N Karasawa and W A Goddard III, 1992b. *The Reduced Cell Multipole Method for Coulomb Interactions in Periodic Systems with Million-Atom Unit Cells*. *Chemical Physics Letters* 196:6-10.)

6.9 Analysing the Results of a Simulation and Estimating Errors

A simulation can generate an enormous amount of data, which should be properly analysed to extract relevant properties and to check that the calculation has behaved properly. The primary reason for undertaking a particular simulation may be to calculate just a single physical or thermodynamic property or to investigate the conformational properties of a molecule. However, it is also advisable to check that other aspects of the simulation have performed as expected. Some properties can be calculated during the simulation itself (such as the energy and the virial), but it is often sensible not to impose too severe a burden on the simulation program itself. In part this is because many properties do not vary significantly from one step to another but can be calculated at less frequent intervals. The configurations (i.e. positions of each atom or molecule in the system) do not change much from one step to another in either a molecular dynamics or Monte Carlo simulation, and so it is usual to store configurations every 5-25 steps, depending upon the nature of the system and the disk space available. It is good practice to visually examine configurations selected from throughout the simulation to ensure that no strange or unexpected behaviour is present. In many simulations of molecular systems the major objective is to investigate the structural behaviour of the system rather than to calculate thermodynamic properties, and so the focus of the analysis will change accordingly.

A computer simulation is subject to error, and this error should be properly calculated and assessed. Of course, computers only do what they are told to by the programmer, and so a program will always give the same results for the same set of initial conditions (if not, some serious fault should be suspected!). The results of a computer simulation may be subject to two kinds of error, just as any other scientific experiment. These are systematic errors and statistical errors. Systematic error results in a constant bias from the 'proper' behaviour. The most obvious effect of a systematic error is to displace the average property from its proper value. Systematic errors are sometimes due to a fault in the simulation algorithm or in the energy model and may be relatively easy to spot, especially if they have an obvious or even catastrophic effect on the simulation. Systematic errors may also arise from approximations inherent in the algorithm, such as truncation (all finite difference methods used in molecular dynamics generate only an approximation to the true integral of the equations of motion) and round-off errors (due to the limited precision with which numbers can be stored in a computer). Such errors can be more difficult to detect. One way to detect systematic error is to compare the distribution of the values of simple thermodynamic properties about their average values. The distribution of such properties about their average values should be normal (i.e. Gaussian), such that the probability of finding a particular value for the property A is given by:

$$p(A) = \frac{1}{\sigma\sqrt{2\pi}} \exp[-(A - \langle A \rangle)^2 / 2\sigma^2] \quad (6.66)$$

where σ^2 is the variance, given by $\sigma^2 = \langle (A - \langle A \rangle)^2 \rangle$. The standard deviation is the square root of the variance. More information on these statistical terms can be found in Section 1.10.7.

The chi-squared test can be used to provide a quantitative estimate of the deviation of a calculated distribution from that expected. Suppose that the value of some property (A) has been calculated from the simulation at regular intervals to give a total of M values. The average value of the property A is determined together with the standard deviation. The data, comprising all of the A values from the simulation, are then divided into bins such that the number of values in each bin (M_i) is approximately the same. The number of values that would be expected in the i th bin is:

$$n_i = \frac{M}{\sigma\sqrt{2\pi}} \int_{A_i - \Delta A/2}^{A_i + \Delta A/2} \exp\left[-\frac{(A_i - \langle A \rangle)^2}{2\sigma^2}\right] \quad (6.67)$$

where A_i is the value of the property in the i th bin and ΔA is the width of each bin. The number of values that would be expected in each bin, n_i , does not have to be integral, though the actual number as determined from the simulation (M_i) will of course be an integer. The chi-squared function is given by:

$$\chi^2 = \sum_i \frac{(M_i - n_i)^2}{n_i} \quad (6.68)$$

If χ^2 is large (bigger than unity) then it is unlikely that the two distributions are the same. Any significant deviations from the expected behaviour should be investigated further to try to eliminate as much of the systematic error as possible. It is good practice to vary as

many of the parameters as possible: using different computers, different compilers, different algorithms and different ways of implementing a given algorithm, and different simulation methods (Monte Carlo and molecular dynamics) not only to test the component parts of the simulation but also the software used to perform the calculation.

If all sources of systematic error can be eliminated, there will still remain statistical errors. These errors are often reported as standard deviations. What we would particularly like to estimate is the error in the average value, $\langle A \rangle$. The standard deviation of the average value is calculated as follows:

$$\sigma_{\langle A \rangle} = \frac{\sigma_A}{\sqrt{M}} = \frac{\sqrt{\sum_{i=1}^M (A(i) - \langle A \rangle)^2}}{\sqrt{M}} \quad (6.69)$$

where $\sigma_{\langle A \rangle}$ is the standard deviation of the average value $\langle A \rangle$ obtained from M data values with respect to the run average, σ_A . Thus the standard deviation of the calculated average is inversely proportional to the square root of the number of data values, and so a longer simulation gives rise to a more accurate value. An important feature of Equation (6.69) is that it applies to *independent* (i.e. random) samples. Thus the number M in the denominator is not simply equal to the number of steps in the simulation. This is because there is a high degree of correlation between successive configurations in either a Monte Carlo or molecular dynamics simulation. What we need to know is the correlation or relaxation 'time' of the simulation; this is the number of steps required for the system to lose its 'memory' of previous configurations. In molecular dynamics, where successive steps are related in a temporal fashion, the correlation 'time' is a true time and will be discussed in more detail in Section 7.6. Usually, the correlation time will be unknown prior to the simulation but it can be estimated as follows. First, the configurations are broken down into a series of blocks. Suppose each block contains t_b successive steps and that there are n_b blocks (so the total simulation contains $t_b n_b$ steps, as shown in Figure 6.26). The average value of the property is calculated for each block:

$$\langle A \rangle_b = \frac{1}{t_b} \sum_{i=1}^{t_b} A_i \quad (6.70)$$

As the number of steps t_b in each block increases, so it would be expected that the block averages become uncorrelated. When this is the case, then the variance of the block averages, $\sigma^2(\langle A \rangle_b)$, will become inversely proportional to t_b . $\sigma^2(\langle A \rangle_b)$ is calculated as follows:

$$\sigma^2(\langle A \rangle_b) = \frac{1}{n_b} \sum_{b=1}^{n_b} (\langle A \rangle_b - \langle A \rangle_{\text{total}})^2 \quad (6.71)$$

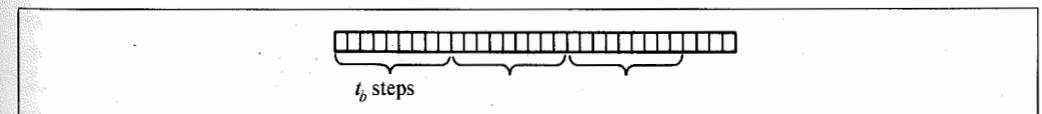


Fig. 6.26: Blocking a simulation to calculate the statistical error.

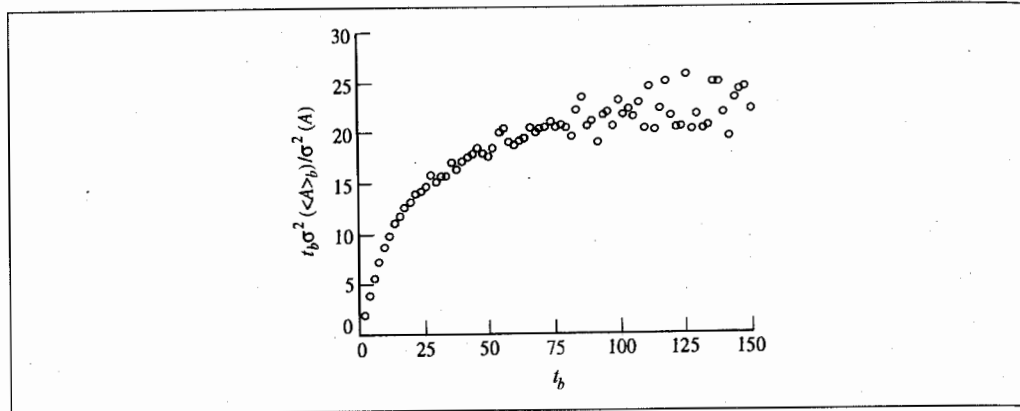


Fig. 6.27: Calculating the statistical efficiency, σ . A plot of $t_b \sigma^2 (\langle A \rangle_b) / \sigma^2(A)$ against t_b shows a steep rise before levelling off. Here the property A corresponds to the pressure calculated from the molecular dynamics simulation of argon.

where $\langle A \rangle_{\text{total}}$ is the average over the entire simulation. The limiting number of steps to obtain uncorrelated configurations (the statistical inefficiency, s) can be calculated using:

$$s = \lim_{t_b \rightarrow \infty} \frac{t_b \sigma^2 (\langle A \rangle_b)}{\sigma^2(A)} \quad (6.72)$$

To determine s , $t_b \sigma^2 (\langle A \rangle_b) / \sigma^2(A)$ is plotted against t_b or $\sqrt{t_b}$. The graph should show a steep rise for low t_b and then level off to give a plateau, as shown in Figure 6.27. The plateau value is the limiting value that gives the correlation time ($s \approx 23$ in this case).

Having determined the value of s , the 'true' standard deviation of the average value is related to the 'true' error for an infinite simulation by:

$$\sigma_{(A)} \approx \sigma \sqrt{\frac{s}{M}} \quad (6.73)$$

M here is the actual number of steps or iterations in the simulation. If the value of s can be reduced, then a more accurate average value can be calculated for a given length of simulation. This should be an important consideration when deciding what simulation protocol to use. For example, it may be more appropriate to use a complex simulation algorithm than a simpler one if the statistical inefficiency is significantly reduced.

If the relaxation time is known, then sample averages are often best calculated using the block method (Figure 6.26). Each block should contain more steps than the relaxation time. The sample average for the whole run can be obtained in a variety of ways:

1. Stratified systematic sampling, in which a single value of the property is taken from each block;
2. Stratified random sampling, in which a single value is taken at random from each block;
3. Coarse graining, in which the average value for each block is determined and then the average for the run is calculated by averaging the coarse-grain averages.

The coarse-graining approach is commonly used for thermodynamic properties whereas the systematic or random sampling methods are appropriate for static structural properties such as the radial distribution function.

Another way to improve the error in a simulation, at least for properties such as the energy and the heat capacity that depend on the size of the system (the extensive properties), is to increase the number of atoms or molecules in the calculation. The standard deviation of the average of such a property is proportional to $1/\sqrt{N}$. Thus, more accurate values can be obtained by running longer simulations on larger systems. In computer simulation it is unfortunately the case that the more effort that is expended the better the results that are obtained. Such is life!

Appendix 6.1 Basic Statistical Mechanics

The Boltzmann distribution is fundamental to statistical mechanics. The Boltzmann distribution is derived by maximising the entropy of the system (in accordance with the second law of thermodynamics) subject to the constraints on the system. Let us consider a system containing N particles (atoms or molecules) such that the energy levels of the particles are $\varepsilon_1, \varepsilon_2, \dots$. If there are n_1 particles in the energy level ε_1 , n_2 particles in ε_2 and so on, then there are W ways in which this distribution can be achieved:

$$W(n_1, n_2, \dots) = N! / n_1! n_2! \dots \quad (6.74)$$

The most favourable distribution is the one with the highest weight, and this corresponds to the configuration with just one particle in each energy level ($W = N!$). However, there are two important constraints on the system. First, the total energy is fixed:

$$\sum_i n_i \varepsilon_i = E \quad (6.75)$$

The second constraint arises from the fact that the total number of particles is fixed:

$$\sum_i n_i = N \quad (6.76)$$

The Boltzmann distribution gives the number of particles n_i in each energy level ε_i as:

$$\frac{n_i}{N} = \frac{\exp(-\varepsilon_i/k_B T)}{\sum_i \exp(-\varepsilon_i/k_B T)} \quad (6.77)$$

The denominator in this expression is the molecular partition function:

$$q = \sum_i \exp(-\varepsilon_i/k_B T) \quad (6.78)$$

For translational, rotational and vibrational motion the partition function can be calculated using standard results obtained by solving the Schrödinger equation:

$$\text{translation: } q^t = \left(\frac{2\pi m k_B T}{h^2} \right)^{3/2} V \quad (6.79)$$

where V is the volume

$$\text{rotation: } q^r \approx \left(\frac{\pi^{1/2}}{\sigma}\right) \left(\frac{2I_A k_B T}{h^2}\right) \left(\frac{2I_B k_B T}{h^2}\right) \left(\frac{2I_C k_B T}{h^2}\right) \quad (6.80)$$

where I_A, I_B, I_C are the moments of inertia and σ is the symmetry number (2 for H_2O , 3 for NH_3 , 12 for benzene)

$$\text{vibration: } r^v = \frac{1}{1 - \exp(-\hbar\omega/k_B T)} \quad (6.81)$$

ω is the angular frequency: $\omega = \sqrt{k/\mu}$, where μ is the reduced mass. This form of the vibrational partition function is measured relative to the zero-point energy.

In computer simulations, we are particularly interested in the properties of a system comprising a number of particles. An ensemble is a collection of such systems, as might be generated using a molecular dynamics or a Monte Carlo simulation. Each member of the ensemble has an energy, and the distribution of the system within the ensemble follows the Boltzmann distribution. This leads to the concept of the ensemble partition function, Q .

Various thermodynamic properties can be calculated from the partition function. Here we simply state some of the most common:

$$\text{internal energy: } U = \frac{k_B T^2}{Q} \left(\frac{\partial Q}{\partial T}\right)_V = k_B T^2 \left(\frac{\partial \ln Q}{\partial T}\right)_V \quad (6.82)$$

$$\text{enthalpy: } H = k_B T^2 \left(\frac{\partial \ln Q}{\partial T}\right)_V + k_B T V \left(\frac{\partial \ln Q}{\partial V}\right)_T \quad (6.83)$$

$$\text{Helmholtz free energy: } A = -k_B T \ln Q \quad (6.84)$$

$$\text{Gibbs free energy: } G = -k_B T \ln Q + k_B T V \left(\frac{\partial \ln Q}{\partial V}\right)_T \quad (6.85)$$

Appendix 6.2 Heat Capacity and Energy Fluctuations

The heat capacity is related to the internal energy U by

$$C_V = \left(\frac{\partial U}{\partial T}\right)_V \quad (6.86)$$

If we differentiate the expression for the internal energy, Equation (6.20), we can obtain the heat capacity in terms of the partition function:

$$C_V = \frac{\partial}{\partial T} \left(\frac{k_B T^2}{Q} \frac{\partial Q}{\partial T}\right)_V = \frac{k_B T^2}{Q} \frac{\partial^2 Q^2}{\partial T^2} + \frac{2k_B T}{Q} \frac{\partial Q}{\partial T} - \frac{k_B T^2}{Q^2} \left(\frac{\partial Q}{\partial T}\right)^2 \quad (6.87)$$

The desired expression is obtained by writing each of these three terms as a function of the average energy, $\langle E \rangle$. The internal energy is just the expectation value of the energy, $\langle E \rangle$, and so:

$$\langle E \rangle = \frac{k_B T^2}{Q} \frac{\partial Q}{\partial T} \quad (6.88)$$

Thus for the second term in Equation (6.87) we have

$$\frac{2k_B T}{Q} \frac{\partial Q}{\partial T} = \frac{2\langle E \rangle}{T} \quad (6.89)$$

We can also rewrite the third term in Equation (6.87):

$$k_B T \left(\frac{1}{Q} \frac{\partial Q}{\partial T}\right)^2 = \frac{\langle E \rangle^2}{k_B T} \quad (6.90)$$

For the first term, we need to do a little more work. The starting point is:

$$\frac{\partial}{\partial T} \left(\frac{\langle E \rangle}{k_B T^2}\right) = \frac{\partial}{\partial T} \left\{ \frac{1}{Q} \left(\frac{\partial Q}{\partial T}\right) \right\} \quad (6.91)$$

or

$$-2 \frac{\langle E \rangle}{k_B T^3} = \frac{1}{Q} \frac{\partial^2 Q}{\partial T^2} + \frac{\partial Q}{\partial T} \frac{\partial}{\partial T} \left(\frac{1}{Q}\right) \quad (6.92)$$

We can use the chain rule as follows:

$$\frac{\partial Q}{\partial T} \frac{\partial}{\partial T} \left(\frac{1}{Q}\right) = \frac{\partial Q}{\partial T} \frac{\partial Q}{\partial T} \frac{\partial}{\partial T} \left(\frac{1}{Q}\right) = -\left(\frac{\partial Q}{\partial T}\right)^2 \left(\frac{1}{Q}\right)^2 \quad (6.93)$$

Thus

$$\frac{k_B T^2}{Q} \frac{\partial^2 Q}{\partial T^2} = -2 \frac{\langle E \rangle}{k_B T^3} + \frac{\langle E^2 \rangle}{k_B^2 T^4} \quad (6.94)$$

So

$$C_V = k_B T^2 \left\{ -2 \frac{\langle E \rangle}{k_B T^3} + \frac{\langle E^2 \rangle}{k_B^2 T^4} \right\} + 2 \frac{\langle E \rangle}{T} - \frac{\langle E \rangle^2}{k_B T^2} \quad (6.95)$$

or

$$C_V = \frac{(\langle E^2 \rangle - \langle E \rangle^2)}{k_B T^2} \quad (6.96)$$

Appendix 6.3 The Real Gas Contribution to the Virial

If the gas particles interact through a pairwise potential, then the contribution to the virial from the intermolecular forces can be derived as follows. Consider two atoms i and j separated by a distance r_{ij} .

$$r_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2} \quad (6.97)$$

The contribution to the virial from the interaction $v(r_{ij})$ between atoms i and j is given by:

$$W_{\text{real}} = \left[x_i \frac{\partial}{\partial x_i} + x_j \frac{\partial}{\partial x_j} + y_i \frac{\partial}{\partial y_i} + y_j \frac{\partial}{\partial y_j} + z_i \frac{\partial}{\partial z_i} + z_j \frac{\partial}{\partial z_j} \right] v(r_{ij}) \quad (6.98)$$

Since

$$x_i \frac{\partial r_{ij}}{\partial x_i} = x_i \frac{(x_i - x_j)}{r_{ij}} \quad \text{and} \quad x_j \frac{\partial r_{ij}}{\partial x_i} = -x_j \frac{(x_i - x_j)}{r_{ij}} \quad (6.99)$$

and similarly for the y and z coordinates, we can apply the chain rule, $\partial/\partial x_i = (\partial/\partial r_{ij})(\partial r_{ij}/\partial x_i)$, as follows:

$$W_{\text{real}} = \left[\frac{(x_i - x_j)^2}{r_{ij}} + \frac{(y_i - y_j)^2}{r_{ij}} + \frac{(z_i - z_j)^2}{r_{ij}} \right] \frac{dv(r_{ij})}{dr_{ij}} = r_{ij} \frac{dv(r_{ij})}{dr_{ij}} \quad (6.100)$$

When we include the contributions from all pairs of atoms, we obtain:

$$W_{\text{real}} = \sum_{i=1}^N \sum_{j=i+1}^N r_{ij} \frac{dv(r_{ij})}{dr_{ij}} \quad (6.101)$$

Appendix 6.4 Translating Particle Back into Central Box for Three Box Shapes

From Smith W 1983. The Periodic Boundary Condition in Non-Cubic MD Cells: Wigner-Seitz Cells with Reflection Symmetry. *CCP5 Quarterly* 10:37–42. This table is expressed using references to several built-in FORTRAN functions. The AINT function returns the integral part of its argument, e.g. AINT(3.4) = 3.0; AINT(4.7) = 4.0; AINT(−0.5) = 0.0 and AINT(−1.7) = −1.0. ANINT returns the nearest integer, so ANINT(0.49) = 0.0 and ANINT(0.51) = 1.0. SIGN(x, y) returns $|x|$ if $y \geq 0$ and $-|x|$ if $y < 0$. ABS(x) returns the absolute value of x , $|x|$. Equivalent functions also exist in most other programming languages.

Rectangular box, side $2a$ (x) by $2b$ (y) by $2c$ (z)

```
x = x - 2 * a * AINT(x/a)
y = y - 2 * b * AINT(y/b)
z = z - 2 * c * AINT(z/c)
A common alternative is:
x = x - a * ANINT(x/a)
y = y - b * ANINT(y/b)
z = z - c * ANINT(z/c)
```

Truncated octahedron derived from cube of side $2a$

```
x = x - 2 * a * AINT(x/a)
y = y - 2 * b * AINT(y/a)
z = z - 2 * c * AINT(z/a)
if (ABS(x) + ABS(y) + ABS(z)) >= 1.5 * A
then
  x = x - SIGN(a, x)
  y = y - SIGN(a, y)
  z = z - SIGN(a, z)
endif
```

Hexagonal prism of length $2a$ (in z direction) and distance between opposite faces of the hexagon $2b$

```
z = z - 2 * a * AINT(z/a)
x = x - 2 * b * AINT(x/b)
if (ABS(x) + sqrt(3) * ABS(y)) >= 2 * B then
  x = x - SIGN(b, x)
  y = y - SIGN(sqrt(3) * b, y)
endif
```

Further Reading

- Allen M P and D J Tildesley 1987. *Computer Simulation of Liquids*. Oxford, Oxford University Press.
- Bradbury T C 1968. *Theoretical Mechanics*. Malabar, FL, Krieger.
- Chandler D 1987. *Introduction to Modern Statistical Mechanics*. New York, Oxford University Press.
- Hansen J P and I R McDonald 1976. *Theory of Simple Liquids*. London, Academic Press.
- Smith P E and van Gunsteren W F 1993. Methods for the Evaluation of Long Range Electrostatic Forces. In van Gunsteren W F, P K Weiner and A J Wilkinson (Editors). *Computer Simulation of Biomolecular Systems*. Leiden, ESCOM.
- van Gunsteren W F and H J C Berendsen 1990. Computer Simulation of Molecular Dynamics: Methodology, Applications and Perspectives in Chemistry. *Angewandte Chemie International Edition in English* 29:992–1023.

References

- Adams D J 1983. Alternatives to the Periodic Cube in Computer Simulation. *CCP5 Quarterly* 10:30–36.
- Alper H E and R M Levy 1989. Computer Simulations of the Dielectric Properties of Water – Studies of the Simple Point-Charge and Transferable Intermolecular Potential Models. *Journal of Chemical Physics* 91:1242–1251.
- Cheatham T E III, J L Miller, T Fox, T A Darden and P A Kollman 1995. Molecular Dynamics Simulations on Solvated Biomolecular Systems: The Particle Mesh Ewald Method Leads to Stable Trajectories of DNA, RNA and Proteins. *Journal of the American Chemical Society* 117:4193–4194.
- Darden T A, L Perera, L Li and L Pedersen 1999. New Tricks for Modelers from the Crystallography Toolkit: The Particle Mesh Ewald Algorithm and Its Use in Nucleic Acid Simulations. *Structure with Folding and Design* 7:R55–R60.
- Darden T A, D York and L Pedersen 1993. Particle-mesh Ewald: An $N \cdot \log(N)$ method for Ewald sums in large systems. *Journal of Chemical Physics* 98:10089–10092.
- Deserno M and C Holm 1998a. How to Mesh Up Ewald Sums. I. A Theoretical and Numerical Comparison of Various Particle Mesh Routines. *Journal of Chemical Physics* 109:7678–7693.
- Deserno M and C Holm 1998b. How to Mesh Up Ewald Sums. II. An Accurate Error Estimate for the Particle-Particle-Particle-Mesh Algorithm. *Journal of Chemical Physics* 109:7694–7701.
- Ding H-Q, N Karasawa and W A Goddard III 1992a. Atomic Level Simulations on a Million Particles: The Cell Multipole Method for Coulomb and London Nonbonding Interactions. *Journal of Chemical Physics* 97:4309–4315.
- Ding H-Q, N Karasawa and W A Goddard III 1992b. The Reduced Cell Multipole Method for Coulomb Interactions in Periodic Systems with Million-Atom Unit Cells. *Chemical Physics Letters* 196:6–10.
- Ewald P 1921. Due Berechnung optischer und elektrostatischer Gitterpotentiale. *Annalen der Physik* 64:253–287.
- Friedman H L 1975. Image Approximation to the Reaction Field. *Molecular Physics* 29:1533–1543.
- Greengard L 1994. Fast Algorithms for Classical Physics. *Science* 265:909–914.
- Greengard L and V I Roklin 1987. A Fast Algorithm for Particle Simulations. *Journal of Computational Physics* 73:325–348.
- Hockney R W and J W Eastwood 1988. *Computer Simulation using Particles*. Bristol, Adam Hilger.
- Luty B A, M E David, I G Tironi and W F van Gunsteren 1994. A Comparison of Particle-Particle, Particle-Mesh and Ewald Methods for Calculating Electrostatics Interactions in Periodic Molecular Systems. *Molecular Simulation* 14:11–20.

- Luty B A, I G Tironi and W F van Gunsteren 1995. Lattice-sum methods for calculating electrostatic interactions in molecular simulations. *Journal of Chemical Physics* **103**:3014–3021.
- Petersen H G, D Soelvaso, J W Perram and E R Smith 1994. The Very Fast Multipole Method. *Journal of Chemical Physics* **101**:8870–8876.
- Thompson S M 1983. Use of Neighbour Lists in Molecular Dynamics. *CCP5 Quarterly* **8**:20–28.
- van Gunsteren W F and H J C Berendsen 1986. *GROMOS User Guide*.
- Verlet L 1967. Computer 'Experiments' on Classical Fluids. II. Equilibrium Correlation Functions. *Physical Review* **165**:201–204.
- Wolf M L, J R Walker and C R A Catlow 1984. A Molecular Dynamics Simulation Study of the Superionic Conductor Lithium Nitride: I. *Journal of Physical Chemistry* **17**:6623–34.
- York D M, A Wlodawer, L G Pedersen and T A Darden 1994. Atomic-level Accuracy in Simulations of Large Protein Crystals. *Proceedings of the National Academy of Sciences USA* **91**:8715–8718.

CHAPTER SEVEN

Molecular Dynamics Simulation Methods

7.1 Introduction

In molecular dynamics, successive configurations of the system are generated by integrating Newton's laws of motion. The result is a trajectory that specifies how the positions and velocities of the particles in the system vary with time. Newton's laws of motion can be stated as follows:

1. A body continues to move in a straight line at constant velocity unless a force acts upon it.
2. Force equals the rate of change of momentum.
3. To every action there is an equal and opposite reaction.

The trajectory is obtained by solving the differential equations embodied in Newton's second law ($F = ma$):

$$\frac{d^2 x_i}{dt^2} = \frac{F_{x_i}}{m_i} \quad (7.1)$$

This equation describes the motion of a particle of mass m_i along one coordinate (x_i) with F_{x_i} being the force on the particle in that direction.

It is helpful to distinguish three different types of problem to which Newton's laws of motion may be applied. In the simplest case, no force acts on each particle between collisions. From one collision to the next, the position of the particle thus changes by $\mathbf{v}_i \delta t$, where \mathbf{v}_i is the (constant) velocity and δt is the time between collisions. In the second situation, the particle experiences a constant force between collisions. An example of this type of motion would be that of a charged particle moving in a uniform electric field. In the third case, the force on the particle depends on its position relative to the other particles. Here the motion is often very difficult, if not impossible, to describe analytically, due to the coupled nature of the particles' motions.

7.2 Molecular Dynamics Using Simple Models

The first molecular dynamics simulation of a condensed phase system was performed by Alder and Wainwright in 1957 using a hard-sphere model [Alder and Wainwright 1957]. In this model, the spheres move at constant velocity in straight lines between collisions. All collisions are perfectly elastic and occur when the separation between the centres of

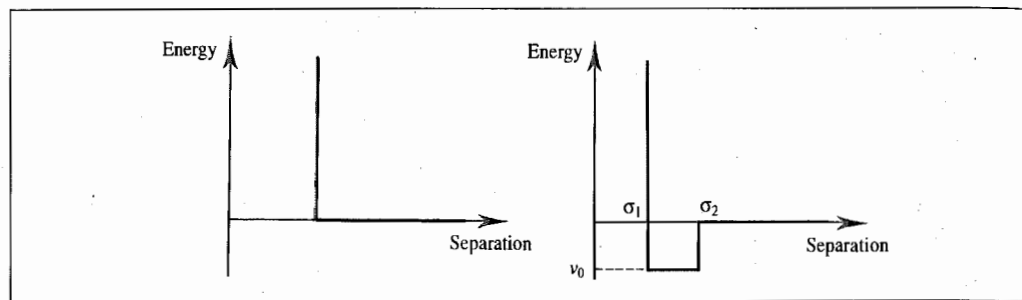


Fig. 7.1: The hard-sphere and square-well potentials.

the spheres equals the sphere diameter. The pair potential thus has the form shown in Figure 7.1. Some early simulations also used the square-well potential, where the interaction energy between two particles is zero beyond a cutoff distance σ_2 ; infinite below a smaller cutoff distance σ_1 ; and equal to ν_0 between the two cutoff values (Figure 7.1). The steps involved in the hard-sphere calculation are as follows:

1. Identify the next pair of spheres to collide and calculate when the collision will occur.
2. Calculate the positions of all the spheres at the collision time.
3. Determine the new velocities of the two colliding spheres after the collision.
4. Repeat from 1 until finished.

The new velocities of the colliding spheres are calculated by applying the principle of conservation of linear momentum.

Simple interaction models such as the hard-sphere potential obviously suffer from many deficiencies but have nevertheless provided many useful insights into the microscopic

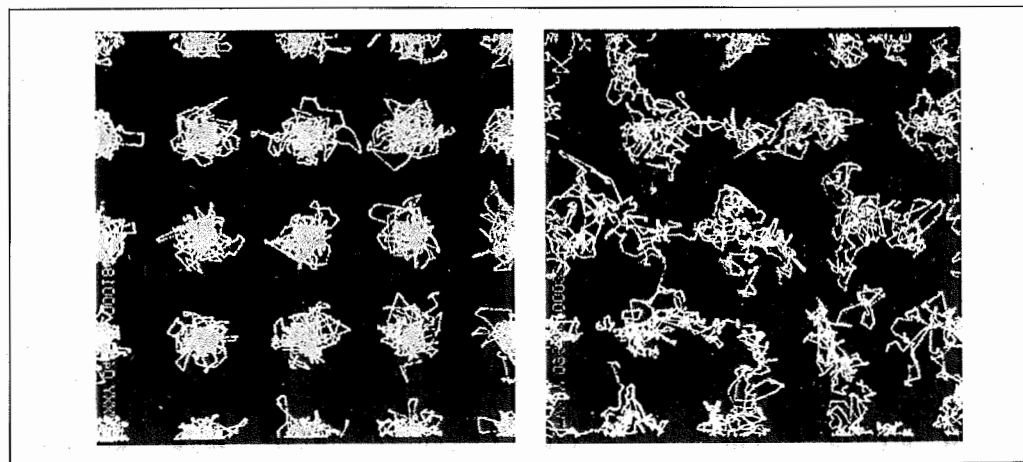


Fig. 7.2: Molecular graphics representation of the paths generated by 32 hard spherical particles in the solid (left) and fluid (right) phase. (Reproduced from Alder B J and T E Wainwright 1959. *Studies in Molecular Dynamics. I. General Method*. Journal of Chemical Physics 31: 459–466.)

nature of fluids. The early workers were particularly keen to quantify the differences between the solid and fluid phases; it is interesting to note that such investigations were facilitated by early molecular graphics systems, which enabled the trajectories of the particles to be represented simultaneously (Figure 7.2).

7.3 Molecular Dynamics with Continuous Potentials

In more realistic models of intermolecular interactions, the force on each particle will change whenever the particle changes its position, or whenever any of the other particles with which it interacts changes position. The first simulation using continuous potentials was of argon by Rahman [Rahman 1964], who also performed the first simulation of a molecular liquid (water) [Rahman and Stillinger 1971]) and made many other important methodological contributions in molecular dynamics. Under the influence of a continuous potential the motions of all the particles are coupled together, giving rise to a many-body problem that cannot be solved analytically. Under such circumstances the equations of motion are integrated using a *finite difference method*.

7.3.1 Finite Difference Methods

Finite difference techniques are used to generate molecular dynamics trajectories with continuous potential models, which we will assume to be pairwise additive. The essential idea is that the integration is broken down into many small stages, each separated in time by a fixed time δt . The total force on each particle in the configuration at a time t is calculated as the vector sum of its interactions with other particles. From the force we can determine the accelerations of the particles, which are then combined with the positions and velocities at a time t to calculate the positions and velocities at a time $t + \delta t$. The force is assumed to be constant during the time step. The forces on the particles in their new positions are then determined, leading to new positions and velocities at time $t + 2\delta t$, and so on.

There are many algorithms for integrating the equations of motion using finite difference methods, several of which are commonly used in molecular dynamics calculations. All algorithms assume that the positions and dynamic properties (velocities, accelerations, etc.) can be approximated as Taylor series expansions:

$$\mathbf{r}(t + \delta t) = \mathbf{r}(t) + \delta t \mathbf{v}(t) + \frac{1}{2} \delta t^2 \mathbf{a}(t) + \frac{1}{6} \delta t^3 \mathbf{b}(t) + \frac{1}{24} \delta t^4 \mathbf{c}(t) + \dots \quad (7.2)$$

$$\mathbf{v}(t + \delta t) = \mathbf{v}(t) + \delta t \mathbf{a}(t) + \frac{1}{2} \delta t^2 \mathbf{b}(t) + \frac{1}{6} \delta t^3 \mathbf{c}(t) + \dots \quad (7.3)$$

$$\mathbf{a}(t + \delta t) = \mathbf{a}(t) + \delta t \mathbf{b}(t) + \frac{1}{2} \delta t^2 \mathbf{c}(t) \dots \quad (7.4)$$

$$\mathbf{b}(t + \delta t) = \mathbf{b}(t) + \delta t \mathbf{c}(t) + \dots \quad (7.5)$$

where \mathbf{v} is the velocity (the first derivative of the positions with respect to time), \mathbf{a} is the acceleration (the second derivative), \mathbf{b} is the third derivative, and so on. The *Verlet algorithm* [Verlet 1967] is probably the most widely used method for integrating the equations of

motion in a molecular dynamics simulation. The Verlet algorithm uses the positions and accelerations at time t , and the positions from the previous step, $\mathbf{r}(t - \delta t)$, to calculate the new positions at $t + \delta t$, $\mathbf{r}(t + \delta t)$. We can write down the following relationships between these quantities and the velocities at time t :

$$\mathbf{r}(t + \delta t) = \mathbf{r}(t) + \delta t \mathbf{v}(t) + \frac{1}{2} \delta t^2 \mathbf{a}(t) + \dots \quad (7.6)$$

$$\mathbf{r}(t - \delta t) = \mathbf{r}(t) - \delta t \mathbf{v}(t) + \frac{1}{2} \delta t^2 \mathbf{a}(t) - \dots \quad (7.7)$$

Adding these two equations gives

$$\mathbf{r}(t + \delta t) = 2\mathbf{r}(t) - \mathbf{r}(t - \delta t) + \delta t^2 \mathbf{a}(t) \quad (7.8)$$

The velocities do not explicitly appear in the Verlet integration algorithm. The velocities can be calculated in a variety of ways; a simple approach is to divide the difference in positions at times $t + \delta t$ and $t - \delta t$ by $2\delta t$:

$$\mathbf{v}(t) = [\mathbf{r}(t + \delta t) - \mathbf{r}(t - \delta t)] / 2\delta t \quad (7.9)$$

Alternatively, the velocities can be estimated at the half-step, $t + \frac{1}{2}\delta t$:

$$\mathbf{v}(t + \frac{1}{2}\delta t) = [\mathbf{r}(t + \delta t) - \mathbf{r}(t)] / \delta t \quad (7.10)$$

Implementation of the Verlet algorithm is straightforward and the storage requirements are modest, comprising two sets of positions ($\mathbf{r}(t)$ and $\mathbf{r}(t - \delta t)$) and the accelerations $\mathbf{a}(t)$. One of its drawbacks is that the positions $\mathbf{r}(t + \delta t)$ are obtained by adding a small term ($\delta t^2 \mathbf{a}(t)$) to the difference of two much larger terms, $2\mathbf{r}(t)$ and $\mathbf{r}(t - \delta t)$. This may lead to a loss of precision. The Verlet algorithm has some other disadvantages. The lack of an explicit velocity term in the equations makes it difficult to obtain the velocities, and indeed the velocities are not available until the positions have been computed at the next step. In addition, it is not a self-starting algorithm; the new positions are obtained from the current positions $\mathbf{r}(t)$ and the positions from the previous time step, $\mathbf{r}(t - \delta t)$. At $t = 0$ there is obviously only one set of positions and so it is necessary to employ some other means to obtain positions at $t - \delta t$. One way to obtain $\mathbf{r}(t - \delta t)$ is to use the Taylor series, Equation (7.2), truncated after the first term. Thus, $\mathbf{r}(t - \delta t) = \mathbf{r}(t) - \delta t \mathbf{v}(t) + \frac{1}{2} \delta t^2 \mathbf{a}(t) - \dots$

Several variations on the Verlet algorithm have been developed. The *leap-frog* algorithm [Hockney 1970] uses the following relationships:

$$\mathbf{r}(t + \delta t) = \mathbf{r}(t) + \delta t \mathbf{v}(t + \frac{1}{2}\delta t) \quad (7.11)$$

$$\mathbf{v}(t + \frac{1}{2}\delta t) = \mathbf{v}(t - \frac{1}{2}\delta t) + \delta t \mathbf{a}(t) \quad (7.12)$$

To implement the leap-frog algorithm, the velocities $\mathbf{v}(t + \frac{1}{2}\delta t)$ are first calculated from the velocities at time $t - \frac{1}{2}\delta t$ and the accelerations at time t . The positions $\mathbf{r}(t + \delta t)$ are then deduced from the velocities just calculated together with the positions at time $\mathbf{r}(t)$ using Equation (7.11). The velocities at time t can be calculated from

$$\mathbf{v}(t) = \frac{1}{2} [\mathbf{v}(t + \frac{1}{2}\delta t) + \mathbf{v}(t - \frac{1}{2}\delta t)] \quad (7.13)$$

The velocities thus 'leap-frog' over the positions to give their values at $t + \frac{1}{2}\delta t$ (hence the name). The positions then leap over the velocities to give their new values at $t + \delta t$, ready for the velocities at $t + \frac{3}{2}\delta t$, and so on. The leap-frog method has two advantages over the

standard Verlet algorithm: it explicitly includes the velocity and also does not require the calculation of the differences of large numbers. However, it has the obvious disadvantage that the positions and velocities are not synchronised. This means that it is not possible to calculate the kinetic energy contribution to the total energy at the same time as the positions are defined (from which the potential energy is determined).

The *velocity Verlet* method [Swope *et al.* 1982] gives positions, velocities and accelerations at the same time and does not compromise precision:

$$\mathbf{r}(t + \delta t) = \mathbf{r}(t) + \delta t \mathbf{v}(t) + \frac{1}{2} \delta t^2 \mathbf{a}(t) \quad (7.14)$$

$$\mathbf{v}(t + \delta t) = \mathbf{v}(t) + \frac{1}{2} \delta t [\mathbf{a}(t) + \mathbf{a}(t + \delta t)] \quad (7.15)$$

The velocity Verlet method is actually implemented as a three-stage procedure because, as can be seen from Equation (7.15), to calculate the new velocities requires the accelerations at both t and $t + \delta t$. Thus in the first step the positions at $t + \delta t$ are calculated according to Equation (7.14) using the velocities and the accelerations at time t . The velocities at time $t + \frac{1}{2}\delta t$ are then determined using:

$$\mathbf{v}(t + \frac{1}{2}\delta t) = \mathbf{v}(t) + \frac{1}{2} \delta t \mathbf{a}(t) \quad (7.16)$$

New forces are next computed from the current positions, thus giving $\mathbf{a}(t + \delta t)$. In the final step, the velocities at time $t + \delta t$ are determined using:

$$\mathbf{v}(t + \delta t) = \mathbf{v}(t + \frac{1}{2}\delta t) + \frac{1}{2} \delta t \mathbf{a}(t + \delta t) \quad (7.17)$$

Beeman's algorithm [Beeman 1976] is also related to the Verlet method:

$$\mathbf{r}(t + \delta t) = \mathbf{r}(t) + \delta t \mathbf{v}(t) + \frac{2}{3} \delta t^2 \mathbf{a}(t) - \frac{1}{6} \delta t^2 \mathbf{a}(t - \delta t) \quad (7.18)$$

$$\mathbf{v}(t + \delta t) = \mathbf{v}(t) + \frac{1}{3} \delta t \mathbf{a}(t) + \frac{5}{6} \delta t \mathbf{a}(t - \delta t) - \frac{1}{6} \delta t \mathbf{a}(t - 2\delta t) \quad (7.19)$$

The Beeman integration scheme uses a more accurate expression for the velocity. As a consequence it often gives better energy conservation, because the kinetic energy is calculated directly from the velocities. However, the expressions used are more complex than those of the Verlet algorithm and so it is computationally more expensive.

We have already encountered four different integration methods, with more to come! Why should we use one method in preference to another? What features characterise a 'good' integration method? As with any other computer algorithm, an ideal integration scheme should be fast, require minimal memory and be easy to program. However, for most molecular dynamics simulations these issues are of secondary importance; most calculations do not make significant memory demands of even a modest workstation, and the time required for the integration is usually trivial compared to the other parts of the calculation. The most demanding part of a molecular dynamics simulation is invariably the calculation of the force on each particle in the system. More important considerations are that the integration algorithm should conserve energy and momentum, be time-reversible, and should permit a long time step, δt , to be used. The size of the time step is particularly relevant to the computational demands as a simulation using a long time step will require fewer iterations to cover a given amount of phase space. A less important requirement is that the integration algorithm should give the same results as an exact, analytical trajectory

(this can be tested using simple problems for which an analytical solution can be derived). We would, in any case, expect the calculated trajectory to deviate from the exact trajectory because the computer can only store numbers to a given precision.

The *order* of an integration method is the degree to which the Taylor series expansion, Equation (7.2), is truncated: it is the lowest term that is not present in the expansion. The order may not always be apparent from the formulae used. For example, the highest-order derivative that appears in the Verlet formulae is the second, $\mathbf{a}(t)$, yet the Verlet algorithm is, in fact, a fourth-order method. This is because the third-order terms, which cancel when Equation (7.6) is added to Equation (7.7), are still implied in the expansion:

$$\mathbf{r}(t + \delta t) = \mathbf{r}(t) + \delta t \mathbf{v}(t) + \frac{1}{2} \delta t^2 \mathbf{a}(t) + \frac{1}{6} \delta t^3 \mathbf{b}(t) + \frac{1}{24} \delta t^4 \mathbf{c}(t) \quad (7.20)$$

$$\mathbf{r}(t - \delta t) = \mathbf{r}(t) - \delta t \mathbf{v}(t) + \frac{1}{2} \delta t^2 \mathbf{a}(t) - \frac{1}{6} \delta t^3 \mathbf{b}(t) + \frac{1}{24} \delta t^4 \mathbf{c}(t) \quad (7.21)$$

7.3.2 Predictor–Corrector Integration Methods

The predictor–corrector methods [Gear 1971] form a general family of integration algorithms from which one can select a scheme that is correct to a given order. These methods have three basic steps. First, new positions, velocities, accelerations and higher-order terms are predicted according to the Taylor expansion, Equations (7.2)–(7.4). In the second stage, the forces are evaluated at the new positions to give accelerations $\mathbf{a}(t + \delta t)$. These accelerations are then compared with the accelerations that are predicted from the Taylor series expansion, $\mathbf{a}^c(t + \delta t)$. The difference between the predicted and calculated accelerations is then used to ‘correct’ the positions, velocities, etc., in the correction step:

$$\Delta \mathbf{a}(t + \delta t) = \mathbf{a}^c(t + \delta t) - \mathbf{a}(t + \delta t) \quad (7.22)$$

Then

$$\mathbf{r}^c(t + \delta t) = \mathbf{r}(t + \delta t) + c_0 \Delta \mathbf{a}(t + \delta t) \quad (7.23)$$

$$\mathbf{v}^c(t + \delta t) = \mathbf{v}(t + \delta t) + c_1 \Delta \mathbf{a}(t + \delta t) \quad (7.24)$$

$$\mathbf{a}^c(t + \delta t)/2 = \mathbf{a}(t + \delta t)/2 + c_2 \Delta \mathbf{a}(t + \delta t) \quad (7.25)$$

$$\mathbf{b}^c(t + \delta t)/6 = \mathbf{b}(t + \delta t)/6 + c_3 \Delta \mathbf{a}(t + \delta t) \quad (7.26)$$

Gear has suggested ‘best’ values of the coefficients c_0, c_1, \dots . The set of coefficients to use depends upon the order of the Taylor series expansion. In Equations (7.23)–(7.26) the expansion has been truncated after the third derivative of the positions (i.e. $\mathbf{b}(t)$). The appropriate set of coefficients to use in this case is $c_0 = \frac{1}{6}$, $c_1 = \frac{5}{6}$, $c_2 = 1$ and $c_3 = \frac{1}{3}$.

The storage required for the Gear predictor–corrector algorithm is $3 \times (O + 1)N$, where O is the highest-order differential used in the Taylor series expansion and N is the number of atoms. Thus the storage required for our example is $15N$, which is rather more than for the Verlet algorithm, which uses $9N$. More importantly, the Gear algorithm requires two time-consuming force evaluations per time step, though this is not necessarily a disadvantage as it may permit a time step more than twice as long as an alternative algorithm.

There are many variants of the ‘predictor–corrector’ theme; of these, we will only mention the algorithm used by Rahman in the first molecular dynamics simulations with continuous potentials [Rahman 1964]. In this method, the first step is to predict new positions as follows:

$$\mathbf{r}(t + \delta t) = \mathbf{r}(t - \delta t) + 2\delta t \mathbf{v}(t) \quad (7.27)$$

New accelerations are calculated at these new positions in the usual way. These accelerations are then used to generate a set of new velocities, and then corrected positions:

$$\mathbf{v}(t + \delta t) = \mathbf{v}(t) + \frac{1}{2} \delta t (\mathbf{a}(t + \delta t) + \mathbf{a}(t)) \quad (7.28)$$

$$\mathbf{r}^c(t + \delta t) = \mathbf{r}(t) + \frac{1}{2} \delta t (\mathbf{v}(t) + \mathbf{v}(t + \delta t)) \quad (7.29)$$

The acceleration can then be recalculated at the new corrected positions to give new velocities. The method then iterates over the two Equations (7.28) and (7.29). Two or three passes are usually required to achieve consistency, with a force evaluation at each step. The computational demands of this scheme mean that it is now rarely used, though it does give accurate solutions of the equations of motion.

7.3.3 Which Integration Algorithm is Most Appropriate?

The wide variety of integration schemes available can make it difficult to decide which is the most appropriate one to use. Various factors may need to be taken into account when deciding which is most appropriate. Clearly, the computational effort required is a major consideration. As we have already indicated, an algorithm that is nominally more expensive (for example, because it requires more than one force evaluation per iteration) may permit a significantly longer time step to be used and so, in fact, be more cost-effective. One of the most important considerations is energy conservation; this can be calculated as the root-mean-square fluctuation and is often plotted against the time step, as shown in Figure 7.3. In Appendix 7.1 we show why energy conservation would be expected in a molecular dynamics simulation. The kinetic and potential energy components would be expected to fluctuate in equal and opposite directions; this is also shown in Figure 7.3.

As the time step increases, so the RMS energy fluctuation also increases. For the argon simulation reported in Figure 7.3, the RMS fluctuation in the total energy is approximately 0.006 kcal/mol and the RMS fluctuations in the kinetic and potential energies are approximately 2.5 kcal/mol. With a time step of 25 fs the RMS fluctuation rises to 0.04 kcal/mol and with a time step of 5 fs the value is 0.002 kcal/mol. A variation of one part in 10^4 is generally considered acceptable. The different algorithms may vary in the rate at which the error varies with the time step. For example, it has been shown that for short time steps the predictor–corrector methods may be more accurate, but for longer time steps the Verlet algorithm may be better [Fincham and Heyes 1982]. Other factors that may be important when choosing an integration algorithm include the memory required; the synchronisation of positions and velocities; whether they are self-starting (some methods require properties at $t - \delta t$, which obviously do not exist); and

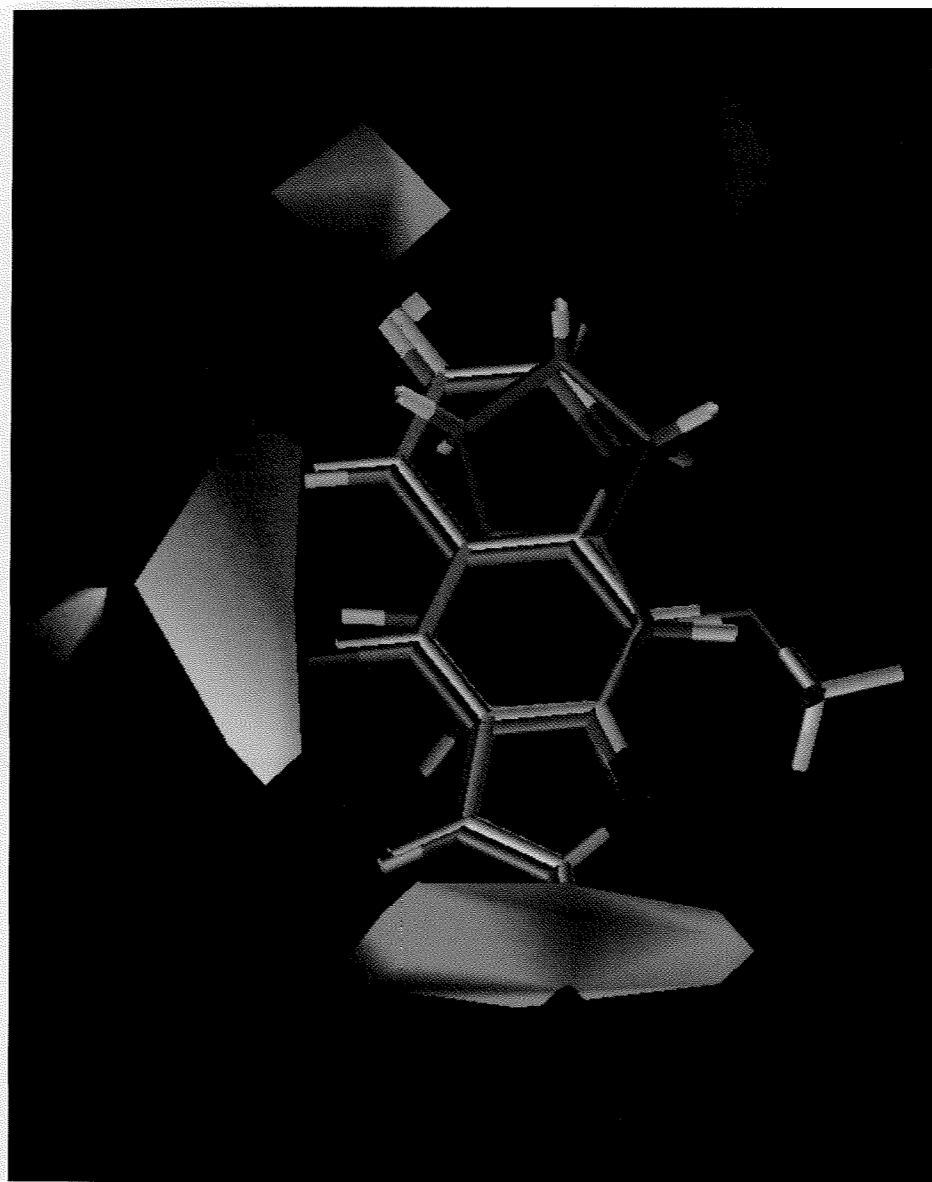


Fig. 12.41: Contour representation of key features from a CoMFA analysis of a series of coumarin substrates and inhibitors of cytochrome P₄₅₀2A5 [Poso et al. 1995]. The red and blue regions indicate positions where it would be favourable and unfavourable respectively to place a negative charge and the green/yellow regions where it would be favourable/unfavourable to locate steric bulk.

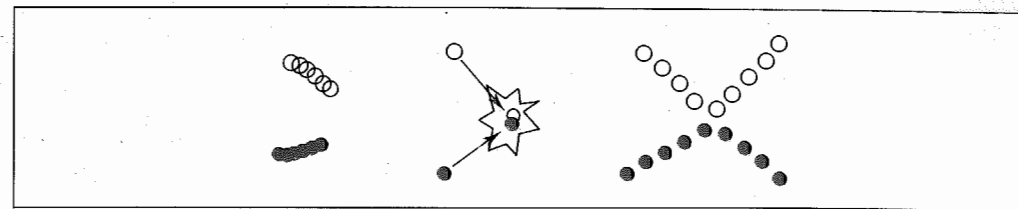


Fig. 7.4: With a very small time step (left) phase space is covered very slowly; a large time step (middle) gives instabilities. With an appropriate time step (right) phase space is covered efficiently and collisions occur smoothly.

consisting of two argon atoms interacting under the Lennard-Jones potential. The behaviour of this system can be determined analytically and so compared with the numerical integration. Suppose the argon atoms are moving towards each other along the x axis with initial velocities of 353 m s^{-1} (this corresponds to the most probable speed of argon at 300 K). We can then plot how the interatomic distance varies with time and compare it to the analytical potential. The result obtained using two time steps (10 fs and 50 fs) are shown in Figure 7.5. In both cases the numerical trajectory initially lags behind the analytical one, but then as the atoms pass through their minimum energy separation and move up the repulsive barrier the atoms 'jump through' the energy barrier. This leads to a gain in energy and the atoms then move apart with velocities that are slightly too high. In both numerical trajectories the total energy rises after the collision. Unfortunately, the atoms move most quickly and take the largest steps in the very region (i.e. near the energy minimum) where it would be best to take the smallest steps. The total error is correlated with the time step, with the largest errors arising for the largest time steps. Of course, with a small time step much more computer time will be required for a given length of calculation; the aim is to find the correct balance between simulating the 'correct' trajectory and covering the phase space. If the time

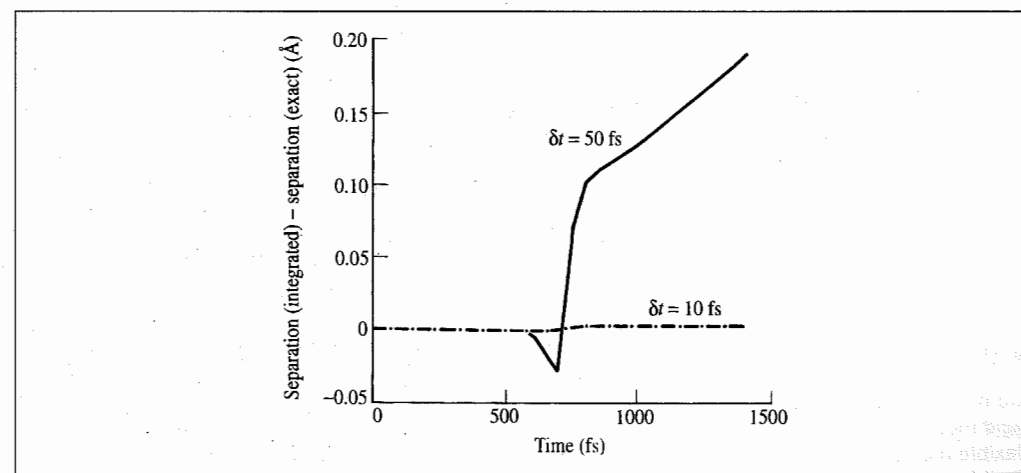


Fig. 7.5: Difference between the exact and numerical trajectories for the approach of two argon atoms with time steps of 10 fs and 50 fs.

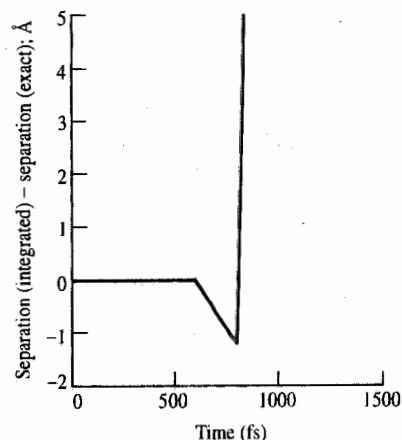


Fig. 7.6: Difference between exact and numerical trajectory for the approach of two argon atoms for a time step of 100 fs. The simulation 'blows up'.

step is too large, then the trajectory will 'blow up', as can be seen for the argon dimer system with a time step of 100 fs (Figure 7.6).

When simulating an atomic fluid the time step should be small compared to the mean time between collisions. When simulating flexible molecules a useful guide is that the time step should be approximately one-tenth the time of the shortest period of motion. In flexible molecules, the highest-frequency vibrations are due to bond stretches, especially those of bonds to hydrogen atoms. A C–H bond vibrates with a repeat period of approximately 10 fs. The timescales of some typical motions together with appropriate time steps are shown in Table 7.1, which can be used to choose an appropriate time step.

The requirement that the time step is approximately one order of magnitude smaller than the shortest motion is clearly a severe restriction, particularly as these high-frequency motions are usually of relatively little interest and have a minimal effect on the overall behaviour of the system. One solution to this problem is to 'freeze out' such vibrations by constraining the appropriate bonds to their equilibrium values while still permitting the rest of the degrees of freedom to vary under the intramolecular and intermolecular forces present. This enables a longer time step to be used. We will consider such constraint dynamics methods in Section 7.5.

System	Types of motion present	Suggested time step (s)
Atoms	Translation	10^{-14}
Rigid molecules	Translation and rotation	5×10^{-15}
Flexible molecules, rigid bonds	Translation, rotation, torsion	2×10^{-15}
Flexible molecules, flexible bonds	Translation, rotation, torsion, vibration	10^{-15} or 5×10^{-16}

Table 7.1 The different types of motion present in various systems together with suggested time steps.

7.3.5 Multiple Time Step Dynamics

Table 7.1 presents us with something of a dilemma. We would obviously desire to explore as much of the phase space as possible but this may be compromised by the need for a small time step. One possible approach is to use a multiple time step method. The underlying rationale is that certain interactions evolve more rapidly with time than other interactions. The twin-range method (Section 6.7.1) is a crude type of multiple time step approach, in that interactions involving atoms between the lower and upper cutoff distance remain constant and change only when the neighbour list is updated. However, this approach can lead to an accumulation of numerical errors in calculated properties. A more sophisticated approach is to approximate the forces due to these atoms using a Taylor series expansion [Streett *et al.* 1978]:

$$\mathbf{f}(t + \tau\delta t) = \mathbf{f}(t) + (\tau\delta t)\mathbf{d}\mathbf{f}(t)/\mathbf{d}t + \frac{1}{2}(\tau\delta t)^2\mathbf{d}^2\mathbf{f}(t)/\mathbf{d}t^2 + \dots \quad (7.30)$$

This series expansion is truncated at a specified order and is probably most easily implemented within a predictor–corrector type of algorithm, where the higher-order terms are already computed. This method has been applied to relatively simple systems such as molecular fluids [Streett *et al.* 1978] and alkane chain liquids [Swindoll and Haile 1984].

An alternative formulation of a multiple time step method is the 'reversible reference system propagation algorithm' (r-RESPA) method [Tuckerman *et al.* 1992]. In this method, the forces within a system are classified into a number of groups according to how rapidly the force varies over time. Each group then has its own time step while maintaining accuracy and numerical stability. The starting point for this algorithm is the Liouville equation, which defines how the state of the system, $\Gamma(t)$, evolves over time:

$$\Gamma(t) = e^{iLt}\Gamma(t=0) \quad (7.31)$$

The exponential $\exp(iLt)$ in Equation (7.31) involves the so-called *Liouville operator*, L , which in the case of a molecular system containing N atoms (and so $3N$ coordinates) can be expressed:

$$iL = \sum_{i=1}^{3N} \left[\frac{\partial x_i}{\partial t} \frac{\partial}{\partial x_i} + F_i(x) \frac{\partial}{\partial p_i} \right] \quad (7.32)$$

In the r-RESPA method this operator is decomposed into two or more parts, for example:

$$L = L_1 + L_2 + L_3 + L_4 \quad (7.33)$$

Each of these parts is then associated with specific terms in the force equation. For example, L_1 may correspond to the bond-stretching terms, L_2 to the angle-bending and torsional terms, L_3 to the short-range non-bonded interactions and L_4 to the long-range non-bonded interactions. Suppose the time step with which we evaluate the bond-stretching terms is δt_1 . Integers n_1 , n_2 and n_3 then define the time steps for the three other forces as follows:

$$\delta t_2 = n_1 \delta t_1; \quad \delta t_3 = n_1 n_2 \delta t_1; \quad \delta t_4 = n_1 n_2 n_3 \delta t_1 \quad (7.34)$$

The underlying theory of r-RESPA is somewhat involved, but the final result and consequent implementation is actually rather straightforward, being very closely related to the velocity Verlet integration scheme. For our four-way decomposition the algorithm would

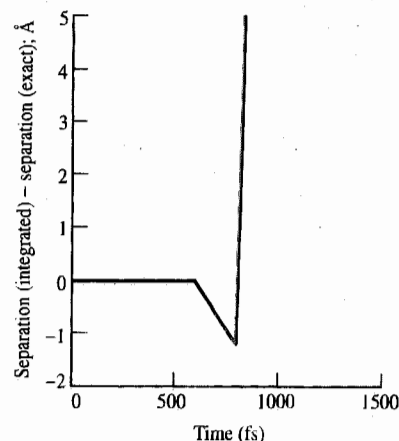


Fig. 7.6: Difference between exact and numerical trajectory for the approach of two argon atoms for a time step of 100 fs. The simulation 'blows up'.

step is too large, then the trajectory will 'blow up', as can be seen for the argon dimer system with a time step of 100 fs (Figure 7.6).

When simulating an atomic fluid the time step should be small compared to the mean time between collisions. When simulating flexible molecules a useful guide is that the time step should be approximately one-tenth the time of the shortest period of motion. In flexible molecules, the highest-frequency vibrations are due to bond stretches, especially those of bonds to hydrogen atoms. A C–H bond vibrates with a repeat period of approximately 10 fs. The timescales of some typical motions together with appropriate time steps are shown in Table 7.1, which can be used to choose an appropriate time step.

The requirement that the time step is approximately one order of magnitude smaller than the shortest motion is clearly a severe restriction, particularly as these high-frequency motions are usually of relatively little interest and have a minimal effect on the overall behaviour of the system. One solution to this problem is to 'freeze out' such vibrations by constraining the appropriate bonds to their equilibrium values while still permitting the rest of the degrees of freedom to vary under the intramolecular and intermolecular forces present. This enables a longer time step to be used. We will consider such constraint dynamics methods in Section 7.5.

System	Types of motion present	Suggested time step (s)
Atoms	Translation	10^{-14}
Rigid molecules	Translation and rotation	5×10^{-15}
Flexible molecules, rigid bonds	Translation, rotation, torsion	2×10^{-15}
Flexible molecules, flexible bonds	Translation, rotation, torsion, vibration	10^{-15} or 5×10^{-16}

Table 7.1 The different types of motion present in various systems together with suggested time steps.

7.3.5 Multiple Time Step Dynamics

Table 7.1 presents us with something of a dilemma. We would obviously desire to explore as much of the phase space as possible but this may be compromised by the need for a small time step. One possible approach is to use a multiple time step method. The underlying rationale is that certain interactions evolve more rapidly with time than other interactions. The twin-range method (Section 6.7.1) is a crude type of multiple time step approach, in that interactions involving atoms between the lower and upper cutoff distance remain constant and change only when the neighbour list is updated. However, this approach can lead to an accumulation of numerical errors in calculated properties. A more sophisticated approach is to approximate the forces due to these atoms using a Taylor series expansion [Streett *et al.* 1978]:

$$\mathbf{f}(t + \tau\delta t) = \mathbf{f}(t) + (\tau\delta t)\mathbf{df}(t)/\mathbf{dt} + \frac{1}{2}(\tau\delta t)^2\mathbf{d}^2\mathbf{f}(t)/\mathbf{dt}^2 + \dots \quad (7.30)$$

This series expansion is truncated at a specified order and is probably most easily implemented within a predictor–corrector type of algorithm, where the higher-order terms are already computed. This method has been applied to relatively simple systems such as molecular fluids [Streett *et al.* 1978] and alkane chain liquids [Swindoll and Haile 1984].

An alternative formulation of a multiple time step method is the 'reversible reference system propagation algorithm' (r-RESPA) method [Tuckerman *et al.* 1992]. In this method, the forces within a system are classified into a number of groups according to how rapidly the force varies over time. Each group then has its own time step while maintaining accuracy and numerical stability. The starting point for this algorithm is the Liouville equation, which defines how the state of the system, $\Gamma(t)$, evolves over time:

$$\Gamma(t) = e^{iLt}\Gamma \quad (t = 0) \quad (7.31)$$

The exponential $\exp(iLt)$ in Equation (7.31) involves the so-called *Liouville operator*, L , which in the case of a molecular system containing N atoms (and so $3N$ coordinates) can be expressed:

$$iL = \sum_{i=1}^{3N} \left[\frac{\partial x_i}{\partial t} \frac{\partial}{\partial x_i} + F_i(x) \frac{\partial}{\partial p_i} \right] \quad (7.32)$$

In the r-RESPA method this operator is decomposed into two or more parts, for example:

$$L = L_1 + L_2 + L_3 + L_4 \quad (7.33)$$

Each of these parts is then associated with specific terms in the force equation. For example, L_1 may correspond to the bond-stretching terms, L_2 to the angle-bending and torsional terms, L_3 to the short-range non-bonded interactions and L_4 to the long-range non-bonded interactions. Suppose the time step with which we evaluate the bond-stretching terms is δt_1 . Integers n_1 , n_2 and n_3 then define the time steps for the three other forces as follows:

$$\delta t_2 = n_1\delta t_1; \quad \delta t_3 = n_1n_2\delta t_1; \quad \delta t_4 = n_1n_2n_3\delta t_1 \quad (7.34)$$

The underlying theory of r-RESPA is somewhat involved, but the final result and consequent implementation is actually rather straightforward, being very closely related to the velocity Verlet integration scheme. For our four-way decomposition the algorithm would

be implemented as follows:

```

Calculate forces-1 (i.e.  $\mathbf{a}_1(t)$ )
Calculate forces-2 (i.e.  $\mathbf{a}_2(t)$ )
Calculate forces-3 (i.e.  $\mathbf{a}_3(t)$ )
Calculate forces-4 (i.e.  $\mathbf{a}_4(t)$ )
do step = 1,  $N_{\text{steps}}$ 
   $\mathbf{v} = \mathbf{v} + \frac{1}{2}n_1n_2n_3\delta t_1\mathbf{a}_4$ 
  do  $i_3 = 1, n_3$ 
     $\mathbf{v} = \mathbf{v} + \frac{1}{2}n_1n_2\delta t_1\mathbf{a}_3$ 
    do  $i_2 = 1, n_2$ 
       $\mathbf{v} = \mathbf{v} + \frac{1}{2}n_1\delta t_1\mathbf{a}_2$ 
      do  $i_1 = 1, n_1$ 
         $\mathbf{v} = \mathbf{v} + \frac{1}{2}\delta t_1\mathbf{a}_1$ 
         $\mathbf{r} = \mathbf{r} + \delta t_1\mathbf{v}$ 
        calculate forces-1 (i.e.  $\mathbf{a}_1$ )
         $\mathbf{v} = \mathbf{v} + \frac{1}{2}\delta t_1\mathbf{a}_1$ 
      enddo
      calculate forces-2 (i.e.  $\mathbf{a}_2$ )
       $\mathbf{v} = \mathbf{v} + \frac{1}{2}n_1\delta t_1\mathbf{a}_2$ 
    enddo
    calculate forces-3 (i.e.  $\mathbf{a}_3$ )
     $\mathbf{v} = \mathbf{v} + \frac{1}{2}n_1n_2\delta t_1\mathbf{a}_3$ 
  enddo
  calculate forces-4 (i.e.  $\mathbf{a}_4$ )
   $\mathbf{v} = \mathbf{v} + \frac{1}{2}n_1n_2n_3\delta t_1\mathbf{a}_4$ 
enddo

```

In this scheme, \mathbf{v} and \mathbf{r} refer to one of the $3N$ velocities or positions, respectively. Note that the different types of force are calculated throughout the algorithm. It can be readily seen that the method reduces to the standard velocity Verlet method if n_1 , n_2 and n_3 are set equal to 1.

The r-RESPA method has been applied to a variety of systems, including simple model systems [Tuckerman *et al.* 1992] but also organic molecules [Watanabe and Karplus 1993], fullerene crystals [Procacci and Berne 1994] and also proteins [Humphreys *et al.* 1994, 1996]. In these studies the reduction in computational time compared with the standard velocity Verlet method varied between 4–5 and 20–40, depending upon the size of the system, without any noticeable loss in accuracy. Other developments of the r-RESPA algorithm include its coupling to the fast multipole method (see Section 6.8.3) [Zhou and Berne 1995].

7.4 Setting Up and Running a Molecular Dynamics Simulation

In this section we will examine some of the steps involved in performing a molecular dynamics simulation in the microcanonical ensemble. First, it is necessary to establish an

initial configuration of the system. As discussed in Section 6.4.2, the initial configuration may be obtained from experimental data, from a theoretical model or from a combination of the two. It is also necessary to assign initial velocities to the atoms. This can be done by randomly selecting from a Maxwell-Boltzmann distribution at the temperature of interest:

$$p(v_{ix}) = \left(\frac{m_i}{2\pi k_B T} \right)^{1/2} \exp \left[-\frac{1}{2} \frac{m_i v_{ix}^2}{k_B T} \right] \quad (7.35)$$

The Maxwell-Boltzmann equation provides the probability that an atom i of mass m_i has a velocity v_{ix} in the x direction at a temperature T . A Maxwell-Boltzmann distribution is a Gaussian distribution, which can be obtained using a random number generator. Most random number generators are designed to produce random numbers that are uniform in the range 0 to 1. However, it is relatively straightforward to convert such a random number generator to sample from a Gaussian distribution (or indeed from one of several other distributions [Rubinstein 1981]). The probability of generating a value from a Gaussian (normal) distribution with mean $\langle x \rangle$ and variance σ^2 ($\sigma^2 = \langle (x - \langle x \rangle)^2 \rangle$) is:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(x - \langle x \rangle)^2}{2\sigma^2} \right] \quad (7.36)$$

One option is to first generate two random numbers ξ_1 and ξ_2 between 0 and 1. The corresponding two numbers from the normal distribution are then calculated using

$$x_1 = \sqrt{-2 \ln \xi_1} \cos(2\pi\xi_2) \quad \text{and} \quad x_2 = \sqrt{-2 \ln \xi_1} \sin(\pi\xi_2) \quad (7.37)$$

An alternative approach is to generate twelve random numbers ξ_1, \dots, ξ_{12} and then calculate:

$$x = \sum_{i=1}^{12} \xi_i - 6 \quad (7.38)$$

These two methods generate random numbers in the normal distribution with zero mean and unit variance. A number (x) generated from this distribution can be related to its counterpart (x') from another Gaussian distribution with mean $\langle x' \rangle$ and variance σ using

$$x' = \langle x' \rangle + \sigma x \quad (7.39)$$

The initial velocities may also be chosen from a uniform distribution or from a simple Gaussian distribution. In either case the Maxwell-Boltzmann distribution of velocities is usually rapidly achieved.

The initial velocities are often adjusted so that the total momentum of the system is zero. Such a system then samples from the constant NVEP ensemble. To set the total linear momentum of the system to zero, the sum of the components of the atomic momenta along the x , y and z axes is calculated. This gives the total momentum of the system in each direction, which, when divided by the total mass, is subtracted from the atomic velocities to give an overall momentum of zero.

Having set up the system and assigned the initial velocities, the simulation proper can commence. At each step the force on each atom must be calculated by differentiating the potential function. The force on an atom may include contributions from the various

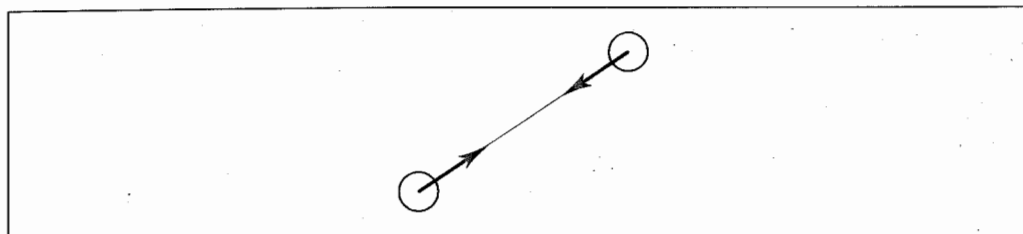


Fig. 7.7: The force between two particles acts along the line joining their centres of mass, in accordance with Newton's third law.

terms in the force field such as bonds, angles, torsional terms and non-bonded interactions. The force is straightforward to calculate for two atoms interacting under the Lennard-Jones potential:

$$\mathbf{f}_{ij} = \frac{\mathbf{r}_{ij}}{|\mathbf{r}_{ij}|} \frac{24\epsilon}{\sigma} \left[2 \left(\frac{\sigma}{r_{ij}} \right)^{13} - \left(\frac{\sigma}{r_{ij}} \right)^7 \right] \quad (7.40)$$

The force between the two atoms is equal in magnitude and opposite in direction and applies along the line connecting the two nuclear centres, in accordance with Newton's third law (Figure 7.7). It is necessary to calculate the force between each atom pair just once. This is most easily achieved by arranging to compute the force between an atom and those atoms with a higher index (i.e. for an atom i the forces are calculated with atoms $i+1, i+2, \dots, N$). Having calculated the force between an atom i and an atom with a higher index j , minus the force is added to the accumulating sum of the forces on j . The force calculation is most easily implemented using two loops, as outlined in the following pseudocode:

```

set elements on force array to zero
while atom1 = 1 to N - 1
  while atom2 = atom1 + 1 to N
    calculate force on atom1 due to interaction with atom2
    add the force to the array element, atom1
    subtract the force from the appropriate force array element atom2
  enddo
enddo

```

At the end of these two loops, the total force on each atom is known. A consequence of the fact that the force between the two atoms is equal and opposite is that the neighbour list for each atom need only contain those atoms with a higher number as the force on an atom due to interactions with lower numbered atoms will be calculated earlier in the loop. This organisation of the neighbour list contrasts with the structure used for Monte Carlo simulations, where all the neighbours of each atom (with both lower and higher indices) must be stored.

Analytical expressions for the forces due to other terms in the molecular mechanics potential function have been published for most of the functional forms encountered in common force

fields. These expressions can seem rather complicated because the intramolecular terms (e.g. bonds, angles, torsions) are calculated in terms of the internal coordinates, whereas molecular dynamics is typically performed using Cartesian coordinates. The chain rule must therefore be used to obtain the desired functional forms. However, the resulting expressions are relatively easy to implement in a computer program.

The first stage of a molecular dynamics simulation is the equilibration phase, the purpose of which is to bring the system to equilibrium from the starting configuration. During equilibration, various parameters are monitored together with the actual configurations. When these parameters achieve stable values then the production phase can commence. It is during the production phase that thermodynamic properties and other data are calculated. The parameters that are used to characterise whether equilibrium has been reached depend to some extent on the system being simulated but invariably include the kinetic, potential and total energies, the velocities, the temperature and the pressure. As we have indicated, the kinetic and potential energies would be expected to fluctuate in a simulation in the microcanonical ensemble but the total energy should remain constant. The components of the velocities should describe a Maxwell-Boltzmann distribution (in all three directions x, y and z) and the kinetic energy should be equally distributed among the three directions x, y and z . It is usually desired to perform a simulation at a specified temperature and so it is common practice to adjust the temperature of the system by scaling the velocities (see Section 7.7.1) during the equilibration phase. During the production phase the temperature is a variable of the system. Order parameters may be calculated to monitor changes in structure, which can supplement visual examination of the evolving trajectory.

When simulating an inhomogeneous system a more detailed equilibration procedure is usually desirable. A typical procedure suitable for a molecular dynamics simulation of a macromolecular solute, such as a protein in solution, would be as follows. First, the solvent alone together with any mobile counterions is subject to energy minimisation with the solute kept fixed in its initial conformation. The solvent and any counterions are then allowed to evolve using either a molecular dynamics (or indeed Monte Carlo) simulation, again keeping the structure of the solute molecule fixed. This solvent equilibration phase should be sufficiently extensive to allow the solvent to completely readjust to the potential field of the solute. For molecular dynamics this implies that the length of this solvent equilibration phase should be longer than the relaxation time of the solvent (the time taken for a molecule to lose any 'memory' of its original orientation, which for water is about 10 ps). Next, the entire system (solute and solvent) is minimised. Only then does the molecular dynamics simulation of the whole system commence.

At the start of the production phase all counters are set to zero and the system is permitted to evolve. In a microcanonical ensemble no velocity scaling is performed during the production phase and so the temperature becomes a calculated property of the system. Various properties are routinely calculated and stored during the production phase for subsequent analysis and processing. Careful monitoring of these properties during the simulation can show whether the simulation is 'well behaved' or not; it may be necessary to restart a simulation if problems are encountered. It is also usual to store the positions, energies

and velocities of configurations at regular intervals (e.g. every five to twenty time steps), from which other properties can be determined once the simulation has finished.

7.4.1 Calculating the Temperature

Many thermodynamic properties can be calculated from a molecular dynamics simulation. Most of these were dealt with in Section 6.2; here we just discuss the calculation of temperature. The instantaneous value of the temperature is related to the kinetic energy via the particles' momenta as follows:

$$\mathcal{K} = \sum_{i=1}^N \frac{|\mathbf{p}_i|^2}{2m_i} = \frac{k_B T}{2} (3N - N_c) \quad (7.41)$$

where N_c is the number of constraints and so $3N - N_c$ is the total number of degrees of freedom. For an isolated system (i.e. for a simulation of a system *in vacuo*) the total translational momentum of the system and the total angular momentum are conserved and can be made equal to zero by an appropriate choice of initial velocities. For a simulation performed using periodic boundary conditions, the total linear momentum is conserved but the total angular momentum is not. It is common practice to choose a set of initial velocities that ensures that the total linear momentum and the total angular momentum are zero. As the system evolves, the linear momentum should remain zero but the angular momentum will not. Molecular dynamics with periodic boundary conditions thus strictly samples from the constant NVEP ensemble where \mathbf{P} is the total linear momentum. This differs trivially from the standard microcanonical ensemble but it should be remembered that the appropriate number of degrees of freedom must be subtracted from the total when calculating the kinetic energy per degree of freedom. Specifically, for a system *in vacuo* where the total linear and angular momenta have been set to zero, six degrees of freedom need to be subtracted. For a simulation using periodic boundary conditions three degrees of freedom need to be subtracted if the centre-of-mass motion of the system is removed. In constraint dynamics, discussed in the next section, rather more degrees of freedom may be fixed and N_c must be calculated accordingly.

7.5 Constraint Dynamics

The earliest molecular dynamics simulations using 'realistic' potentials were of atoms interacting under the Lennard-Jones potential. In such calculations the only forces on the atoms are those due to non-bonded interactions. It is rather more difficult to simulate molecules because the interaction between two non-spherical molecules depends upon their relative orientation as well as the distance between them. If the molecules are flexible then there will also be intramolecular interactions, which give rise to changes in conformation. Clearly, the simplest model is to treat the species present as rigid bodies with no intramolecular conformational freedom. In such cases the dynamics of each molecule can often be considered in terms of translations of its centre of mass and rotations about its centre of mass. The force on the molecule equals the vector sum of all the forces acting at the

centre of mass, and the rotational motion is determined by the torque about the centre of mass. To deal with these rotational motions is considerably more complicated than for the translational motions, but in favourable cases they can be programmed quite efficiently.

When the simulation involves conformationally flexible molecules then the motion is inevitably described in terms of the atomic Cartesian coordinates. The conformational behaviour of a flexible molecule is usually a complex superposition of different motions. The high frequency motions (e.g. bond vibrations) are usually of less interest than the lower frequency modes, which often correspond to major conformational changes. Unfortunately, the time step of a molecular dynamics simulation is dictated by the highest frequency motion present in the system. It would therefore be of considerable benefit to be able to increase the time step without prejudicing the accuracy of the simulation. Constraint dynamics enables individual internal coordinates or combinations of specified coordinates to be constrained, or 'fixed' during the simulation without affecting the other internal degrees of freedom.

Before we consider in detail the use of constraint dynamics, it is helpful to establish the difference between *constraints* and *restraints*; we shall discuss the method of restrained molecular dynamics in a later chapter (see Section 9.10). A constraint is a requirement that the system is forced to satisfy. As we shall see, in constraint dynamics bonds or angles are forced to adopt specific values throughout a simulation. When a bond or angle is restrained then it is able to deviate from the desired value; the restraint only acts to 'encourage' a particular value. Restraints are most easily incorporated using additional terms in the force field which impose a penalty for deviations from the reference value. An additional difference is that restrained degrees of freedom still have an energy $k_B T/2$ associated with them, whereas constrained degrees of freedom do not.

The most commonly used method for applying constraints, particularly in molecular dynamics, is the SHAKE procedure of Ryckaert, Ciccotti and Berendsen [Ryckaert *et al.* 1977]. In constraint dynamics the equations of motion are solved while simultaneously satisfying the imposed constraints. Constrained systems have been much studied in classical mechanics; we shall illustrate the general principles using a simple system comprising a box sliding down a frictionless slope in two dimensions (Figure 7.8). The box is constrained to remain on the slope and so the box's x and y coordinates must always satisfy the equation of the slope (which we shall write as $y = mx + c$). If the slope were not present then the box

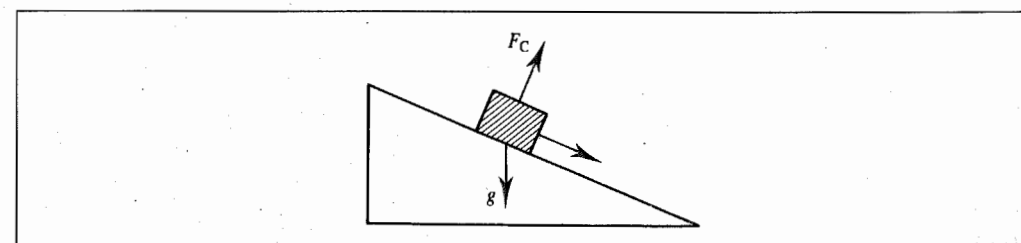


Fig. 7.8: A box sliding down a slope under the influence of gravity is subject to the constraint that it must remain on the slope. The constraint force F_c acts perpendicular to the direction of motion.

would fall vertically downwards. Constraints are often categorised as *holonomic* or *non-holonomic*. Holonomic constraints can be expressed in the form

$$f(q_1, q_2, q_3, \dots, t) = 0 \quad (7.42)$$

q_1, q_2, \dots , are the coordinates of the particles. Non-holonomic constraints cannot be expressed in this way. For example, the motion of a particle constrained to lie on the surface of a sphere is subject to a holonomic constraint, but if the particle is able to fall off the sphere under the influence of gravity then the constraint becomes non-holonomic. A holonomic constraint that keeps a particle on the surface of a sphere can be written:

$$r^2 - a^2 = 0 \quad (7.43)$$

r is the distance of the particle from the origin where the sphere of radius a is centred. The equivalent non-holonomic constraint is written as an inequality:

$$r^2 - a^2 \geq 0 \quad (7.44)$$

SHAKE uses holonomic constraints. In a constrained system the coordinates of the particles are not independent and the equations of motion in each of the coordinate directions are connected. A second difficulty is that the magnitude of the constraint forces is unknown. Thus in the case of the box on the slope, the gravitational force acting on the box is in the y direction whereas the motion is down the slope. The motion is thus not in the same direction as the gravitational force. As such, the total force on the box can be considered to arise from two sources: one due to gravity and the other a constraint force that is perpendicular to the motion of the box (Figure 7.8). As there is no motion perpendicular to the surface of the slope, the constraint force does no work.

As we know, the motion of a system of N particles can be described in terms of $3N$ independent coordinates or degrees of freedom. If there are k holonomic constraints then the number of degrees of freedom is reduced to $3N - k$. It is possible, in principle at least, to find $3N - k$ independent coordinates (the *generalised coordinates*), which can then be used to solve the problem directly. For example, the motion of the box can be described using the single coordinate, q , along the direction of the slope. The component of the gravitational force that acts along the slope is $Mg \sin \theta$ and so the acceleration down the slope is $g \sin \theta$. The position at any time t can thus be obtained by integrating the following equation of motion:

$$\frac{d^2 q}{dt^2} = g \sin \theta \quad (7.45)$$

The solution to this equation is:

$$q(t) = q(0) + t\dot{q}(0) + \frac{t^2}{2}g \sin \theta \quad (7.46)$$

where $q(0)$ is the value of q at time $t = 0$ and $\dot{q}(0)$ is the initial velocity of the box along the slope. In this simple example it is quite easy to identify the single generalised coordinate that can be used to describe the motion in the constrained system. When there are many constraints, it can be difficult to determine the generalised coordinates. In any case, it is usually desirable to work with the atomic Cartesian coordinates. The motion of the box can be more generally described in terms of the Cartesian (x, y) coordinates of the box as follows.

Newton's equations in the x and y directions are:

$$M \frac{d^2 x}{dt^2} = F_{cx} \quad (7.47)$$

$$M \frac{d^2 y}{dt^2} = -Mg + F_{cy} \quad (7.48)$$

where F_{cx} and F_{cy} are the components of the as yet unknown constraint force in the x and the y directions, respectively. We know that the constraint force acts perpendicular to the slope, and so the ratio of its x and y components must be:

$$\frac{F_{cx}}{F_{cy}} = -m \quad (7.49)$$

The constraint force can be introduced into Newton's equations as a Lagrange multiplier (see Section 1.10.5). To achieve consistency with the usual Lagrangian notation, we write F_{cy} as $-\lambda$ and so F_{cx} equals λm . Thus:

$$M \frac{d^2 x}{dt^2} = \lambda m \quad (7.50)$$

$$M \frac{d^2 y}{dt^2} = -Mg - \lambda \quad (7.51)$$

Equations (7.50) and (7.51) contain three unknowns (d^2x/dt^2 , d^2y/dt^2 and λ). A third equation that links x and y is the equation of the slope, which can be written in the following form:

$$\sigma = mx - y + c = 0 \quad (7.52)$$

This constraint equation is expressed in terms of x and y rather than their second derivatives. However, as $\sigma(x, y) = 0$ holds for all x, y , it follows that $d\sigma = 0$ and $d^2\sigma = 0$ also. Consequently, the constraint equation can be written:

$$m \frac{d^2 x}{dt^2} - \frac{d^2 y}{dt^2} = 0 \quad (7.53)$$

Solving the three equations gives:

$$\frac{d^2 x}{dt^2} = -g \frac{m}{1 + m^2} \quad (7.54)$$

$$\frac{d^2 y}{dt^2} = -g \frac{m^2}{1 + m^2} \quad (7.55)$$

The x and y coordinates at time t are thus given by:

$$x(t) = x(0) + t \frac{dx(0)}{dt} - g \frac{t^2}{2} \frac{m}{(1 + m^2)} \quad (7.56)$$

$$y(t) = y(0) + t \frac{dy(0)}{dt} - g \frac{t^2}{2} \frac{m^2}{(1 + m^2)} \quad (7.57)$$

In the general case, the equations of motion for a constrained system involve two types of force: the 'normal' forces arising from the intra- and intermolecular interactions, and the forces due to the constraints. We are particularly interested in the case where the constraint σ_k requires the bond between atoms i and j to remain fixed. The constraint influences the Cartesian coordinates of atoms i and j . The force due to this constraint can be written as follows:

$$F_{c_{kx}} = \lambda_k \frac{\partial \sigma_k}{\partial x} \quad (7.58)$$

where λ_k is the Lagrange multiplier and x represents one of the Cartesian coordinates of the two atoms. Applying Equation (7.58) to the above example, we would write $F_{cx} = \lambda \partial \sigma / \partial x = \lambda m$ and $F_{cy} = \lambda \partial \sigma / \partial y = -\lambda$. If an atom is involved in a number of constraints (because it is involved in more than one constrained bond) then the total constraint force equals the sum of all such terms. The nature of the constraint for a bond between atoms i and j is:

$$\sigma_{ij} = (\mathbf{r}_i - \mathbf{r}_j)^2 - d_{ij}^2 = 0 \quad (7.59)$$

The constraint force lies along the bond at all times. For each constrained bond, there is an equal and opposite force on the two atoms that comprise the bond. The overall effect is that the constraint forces do no work. Suppose the constraint k corresponds to the bond length between atoms i and j . The constraint forces are obtained by differentiating the constraint with respect to the coordinates of atoms i and j and multiplying by an (as yet) undetermined multiplier:

$$\partial \sigma_k / \partial \mathbf{r}_i = 2(\mathbf{r}_i - \mathbf{r}_j) \quad \text{so} \quad F_{ci} = \lambda(\mathbf{r}_i - \mathbf{r}_j) \quad (7.60)$$

$$\partial \sigma_k / \partial \mathbf{r}_j = -2(\mathbf{r}_i - \mathbf{r}_j) \quad \text{and} \quad F_{cj} = -\lambda(\mathbf{r}_i - \mathbf{r}_j) \quad (7.61)$$

The factor of 2 that arises when we differentiate the square term has been incorporated into the Lagrange multiplier λ . The above expression for the forces can be incorporated into the Verlet algorithm as follows:

$$\mathbf{r}_i(t + \delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \delta t) + \frac{\delta t^2}{m_i} \mathbf{F}_i(t) + \sum_k \frac{\lambda_k \delta t^2}{m_i} \mathbf{r}_{ij}(t) \quad (7.62)$$

Recall that the positions that would be obtained from the Verlet algorithm without constraints are $\mathbf{r}'_i(t + \delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \delta t) + \delta t^2 \mathbf{F}_i(t) / m_i$. The summation in Equation (7.62) is over all constraints k that affect atom i . These constraints perturb the positions that would otherwise have been obtained from the integration algorithm, and so the above expression can be written:

$$\mathbf{r}_i(t + \delta t) = \mathbf{r}'_i(t + \delta t) + \sum_k \frac{\lambda_k \delta t^2}{m_i} \mathbf{r}_{ij}(t) \quad (7.63)$$

The next problem is to determine the multipliers λ_k that enable all the constraints to be satisfied simultaneously. This can be done algebraically in simple cases. Suppose we wish to fix the bond in a diatomic molecule. There is just one constraint in this case, and if the

atoms are labelled 1 and 2 we can write:

$$\mathbf{r}_1(t + \delta t) = \mathbf{r}'_1(t + \delta t) + \lambda_{12}(\delta t^2 / m_1)(\mathbf{r}_1(t) - \mathbf{r}_2(t)) \quad (7.64)$$

$$\mathbf{r}_2(t + \delta t) = \mathbf{r}'_2(t + \delta t) - \lambda_{12}(\delta t^2 / m_2)(\mathbf{r}_1(t) - \mathbf{r}_2(t)) \quad (7.65)$$

A third equation is derived from the requirement that the new positions keep the bond at the required distance:

$$|\mathbf{r}_1(t + \delta t) - \mathbf{r}_2(t + \delta t)|^2 = |\mathbf{r}_1(t) - \mathbf{r}_2(t)|^2 = d_{12}^2 \quad (7.66)$$

We now have three equations and three unknowns ($\mathbf{r}_1(t + \delta t)$, $\mathbf{r}_2(t + \delta t)$ and λ_{12}). Subtracting, and putting $\mathbf{r}_{12}(t) = (\mathbf{r}_1(t) - \mathbf{r}_2(t))$ and $\mathbf{r}'_{12}(t + \delta t) = \mathbf{r}'_1(t + \delta t) - \mathbf{r}'_2(t + \delta t)$ gives:

$$\mathbf{r}_1(t + \delta t) - \mathbf{r}_2(t + \delta t) = \mathbf{r}'_{12}(t + \delta t) + \lambda_{12} \delta t^2 (1/m_1 + 1/m_2) \mathbf{r}_{12}(t) \quad (7.67)$$

Squaring both sides and imposing the constraint gives:

$$\mathbf{r}'_{12}(t + \delta t)^2 + 2\lambda_{12} \delta t^2 (1/m_1 + 1/m_2) \mathbf{r}_{12}(t) + \lambda_{12}^2 \delta t^4 (1/m_1 + 1/m_2)^2 \mathbf{r}_{12}(t)^2 = d_{12}^2 \quad (7.68)$$

Solving this quadratic equation for λ_{12} enables the new positions $\mathbf{r}_1(t + \delta t)$ and $\mathbf{r}_2(t + \delta t)$ to be determined.

In the case of a triatomic molecule with two bonds (between atoms 1,2 and 2,3), two constraint equations are obtained:

$$\mathbf{r}_{12}(t + \delta t) = \mathbf{r}'_{12}(t + \delta t) + \delta t^2 (1/m_1 + 1/m_2) \lambda_{12} \mathbf{r}_{12}(t) - (\delta t^2 / m_2) \lambda_{23} \mathbf{r}_{23}(t) \quad (7.69)$$

$$\mathbf{r}_{23}(t + \delta t) = \mathbf{r}'_{23}(t + \delta t) + \delta t^2 (1/m_2 + 1/m_3) \lambda_{23} \mathbf{r}_{23}(t) - (\delta t^2 / m_2) \lambda_{12} \mathbf{r}_{12}(t) \quad (7.70)$$

These expressions could be solved algebraically but even in this simple case the algebra becomes rather complicated. A solution can be obtained by ignoring the terms that are quadratic in λ as this produces equations which are linear in the Lagrange multipliers λ . When there are many constraints, the problem is equivalent to inverting a $k \times k$ matrix, even when the quadratic terms are ignored. The SHAKE method uses an alternative approach in which each constraint is considered in turn and solved. Satisfying one constraint may cause another constraint to be violated, and so it is necessary to iterate around the constraints until they are all satisfied to within some tolerance. The tolerance should be tight enough to ensure that the fluctuations in the simulation due to the SHAKE algorithm are much smaller than the fluctuations due to other sources, such as the use of cutoffs. Another important requirement is that the constrained degrees of freedom should be only weakly coupled to the remaining degrees of freedom, so that the motion of the molecule is not affected by the application of the constraints. The sampling of unconstrained degrees of freedom should also be unaffected. For example, if the bond lengths and angles are constrained in butane then the only degree of freedom remaining is the torsion angle. It is important that this torsion is able to explore its entire range of values in a way that is not biased because of the SHAKE procedure.

Our discussion so far has considered the use of SHAKE with the Verlet algorithm. Versions have also been derived for other integration schemes, such as the leap-frog algorithm, the predictor-corrector methods and the velocity Verlet algorithm. In the case of the velocity Verlet algorithm, the method has been named RATTLE [Anderson 1983].

When velocities appear in the integration algorithm these must be corrected as well as the positions.

Angle constraints can be easily accommodated in the SHAKE scheme by recognising that an angle constraint simply corresponds to an additional distance constraint. The angle in a triatomic molecule could thus be maintained at the desired value by requiring the distance between the two end atoms to adopt the appropriate value. This is how some small molecules (e.g. water) are maintained in a rigid geometry. For example, the simple point-charge (SPC) model of water uses three distance constraints. However, it is generally accepted that to constrain the bond angles in simulations of conformationally flexible molecules can have a deleterious effect on the efficiency with which the system explores configurational space. This is because many conformational transitions involve some opening or closing of angles as well as rotations about bonds. The most common use of SHAKE is for constraining bonds involving hydrogen atoms due to their much higher vibrational frequencies. This can enable the time step in a molecular dynamics simulation to be increased (e.g. from 1 fs to 2 fs).

The SHAKE method has been extended by Tobias and Brooks [Tobias and Brooks 1988] to enable constraints to be applied to arbitrary internal coordinates. This enables the torsion angle of a rotatable bond to be constrained to a particular value during a molecular dynamics simulation, which is particularly useful when used in conjunction with methods for calculating free energies (see Section 11.7).

7.6 Time-dependent Properties

Molecular dynamics generates configurations of the system that are connected in time and so an MD simulation can be used to calculate time-dependent properties. This is a major advantage of molecular dynamics over the Monte Carlo method. Time-dependent properties are often calculated as *time correlation coefficients*.

7.6.1 Correlation Functions

Suppose we have two sets of data values, x and y , and we wish to determine what correlation (if any) exists between them. For example, imagine that we are performing a simulation of a fluid in a capillary, and that we wish to determine the correlation between the absolute velocity of an atom and its distance from the wall of the tube. One way to do this would be to plot the sets of values as a graph. A correlation function (also known as a correlation coefficient) provides a numerical value that encapsulates the data and quantifies the strength of the correlation. A series of simulations with different capillary diameters could then be compared by examining the correlation coefficients. A variety of correlation functions can be defined, a commonly used one being:

$$C_{xy} = \frac{1}{M} \sum_{i=1}^M x_i y_i \equiv \langle x_i y_i \rangle \quad (7.71)$$

We have assumed that there are M values of x_i and y_i in the data sets. This correlation function can be normalised to a value between -1 and $+1$ by dividing by the root-mean-square values of x and y :

$$c_{xy} = \frac{\frac{1}{M} \sum_{i=1}^M x_i y_i}{\sqrt{\left(\frac{1}{M} \sum_{i=1}^M x_i^2\right) \left(\frac{1}{M} \sum_{i=1}^M y_i^2\right)}} = \frac{\langle x_i y_i \rangle}{\sqrt{\langle x_i^2 \rangle \langle y_i^2 \rangle}} \quad (7.72)$$

A value of 0 indicates no correlation and an absolute value of 1 indicates a high degree of correlation (which may be positive or negative). We will use a lowercase c to indicate a normalised correlation coefficient.

Sometimes the quantities x and y will fluctuate about non-zero mean values $\langle x \rangle$ and $\langle y \rangle$. Under such circumstances it is typical to consider just the fluctuating part and to define the correlation function as:

$$c_{xy} = \frac{\frac{1}{M} \sum_{i=1}^M (x_i - \langle x \rangle)(y_i - \langle y \rangle)}{\sqrt{\left(\frac{1}{M} \sum_{i=1}^M (x_i - \langle x \rangle)^2\right) \left(\frac{1}{M} \sum_{i=1}^M (y_i - \langle y \rangle)^2\right)}} = \frac{\langle (x_i - \langle x \rangle)(y_i - \langle y \rangle) \rangle}{\sqrt{\langle (x_i - \langle x \rangle)^2 \rangle \langle (y_i - \langle y \rangle)^2 \rangle}} \quad (7.73)$$

c_{xy} can also be written in the following useful way:

$$c_{xy} = \frac{\sum_{i=1}^M x_i y_i - \frac{1}{M} \left(\sum_{i=1}^M x_i\right) \left(\sum_{i=1}^M y_i\right)}{\sqrt{\left[\sum_{i=1}^M x_i^2 - \frac{1}{M} \left(\sum_{i=1}^M x_i\right)^2\right] \left[\sum_{i=1}^M y_i^2 - \frac{1}{M} \left(\sum_{i=1}^M y_i\right)^2\right]}} \quad (7.74)$$

Equation (7.74) does not require the mean values $\langle x \rangle$ and $\langle y \rangle$ to be determined before the correlation coefficient can be calculated and so values can be accumulated as the simulation proceeds.

A molecular dynamics simulation provides data values at specific times. This enables the value of some property at some instant to be correlated with the value of the same or another property at a later time t . The resulting values are known as *time correlation coefficients*. The correlation function is then written:

$$C_{xy}(t) = \langle x(t)y(0) \rangle \quad (7.75)$$

The following two results are useful:

$$\lim_{t \rightarrow 0} C_{xy}(t) = \langle xy \rangle \quad (7.76)$$

$$\lim_{t \rightarrow \infty} C_{xy}(t) = \langle x \rangle \langle y \rangle \quad (7.77)$$

If the quantities x and y are different, then the correlation function is sometimes referred to as a *cross-correlation function*. When x and y are the same then the function is usually called an *autocorrelation function*. An autocorrelation function indicates the extent to which the system retains a 'memory' of its previous values (or, conversely, how long it takes the system to 'lose' its memory). A simple example is the velocity autocorrelation coefficient whose value indicates how closely the velocity at a time t is correlated with the velocity at time 0. Some correlation functions can be averaged over all the particles in the system (as can the velocity autocorrelation function) whereas other functions are a property of the entire system (e.g. the dipole moment of the sample). The value of the velocity autocorrelation coefficient can be calculated by averaging over the N atoms in the simulation:

$$C_{vv}(t) = \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i(t) \cdot \mathbf{v}_i(0) \quad (7.78)$$

To normalise the function, we divide by $\langle \mathbf{v}_i(0) \cdot \mathbf{v}_i(0) \rangle$:

$$c_{vv}(t) = \frac{1}{N} \sum_{i=1}^N \frac{\langle \mathbf{v}_i(t) \cdot \mathbf{v}_i(0) \rangle}{\langle \mathbf{v}_i(0) \cdot \mathbf{v}_i(0) \rangle} \quad (7.79)$$

In general, an autocorrelation function such as the velocity autocorrelation coefficient has an initial value of 1 and at long times has the value 0. The time taken to lose the correlation is often called the *correlation time*, or the *relaxation time*. If the duration of the simulation is significantly longer than the relaxation time (as it should be), then many sets of data can be extracted from the simulation to calculate the correlation function and to reduce the uncertainty in the calculation. If P steps of molecular dynamics are required for complete relaxation, and the simulation has been run for a total of Q steps, then $(Q - P)$ different sets of values could be used to calculate a value for the correlation function. The first set would run from step 1 to step N ; the second from step 2 to step $N + 1$, and so on (Figure 7.9). In fact, as we saw in Section 6.9 the high degree of correlation between successive time steps means that it is common to use time origins that are separated by several time steps, as shown in Figure 7.9. If we use M time origins (t_j) then the velocity autocorrelation function is given by:

$$C_{vv}(t) = \frac{1}{MN} \sum_{j=1}^M \sum_{i=1}^N \mathbf{v}_i(t_j) \cdot \mathbf{v}_i(t_j + t) \quad (7.80)$$

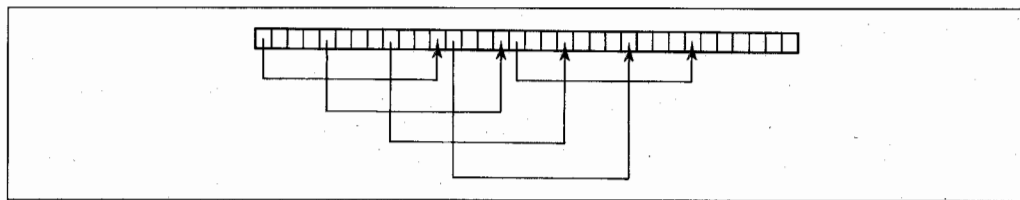


Fig. 7.9: The use of different time origins improves the accuracy with which time correlation functions can be calculated.

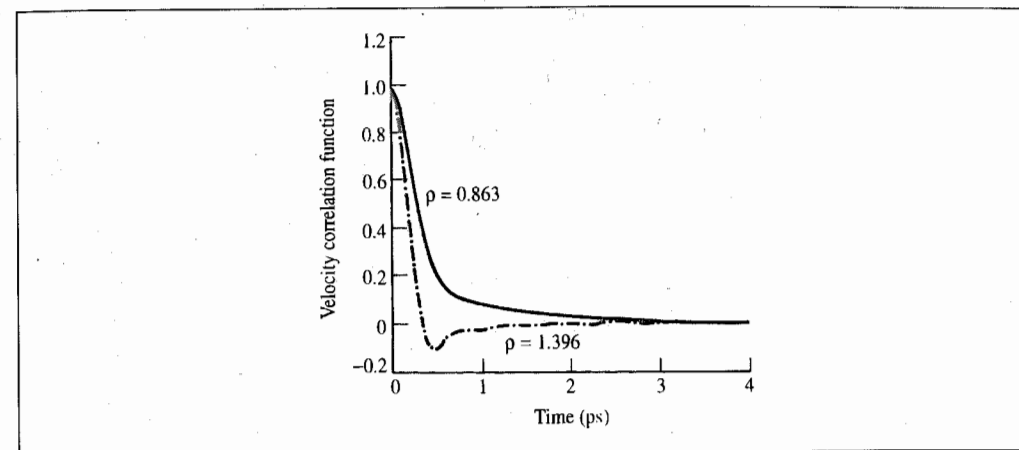


Fig. 7.10: Velocity autocorrelation functions for liquid argon at densities of 1.396 g cm^{-3} and 0.863 g cm^{-3} .

Quantities with small relaxation times can thus be determined with greater statistical precision, as it will be possible to include a greater number of data sets from a given simulation. Moreover, no quantity with a relaxation time greater than the length of the simulation can be determined accurately.

The velocity autocorrelation functions obtained from molecular dynamics simulations of argon at two different densities are shown in Figure 7.10. The correlation coefficient has an initial value of 1 and is then quadratic at short times, a result that is predicted theoretically. The subsequent behaviour depends upon the density of the fluid. For a low-density fluid, the velocity autocorrelation coefficient gradually decreases to zero. At high densities $c_{vv}(t)$ crosses the axis and then becomes negative. A negative correlation coefficient simply means that the particle is now moving in the direction opposite to that at $t = 0$. This result has been interpreted in terms of a 'cage' structure of the liquid; the atom hits the side of the cage formed by its nearest neighbours and rebounds, reversing the direction of its motion. At both low density and high density the decay towards zero is significantly slower than the exponential decay predicted by kinetic theory. In fact, the decay varies as $t^{-3/2}$. This was one of the most interesting results obtained from the early molecular dynamics simulations and can be observed even with a hard-sphere model [Alder and Wainwright 1970]. The phenomenon is ascribed to the formation of a 'hydrodynamic vortex'. As the atom moves through the fluid it pushes other atoms out of the way. These atoms circle around and eventually give it a final 'push' so resulting in a less rapid decrease to zero (Figure 7.11).

The slow decay of the velocity autocorrelation function can present practical problems when deriving other properties, such as the transport coefficients, that require the correlation function to be integrated between $t = 0$ and $t = \infty$. The so-called 'long time-tail' of the autocorrelation function makes a significant contribution to the integral, but unfortunately the statistical uncertainty with which this part of the function can be calculated is greater as fewer segments of the appropriate length can be extracted from the simulation.

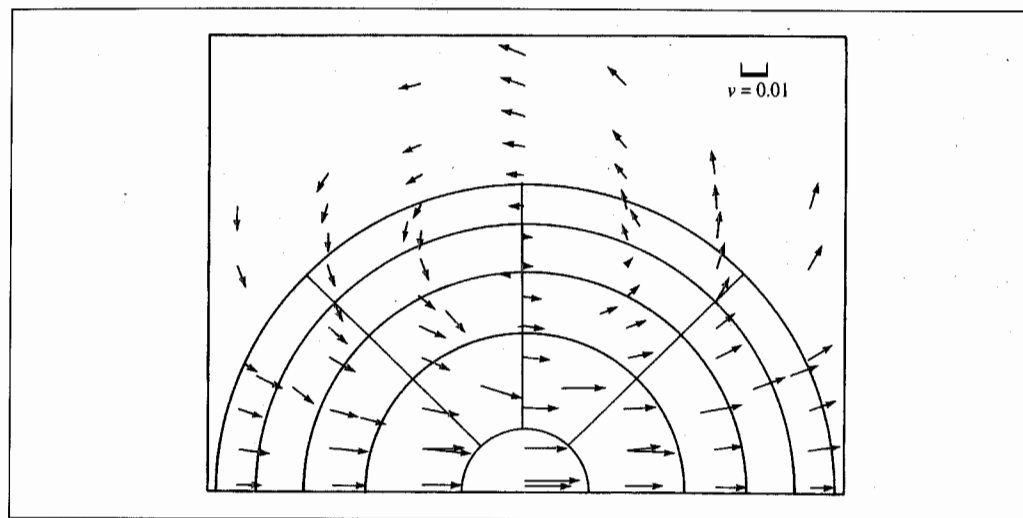


Fig. 7.11: The slow decay of the velocity autocorrelation function towards zero can be explained in terms of the formation of a hydrodynamic vortex. (Figure adapted from Alder B J and T E Wainwright 1970. *Decay of the Velocity Autocorrelation Function*. *Physical Review A* 1:18–21.)

The velocity autocorrelation function is an example of a single-particle correlation function, in which the average is calculated not only over time origins but also over all the atoms. Some properties are calculated for the entire system. One such property is the net dipole moment of the system, which is the vector sum of all the individual dipoles of the molecules in the system (clearly the dipole moment of the system can be non-zero only if each individual molecule has a dipole). The magnitude and orientation of the net dipole will change with time and is given by:

$$\boldsymbol{\mu}_{\text{tot}}(t) = \sum_{i=1}^N \boldsymbol{\mu}_i(t) \quad (7.81)$$

$\boldsymbol{\mu}_i(t)$ is the dipole moment of molecule i at time t . The total dipolar correlation function is given by:

$$c_{\text{dipole}}(t) = \frac{\langle \boldsymbol{\mu}_{\text{tot}}(t) \cdot \boldsymbol{\mu}_{\text{tot}}(0) \rangle}{\langle \boldsymbol{\mu}_{\text{tot}}(0) \cdot \boldsymbol{\mu}_{\text{tot}}(0) \rangle} \quad (7.82)$$

The dipole correlation time of the system is a valuable quantity to calculate as it is related to the sample's absorption spectrum. Liquids usually absorb in the infrared region of the electromagnetic spectrum, a typical spectrum being shown in Figure 7.12. As can be seen, the spectrum is very broad with none of the sharp peaks characteristic of a well-resolved spectrum for a species in the gas phase. This is because the overall dipole of a liquid does not change at a constant rate but, rather, there is a distribution of frequencies. The intensity of absorption at any frequency depends upon the relative contribution of that frequency to the overall distribution. If, on average, the overall dipole changes very rapidly (i.e. the relaxation time is short) then the maximum in the absorption spectrum will occur at a

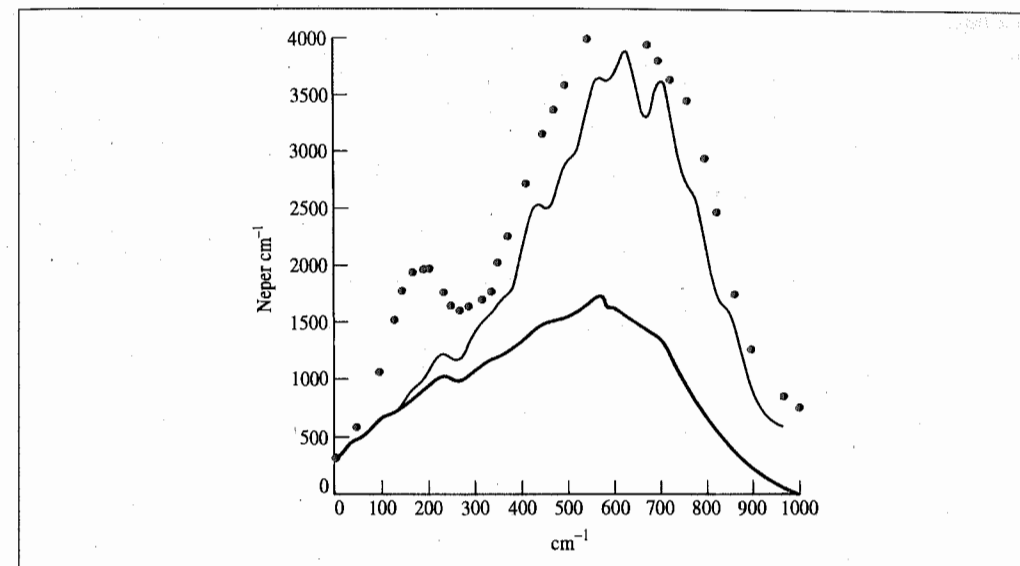


Fig. 7.12: Experimental and calculated infrared spectra for liquid water. The black dots are the experimental values. The thick curve is the classical profile produced by the molecular dynamics simulation. The thin curve is obtained by applying quantum corrections. (Figure redrawn from Guillot B 1991. *A Molecular Dynamics Study of the Infrared Spectrum of Water*. *Journal of Chemical Physics* 95:1543–1551.)

higher frequency than if the relaxation time is long. To predict the spectrum from the correlation function it is therefore necessary to extract the relative contribution of each dipole fluctuation. This is done using Fourier analysis techniques, in which the correlation function is transformed from the time domain into the frequency domain (an introduction to Fourier analysis is provided in Section 1.10.8). The Fourier analysis picks out the intensity of dipole fluctuation at each frequency using the following relationship:

$$\hat{c}_{\text{dipole}}(\nu) = \int_{-\infty}^{\infty} c_{\text{dipole}}(t) \exp(-i2\pi\nu t) dt \quad (7.83)$$

Having calculated the Fourier transform it is then possible to plot the simulated spectrum and compare it to that obtained by experiment, as in Figure 7.12.

7.6.2 Orientational Correlation Functions

Other orientational correlation coefficients can be calculated in the same way as the correlation coefficients that we have discussed already. Thus, the reorientational correlation coefficient of a single rigid molecule indicates the degree to which the orientation of a molecule at a time t is related to its orientation at time 0. The angular velocity autocorrelation function is the rotational equivalent of the velocity correlation function:

$$c_{\omega\omega}(t) = \frac{\langle \boldsymbol{\omega}_i(t) \cdot \boldsymbol{\omega}_i(0) \rangle}{\langle \boldsymbol{\omega}_i(0) \cdot \boldsymbol{\omega}_i(0) \rangle} \quad (7.84)$$

In a liquid, the rotation of a molecule is influenced by neighbouring molecules and over time the correlation will decay to zero. The information embodied in the orientational correlation functions can be compared to a variety of spectroscopic experiments, including infrared, Raman and NMR spectra. For non-spherical molecules it can be useful to derive separate auto-correlation functions for the angular velocity along each of the principal axes of rotation. For example, for a spherical molecule such as CBr_4 neighbouring molecules have a relatively small influence on the loss in correlation in the angular velocity. By contrast, a linear molecule such as CS_2 experiences significant torques as it rotates. This has the effect of damping the rotational motion more severely than for the spherical case, and indeed the correlation function can change sign, indicating that the molecule is now rotating in the opposite direction. For some molecules such as water the presence of specific interactions between molecules (for example, due to hydrogen bonding) can give rise to very rapid damping and several minima in $c_{\omega\omega}(t)$.

7.6.3 Transport Properties

Transport refers to a phenomenon that gives rise to a flow of material from one region to another. For example, if a solution is prepared with a non-equilibrium solute distribution, then the solute diffuses until the concentration is equal throughout. If a thermal gradient is created, energy flows until the temperature is equalised. A momentum gradient gives rise to viscosity. The very existence of transport implies that the system is not in equilibrium. Techniques have been developed to perform non-equilibrium molecular dynamics simulations from which transport properties can be calculated, but we will not consider these here. We are thus faced with the problem of calculating non-equilibrium properties from equilibrium simulations. This may seem an impossible task but can be achieved by considering the microscopic local fluctuations that occur even in systems at equilibrium. We should be aware, however, that non-equilibrium molecular dynamics simulations can be a more efficient way to calculate transport properties and other quantities [Allen and Tildesley 1987].

To a first approximation the rate at which transport of the relevant quantity occurs (called the *flux*) is proportional to the gradient of the property with the constant of proportionality being the relevant transport property coefficient. For example, the flux of matter J_z (i.e. the amount passing through unit area in unit time) equals the diffusion coefficient (D) multiplied by the concentration gradient; this is Fick's first law of diffusion:

$$J_z = -D(d\mathcal{N}/dz) \quad (7.85)$$

\mathcal{N} is the number density (the number of atoms per unit volume). Equation (7.85) refers to diffusion in the z direction. The minus sign indicates that flux increases in the direction of negative concentration gradient. The time dependence of diffusive behaviour (which applies if a distribution is established at some time and is then allowed to evolve) is governed by Fick's second law, which gives the rate of change of concentration with time:

$$\frac{\partial \mathcal{N}(z, t)}{\partial t} = D \frac{\partial^2 \mathcal{N}(z, t)}{\partial z^2} \quad (7.86)$$

To solve Fick's second law equation it is necessary to impose two boundary conditions for the spatial dependence and one boundary condition for the temporal dependence (the

equation is second order in space and first order in time). For example, we might require that at time zero all N_0 particles have $z = 0$. The solution to the equation is then:

$$\mathcal{N}(z, t) = \frac{N_0}{A\sqrt{\pi Dt}} \exp\left[-\frac{z^2}{4Dt}\right] \quad (7.87)$$

where A is the cross-sectional area of the sample. Equation (7.87) is a Gaussian function which initially has a sharp peak at $z = 0$ but which gradually becomes more spread out as time progresses. When the material being simulated is a pure liquid then the coefficient D is often referred to as a *self-diffusion coefficient*. The diffusion coefficient is related to the mean square distance, $\langle |\mathbf{r}(t) - \mathbf{r}(0)|^2 \rangle$, which Einstein showed was equal to $2Dt$. In three dimensions, the mean square displacement is given by:

$$3D = \lim_{t \rightarrow \infty} \frac{\langle |\mathbf{r}(t) - \mathbf{r}(0)|^2 \rangle}{2t} \quad (7.88)$$

As indicated, the relationship strictly holds only in the limit as $t \rightarrow \infty$.

The Einstein relationship can thus be used to calculate the diffusion coefficient from an equilibrium simulation, by plotting the mean square displacement as a function of time and then attempting to obtain the limiting behaviour as $t \rightarrow \infty$ (Fick's law is inapplicable at short times). The quantity $|\mathbf{r}(t) - \mathbf{r}(0)|$ can be averaged over the particles in the system to reduce the statistical error. It is also usual practice to average over time origins where possible. When using this method for calculating the diffusion coefficient the mean-squared distances should not be limited by the edges of the periodic box. In other words, we require a set of positions that have not been translated back into the central simulation cell. This can be achieved either by storing a set of 'uncorrected' positions or indeed by not correcting any of the positions during the simulation and simply generating the appropriate minimum image positions as required for the calculation of the energy or forces.

Einstein relationships hold for other transport properties, e.g. the shear viscosity, the bulk viscosity and the thermal conductivity. For example, the shear viscosity η is given by:

$$\eta_{xy} = \frac{1}{Vk_B T} \lim_{t \rightarrow \infty} \frac{\langle (\sum_{i=1}^N m\dot{x}_i(t)y_i(t) - \sum_{i=1}^N m\dot{y}_i(t)x_i(t))^2 \rangle}{2t} \quad (7.89)$$

The shear viscosity is a tensor quantity, with components η_{xy} , η_{xz} , η_{yx} , η_{yz} , η_{zx} , η_{zy} . It is a property of the whole sample rather than of individual atoms and so cannot be calculated with the same accuracy as the self-diffusion coefficient. For a homogeneous fluid the components of the shear viscosity should all be equal and so the statistical error can be reduced by averaging over the six components. An estimate of the precision of the calculation can then be determined by evaluating the standard deviation of these components from the average. Unfortunately, Equation (7.89) cannot be directly used in periodic systems, even if the positions have been unfolded, because the 'unfolded' distance between two particles may not correspond to the distance of the minimum image that is used to calculate the force. For this reason alternative approaches are required.

One alternative approach to the calculation of the diffusion and other transport coefficients is via an appropriate autocorrelation function. For example, the diffusion coefficient

depends upon the way in which the atomic position $\mathbf{r}(t)$ changes with time. At a time t the difference between $\mathbf{r}(t)$ and $\mathbf{r}(0)$ is given by:

$$|\mathbf{r}(t) - \mathbf{r}(0)| = \int_0^t \mathbf{v}(t') dt' \quad (7.90)$$

If both sides of Equation (7.90) are now squared then we obtain the following for the mean-square value:

$$\langle |\mathbf{r}(t) - \mathbf{r}(0)|^2 \rangle = \int_0^t dt' \int_0^{t'} dt'' \langle \mathbf{v}(t') \cdot \mathbf{v}(t'') \rangle \quad (7.91)$$

The crucial feature to recognise is that the relevant correlation functions are unaffected by changing the origin, which means that the following holds:

$$\langle \mathbf{v}(t') \cdot \mathbf{v}(t'') \rangle = \langle \mathbf{v}(t'' - t') \cdot \mathbf{v}(0) \rangle \quad (7.92)$$

Integration of the double integral, Equation (7.91) leads to the *Green-Kubo* formula:

$$\frac{\langle |\mathbf{r}(t) - \mathbf{r}(0)|^2 \rangle}{2t} = \int_0^t \langle \mathbf{v}(\tau) \cdot \mathbf{v}(0) \rangle \left(1 - \frac{\tau}{t}\right) d\tau \quad (7.93)$$

In the limit:

$$\int_0^\infty \langle \mathbf{v}(\tau) \cdot \mathbf{v}(0) \rangle d\tau = \lim_{t \rightarrow \infty} \frac{\langle |\mathbf{r}(t) - \mathbf{r}(0)|^2 \rangle}{2t} = 3D \quad (7.94)$$

We can now see why long time-tails in the autocorrelation functions are so important. The area under the curve during the slow decay towards zero may be a significant part of the integral in the Green-Kubo formula. In practice, these integrals are determined numerically. The long time-tail may be dealt with by fitting a function to the curve and then attempting to integrate to infinity.

7.7 Molecular Dynamics at Constant Temperature and Pressure

Molecular dynamics is traditionally performed in the constant *NVE* (or *NVEP*) ensemble. Although thermodynamic results can be transformed between ensembles, this is strictly only possible in the limit of infinite system size ('the thermodynamic limit'). It may therefore be desired to perform the simulation in a different ensemble. The two most common alternative ensembles are the constant *NVT* and the constant *NPT* ensembles. In this section we will therefore consider how molecular dynamics simulations can be performed under conditions of constant temperature and/or constant pressure.

7.7.1 Constant Temperature Dynamics

There are several reasons why we might want to maintain or otherwise control the temperature during a molecular dynamics simulation. Even in a constant *NVE* simulation it is common practice to adjust the temperature to the desired value during the equilibration phase. A constant temperature simulation may be required if we wish to determine how

the behaviour of the system changes with temperature, such as the unfolding of a protein or glass formation. Simulated annealing algorithms require the temperature of the system to be reduced gradually while the system explores its degrees of freedom. Simulated annealing is used in searching conformational space and in the elucidation of macromolecular structure from NMR and X-ray data and is discussed in Section 9.9.2.

The temperature of the system is related to the time average of the kinetic energy, which for an unconstrained system is given by:

$$\langle \mathcal{K} \rangle_{NVT} = \frac{3}{2} Nk_B T \quad (7.95)$$

An obvious way to alter the temperature of the system is thus to scale the velocities [Woodcock 1971]. If the temperature at time t is $T(t)$ and the velocities are multiplied by a factor λ , then the associated temperature change can be calculated as follows:

$$\Delta T = \frac{1}{2} \sum_{i=1}^N \frac{2}{3} \frac{m_i (\lambda v_i)^2}{Nk_B} - \frac{1}{2} \sum_{i=1}^N \frac{2}{3} \frac{m_i v_i^2}{Nk_B} \quad (7.96)$$

$$\Delta T = (\lambda^2 - 1)T(t) \quad (7.97)$$

$$\lambda = \sqrt{T_{\text{new}}/T(t)} \quad (7.98)$$

The simplest way to control the temperature is thus to multiply the velocities at each time step by the factor $\lambda = \sqrt{T_{\text{req}}/T_{\text{curr}}}$, where T_{curr} is the current temperature as calculated from the kinetic energy and T_{req} is the desired temperature.

An alternative way to maintain the temperature is to couple the system to an external heat bath that is fixed at the desired temperature [Berendsen *et al.* 1984]. The bath acts as a source of thermal energy, supplying or removing heat from the system as appropriate. The velocities are scaled at each step, such that the rate of change of temperature is proportional to the difference in temperature between the bath and the system:

$$\frac{dT(t)}{dt} = \frac{1}{\tau} (T_{\text{bath}} - T(t)) \quad (7.99)$$

τ is a coupling parameter whose magnitude determines how tightly the bath and the system are coupled together. This method gives an exponential decay of the system towards the desired temperature. The change in temperature between successive time steps is:

$$\Delta T = \frac{\delta t}{\tau} (T_{\text{bath}} - T(t)) \quad (7.100)$$

The scaling factor for the velocities is thus:

$$\lambda^2 = 1 + \frac{\delta t}{\tau} \left(\frac{T_{\text{bath}}}{T(t)} - 1 \right) \quad (7.101)$$

If τ is large, then the coupling will be weak. If τ is small, the coupling will be strong and when the coupling parameter equals the time step ($\tau = \delta t$) then the algorithm is equivalent to the simple velocity scaling method. A coupling constant of approximately 0.4 ps has been suggested as an appropriate value to use when the time step is 1 fs, giving $\delta t/\tau \approx 0.0025$. The advantage of this approach is that it does permit the system to fluctuate about the desired temperature.

These two relatively simple temperature scaling methods do not generate rigorous canonical averages. Velocity scaling artificially prolongs any temperature differences among the components of the system, which can lead to the phenomenon of 'hot solvent, cold solute', in which the 'temperature' of the solute is lower than that of the solvent, even though the overall temperature of the system is at the desired value. One 'solution' to this problem is to apply temperature coupling separately to the solute and to the solvent but the problem of unequal distribution of energy between the various components (and between the various modes of motion) may still remain. Two methods that do generate rigorous canonical ensembles if properly implemented are the *stochastic collisions* method and the *extended system* method.

In the stochastic collisions method a particle is randomly chosen at intervals and its velocity is reassigned by random selection from the Maxwell-Boltzmann distribution [Anderson 1980]. This is equivalent to the system being in contact with a heat bath that randomly emits 'thermal particles' which collide with the atoms in the system. Between each collision the system is simulated at constant energy and so the overall effect is equivalent to a series of microcanonical simulations, each performed at a slightly different energy. The distribution of the energies of these 'mini microcanonical' simulations will be a Gaussian function. The stochastic collisions method does not, of course, generate a smooth trajectory, which may be a drawback. By calculating the energy change due to a collision, Anderson showed that the mean rate (ν) at which each particle should suffer a stochastic collision is given by:

$$\nu = \frac{2a\kappa}{3k_B \mathcal{N}^{1/3} N^{2/3}} \quad (7.102)$$

a is a dimensionless constant, κ is the thermal conductivity and \mathcal{N} is the number density of the particles. If the thermal conductivity is not known then a suitable value of ν can be obtained from the intermolecular collision frequency ν_c :

$$\nu = \nu_c / N^{2/3} \quad (7.103)$$

If the collision rate is too low then the system does not sample from a canonical distribution of energies. If it is too high then the temperature control algorithm dominates and the system does not show the expected fluctuations in kinetic energy. The velocity of more than one particle can be changed in the stochastic collision method; in the limit the velocities of all the particles are changed simultaneously, though it is preferable to do this at quite long intervals. A distinction can thus be made between 'minor' collisions, in which only one (or a few) particles are affected, and 'major' (or 'massive') collisions, where the velocities of all particles are changed. It is also possible to use a combined approach, with minor collisions occurring relatively frequently and major collisions at longer intervals.

Extended system methods, originally introduced for performing constant temperature molecular dynamics by Nosé [Nosé 1984] and subsequently developed by Hoover [Hoover 1985], consider the thermal reservoir to be an integral part of the system. The reservoir is represented by an additional degree of freedom, labelled s . The reservoir has potential energy $(f+1)k_B T \ln s$, where f is the number of degrees of freedom in the physical system and T is the desired temperature. The reservoir also has kinetic energy $(Q/2)(ds/dt)^2$. Q is a parameter with the dimensions of energy \times (time)² and can be considered the

(fictitious) mass of the extra degree of freedom. The magnitude of Q determines the coupling between the reservoir and the real system and so influences the temperature fluctuations.

Each state of the extended system that is generated by the molecular dynamics simulation corresponds to a unique state of the real system. There is not, however, a direct correspondence between the velocities and the time in the real and the extended systems. The velocities of the atoms in the real system are given by:

$$\mathbf{v}_i = s \frac{d\mathbf{r}_i}{dt} \quad (7.104)$$

\mathbf{r}_i is the position of particle i in the simulation and \mathbf{v}_i is considered to be the real velocity of the particle. The time step $\delta t'$ is related to the time step in 'real time' δt by

$$\delta t = s \delta t' \quad (7.105)$$

The value of the additional degree of freedom s can change and so the time step in real time can fluctuate. Thus regular time intervals in the extended system correspond to a trajectory of the real system which is unevenly spaced in time.

The parameter Q controls the energy flow between the system and the reservoir. If Q is large then the energy flow is slow; in the limit of infinite Q , conventional molecular dynamics is regained and there is no energy exchange between the reservoir and the real system. However, if Q is too small then the energy oscillates, resulting in equilibration problems. Nosé has suggested that Q should be proportional to $f k_B T$; the constant of proportionality can then be obtained by performing a series of trial simulations for a test system and observing how well the system maintains the desired temperature.

7.7.2 Constant Pressure Dynamics

Just as one may wish to specify the temperature in a molecular dynamics simulation, so it may be desired to maintain the system at a constant pressure. This enables the behaviour of the system to be explored as a function of the pressure, enabling one to study phenomena such as the onset of pressure-induced phase transitions. Many experimental measurements are made under conditions of constant temperature and pressure, and so simulations in the isothermal-isobaric ensemble are most directly relevant to experimental data. Certain structural rearrangements may be achieved more easily in an isobaric simulation than in a simulation at constant volume. Constant pressure conditions may also be important when the number of particles in the system changes (as in some of the 'test particle' methods for calculating free energies and chemical potentials; see Section 8.9).

The pressure often fluctuates much more than quantities such as the total energy in a constant NVE molecular dynamics simulation. This is as expected because the pressure is related to the virial, which is obtained as the product of the positions and the derivative of the potential energy function. This product, $r_{ij} dV(r_{ij})/dr_{ij}$, changes more quickly with r than does the internal energy, hence the greater fluctuation in the pressure.

A macroscopic system maintains constant pressure by changing its volume. A simulation in the isothermal-isobaric ensemble also maintains constant pressure by changing the volume

of the simulation cell. The amount of volume fluctuation is related to the isothermal compressibility, κ :

$$\kappa = -\frac{1}{V} \left(\frac{\partial V}{\partial P} \right)_T \quad (7.106)$$

An easily compressible substance has a larger value of κ , so larger volume fluctuations occur at a given pressure than in a more incompressible substance. Conversely, in a constant volume simulation a less compressible substance shows larger fluctuations in the pressure. The isothermal compressibility is the pressure analogue of the heat capacity, which is related to the energy fluctuations.

A volume change in an isobaric simulation can be achieved by changing the volume in all directions, or in just one direction. It is instructive to consider the range of volume changes that one might expect to observe in a constant pressure simulation of a 'typical' system. The isothermal compressibility is related to the mean square volume displacement by:

$$\kappa = \frac{1}{k_B T} \frac{\langle V^2 \rangle - \langle V \rangle^2}{\langle V^2 \rangle} \quad (7.107)$$

The isothermal compressibility of an ideal gas is approximately 1 atm^{-1} . So for a simulation in a box of side 20 \AA (volume 8000 \AA^3) at 300 K , the root mean square change in the volume is approximately 18100 \AA^3 . This is larger than the initial size of the box! For a relatively incompressible substance such as water ($\kappa = 44.75 \times 10^{-6} \text{ atm}^{-1}$) the fluctuation is 121 \AA^3 , which corresponds to the box only changing by about 0.1 \AA in each direction. These values have clear implications for the appropriate size of the simulation system.

Many of the methods used for pressure control are analogous to those used for temperature control. Thus, the pressure can be maintained at a constant value by simply scaling the volume. An alternative is to couple the system to a 'pressure bath', analogous to a temperature bath [Berendsen *et al.* 1984]. The rate of change of pressure is given by:

$$\frac{dP(t)}{dt} = \frac{1}{\tau_p} (P_{\text{bath}} - P(t)) \quad (7.108)$$

τ_p is the coupling constant, P_{bath} is the pressure of the 'bath', and $P(t)$ is the actual pressure at time t . The volume of the simulation box is scaled by a factor λ , which is equivalent to scaling the atomic coordinates by a factor $\lambda^{1/3}$. Thus:

$$\lambda = 1 - \kappa \frac{\delta t}{\tau_p} (P - P_{\text{bath}}) \quad (7.109)$$

The new positions are given by:

$$\mathbf{r}'_i = \lambda^{1/3} \mathbf{r}_i \quad (7.110)$$

The constant κ can be combined with the relaxation constant τ_p as a single constant. This expression can be applied isotropically (i.e. such that the scaling factor is equal for all three directions) or anisotropically (where the scaling factor is calculated independently for each of the three axes). In general, it is best to use the anisotropic approach as this enables the box dimensions to change independently. Unfortunately, it has not been possible to determine from which ensemble this method samples.

In the extended pressure-coupling system methods, first introduced by Anderson [Anderson 1980], an extra degree of freedom, corresponding to the volume of the box, is added to the system. The kinetic energy associated with this degree of freedom (which can be considered to be equivalent to a piston acting on the system) is $\frac{1}{2}Q(dV/dt)^2$, where Q is the 'mass' of the piston. The piston also has potential energy PV , where P is the desired pressure and V is the volume of the system. A piston of small mass gives rise to rapid oscillations in the box, whereas a large mass has the opposite effect. An infinite mass returns normal molecular dynamics behaviour. The volume can vary during the simulation, with the average volume being determined by the balance between the internal pressure of the system and the desired external pressure. The extended-system temperature-scaling method of Nosé uses a scaled time; in the extended pressure method the coordinates of the extended system are related to the 'real' coordinates by:

$$\mathbf{r}'_i = V^{-1/3} \mathbf{r}_i \quad (7.111)$$

7.8 Incorporating Solvent Effects into Molecular Dynamics: Potentials of Mean Force and Stochastic Dynamics

In many simulations of solute-solvent systems the primary focus is the behaviour of the solute; the solvent is of relatively little interest, particularly in regions far from the solute molecule. The use of non-rectangular periodic boundary conditions, stochastic boundaries and 'solvent shells' can all help to reduce the number of solvent molecules required and enable a larger proportion of the computing time to be spent simulating the solute. In this section we consider a group of techniques that incorporate the effects of solvent without requiring any explicit specific solvent molecules to be present.

One approach to this problem is to use a *potential of mean force* (PMF), which describes how the free energy changes as a particular coordinate (such as the separation of two atoms or the torsion angle of a bond) is varied. The free energy change described by the potential of mean force includes the averaged effects of the solvent.

Potentials of mean force may be determined using a molecular dynamics or Monte Carlo simulation using the techniques of umbrella sampling or free energy perturbation, which will be discussed in Chapter 11. Here we illustrate the concept using an example. The energy difference between the *trans* and *gauche* conformations for an isolated molecule of 1,2-dichloroethane (i.e. in the gas phase) is approximately $1.14 \text{ kcal mol}^{-1}$ with a population containing 77% *trans* and 23% *gauche* conformers. In liquid 1,2-dichloroethane, however, the relative population of the *gauche* conformer is significantly increased relative to the *trans* conformer by comparison with the isolated molecule, with 44% *trans* and 56% *gauche*. These experimental results were reproduced by Jorgensen (see Figure 7.13) using Monte Carlo simulations [Jorgensen *et al.* 1981]. The potential of mean force would be designed to reproduce this new population and so enable a single 1,2-dichloroethane molecule to be simulated as if it were present in the liquid.

A simulation performed using a potential of mean force enables the modulating effects of the solvent to be taken into account. The solvent also influences the dynamic behaviour of the

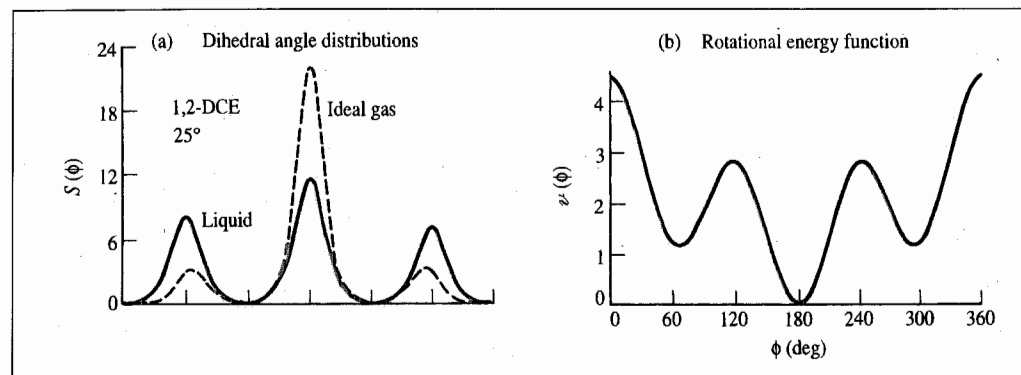


Fig. 7.13: Population distribution for 1,2-dichloroethane in the gas and liquid phases. (Figure redrawn from Jorgensen W L, R C Binning Jr and B Bigot 1981. *Structures and Properties of Organic Liquids: n-Butane and 1,2-Dichloroethane and Their Conformational Equilibria*. *Journal of the American Chemical Society* 103:4393-4399.)

solute via random collisions, and by imposing a frictional drag on the motion of the solute through the solvent. The Langevin equation of motion is the starting point for the *stochastic dynamics* models, which also incorporate these two effects. In stochastic dynamics the force on a particle is considered to arise from three sources. The first component is due to interactions between the particle and other particles. This force (F_i) depends upon the position of the particle relative to the other particles and is modelled using a potential of mean force. The second force arises from the motion of the particle through the solvent and is equivalent to the frictional drag on the particle due to the solvent. This frictional force is proportional to the speed of the particle with the constant of proportionality being the friction coefficient:

$$\mathbf{F}_{\text{frictional}} = -\xi \mathbf{v} \quad (7.112)$$

where \mathbf{v} is the velocity and ξ is the friction coefficient. The friction coefficient is related to the collision frequency (γ) by $\gamma = \xi/m$ (m is the mass of the particle). γ^{-1} can be considered as the time taken for the particle to lose memory of its initial velocity (the velocity relaxation time). For a spherical particle the friction coefficient is related to the diffusion constant D by:

$$\xi = k_B T/D \quad (7.113)$$

If the radius of the spherical particle is a then the frictional force is given by Stokes' law:

$$\mathbf{F}_{\text{frictional}} = 6\pi a \eta \mathbf{v} \quad (7.114)$$

where η is the viscosity of the fluid.

The third contribution to the force on the particle is due to random fluctuations caused by interactions with solvent molecules. We will write this force as $\mathbf{R}(t)$. The Langevin equation of motion for a particle i can therefore be written*:

$$m_i \frac{d^2 x_i(t)}{dt^2} = \mathbf{F}_i\{x_i(t)\} - \gamma_i \frac{dx_i(t)}{dt} m_i + \mathbf{R}_i(t) \quad (7.115)$$

* γ_i in Equation (7.115) is often referred to as the friction coefficient in the literature.

A number of simulation methods based on Equation (7.115) have been described. These differ in the assumptions that are made about the nature of frictional and random forces. A common simplifying assumption is that the collision frequency γ is independent of time and position. The random force $\mathbf{R}(t)$ is often assumed to be uncorrelated with the particle velocities, positions and the forces acting on them, and to obey a Gaussian distribution with zero mean. The force F_i is assumed to be constant over the time step of the integration.

Three different situations can be considered, depending upon the relative magnitudes of the integration time step and the velocity relaxation time. The first category corresponds to timescales that are short relative to the velocity relaxation time ($\gamma \delta t \ll 1$). Under such circumstances the solvent does not activate or deactivate the particle to any significant extent. In the limit of zero γ (when there are no effects due to solvent) then the Langevin Equation (7.115) reduces to that obtained from Newton's laws of motion. At the other extreme the velocity relaxation time is much smaller than the time step. This corresponds to the diffusive regime, where the motion is rapidly damped by the solvent. The third situation is intermediate between these two extremes. Various methods have been proposed for integrating the Langevin equation of motion in these three regions.

In the region where $\gamma \delta t \ll 1$ the following is a simple integration algorithm [van Gunsteren *et al.* 1981]:

$$x_{i+1} = x_i + v_i \delta t + \frac{1}{2} (\delta t)^2 \{-\gamma v_i + m^{-1}(F_i + R_i)\} \quad (7.116)$$

$$v_{i+1} = v_i + (\delta t) \{-\gamma v_i + m^{-1}(F_i + R_i)\} \quad (7.117)$$

The average random force over the time step is taken from a Gaussian with a variance $2mk_B T \gamma (\delta t)^{-1}$. x_i is one of the $3N$ coordinates at time step i ; F_i and R_i are the relevant components of the frictional and random forces at that time; v_i is the velocity component.

An alternative expression is based on the following finite difference approximations [Brunger *et al.* 1984]:

$$d^2 x/dt^2 \approx (x_{i+1} - 2x_i + x_{i-1})/\delta t^2 \quad (7.118)$$

$$dx/dt \approx (x_{i+1} - x_{i-1})/2\delta t \quad (7.119)$$

This leads to the following expressions for the coordinates x_{i+1} :

$$x_{i+1} = x_i + (x_i - x_{i-1}) \frac{1 - \frac{1}{2}\gamma \delta t}{1 + \frac{1}{2}\gamma \delta t} + \left(\frac{\delta t^2}{m} \right) \frac{F_i + R_i}{1 + \frac{1}{2}\gamma \delta t} \quad (7.120)$$

In the region where $\gamma \delta t \gg 1$ then if the interparticle force is assumed to be constant over the integration time step the following result is obtained [van Gunsteren *et al.* 1981]:

$$x_{i+1} = x_i + F_i (m\gamma)^{-1} \delta t + X_i(\delta t) \quad (7.121)$$

where X_i is a Gaussian distribution with zero mean and a variance of $2k_B T (m\gamma)^{-1} = 2D\delta t$. An extension of this treatment is to permit force F_i to vary linearly over the time step, giving:

$$x_{i+1} = x_i + \frac{\delta t}{m\gamma} (F_i + \frac{1}{2} \dot{F}_i \delta t) + X_i \quad (7.122)$$

\dot{F}_i is the derivative of the force at the time step i and is obtained numerically:

$$\dot{F}_i = (F_i - F_{i-1})/\delta t \quad (7.123)$$

In the intermediate region, where there are no restrictions on $\gamma\delta t$, then integration of the equations of motion gives the following rather complicated result [van Gunsteren and Berendsen 1982]:

$$x_{i+1} = x_i + v_i\gamma^{-1}(1 - \exp(-\gamma\delta t)) + F_i(m\gamma)^{-1}[\delta t - \gamma^{-1}(1 - \exp(-\gamma\delta t))] + (m\gamma)^{-1} \int_{t_i}^{t_{i+1}} [1 - \exp(-\gamma(t_{i+1} - t'))]R(t') dt' \quad (7.124)$$

$$v_{i+1} = v_i \exp(-\gamma\delta t) + F_i(m\gamma)^{-1}(1 - \exp(-\gamma\delta t)) + (m)^{-1} \int_{t_i}^{t_{i+1}} \exp(-\gamma(t_{i+1} - t'))R(t') dt' \quad (7.125)$$

The important feature of these two equations is that the new positions and the new velocities both depend upon an integral over the random force, $R(t)$ (the final terms in Equations (7.124) and (7.125)). As both of these integrals depend upon $R_i(t)$ they are correlated. Specifically, they obey a *bivariate* Gaussian distribution. Such a distribution provides the probability that a particle located at x_i at time t with velocity v_i and experiencing a force F_i will be at x_{i+1} at time $t + \delta t$ with velocity v_{i+1} . In practice, this means that the distribution for the second variable depends upon the value selected for the first variable. It can be difficult to properly sample from such distributions, but van Gunsteren and Berendsen showed that the equations can be reformulated in terms of sampling from two independent Gaussian functions.

More complex stochastic dynamics treatments are possible; our treatment has only provided a rather simple treatment of solvent effects. For example, we have assumed that the frictional force at a given instant is proportional only to its velocity at the same time. A more realistic model assumes that the frictional forces are correlated; they have a 'memory' of previous values. The friction coefficient can also be made to depend on the coordinates of the other particles.

7.8.1 Practical Aspects of Stochastic Dynamics Simulations

A stochastic dynamics simulation requires a value to be assigned to the collision frequency friction coefficient γ . For simple particles such as spheres this can be related to the diffusion constant in the fluid. For the simulation of a rigid molecule it may be possible to derive γ via the diffusion coefficient from a standard molecular dynamics situation. In the more general case we require the friction coefficient of each atom. For simple molecules such as butane the friction coefficient can be considered to be the same for all atoms. The optimal value for γ can be determined by trial and error, performing a stochastic dynamics simulation for different values of γ and comparing the results with those from experiment (where available) or from standard molecular dynamics simulations. For large molecules the atomic friction coefficient is considered to depend upon the degree to which each atom is in contact with the solvent and is usually taken to be proportional to the accessible surface area of the atom (as defined in Section 1.5).

One of the main advantages of the stochastic dynamics methods is that dramatic time savings can be achieved, which enables much longer stimulations to be performed. For example, Widmalm and Pastor performed 1 ns molecular dynamics and stochastic dynamics simulations of an ethylene glycol molecule in aqueous solution of the solute and 259 water molecules [Widmalm and Pastor 1992]. The molecular dynamics simulation required 300 hours whereas the stochastic dynamics simulation of the solute alone required just 24 minutes. The dramatic reduction in time for the stochastic dynamics calculation is due not only to the very much smaller number of molecules present but also to the fact that longer time steps can often be used in stochastic dynamics simulations.

Stochastic dynamics has been widely used to study the behaviour of long-chain molecules and polymers. The advantages of stochastic dynamics are especially important for polymers [Helfand 1984], where many interesting phenomena occur over relatively long time periods, so putting them beyond the scope of conventional molecular dynamics. However, one must take care with the Langevin method when simulating systems in which specific solute-solvent interactions are present. For example, Yun-Yu, Lu and van Gunsteren used both stochastic dynamics and molecular dynamics to study the immunosuppressant drug cyclosporin (Figure 7.14) in two solvents: carbon tetrachloride and water [Yun-Yu *et al.* 1988]. The time-averaged structures obtained from each method were compared to determine the similarity between the average structure obtained for each simulation. Fluctuations in torsion angles were also compared. The analysis showed that the structures obtained from the molecular dynamics and stochastic dynamics simulations of cyclosporin in carbon tetrachloride were very similar, but that the results were very different for the Langevin and molecular dynamics simulations performed in water. This was due to an excessive degree of internal hydrogen bonding in the stochastic dynamics simulation; the equivalent

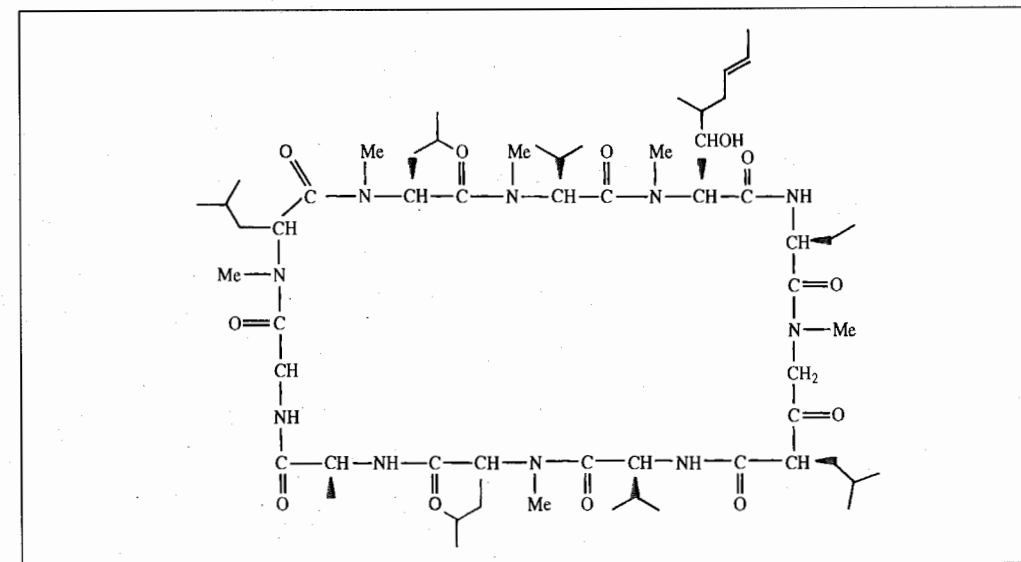


Fig. 7.14: Cyclosporin.

molecular dynamics simulation contained much more hydrogen bonding between cyclosporin and the solvent.

7.9 Conformational Changes from Molecular Dynamics Simulations

Molecular dynamics can provide information about the conformational properties of molecular systems and the way in which the conformation changes with time. Molecular graphics programs can facilitate the analysis of such simulations by displaying the structural parameters of interest in a manner that enables the time dimension to be taken into account. Perhaps the most direct way to demonstrate the conformational behaviour of the system is as a movie, where coordinate sets saved at regular intervals are displayed in sequence. For publication purposes, time-dependent data can be displayed graphically, with one of the axes corresponding to the time, such as the plots of energy or autocorrelation function versus time (Figures 7.3 and 7.10). The representation of bond rotations is difficult using x/y plots due to the 2π periodicity of a torsion angle. Lavery and Sklenar have developed a method to represent torsion data as a polar plot [Lavery and Sklenar 1988], where the distance from the origin corresponds to the time (Figure 7.15). Such 'dials' are very useful for detecting the presence of correlated conformational changes.

When viewing a movie of a molecular dynamics simulation of a complex molecule one is often struck by the chaotic nature of the motion. This should be expected; the motion of complex molecules is chaotic, but there are often underlying low-frequency motions which correspond to more significant and more interesting conformational changes. Fourier analysis techniques can be used to filter out the unwanted high-frequency motions, enabling the important low-frequency changes to be observed unhindered. Here we describe the filtering method of Dauber-Osguthorpe and Osguthorpe [Dauber-Osguthorpe and Osguthorpe 1990, 1993].

A Fourier transform enables one to convert the variation of some quantity as a function of time into a function of frequency, and vice versa. Thus, if we represent the quantity that varies in time as $x(t)$, then Fourier analysis enables us to also represent that quantity as a function $X(\nu)$, where ν is the frequency ($-\infty < \nu < \infty$). Fourier analysis is usually introduced by considering functions that vary in a periodic manner with time which can be written as a superposition of sine and cosine functions (a Fourier series; see Section 1.10.8). If the period of the function $x(t)$ is τ then the cosine and sine terms in the Fourier series are functions of frequencies $2\pi n/\tau$, where n can take integer values 1, 2, 3, ...

A Fourier series is rarely relevant to the interpretation of a molecular dynamics simulation as the movement of the atoms is not periodic but chaotic. The Fourier transform enables a non-periodic function to be converted into the equivalent frequency function (and vice versa). The Fourier transform can be developed from the Fourier series simply by considering the effect of increasing the period of a periodic function to infinity. The frequency function obtained from a Fourier transform is a continuous function rather than one written as a series of discrete frequencies. Further details are provided in Section 1.10.8.

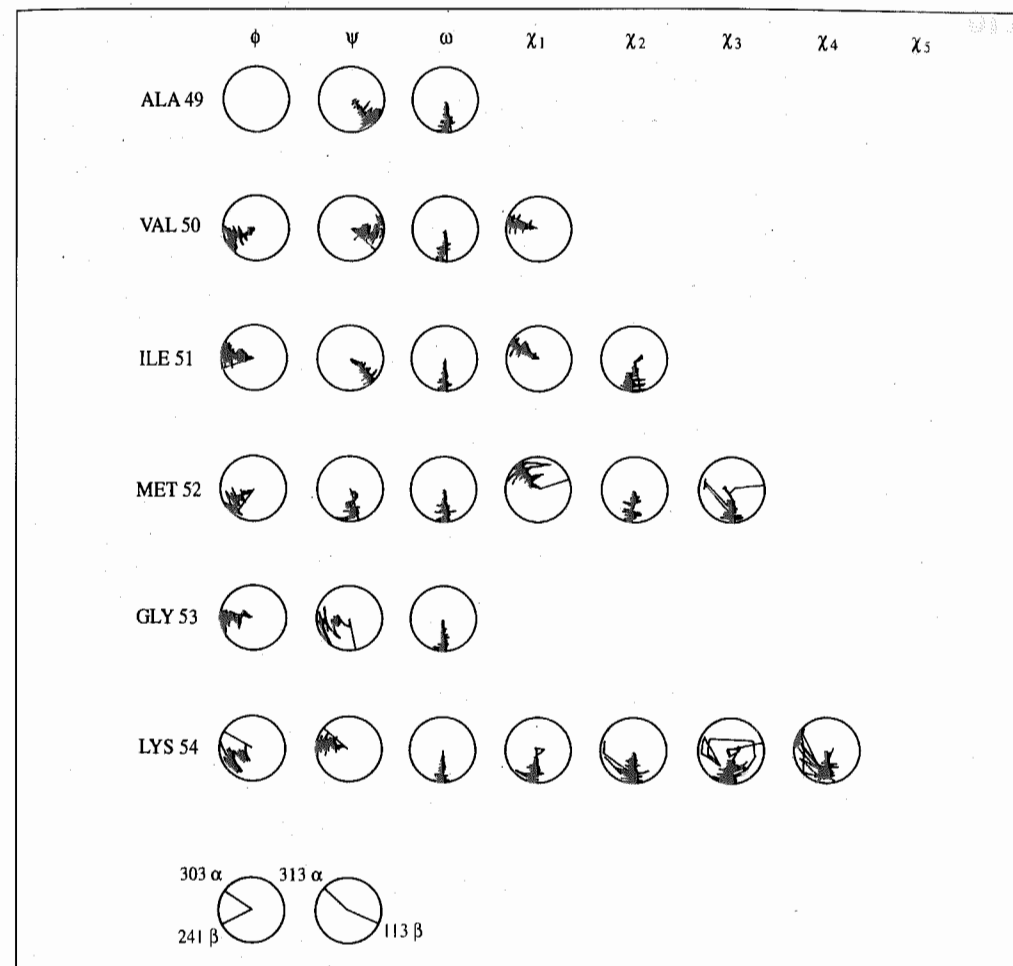


Fig. 7.15: The variation in torsion angles can be effectively represented as a series of 'dials', where the time corresponds to the distance from the centre of the dial. Data from a molecular dynamics simulation of an intermolecular complex between the enzyme dihydrofolate reductase and a triazine inhibitor [Leach and Klein 1995].

At each step of the Fourier analysis of a molecular dynamics simulation the variation with time of one of the Cartesian coordinates of one of the atoms in the system is converted into the corresponding frequency function. Fast Fourier techniques are usually employed for this step. The frequency spectrum can then be filtered to remove high frequencies. This is achieved simply by setting the coefficients of the unwanted frequencies in the frequency function to zero. The resulting spectrum is then converted back to the time domain to give a new set of coordinate values at each of the time steps in the trajectory. This new coordinate set includes only the selected frequencies. This process can be repeated for the three coordinates of each atom to give a filtered trajectory for the entire system. It is also possible to select just a single frequency (i.e. a single normal mode) from the frequency spectrum and view this in isolation.

7.10 Molecular Dynamics Simulations of Chain Amphiphiles

The molecular dynamics technique is widely used for simulating large molecular systems, some of which have many degrees of conformational freedom. In this section we will examine the application of molecular dynamics to chain amphiphiles, a class of molecules of interest to both the 'biological' and 'materials science' communities. These molecules have a polar head group attached to one or more hydrocarbon chains. Some examples are shown in Figure 7.16. The head group has a high affinity for water, whereas the hydrocarbon tail prefers to exist in a hydrophobic environment. The molecules therefore exist in both phases at a water/oil interface. A characteristic feature of these molecules is their ability

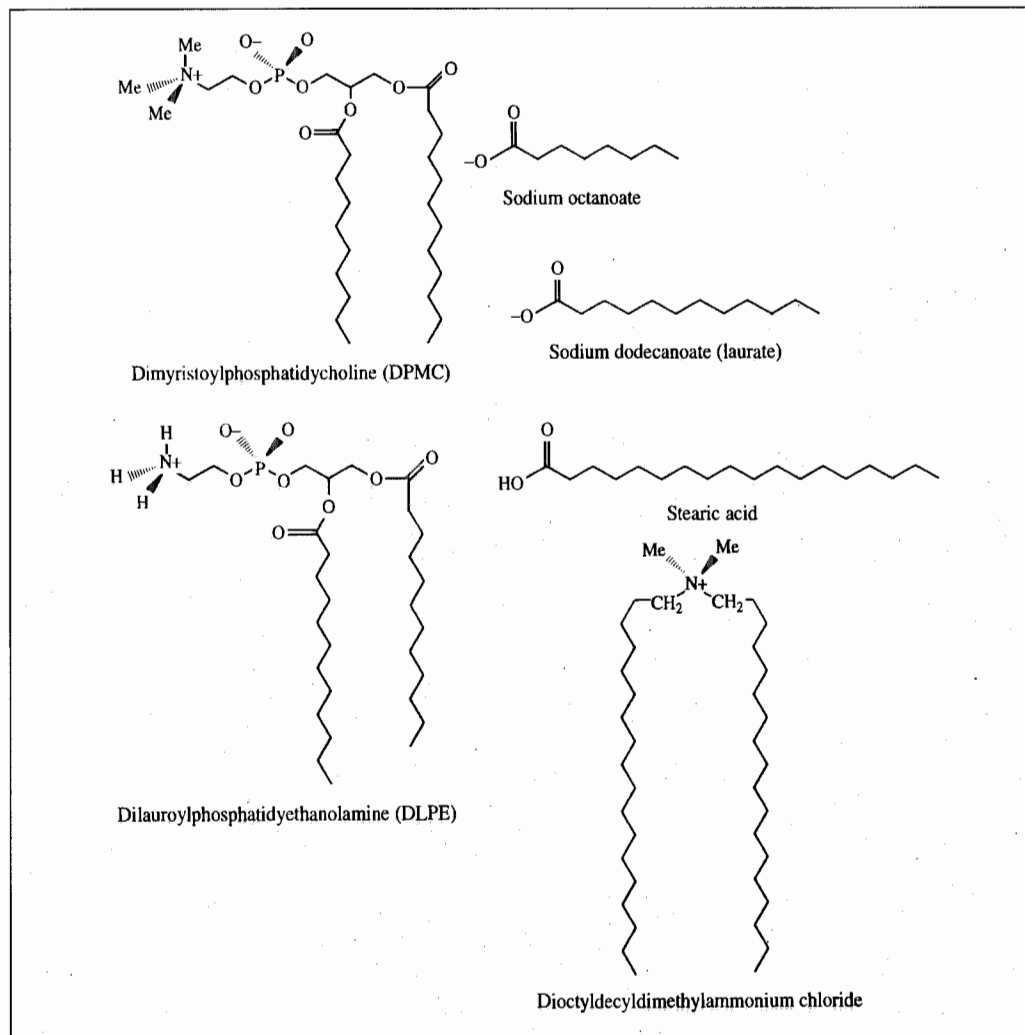


Fig. 7.16: Some typical amphiphiles.

to form extended layer structures. Monolayers, bilayers and multiple layers are all possible. A monolayer at the water/air interface is known as a Langmuir film; when this is transferred to a solid substrate it is known as a Langmuir-Blodgett film. Langmuir-Blodgett films with many layers can be constructed in the laboratory but most simulation studies of these systems have been restricted to monolayers or bilayers. The ability to control the thickness of a Langmuir-Blodgett film and their high degree of order means that they are intensively investigated as insulators in semiconductors, filtration devices and as anti-reflective coatings. Amphiphiles are important in biology as cell membranes are formed from lipid bilayers. At a high enough concentration some amphiphiles can form micelles, which are globular structures that have the head groups all pointing into solution and the tails inside (Figure 7.17).

Amphiphiles often have a complex phase behaviour with several liquid crystalline phases. These liquid crystalline phases are often characterised by long-range order in one direction together with the formation of a layer structure. The molecules may nevertheless be able to move laterally within the layer and perpendicular to the surface of the layer. Structural information can be obtained using spectroscopic techniques including X-ray and neutron diffraction and NMR. The quadrupolar splitting in the deuterium NMR spectrum can be

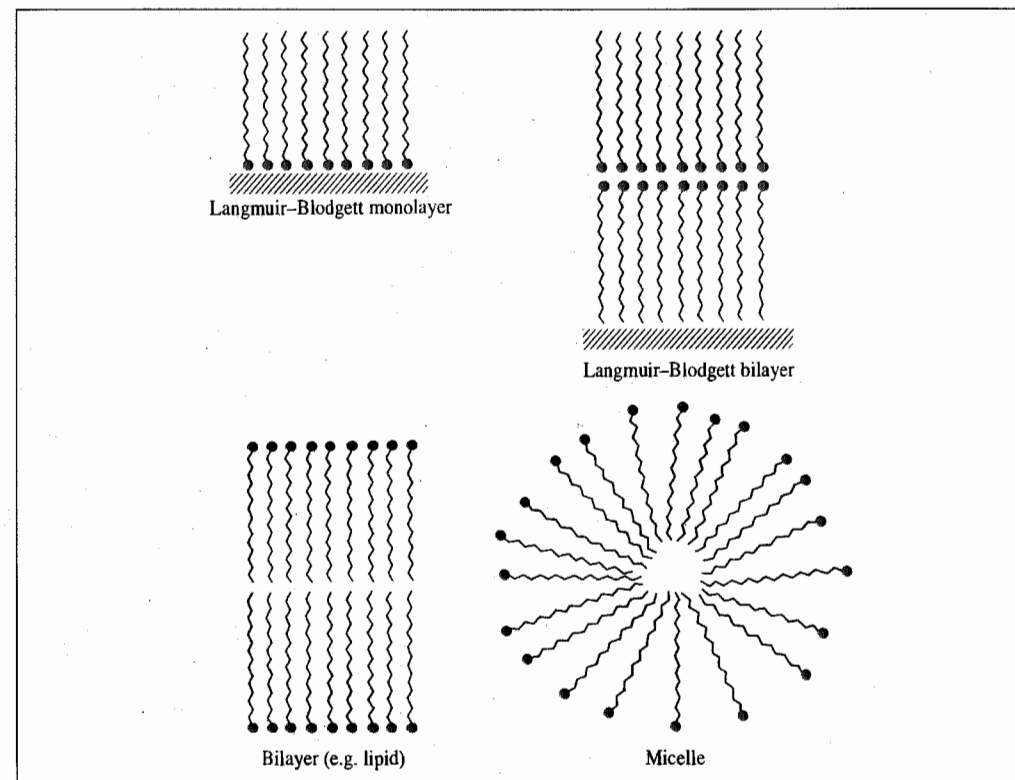


Fig. 7.17: Some of the various phases that amphiphiles may form.

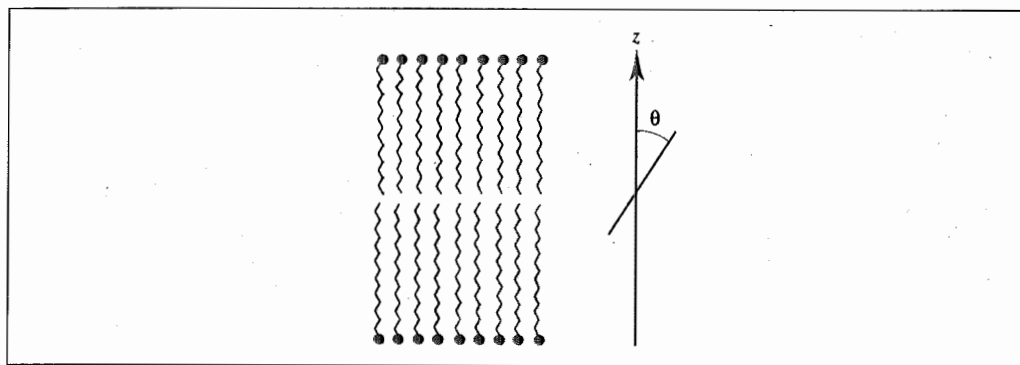


Fig. 7.18: Definition of the order parameter.

used to determine an *order parameter* for the carbon atoms on the hydrocarbon tail. The order parameter is defined as:

$$S = 0.5(3 \cos \theta_i \cos \theta_j - \delta_{ij}) \quad (7.126)$$

θ_i is the angle between the i th molecular axis and the *director*, which is the average of the molecular axes over the sample. For a bilayer in the $L\alpha$ phase (the one present in cell membranes) the director is the same as the bilayer normal and is conventionally taken to be the z axis; see Figure 7.18. δ_{ij} is the Kronecker delta function ($\delta_{ij} = 1$ if $i = j$; $\delta_{ij} = 0$ if $i \neq j$). The expression for S is averaged over time and over molecules. The deuterium NMR experiment provides the order parameter S_{CD} , which indicates the average orientation of the C–D bond vector with respect to the bilayer normal. The experimental order parameters S_{CD} can range from 1.0 (indicating full order along the bilayer normal) to -0.5 (full order perpendicular to the bilayer normal) [Seelig and Seelig 1974]. A value of zero is considered to indicate full isotropic motion of the group. Experimental values are determined using molecules with deuterium-substituted methylene groups at positions along the hydrocarbon chain. Many simulations of amphiphiles are performed using united atom models for the hydrocarbon chains and it is therefore necessary to be able to relate the experimental order parameters to values that can be calculated from a simulation. This is done as follows [Essex *et al.* 1994]. Molecular axes are defined for each CH_2 unit in the chain as shown in Figure 7.19. These molecular axes are defined for the n th CH_2 unit as follows:

z : vector from C_{n-1} to C_{n+1}

y : vector perpendicular to z and in the plane through C_{n-1} , C_n and C_{n+1}

x : perpendicular to y and z

Using these definitions, components of the molecular order parameter tensor can be determined (for example, S_{zz} is determined by measuring the angle between the molecular z axis and the bilayer normal). The experimental order parameter can be related to the molecular order parameter using the equation:

$$S_{CD} = 2S_{xx}/3 + S_{yy}/3 \quad (7.127)$$

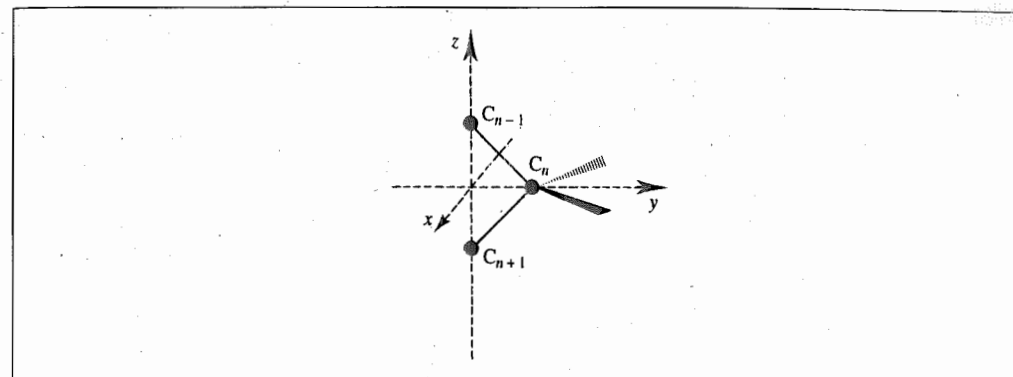


Fig. 7.19: Calculation of the order parameter for united atom simulations.

With all-atom simulations the locations of the hydrogen atoms are known and so the order parameters can be calculated directly. Another structural property of interest is the ratio of *trans* conformations to *gauche* conformations for the CH_2 – CH_2 bonds in the hydrocarbon tail. The *trans* : *gauche* ratio can be estimated using a variety of experimental techniques such as Raman, infrared and NMR spectroscopy.

7.10.1 Simulation of Lipids

There has been considerable interest in the simulation of lipid bilayers due to their biological importance. Early calculations on amphiphilic assemblies were limited by the computing power available, and so relatively simple models were employed. One of the most important of these is the mean field approach of Marcelja [Marcelja 1973, 1974], in which the interaction of a single hydrocarbon chain with its neighbours is represented by two additional contributions to the energy function. The energy of a chain in the mean field is given by:

$$\mathcal{V}_{\text{tot}} = \mathcal{V}_{\text{int}} + \mathcal{V}_{\text{disp}} + \mathcal{V}_{\text{rep}} \quad (7.128)$$

where \mathcal{V}_{int} is the internal energy of a chain, which can be calculated using standard force field methods. $\mathcal{V}_{\text{disp}}$ simulates the van der Waals interactions with the neighbouring molecules. It is often modelled using a Maier-Saupe potential:

$$\mathcal{V}_{\text{disp}} = -\Phi \sum_{i=1}^{\text{carbons}} \frac{1}{2} (3 \cos^2 \theta_i - 1) \quad (7.129)$$

The summation runs over all carbon atoms in the chain. θ_i is the angle between the bilayer normal and the molecular axis, as discussed above. Φ is the field strength; this may be parametrised to reproduce appropriate experimental data such as the deuterium NMR order parameters or it may be obtained by a self-consistent protocol, as described below. In his work on lipid bilayers Marcelja used a slightly different expression for $\mathcal{V}_{\text{disp}}$ which

involved the fraction of *trans* bonds in the system:

$$\mathcal{V}_{\text{disp}} = -\Phi \frac{n_{\text{trans}}}{n} \sum_{i=1}^{\text{carbons}} \frac{1}{2} (3 \cos^2 \theta_i - 1) \quad (7.130)$$

This additional factor was introduced to ensure the proper behaviour over both liquid crystalline and solid phases. In simulations of the liquid crystalline phase alone this term may be omitted for computational efficiency.

The repulsive contribution, \mathcal{V}_{rep} , is due to lateral pressure on each chain. In Marcelja's original treatment, this was set equal to the product of the lateral pressure, γ , and the cross-sectional area of the chain. The cross-sectional area was approximated by:

$$A = A_0 l_0 / l \quad (7.131)$$

where l_0 and A_0 are the length and cross-sectional area, respectively, of the hydrocarbon chain in a fully extended conformation. l is the length of the chain in the current conformation, projected onto the bilayer normal. If the bilayer normal is along the z axis then l is taken to be the z coordinate of the last carbon atom in the hydrocarbon chain. In other mean field models [Pastor *et al.* 1988] the product $\gamma A_0 / l_0$ is replaced with a single adjustable parameter Γ and so \mathcal{V}_{rep} is given by:

$$\mathcal{V}_{\text{rep}} = \sum_{\text{chains}} \frac{\Gamma}{(z_n - z_0)} \quad (7.132)$$

where z_n is the z coordinate of the last carbon in the chain and z_0 is the coordinate of the surface of the monolayer or bilayer. This force acts to keep the last carbon away from the surface; the closer it gets the larger the force pulling it away.

In his calculations, Marcelja generated all possible conformations of the hydrocarbon chain, restricting each carbon-carbon bond to the *trans* and *gauche* conformations. The energy of each conformation was evaluated. From the ensemble of conformations a partition function can be computed:

$$Z = \sum_{\text{all conformations}} \exp[-\mathcal{V}_{\text{tot}}/k_B T] \quad (7.133)$$

The molecular field is related to the partition function:

$$\Phi = \sum_{\text{all conformations}} \left\{ \frac{\frac{n_{\text{trans}}}{n} \sum_{i=1}^{\text{carbons}} \frac{1}{2} (3 \cos^2 \theta_i - 1) \exp[-\mathcal{V}_{\text{tot}}/k_B T]}{Z} \right\} \quad (7.134)$$

The molecular field is thus related to the partition function and so it is possible to generate a self-consistent value of the molecular field, Φ . Thermodynamic properties can then be calculated from the partition function. For example, Marcelja calculated the pressure as a function of the area per polar head group for surface monolayers at a variety of temperatures. His results showed good qualitative agreement with experimental results for such systems.

The mean field approach can be incorporated into a molecular dynamics simulation. It is particularly useful when used in conjunction with Langevin dynamics, as very long simulations can be performed. For example, Pearce and Harvey were able to perform simulations of three unsaturated phospholipids for 100 ns (i.e. 0.1 μ s) in single-molecule Langevin dynamics calculations [Pearce and Harvey 1993]. An extension of this strategy is to use a central 'core' containing one or more molecules that are simulated using molecular dynamics. This core is surrounded by a shell of molecules that are simulated using Langevin dynamics with the mean field. In this way one attempts to simulate a more 'realistic' system without incurring the computational penalty of a full molecular dynamics simulation of the entire system [De Loof *et al.* 1991].

The first molecular dynamics simulations of a lipid bilayer which used an explicit representation of all the molecules was performed by van der Ploeg and Berendsen in 1982 [van der Ploeg and Berendsen 1982]. Their simulation contained 32 decanoate molecules arranged in two layers of sixteen molecules each. Periodic boundary conditions were employed and a united atom force potential was used to model the interactions. The head groups were restrained using a harmonic potential of the form:

$$v(z) = \frac{k_h}{2} (z - \langle z \rangle)^2 \quad (7.135)$$

By writing the restraint in terms of the average z coordinates of the head groups ($\langle z \rangle$) van der Ploeg and Berendsen ensured that the bilayer was able to change its thickness to reach its equilibrium value. This restraining potential was designed to reproduce the interactions between the head groups and the water layer, neither of which was explicitly included in the calculation. A key feature of the simulation was the long equilibration time required. By explicitly representing all the molecules in the system it was possible to determine the collective motion of the system as a whole. One distinct feature was a slowly fluctuating collective tilt of the molecules away from the normal to the bilayer surface (Figure 7.20). The degree to which the molecules were aligned with each other was also correlated with the tilt angle. When the average tilt angle reached a maximum the chains were much more likely to be well aligned, but when the average tilt angle was close to zero (i.e. such that the average orientation of the chains was almost normal to the bilayer surface) much less order was observed. In their original simulations this collective tilt phenomenon was observed to extend over the entire simulation cell, suggesting that the cell dimensions were too small and that the use of periodic boundary conditions was enhancing the long-range correlations. Simulations using a larger system subsequently showed that this collective tilt could be observed for subsets of the molecules.

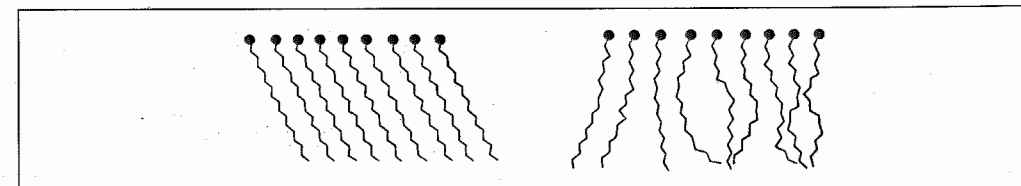


Fig. 7.20: Variation in alignment of chains in lipid simulation with tilt angle [van der Ploeg and Berendsen 1982].

Faster computers have enabled more realistic simulations of lipid bilayers to be performed, of larger systems, with more accurate models and for longer times [Stouch 1993; Tobias *et al.* 1997]. The trend is very much towards simulations that use full representations of all species present (i.e. 'all-atom' models with explicit solvent and counterions). The charged and highly polar nature of lipid head groups means that a proper representation of the long-range electrostatic forces can be critical, using a method such as the Ewald summation. Equilibration of such systems often requires hundreds of picoseconds and certain phenomena are only observed on a nanosecond timescale. In addition, molecules such as cholesterol and proteins may be included within the membrane. These simulations have revealed many hitherto unknown features of the behaviour of such systems. For example, considerable conformational mobility of the hydrocarbon chains is often observed in the liquid crystalline phases. This is illustrated in Figure 7.21 (colour plate section), which shows a snapshot of a lipid bilayer after a molecular dynamics simulation of several hundred picoseconds. The considerable degree of disorder in the hydrocarbon chains near the middle of the bilayer is clear from this figure and is very different to the idealised, 'textbook' pictures in which the chains are perfectly aligned in completely extended conformations. The distribution of *gauche* conformations tends to be higher towards the end of the chain, though in some systems a *gauche* link is required near the head group to enable the chain to lie perpendicular to the interface. 'Kinks' are often observed in the chains; these are arrangements of three successive bonds with *gauche*(+)-*trans*-*gauche*(-) torsion angle, which enable the chain to remain perpendicular to the surface.

7.10.2 Simulations of Langmuir-Blodgett films

The simulations of Langmuir-Blodgett systems can be difficult due to the need to correctly model the solid support. To illustrate the procedure we will describe the calculations of Kim, Moller, Tildesley and Quirke [Kim *et al.* 1994a] who simulated stearic acid ($\text{CH}_3(\text{CH}_2)_{16}\text{COOH}$) adsorbed onto graphite. The surface was modelled using a Lennard-Jones 9-3 potential that depends upon the height of the atom (α) above the surface (z_α):

$$v_{\text{as}}(z_\alpha) = \frac{2\pi\rho}{3} \epsilon_{\text{ss}} \left[\frac{2}{15} \left(\frac{\sigma_{\text{as}}}{z_\alpha} \right)^9 - \left(\frac{\sigma_{\text{as}}}{z_\alpha} \right)^3 \right] \quad (7.136)$$

where ρ is the density of the solid and ϵ_{xx} and δ_{ss} are its Lennard-Jones parameters. An image-charge method was also applied to the acid head group with the interaction between a charge and its image being:

$$v_{\text{ic}}(z) = \frac{1}{2} \frac{(\epsilon - \epsilon')}{(\epsilon + \epsilon')} \left[\frac{q_\alpha^2}{8\pi\epsilon_0(z - z_{\text{ip}})} \right] \quad (7.137)$$

where ϵ' is the relative permittivity of the solid (taken to be 4.0) and ϵ is the permittivity above the surface ($\epsilon = 1.0$). The image plane is located at $z_{\text{ip}} = \sigma_{\text{ss}}/2$. Each charge interacts with its own image and with the images of other charges, but there are no interactions between the image charges themselves. The hydrocarbon chain of the stearic acid was modelled using an all-atom model, with explicit hydrogen atoms.

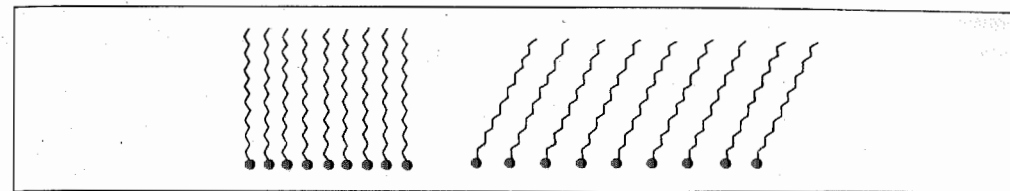


Fig. 7.22: Simulations of a Langmuir-Blodgett film [Kim *et al.* 1994a]: as the area per head group increases the chains tilt away from the normal.

A molecular dynamics simulation of 64 molecules with periodic boundary conditions confirmed the presence of a transition in which the collective tilt of the chains changed from being upright (i.e. perpendicular to the surface) to having an angle of around 20° (Figure 7.22). This transition was induced by increasing the area per head group. The proportion of molecules in the all-*trans* conformation decreased significantly as the head group area was increased (97.7% of molecules were fully extended for a head group area of 20.6 \AA^2 but only 66.9% for an area of 21.2 \AA^2). The bond linking the chain to the acid head growth showed a considerable degree of rotational disorder.

Bilayers of stearic acid were also simulated on a hydrophobic surface [Kim *et al.* 1994b]. In the bilayer the molecules are arranged head to head, with the hydrocarbon tail on the surface. In this arrangement hydrogen bonds form between the head groups (Figure 7.23). The bilayer also showed the tilt angle transition that was observed for the monolayer, though the degree of tilt was considerably less for the bilayer, suggesting that hydrogen bonding between the head groups was important in controlling the orientation of the molecules.

An extension of these calculations to cationic dialkylamide salts required an even more complex model [Adolf *et al.* 1995]. These molecules have the general formula $(\text{CH}_3)_2\text{N}^+[(\text{CH}_2)_{n-1}\text{CH}_3][(\text{CH}_2)_{m-1}\text{CH}_3]\text{Cl}^-$ and the isomer with $m = n = 18$ is one of the main active ingredients in commercial fabric softeners. The presence of two long alkyl

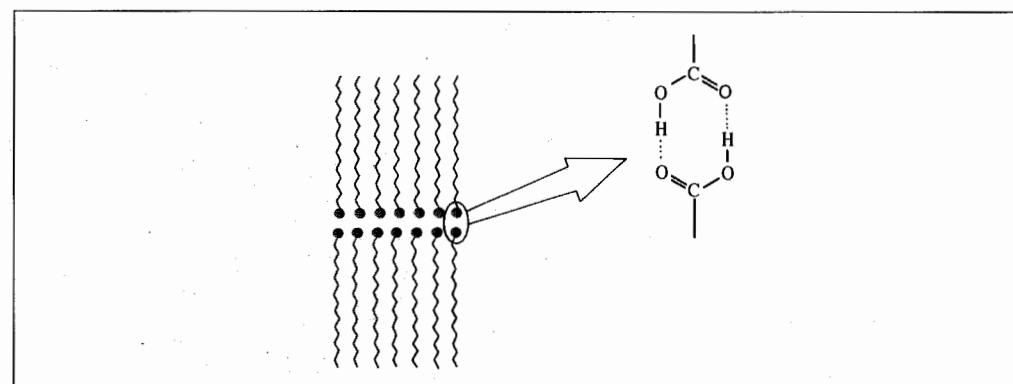


Fig. 7.23: In simulations of stearic acid on a hydrophobic surface hydrogen bonding between the head groups is important in controlling the orientation of the molecules [Kim *et al.* 1994b].

chains and an ionic head group means that these molecules are also structurally similar to phospholipids. A modified Ewald method was used to calculate electrostatic interactions in the two dimensions parallel to the surface, and the anisotropic potential model of Toxvaerd (see Section 4.15) was employed to retain the computational savings of a united-atom model. This system also showed a variation in the tilt with head group area, though the results at the highest head group densities were not as 'solid-like' as was suggested by the experimental data. Nevertheless, there were some areas where the model could be improved, including the need to incorporate water molecules and use a more appropriate representation of the chloride anion.

7.10.3 Mesoscale Modelling: Dissipative Particle Dynamics

The molecular dynamics methods that we have discussed in this chapter, and the examples that have been used to illustrate them, fall into the category of 'atomistic simulations', in that all of the actual atoms (or at least the non-hydrogen atoms) in the core system are represented explicitly. Atomistic simulations can provide very detailed information about the behaviour of the system, but as we have discussed this typically limits a simulation to the nanosecond timescale. Many processes of interest occur over a longer timescale. In the case of processes which occur on a 'macroscopic' timescale (i.e. of the order of seconds) then rather simple models may often be applicable. Between these two extremes are phenomena that occur on an intermediate scale (of the order of microseconds). This is the realm of the *mesoscale*. Dissipative particle dynamics (DPD) is particularly useful in this region; examples include complex fluids such as surfactants and polymer melts.

These three general regions (atomistic, mesoscopic and macroscopic) are not only characterised by different timescales but also varying length scales. Indeed, there is a general inverse relationship between the time and the length. In the case of the dissipative particle dynamics method the fast motion of the atoms is integrated out, leaving as the fundamental 'unit' a set of beads that interact with other beads via an appropriate potential [Koelman and Hoogerbrugge 1993]. Each bead represents a small 'droplet' of the fluid. The total force on each bead is due to a combination of direct interactions with other beads together with random and dissipative forces. The trajectory of the system is calculated by integrating Newton's laws of motion in the usual way, from which properties can be derived.

The underlying model in dissipative particle dynamics is usually developed in such a way that the mass, length and timescales are all unity. This is similar to the use of reduced units for the Lennard-Jones potential (Section 4.10.5). A particular advantage of such an approach is that a single simulation may often be able to explain the behaviour of many different systems. With a mass of 1 the force acting on a particle is equal to its acceleration. In DPD there are three forces on each bead [Groot and Warren 1997]:

$$\mathbf{f}_i = \sum_{j=1; j \neq i}^N (\mathbf{F}_{ij}^C + \mathbf{F}_{ij}^D + \mathbf{F}_{ij}^R) \quad (7.138)$$

The summation is over all other particles j which are within a certain cutoff radius r_c of i . This cutoff radius becomes the unit of length in the subsequent treatment (i.e. $r_c = 1$). The

first of these forces, the conservative force, \mathbf{F}_{ij}^C , is a soft repulsion that acts along the line joining i to j :

$$\mathbf{F}_{ij}^C = \begin{cases} a_{ij}(1 - r_{ij})\hat{\mathbf{r}}_{ij} & r_{ij} < 1 \\ 0 & r_{ij} > 1 \end{cases} \quad (7.139)$$

r_{ij} is the distance between beads i and j and $\hat{\mathbf{r}}_{ij}$ is the corresponding unit vector. The second force is a dissipative (or drag) force, which is given by:

$$\mathbf{F}_{ij}^D = \begin{cases} -\gamma w^D(r_{ij})(\hat{\mathbf{r}}_{ij} \cdot \mathbf{v}_{ij})\hat{\mathbf{r}}_{ij} & r_{ij} < 1 \\ 0 & r_{ij} > 1 \end{cases} \quad (7.140)$$

This dissipative force is proportional to the relative velocity of the two beads and acts so as to reduce their relative momentum. \mathbf{v}_{ij} is the difference between the two velocities ($\mathbf{v}_{ij} = \mathbf{v}_i - \mathbf{v}_j$) and $w^D(r_{ij})$ is a weight function that depends upon the distance r_{ij} and disappears for inter-bead distances greater than unity (i.e. r_c).

The third and final force acting between any pair of beads is a random force:

$$\mathbf{F}_{ij}^R = \begin{cases} \sigma w^R(r_{ij})\theta_{ij}\hat{\mathbf{r}}_{ij} & r_{ij} < 1 \\ 0 & r_{ij} > 1 \end{cases} \quad (7.141)$$

$w^R(r_{ij})$ is a distance-dependent weight function similar to that for the dissipative force. θ_{ij} is a function which ensures that the random force between each pair of particles averages to zero over time and is independent of the force between every other pair of particles. The random force can be more usefully expressed in terms of the timestep in the integration scheme:

$$\mathbf{F}_{ij}^R = \frac{\sigma w^R(r_{ij})\zeta_{ij}\hat{\mathbf{r}}_{ij}}{\sqrt{\delta t}} \quad (7.142)$$

ζ_{ij} is a random number with zero mean and unit variance, chosen independently for each pair of particles and at each time step in the integration.

Both the dissipative force and the random force act along the line joining the pair of beads and also conserve linear and angular momentum. The model thus has two unknown functions $w^D(r_{ij})$ and $w^R(r_{ij})$ and two unknown constants γ and σ . In fact, only one of the two weight functions can be chosen arbitrarily as they are related [Espanol and Warren 1995]. Moreover, the temperature of the system relates the two constants:

$$w^D(r) = [w^R(r)]^2 \quad (7.143)$$

$$\sigma^2 = 2\gamma k_B T \quad (7.144)$$

The usual choice for the weight functions is to make the random force the same as the conservative force:

$$w^D(r) = [w^R(r)]^2 = \begin{cases} (1 - r)^2 & r < 1 \\ 0 & r > 1 \end{cases} \quad (7.145)$$

The equations of motion are integrated using a modified velocity Verlet algorithm. The modification is required because the force depends upon the velocity; the extra step involves

a prediction followed by a correction. If, in addition to the use of units of mass and length, we assume that $k_B T$ is equal to 1, then the unit of time is:

$$\tau = r_c \sqrt{m/k_B T} \quad (7.146)$$

It remains to assign values to the noise amplitude, σ , the time step for the integration, δt , and the repulsion parameter, a_{ij} . The effects of the first two of these upon the stability of the simulation are also related to the integration method. Groot and Warren determined that when the noise amplitude was larger than 8 the integration scheme became unstable and that a value of 3 gave good results over a range of temperatures [Groot and Warren 1997]. The integration time step with the modified Verlet algorithm should have a value between 0.04 and 0.06; any larger and the temperature would artificially increase by an unacceptable amount. The repulsion parameter is the key determinant of the interactions between the beads. This can be achieved by relating the DPD model to bulk properties. For example, to model the compressibility of water at room temperature the repulsion parameter is related to the density, ρ , by:

$$a_{ii} \rho = 75 k_B T \quad (7.147)$$

These interaction parameters between particles of the same type can be used to derive the values of a_{ij} between unlike beads. For polymers which involve beads of different types the repulsion between unlike beads is made larger than between like beads.

A good example of the use of DPD is the study by Groot and Madden of the microphase separation of diblock copolymer melts [Groot and Madden 1998; Groot *et al.* 1999]. Block copolymers are surfactants which are present in many consumer products such as foods (e.g. ice cream and margarine), detergents and personal care products (e.g. shampoo). The properties of these materials are strongly dependent upon their bulk organisation (or *morphology*), which in turn depends upon the relative sizes of the head and the tail groups and how they interact. The diblock copolymers of interest can be represented by the general formula $A_m B_n$ where A and B represent amalgamations of the smaller building blocks from which the polymer is constructed. Of particular interest was the way in which the behaviour of the system varied as the ratio of A to B was changed, for a fixed polymer length. In this particular case the length was fixed at 10 beads and the entire simulation contained a total of 40 000 particles. A variety of systems were investigated, such as $A_2 B_8$, $A_3 B_7$ and $A_5 B_5$. The beads in each polymer chain were kept together by adding an extra term to the force (Equation (7.139)) of the form $C r_{ij}$ if i is connected to j .

Due to the greater degree of repulsion between unlike beads the final configuration of the system contains domains which are rich in either the A or the B type of bead. Some regions are rich in A; others are rich in B. The organisation of the A-rich and B-rich domains can be visualised by plotting a three-dimensional contour that connects regions where the density is intermediate between purely A and purely B.

For the 1:1 polymer ($A_5 B_5$) a lamellar phase was obtained, in which the A- and B-rich domains form parallel planes of alternating A and B beads (Figure 7.24(a), colour plate section). For other configurations, however, different structures were observed. The $A_3 B_7$ system evolved to a hexagonal phase (Figure 7.24(b)) and the $A_2 B_8$ structure produces a set of peanut-shaped micelles (Figure 7.24(c)).

Appendix 7.1 Energy Conservation in Molecular Dynamics

The total energy is the sum of the kinetic $\mathcal{K}(t)$ and potential energies $\mathcal{V}(t)$:

$$E(t) = \mathcal{K}(t) + \mathcal{V}(t) \quad (7.148)$$

We want to derive an expression for the rate of change of the energy with time, dE/dt . First, we differentiate the kinetic energy term with respect to time:

$$\frac{d\mathcal{K}}{dt} = \sum_{i=1}^N \frac{d}{dt} \left(\frac{1}{2} m_i v_i^2 \right) = \sum_{i=1}^N m_i v_i \frac{dv_i}{dt} \quad (7.149)$$

As $m_i dv_i/dt$ is equal to the force on the atom i , the result can be written:

$$\frac{d\mathcal{K}}{dt} = \sum_{i=1}^N v_i f_i \quad (7.150)$$

f_i is the force on atom i .

The potential energy is written as a series of pairwise interaction terms:

$$\mathcal{V}(t) = \sum_{i=1}^N \sum_{j=i+1}^N v(r_{ij}(t)) \quad (7.151)$$

The derivative of the potential energy with respect to time can be written:

$$\frac{d\mathcal{V}}{dt} = \sum_{i=1}^N \sum_{j=i+1}^N \frac{\partial v}{\partial r_{ij}} \frac{dr_{ij}}{dt} \quad (7.152)$$

$\partial v / \partial r_{ij}$ equals 1 for each pairwise combination i and j . Each term $v(r_{ij})$ is a function of the positions of atom i and j (\mathbf{r}_i and \mathbf{r}_j) and we can then write:

$$\frac{dv(r_{ij})}{dt} = \frac{dv(r_{ij})}{dr_i} \frac{dr_i}{dt} + \frac{dv(r_{ij})}{dr_j} \frac{dr_j}{dt} \quad (7.153)$$

For a given atom i , there will be a total of $N - 1$ terms of the form $v(r_{ij})$ in the expression for the potential energy due to the interactions between i and all other atoms j . Hence we can write $d\mathcal{V}/dt$ as follows:

$$\frac{d\mathcal{V}}{dt} = \sum_{i=1}^N \sum_{j=1; j \neq i}^N \frac{\partial v(r_{ij})}{\partial \mathbf{r}_i} \frac{d\mathbf{r}_i}{dt} = \sum_{i=1}^N \frac{d\mathbf{r}_i}{dt} \sum_{j=1; j \neq i}^N \frac{\partial v(r_{ij})}{\partial \mathbf{r}_i} \quad (7.154)$$

The force on atom i due to its interaction with atom j equals minus the gradient with respect to \mathbf{r}_i , or $-dv(r_{ij})/d\mathbf{r}_i$. Thus the total force on the atom is equal to

$$- \sum_{j=1; j \neq i}^N \frac{\partial v(r_{ij})}{\partial \mathbf{r}_i} \quad (7.155)$$

and so we have:

$$\frac{d\mathcal{V}}{dt} = - \sum_{i=1}^N \frac{d\mathbf{r}_i}{dt} \cdot \mathbf{f}_i = - \sum_{i=1}^N v_i f_i \quad (7.156)$$

Thus $(d\mathcal{V}/dt) + (d\mathcal{K}/dt) = dE/dt = 0$, which implies that the energy is constant. In practice, the total energy fluctuates about a constant value.

Further Reading

- Allen M P and D J Tildesley 1987. *Computer Simulation of Liquids*. Oxford, Oxford University Press.
- Berendsen H C and W F van Gunsteren 1984. Molecular Dynamics Simulations: Techniques and Approaches. In Barnes A J, W J Orville-Thomas and J Yarwood (Editors). *Molecular Liquids, Dynamics and Interactions*. NATO ASI Series C135, New York, Reidel, pp. 475–600.
- Berendsen H C and W F van Gunsteren 1986. Practical Algorithms for Dynamic Simulations. Molecular Dynamics Simulation of Statistical Mechanical Systems. *Proceedings of the Enrico Fermi Summer School Varenna Soc. Italiana di Fisica*. Bologna, pp. 43–65.
- Brooks C L III, M Karplus and B M Pettitt 1988. Proteins. A Theoretical Perspective of Dynamics, Structure and Thermodynamics. *Advances in Chemical Physics* Volume LXXI. New York, John Wiley & Sons.
- Goldstein H 1980. *Classical Mechanics* (2nd Edition). Reading, MA, Addison-Wesley.
- Haile J M 1992. *Molecular Dynamics Simulation. Elementary Methods*. New York, John Wiley & Sons.
- McCammon J A and S C Harvey 1987. *Dynamics of Proteins and Nucleic Acids*. Cambridge, Cambridge University Press.
- van Gunsteren W F 1994. Molecular Dynamics and Stochastic Dynamics Simulations: A Primer. In van Gunsteren W F, P K Weiner and A J Wilkinson (Editors). *Computer Simulations of Biomolecular Systems* Volume 2. Leiden, ESCOM.
- van Gunsteren W F and H J C Berendsen 1990. Computer Simulation of Molecular Dynamics: Methodology, Applications and Perspectives in Chemistry. *Angewandte Chemie International Edition in English* 29:992–1023.

References

- Adolf D B, D J Tildesley, M R S Pinches, J B Kingdon, T Madden and A Clark 1995. Molecular Dynamics Simulations of Dioctadecyldimethylammonium Chloride Monolayers. *Langmuir* 11:237–246.
- Alder B J and T E Wainwright 1957. Phase Transition for a Hard-sphere System. *Journal of Chemical Physics* 27:1208–1209.
- Alder B J and T E Wainwright 1970. Decay of the Velocity Autocorrelation Function. *Physical Review A* 1:18–21.
- Allen M P and D J Tildesley 1987. *Computer Simulation of Liquids*. Oxford, Oxford University Press.
- Anderson H C 1980. Molecular Dynamics Simulations at Constant Pressure and/or Temperature. *Journal of Chemical Physics* 72:2384–2393.
- Anderson H C 1983. Rattle: A 'Velocity' Version of the Shake Algorithm for Molecular Dynamics Calculations. *Journal of Computational Physics* 54:24–34.
- Beeman D 1976. Some Multistep Methods for Use in Molecular Dynamics Calculations. *Journal of Computational Physics* 20:130–139.

- Berendsen H J C, J P M Postma, W F van Gunsteren, A Di Nola and J R Haak 1984. Molecular Dynamics with Coupling to an External Bath. *Journal of Chemical Physics* 81:3684–3690.
- Brunger A, C B Brooks and M Karplus 1984. Stochastic Boundary Conditions for Molecular Dynamics Simulations of ST2 Water. *Chemical Physics Letters* 105:495–500.
- Dauber-Osguthorpe P and D J Osguthorpe 1990. Analysis of Intramolecular Motions by Filtering Molecular Dynamics Trajectories. *Journal of the American Chemical Society* 112:7921–7935.
- Dauber-Osguthorpe P and D J Osguthorpe 1993. Partitioning the Motion in Molecular Dynamics Simulations into Characteristic Modes of Motion. *Journal of Computational Chemistry* 14:1259–1271.
- De Loof H, S C Harvey, J P Segrest and R W Pastor 1991. Mean Field Stochastic Boundary Molecular Dynamics Simulation of a Phospholipid in a Membrane. *Biochemistry* 30:2099–2113.
- Espanol P and P B Warren 1995. Statistical Mechanics of Dissipative Particle Dynamics. *Europhysics Letters* 30:191–196.
- Essex J W, M M Hann and W G Richards 1994. Molecular Dynamics of a Hydrated Phospholipid Bilayer. *Philosophical Transactions of the Royal Society of London* B344:239–260.
- Fincham D and Heyes D M 1982. Integration Algorithms in Molecular Dynamics. *CCP5 Quarterly* 6:4–10.
- Gear C W 1971. *Numerical Initial Value Problems in Ordinary Differential Equations*. Englewood Cliffs, NJ, Prentice Hall.
- Groot R D and T J Madden 1998. Dynamic Simulation of Diblock Copolymer Microphase Separation. *Journal of Chemical Physics* 108:8713–8724.
- Groot R D, T J Madden and D J Tildesley 1999. On the Role of Hydrodynamic Interactions in Block Copolymer Microphase Separation. *Journal of Chemical Physics* 110:9739–9749.
- Groot R D and P B Warren 1997. Dissipative Particle Dynamics: Bridging the Gap Between Atomistic and Mesoscopic Simulation. *Journal of Chemical Physics* 107:4423–4435.
- Helfand E 1984. Dynamics of Conformational Transitions in Polymers. *Science* 226:647–650.
- Hockney R W 1970. The Potential Calculation and Some Applications. *Methods in Computational Physics* 9:136–211.
- Hoover W G 1985. Canonical Dynamics: Equilibrium Phase-space Distributions. *Physical Review A* 31:1695–1697.
- Humphreys D D, R A Friesner and B J Berne 1994. A Multiple Time-step Molecular Dynamics Algorithm for Macromolecules. *Journal of Physical Chemistry* 98:6885–6892.
- Humphreys D D, R A Friesner and B J Berne 1995. Simulated Annealing of a Protein in a Continuum Solvent by Multiple Time-step Molecular Dynamics. *Journal of Physical Chemistry* 99:10674–10685.
- Humphreys, D D, R A Friesner and B J Berne 1996. A Multiple Time-step Molecular Dynamics Algorithm for Macromolecules. *Journal of Physical Chemistry* 98:6885–6892.
- Jorgensen W L, R C Binning Jr and B Bigot 1981. Structures and Properties of Organic Liquids: *n*-Butane and 1,2-Dichloroethane and Their Conformational Equilibria. *Journal of the American Chemical Society* 103:4393–4399.
- Kim K S, M A Moller, D J Tildesley and N Quirke 1994a. Molecular Dynamics Simulations of Langmuir–Blodgett Monolayers with Explicit Head-group Interactions. *Molecular Simulation* 13:77–99.
- Kim K S, D J Tildesley and N Quirke 1994b. Molecular Dynamics of Langmuir–Blodgett Films: II. Bilayers. *Molecular Simulation* 13:101–114.
- Koelman J M V A and P J Hoogerbrugge 1993. Dynamic Simulations of Hard-sphere Suspensions Under Steady Shear. *Europhysics Letters* 21:363–368.
- Lavery R and H Sklenar 1988. The Definition of Generalized Helicoidal Parameters and of Axis Curvature for Irregular Nucleic Acids. *Journal of Biomolecular Structure and Dynamics* 6:63–91.
- Leach A R and T E Klein 1995. A Molecular Dynamics Study of the Inhibitors of Dihydrofolate Reductase by a Phenyl Triazine. *Journal of Computational Chemistry* 16:1378–1393.

- Marcelja S 1973. Molecular Model for Phase Transition in Biological Membranes. *Nature* 241:451-453.
- Marcelja S 1974. Chain Ordering in Liquid Crystals. II. Structure of Bilayer Membranes. *Biochimica et Biophysica Acta* 367:165-176.
- Nosé S 1984. A Molecular Dynamics Method for Simulations in the Canonical Ensemble. *Molecular Physics* 53:255-268.
- Pastor R W, R M Venable and M Karplus 1988. Brownian Dynamics Simulation of a Lipid Chain in a Membrane Bilayer. *Journal of Chemical Physics* 89:1112-1127.
- Pearce L L and S C Harvey 1993. Langevin Dynamics Studies of Unsaturated Phospholipids in a Membrane Environment. *Biophysical Journal* 65:1084-1092.
- Procacci P and B Berne 1994. Computer Simulation of Solid C₆₀ Using Multiple Time-step Algorithms. *Journal of Chemical Physics* 101:2421-2431.
- Rahman A 1964. Correlations in the Motion of Atoms in Liquid Argon. *Physical Review* A136:405-411.
- Rahman A and F H Stillinger 1971. Molecular Dynamics Study of Liquid Water. *Journal of Chemical Physics* 55:3336-3359.
- Robinson A J, W G Richards, P J Thomas and M M Hann 1994. Head Group and Chain Behaviour in Biological Membranes - A Molecular Dynamics Simulation. *Biophysical Journal* 67:2345-2354.
- Rubinstein R Y 1981. *Simulation and Monte Carlo Methods*. New York, John Wiley & Sons.
- Ryckaert J P, G Cicotti and H J C Berendsen 1977. Numerical Integration of the Cartesian Equations of Motion of a System with Constraints: Molecular Dynamics of *n*-Alkanes. *Journal of Computational Physics* 23:327-341.
- Seelig A and J Seelig 1974. The Dynamics Structure of Fatty Acyl Chains in a Phospholipid Bilayer Measured by Deuterium Magnetic Resonance. *Biochemistry* 13:4839-4845.
- Stouch T R 1993. Lipid Membrane Structure and Dynamics Studied by All-atom Molecular Dynamics Simulations of Hydrated Phospholipid Bilayers. *Molecular Simulation* 10:335-362.
- Streit W B, D Tildesley and G Saville 1978. Multiple Time-step Methods in Molecular Dynamics. *Molecular Physics* 35:639-648.
- Swindoll R D and J M Haile 1984. A Multiple Time-step Method for Molecular Dynamics Simulations of Fluids of Chain Molecules. *Journal of Computational Physics* 53:289-298.
- Swope W C, H C Anderson, P H Berens and K R Wilson 1982. A Computer Simulation Method for the Calculation of Equilibrium Constants for the Formation of Physical Clusters of Molecules: Application to Small Water Clusters. *Journal of Chemical Physics* 76:637-649.
- Tobias D J and C L Brooks III 1988. Molecular Dynamics with Internal Coordinate Constraints. *Journal of Chemical Physics* 89:5115-5126.
- Tobias D J, K Tu and M L Klein 1997. Atomic-scale Molecular Dynamics Simulations of Lipid Membranes. *Current Opinion in Colloid and Interface Science* 2:15-26.
- Tuckerman M, B J Berne and G J Martyna 1992. Reversible Multiple Time Scale Molecular Dynamics. *Journal of Chemical Physics* 97:1990-2001.
- van der Ploeg P and H J C Berendsen 1982. Molecular Dynamics Simulation of a Bilayer Membrane. *Journal of Chemical Physics* 76:3271-3276.
- van Gunsteren W F and H J C Berendsen 1982. Algorithms for Brownian Dynamics. *Molecular Physics* 45:637-547.
- van Gunsteren W F, H J C Berendsen and J A C Rullmann 1981. Stochastic Dynamics for Molecules with Constraints. Brownian Dynamics of *n*-Alkanes. *Molecular Physics* 44:69-95.
- Verlet L 1967. Computer 'Experiments' on Classical Fluids. I. Thermodynamical Properties of Lennard-Jones Molecules. *Physical Review* 159:98-103.
- Watanabe M and M Karplus 1993. Dynamics of Molecules with Internal Degrees of Freedom by Multiple Time-step Methods. *Journal of Chemical Physics* 99:8063-8074.
- Widmalm G and R W Pastor 1992. Comparison of Langevin and Molecular Dynamics Simulations. *Journal of the Chemical Society Faraday Transactions* 88:1747-1754.

- Woodcock L V 1971. Isothermal Molecular Dynamics Calculations for Liquid Salts. *Chemical Physics Letters* 10:257-261.
- Yun-Yu S, W Lu and W F van Gunsteren 1988. On the Approximation of Solvent Effects on the Conformation and Dynamics of Cyclosporin A by Stochastic Dynamics Simulation Techniques. *Molecular Simulation* 1:369-383.
- Zhou R and B J Berne 1995. A New Molecular Dynamics Method Combining the Reference System Propagator Algorithm with a Fast Multipole Method for Simulating Proteins and Other Complex Systems. *Journal of Chemical Physics* 103:9444-9459.

CHAPTER EIGHT

Monte Carlo Simulation Methods

8.1 Introduction

The Monte Carlo simulation method occupies a special place in the history of molecular modelling, as it was the technique used to perform the first computer simulation of a molecular system. A Monte Carlo simulation generates configurations of a system by making random changes to the positions of the species present, together with their orientations and conformations where appropriate. Many computer algorithms are said to use a 'Monte Carlo' method, meaning that some kind of random sampling is employed. In molecular simulations 'Monte Carlo' is almost always used to refer to methods that use a technique called *importance sampling*. Importance sampling methods are able to generate states of low energy, as this enables properties to be calculated accurately. We can calculate the potential energy of each configuration of the system, together with the values of other properties, from the positions of the atoms. The Monte Carlo method thus samples from a $3N$ -dimensional space of the positions of the particles. There is no momentum contribution in a Monte Carlo simulation, in contrast to a molecular dynamics simulation. How then can a Monte Carlo simulation be used to calculate thermodynamic quantities, given that phase space is $6N$ -dimensional?

To resolve this difficulty, let us return to the canonical ensemble partition, Q , which for a system of N identical particles of mass m can be written:

$$Q_{NVT} = \frac{1}{N!} \frac{1}{h^{3N}} \iint d\mathbf{p}^N d\mathbf{r}^N \exp \left[-\frac{\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{k_B T} \right] \quad (8.1)$$

The factor $N!$ disappears when the particles are no longer indistinguishable. $\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)$ is the Hamiltonian that corresponds to the total energy of the system. The value of the Hamiltonian depends upon the $3N$ positions and $3N$ momenta of the particles in the system (one position and one momentum for each of the three coordinates of each particle). The Hamiltonian can be written as the sum of the kinetic and potential energies of the system:

$$\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N) = \sum_{i=1}^N \frac{|\mathbf{p}_i|^2}{2m} + \mathcal{V}(\mathbf{r}^N) \quad (8.2)$$

The crucial point to recognise is that the double integral in Equation (8.1) can be separated into two separate integrals, one over positions and the other over the momenta:

$$Q_{NVT} = \frac{1}{N!} \frac{1}{h^{3N}} \int d\mathbf{p}^N \exp \left[-\frac{|\mathbf{p}|^2}{2mk_B T} \right] \int d\mathbf{r}^N \exp \left[-\frac{\mathcal{V}(\mathbf{r}^N)}{k_B T} \right] \quad (8.3)$$

This separation is possible only if the potential energy function, $\mathcal{V}(\mathbf{r}^N)$, is not dependent upon the velocities (this is a safe assumption for almost all potential functions in common use). The integral over the momenta can now be performed analytically, the result being:

$$\int d\mathbf{p}^N \exp \left[-\frac{|\mathbf{p}|^2}{2mk_B T} \right] = (2\pi mk_B T)^{3N/2} \quad (8.4)$$

The partition function can thus be written:

$$Q_{NVT} = \frac{1}{N!} \left(\frac{2\pi mk_B T}{h^2} \right)^{3N/2} \int d\mathbf{r}^N \exp \left(-\frac{\mathcal{V}(\mathbf{r}^N)}{k_B T} \right) \quad (8.5)$$

The integral over the positions is often referred to as the *configurational integral*, Z_{NVT} :

$$Z_{NVT} = \int d\mathbf{r}^N \exp \left(-\frac{\mathcal{V}(\mathbf{r}^N)}{k_B T} \right) \quad (8.6)$$

In an ideal gas there are no interactions between the particles and so the potential energy function, $\mathcal{V}(\mathbf{r}^N)$, equals zero. $\exp(-\mathcal{V}(\mathbf{r}^N)/k_B T)$ is therefore equal to 1 for every gas particle in the system. The integral of 1 over the coordinates of each atom is equal to the volume, and so for N ideal gas particles the configurational integral is given by V^N ($V \equiv$ volume). This leads to the following result for the canonical partition function of an ideal gas:

$$Q_{NVT} = \frac{V^N}{N!} \left(\frac{2\pi k_B T m}{h^2} \right)^{3N/2} \quad (8.7)$$

This is often written in terms of the *de Broglie thermal wavelength*, Λ :

$$Q_{NVT} = \frac{V^N}{N! \Lambda^{3N}} \quad (8.8)$$

where $\Lambda = \sqrt{h^2/2\pi k_B T m}$.

By combining Equations (8.4) and (8.6) we can see that the partition function for a 'real' system has a contribution due to ideal gas behaviour (the momenta) and a contribution due to the interactions between the particles. Any deviations from ideal gas behaviour are due to interactions within the system as a consequence of these interactions. This enables us to write the partition function as:

$$Q_{NVT} = Q_{NVT}^{\text{ideal}} Q_{NVT}^{\text{excess}} \quad (8.9)$$

The excess part of the partition function is given by:

$$Q_{NVT}^{\text{excess}} = \frac{1}{V^N} \int d\mathbf{r}^N \exp \left[-\frac{\mathcal{V}(\mathbf{r}^N)}{k_B T} \right] \quad (8.10)$$

A consequence of writing the partition function as a product of a real gas and an ideal gas part is that thermodynamic properties can be written in terms of an ideal gas value and an excess value. The ideal gas contributions can be determined analytically by integrating over the momenta. For example, the Helmholtz free energy is related to the canonical partition function by:

$$A = -k_B T \ln Q_{NVT} \quad (8.11)$$

Writing the partition function as the product, Equation (8.9), leads to:

$$A = A^{\text{ideal}} + A^{\text{excess}} \quad (8.12)$$

The important conclusion is that all of the deviations from ideal gas behaviour are due to the presence of interactions between the atoms in the system, as calculated using the potential energy function. This energy function is dependent only upon the positions of the atoms and not their momenta, and so a Monte Carlo simulation is able to calculate the excess contributions that give rise to deviations from ideal gas behaviour.

8.2 Calculating Properties by Integration

Having established that we can indeed explore configurational phase space and derive useful thermodynamic properties, let us consider how we might achieve this in practice. For example, the average potential energy can, in principle at least, be determined by evaluating the integral:

$$\langle \mathcal{V}(\mathbf{r}^N) \rangle = \int d\mathbf{r}^N \mathcal{V}(\mathbf{r}^N) \rho(\mathbf{r}^N) \quad (8.13)$$

This is a multidimensional integral over the $3N$ degrees of freedom of the N particles in the system. $\rho(\mathbf{r}^N)$ is the probability of obtaining the configuration \mathbf{r}^N and is given by

$$\rho(\mathbf{r}^N) = \frac{\exp[-\mathcal{V}(\mathbf{r}^N)/k_B T]}{Z} \quad (8.14)$$

The denominator, Z , is the configurational integral (Equation (8.6)). For the potential functions commonly used in molecular modelling, it is not possible to evaluate these integrals analytically. However, we could attempt to obtain values for the integrals using numerical methods. One simple numerical integration method is the trapezium rule. This approximates the integral as a series of trapeziums between the two limits, as illustrated for a one-dimensional problem in Figure 8.1. In this case we have divided the integral into ten trapeziums, which requires eleven function evaluations. Simpson's rule involves a similar procedure and may provide a more accurate value of the integral [Stephenson 1973]. For a function of two variables ($f(x, y)$), it is necessary to square the number of function

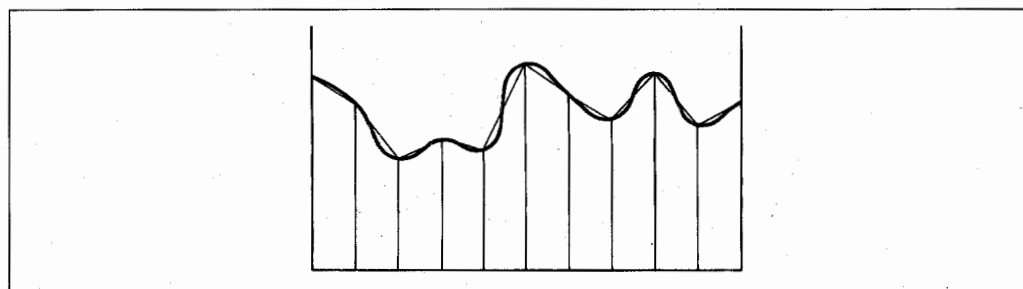


Fig. 8.1: Evaluation of a one-dimensional integral using the trapezium rule. The area under the curve is approximated as the sum of the areas of the trapeziums.

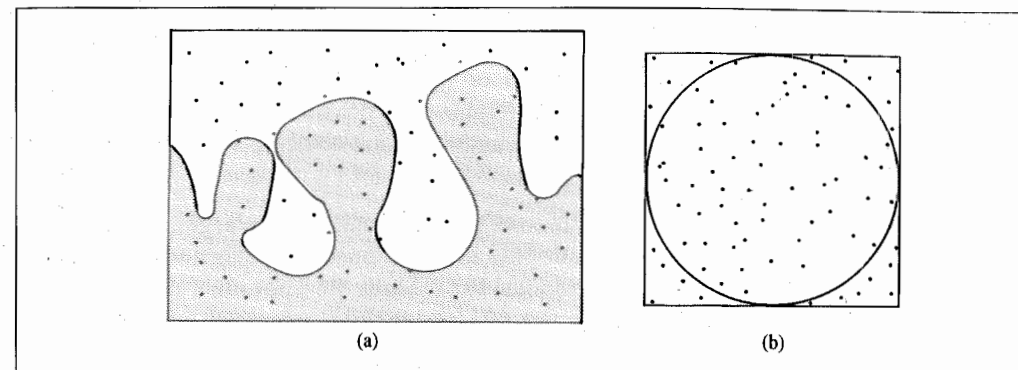


Fig. 8.2: Simple Monte Carlo integration. (a) The shaded area under the irregular curve equals the ratio of the number of random points under the curve to the total number of points, multiplied by the area of the bounding area. (b) An estimate of π can be obtained by generating random numbers within the square. π then equals the number of points within the circle divided by the total number of points within the square, multiplied by 4.

evaluations required. For a $3N$ -dimensional integral the total number of function evaluations required to determine the integral would be m^{3N} , where m is the number of points needed to determine the integral in each dimension. This number is enormous even for very small numbers of particles. For example, with just 50 particles and three points per dimension, a total of 3^{150} ($\sim 10^{71}$) evaluations would be required. Integration using the trapezium rule or Simpson's rule is clearly not a feasible approach.

We could consider a random method as a possible alternative. The general principle can be illustrated using the function shown in Figure 8.2. To determine the area under the curve in Figure 8.2 a series of random points would be generated within the bounding area. The area under the curve is then calculated by multiplying the bounding area A by the ratio of the number of trial points that lie under the curve to the total number of points generated. An estimate of π can be determined in this way, as illustrated in Figure 8.2.

To calculate the partition function for a system of N atoms using this simple Monte Carlo integration method would involve the following steps:

1. Obtain a configuration of the system by randomly generating $3N$ Cartesian coordinates, which are assigned to the particles.
2. Calculate the potential energy of the configuration, $\mathcal{V}(\mathbf{r}^N)$.
3. From the potential energy, calculate the Boltzmann factor, $\exp(-\mathcal{V}(\mathbf{r}^N)/k_B T)$.
4. Add the Boltzmann factor to the accumulated sum of Boltzmann factors and the potential energy contribution to its accumulated sum and return to step 1.
5. After a number, N_{trial} , of iterations, the mean value of the potential energy would be calculated using:

$$\langle \mathcal{V}(\mathbf{r}^N) \rangle = \frac{\sum_{i=1}^{N_{\text{trial}}} \mathcal{V}_i(\mathbf{r}^N) \exp[-\mathcal{V}_i(\mathbf{r}^N)/k_B T]}{\sum_{i=1}^{N_{\text{trial}}} \exp[-\mathcal{V}_i(\mathbf{r}^N)/k_B T]} \quad (8.15)$$

Unfortunately, this is not a feasible approach for calculating thermodynamic properties due to the large number of configurations that have extremely small (effectively zero) Boltzmann

factors caused by high-energy overlaps between the particles. This reflects the nature of the phase space, most of which corresponds to non-physical configurations with very high energies. Only a very small proportion of the phase space corresponds to low-energy configurations where there are no overlapping particles and where the Boltzmann factor has an appreciable value. These low-energy regions coincide with the physically observed phases such as solid, liquid, etc.

One way around this impasse is to generate configurations that make a large contribution to the integral (8.15), which is the strategy adopted in importance sampling and which is the essence of the method described by Metropolis, Rosenbluth, Rosenbluth, Teller and Teller in 1953 [Metropolis *et al.* 1953]. For many thermodynamic properties of a molecular system, those states with a high probability ρ are also the ones that make a significant contribution to the integral (there are some notable exceptions to this, such as the free energy). The Metropolis method has become so widely adopted that in the simulation and molecular modelling communities it is usually referred to as 'the Monte Carlo method'. Fortunately, there is rarely any confusion with the simple Monte Carlo methods. The crucial feature of the Metropolis approach is that it biases the generation of configurations towards those that make the most significant contribution to the integral. Specifically, it generates states with a probability $\exp(-\mathcal{V}(\mathbf{r}^N)/k_B T)$ and then counts each of them equally. By contrast, the simple Monte Carlo integration method generates states with equal probability (both high- and low-energy) and then assigns them a weight $\exp(-\mathcal{V}(\mathbf{r}^N)/k_B T)$.

8.3 Some Theoretical Background to the Metropolis Method

The Metropolis algorithm generates a *Markov chain* of states. A Markov chain satisfies the following two conditions:

1. The outcome of each trial depends only upon the preceding trial and not upon any previous trials.
2. Each trial belongs to a finite set of possible outcomes.

Condition (1) provides a clear distinction between the molecular dynamics and Monte Carlo methods, for in a molecular dynamics simulation all of the states are connected in time. Suppose the system is in a state m . We denote the probability of moving to state n as π_{mn} . The various π_{mn} can be considered to constitute an $N \times N$ matrix π (the transition matrix), where N is the number of possible states. Each row of the transition matrix sums to 1 (i.e. the sum of the probabilities π_{mn} for a given m equals 1). The probability that the system is in a particular state is represented by a probability vector ρ :

$$\rho = (\rho_1, \rho_2, \dots, \rho_m, \rho_n, \dots, \rho_N) \quad (8.16)$$

Thus ρ_1 is the probability that the system is in state 1 and ρ_m the probability that the system is in state m . If $\rho(1)$ represents the initial (randomly chosen) configuration, then the probability of the second state is given by:

$$\rho(2) = \rho(1)\pi \quad (8.17)$$

The probability of the third state is:

$$\rho(3) = \rho(2)\pi = \rho(1)\pi\pi \quad (8.18)$$

The equilibrium distribution of the system can be determined by considering the result of applying the transition matrix an infinite number of times. This limiting distribution of the Markov chain is given by $\rho_{\text{limit}} = \lim_{N \rightarrow \infty} \rho(1)\pi^N$.

One feature of the limiting distribution is that it is independent of the initial guess $\rho(1)$. The limiting or equilibrium distribution for a molecular or atomic system is one in which the probabilities of each state are proportional to the Boltzmann factor. We can illustrate the use of the probability distribution and the transition matrix by considering a two-level system in which the energy levels are such that the ratio of the Boltzmann factors is 2:1. The expected limiting distribution thus corresponds to a configuration vector $(\frac{2}{3}, \frac{1}{3})$. The following transition matrix enables the limiting distribution to be achieved:

$$\pi = \begin{pmatrix} 0.5 & 0.5 \\ 1 & 0 \end{pmatrix} \quad (8.19)$$

We can illustrate the use of this transition matrix as follows. Suppose the initial probability vector is $(1, 0)$ and so the system starts with a 100% probability of being in state 1 and no probability of being in state 2. Then the second state is given by:

$$\rho(2) = (1 \ 0) \begin{pmatrix} 0.5 & 0.5 \\ 1 & 0 \end{pmatrix} = (0.5 \ 0.5) \quad (8.20)$$

The third state is $\rho(3) = (0.75, 0.25)$. Successive applications of the transition matrix give the limiting distribution $(2/3, 1/3)$.

When the limiting distribution is reached then application of the transition matrix must return the same distribution back:

$$\rho_{\text{limit}} = \rho_{\text{limit}}\pi \quad (8.21)$$

Thus, if an ensemble can be prepared that is at equilibrium, then one Metropolis Monte Carlo step should return an ensemble that is still at equilibrium. A consequence of this is that the elements of the probability vector for the limiting distribution must satisfy:

$$\sum_m \rho_m \pi_{mn} = \rho_n \quad (8.22)$$

This can be seen to hold for our simple two-level example:

$$(2/3 \ 1/3) \begin{pmatrix} 1/2 & 1/2 \\ 1 & 0 \end{pmatrix} = (2/3 \ 1/3) \quad (8.23)$$

We will henceforth use the symbol ρ to refer to the limiting distribution.

Closely related to the transition matrix is the *stochastic matrix*, whose elements are labelled α_{mn} . This matrix gives the probability of choosing the two states m and n between which the move is to be made. It is often known as the *underlying matrix* of the Markov chain. If the probability of accepting a trial move from m to n is p_{mn} then the probability of making a transition from m to n (π_{mn}) is given by multiplying the probability of choosing states m

and n (α_{mn}) by the probability of accepting the trial move (p_{mn}):

$$\pi_{mn} = \alpha_{mn} p_{mn} \quad (8.24)$$

It is often assumed that the stochastic matrix \mathbf{a} is symmetrical (i.e. the probability of choosing the states m and n is the same whether the move is made from m to n or from n to m). If the probability of state n is greater than that of state m in the limiting distribution (i.e. if the Boltzmann factor of n is greater than that of m because the energy of n is lower than the energy of m) then in the Metropolis recipe, the transition matrix element π_{mn} for progressing from m to n equals the probability of selecting the two states in the first place (i.e. $\pi_{mn} = \alpha_{mn}$ if $\rho_n \geq \rho_m$). If the Boltzmann weight of the state n is less than that of state m , then the probability of permitting the transition is given by multiplying the stochastic matrix element α_{mn} by the ratio of the probabilities of the state n to the previous state m . This can be written:

$$\pi_{mn} = \alpha_{mn} \quad (\rho_n \geq \rho_m) \quad (8.25)$$

$$\pi_{mn} = \alpha_{mn} (\rho_n / \rho_m) \quad (\rho_n < \rho_m) \quad (8.26)$$

These two conditions apply if the initial and final states m and n are different. If m and n are the same state, then the transition matrix element is calculated from the fact that the rows of the stochastic matrix sum to 1:

$$\pi_{mm} = 1 - \sum_{m \neq n} \pi_{mn} \quad (8.27)$$

Let us now try to reconcile the Metropolis algorithm as outlined in Section 6.1.3 with the more formal approach that we have just developed. We recall that in the Metropolis method a new configuration n is accepted if its energy is lower than the original state m . If the energy is higher, however, then we would like to choose the move with a probability according to Equation (8.24). This is achieved by comparing the Boltzmann factor $\exp(-\Delta\mathcal{V}(\mathbf{r}^N)/k_B T)$ ($\Delta\mathcal{V}(\mathbf{r}^N) = [\mathcal{V}(\mathbf{r}^N)_n - \mathcal{V}(\mathbf{r}^N)_m]$) to a random number between 0 and 1. If the Boltzmann factor is greater than the random number then the new state is accepted. If it is smaller then the new state is rejected. Thus if the energy of the new state (n) is very close to that of the old state (m) then the Boltzmann factor of the energy difference will be very close to 1, and so the move is likely to be accepted. If the energy difference is very large, however, then the Boltzmann factor will be close to zero and the move is unlikely to be accepted.

The Metropolis method is derived by imposing the condition of microscopic reversibility: at equilibrium the transition between two states occurs at the same rate. The rate of transition from a state m to a state n equals the product of the population (ρ_m) and the appropriate element of the transition matrix (π_{mn}). Thus, at equilibrium we can write:

$$\pi_{nm} \rho_m = \pi_{mn} \rho_n \quad (8.28)$$

The ratio of the transition matrix elements thus equals the ratio of the Boltzmann factors of the two states:

$$\frac{\pi_{mn}}{\pi_{nm}} = \exp[-(\mathcal{V}(\mathbf{r}^N)_n - \mathcal{V}(\mathbf{r}^N)_m)/k_B T] \quad (8.29)$$

8.4 Implementation of the Metropolis Monte Carlo Method

A Monte Carlo program to simulate an atomic fluid is quite simple to construct. At each iteration of the simulation a new configuration is generated. This is usually done by making a random change to the Cartesian coordinates of a single randomly chosen particle using a random number generator. If the random number generator produces numbers (ξ) in the range 0 to 1, moves in both positive and negative directions are possible if the coordinates are changed as follows:

$$x_{\text{new}} = x_{\text{old}} + (2\xi - 1)\delta r_{\text{max}} \quad (8.30)$$

$$y_{\text{new}} = y_{\text{old}} + (2\xi - 1)\delta r_{\text{max}} \quad (8.31)$$

$$z_{\text{new}} = z_{\text{old}} + (2\xi - 1)\delta r_{\text{max}} \quad (8.32)$$

A unique random number is generated for each of the three directions x , y and z . δr_{max} is the maximum possible displacement in any direction. The energy of the new configuration is then calculated; this need not require a complete recalculation of the energy of the entire system but only those contributions involving the particle that has just been moved. As a consequence, the neighbour list used by a Monte Carlo simulation must contain *all* the neighbours of each atom, because it is necessary to identify all the atoms which interact with the moving atom (recall that in molecular dynamics the neighbour list for each atom contains only neighbours with a higher index). Proper account should be taken of periodic boundary conditions and the minimum image convention when generating new configurations and calculating their energies. If the new configuration is lower in energy than its predecessor then the new configuration is retained as the starting point for the next iteration. If the new configuration is higher in energy than its predecessor then the Boltzmann factor, $\exp(-\Delta\mathcal{V}/k_B T)$, is compared to a random number between 0 and 1. If the Boltzmann factor is greater than the random number then the new configuration is accepted; if not then it is rejected and the initial configuration is retained for the next move. This acceptance condition can be written in the following concise fashion:

$$\text{rand}(0, 1) \leq \exp(-\Delta\mathcal{V}(\mathbf{r}^N)/k_B T) \quad (8.33)$$

The size of the move at each iteration is governed by the maximum displacement, δr_{max} . This is an adjustable parameter whose value is usually chosen so that approximately 50% of the trial moves are accepted. If the maximum displacement is too small then many moves will be accepted but the states will be very similar and the phase space will only be explored very slowly. Too large a value of δr_{max} and many trial moves will be rejected because they lead to unfavourable overlaps. The maximum displacement can be adjusted automatically while the program is running to achieve the desired acceptance ratio by keeping a running score of the proportion of moves that are accepted. Every so often the maximum displacement is then scaled by a few percent: if too many moves have been accepted then the maximum displacement is increased; too few and δr_{max} is reduced.

As an alternative to the random selection of particles it is possible to move the atoms sequentially (this requires one fewer call to the random number generator per iteration). Alternatively, several atoms can be moved at once; if an appropriate value for the maximum displacement is chosen then this may enable phase space to be covered more efficiently.

As with a molecular dynamics simulation, a Monte Carlo simulation comprises an equilibration phase followed by a production phase. During equilibration, appropriate thermodynamic and structural quantities such as the total energy (and the partitioning of the energy among the various components), mean square displacement and order parameters (as appropriate) are monitored until they achieve stable values, whereupon the production phase can commence. In a Monte Carlo simulation from the canonical ensemble, the temperature and volume are, of course, fixed. In a constant pressure simulation the volume will change and should therefore also be monitored to ensure that a stable system density is achieved.

8.4.1 Random Number Generators

The random number generator at the heart of every Monte Carlo simulation program is accessed a very large number of times, not only to generate new configurations but also to decide whether a given move should be accepted or not. Random number generators are also used in other modelling applications; for example, in a molecular dynamics simulation the initial velocities are normally assigned using a random number generator. The numbers produced by a random number generator are not, in fact, truly random; the same sequence of numbers should always be generated when the program is run with the same initial conditions (if not, then a serious error in the hardware or software must be suspected!). The sequences of numbers are thus often referred to as 'pseudo-random' numbers as they possess the statistical properties of 'true' sequences of random numbers. Most random number generators are designed to generate different sequences of numbers if a different 'seed' is provided. In this way, several independent runs can be performed using different seeds. One simple strategy is to use the time and/or date as the seed; this is information that can often be obtained automatically by the program from the computer's operating system.

The numbers produced by a random number generator should satisfy certain statistical properties. This requirement usually supersedes the need for a computationally very fast algorithm as other parts of a Monte Carlo simulation take much more time (such as calculating the change in energy). One useful and simple test of a random number generator is to break a sequence of random numbers into blocks of k numbers, which are taken to be coordinates in a k -dimensional space. A good random number should give a random distribution of points. Many of the common generators do not satisfy this test because the points lie on a plane or because they show clear correlations [Sharp and Bays 1992].

The *linear congruential* method is widely used for generating random numbers. Each number in the sequence is generated by taking the previous number, multiplying by a constant (the multiplier, a), adding a second constant (the increment, b), and taking the remainder when dividing by a third constant (the modulus, m). The first value is the seed, supplied by the user. Thus:

$$\xi[1] = \text{seed} \quad (8.34)$$

$$\xi[i] = \text{MOD}\{(\xi[i-1] \times a + b), m\} \quad (8.35)$$

The MOD function returns the remainder when the first argument is divided by the second (for example, MOD(14,5) equals 4). If the constants are chosen carefully, the linear congruential method generates all possible integers between 0 and $m-1$, and the period (i.e. the number of iterations before the sequence starts to repeat itself) will be equal to

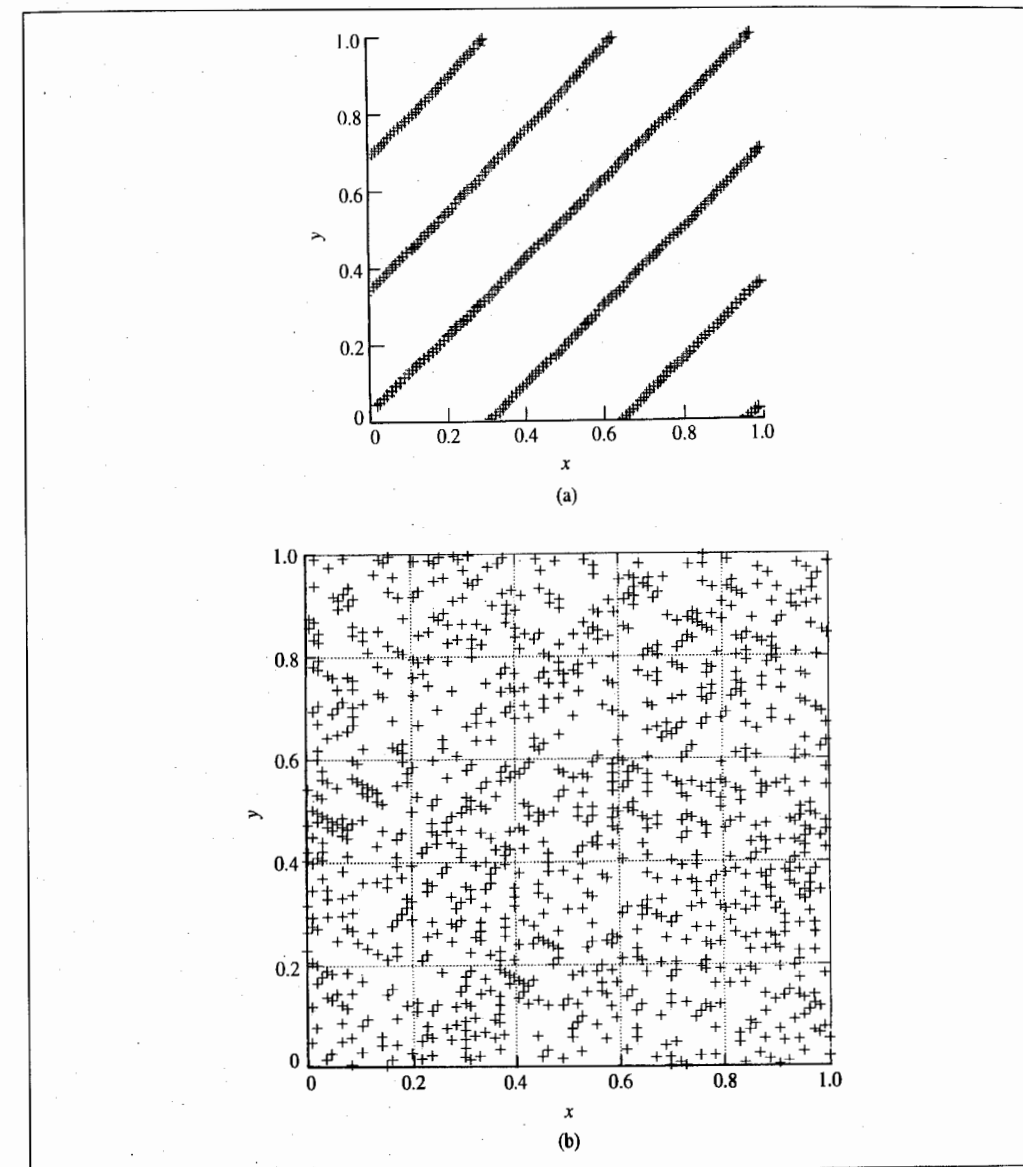


Fig. 8.3: Two 'random' distributions obtained by plotting pairs of values from a linear congruential random generator. The distribution (a) was obtained using $m = 32\,769$, $a = 10\,924$, $b = 11\,830$. The distribution (b) was obtained using $m = 6075$, $a = 106$, $b = 1283$. Data from [Sharp and Bays 1992].

the modulus. The period cannot of course be greater than m . The linear congruential method generates integral values, which can be converted to real numbers between 0 and 1 by dividing by m . The modulus is often chosen to be the largest prime number that can be represented in a given number of bits (usually chosen to be the number of bits per word; $2^{31} - 1$ is thus a common choice on a 32-bit machine).

Although popular, by virtue of the ease with which it can be programmed, the linear congruential method does not satisfy all of the requirements that are now regarded as important in a random number generator. For example, the points obtained from a linear congruential generator lie on $(k - 1)$ -dimensional planes rather than uniformly filling up the space. Indeed, if the constants a , b and m are chosen inappropriately then the linear congruential method can give truly terrible results, as shown in Figure 8.3. One random number generator that is claimed to perform well in all of the standard tests is that of G Marsaglia, which is described in Appendix 8.1.

8.5 Monte Carlo Simulation of Molecules

The Monte Carlo method is most easily implemented for atomic systems because it is only necessary to consider the translational degrees of freedom. The algorithm is easy to implement and accurate results can be obtained from relatively short simulations of a few tens of thousands of steps. There can be practical problems in applying the method to molecular systems, and especially to molecules which have a significant degree of conformational flexibility. This is because, in such systems, it is necessary to permit the internal degrees of freedom to vary. Unfortunately, such changes often lead to high-energy overlaps either within the molecule or between the molecule and its neighbours and thus a high rejection rate.

8.5.1 Rigid Molecules

For rigid, non-spherical molecules, the orientations of the molecules must be varied as well as their positions in space. It is usual to translate and rotate one molecule during each Monte Carlo step. Translations are usually described in terms of the position of the centre of mass. There are various ways to generate a new orientation of a molecule. The simplest approach is to choose one of the three Cartesian axes (x , y or z) and to rotate about the chosen axis by a randomly chosen angle $\delta\omega$, chosen to lie within the maximum angle variation, $\delta\omega_{\max}$ [Barker and Watts 1969]. The rotation is achieved by applying routine trigonometric relationships. For example, if the vector (xi, yj, zk) describes the orientation of a molecule then the new vector $(x'i, y'j, z'k)$ that corresponds to rotation by $\delta\omega$ about the x axis is calculated as follows:

$$\begin{pmatrix} x' \\ y' \\ z' \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos \delta\omega & \sin \delta\omega \\ 0 & -\sin \delta\omega & \cos \delta\omega \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} \quad (8.36)$$

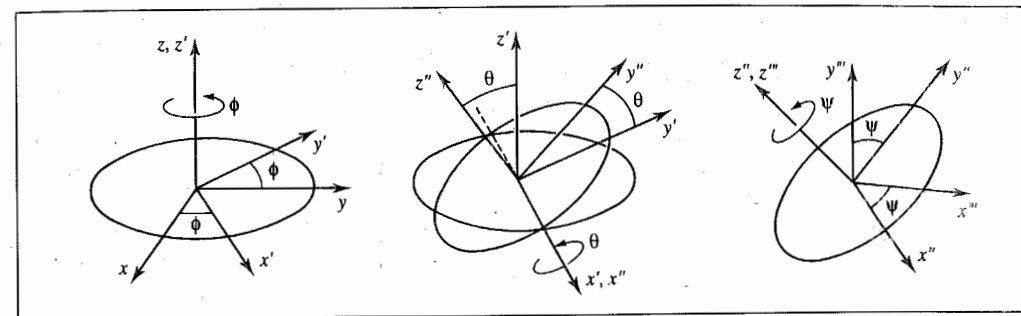


Fig. 8.4: The Euler angles ϕ , θ and ψ .

The Euler angles are often used to describe the orientations of a molecule. There are three Euler angles; ϕ , θ and ψ . ϕ is a rotation about the Cartesian z axis; this has the effect of moving the x and y axes. θ is a rotation about the new x axis. Finally, ψ is a rotation about the new z axis (Figure 8.4). If the Euler angles are randomly changed by small amounts $\delta\phi$, $\delta\theta$ and $\delta\psi$ then a vector \mathbf{v}_{old} is moved according to the following matrix equation:

$$\mathbf{v}_{\text{new}} = \mathbf{A}\mathbf{v}_{\text{old}} \quad (8.37)$$

where the matrix \mathbf{A} is

$$\begin{pmatrix} \cos \delta\phi \cos \delta\psi - \sin \delta\phi \cos \delta\theta \sin \delta\psi & \sin \delta\phi \cos \delta\psi + \cos \delta\phi \cos \delta\theta \sin \delta\psi & \sin \delta\theta \sin \delta\psi \\ -\cos \delta\phi \sin \delta\psi - \sin \delta\phi \cos \delta\theta \cos \delta\psi & -\sin \delta\phi \sin \delta\psi + \cos \delta\phi \cos \delta\theta \cos \delta\psi & \sin \delta\theta \cos \delta\psi \\ \sin \delta\phi \sin \delta\theta & -\cos \delta\phi \sin \delta\theta & \cos \delta\theta \end{pmatrix} \quad (8.38)$$

It is important to note that simply sampling displacements of the three Euler angles does not lead to a uniform distribution; it is necessary to sample from $\cos \theta$ rather than θ (Figure 8.5).

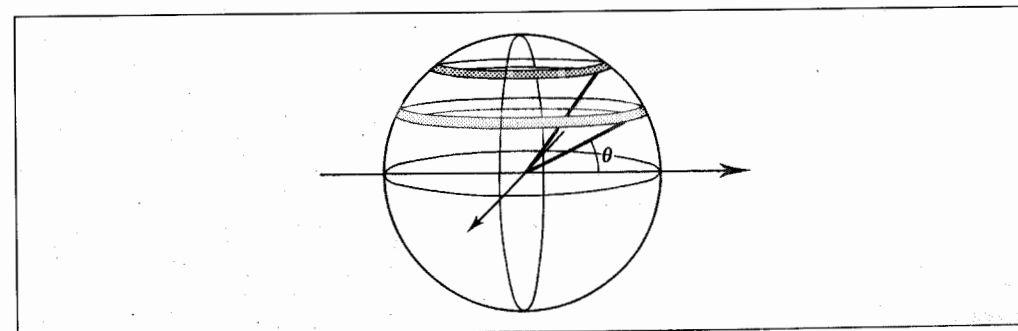


Fig. 8.5: To achieve a uniform distribution of points over the surface of a sphere it is necessary to sample from $\cos \theta$ rather than θ . If the sampling is uniform in θ then the number of points per unit area increases with θ , leading to an uneven distribution over the sphere.

The preferred approach is to sample directly in $\cos \theta$ as follows:

$$\phi_{\text{new}} = \phi_{\text{old}} + 2(\xi - 1)\delta\phi_{\text{max}} \quad (8.39)$$

$$\cos \theta_{\text{new}} = \cos \theta_{\text{old}} + (2\xi - 1)\delta(\cos \theta)_{\text{max}} \quad (8.40)$$

$$\psi_{\text{new}} = \psi_{\text{old}} + 2(\xi - 1)\delta\psi_{\text{max}} \quad (8.41)$$

The alternative is to sample in θ and to modify the acceptance or rejection criteria as follows:

$$\theta_{\text{new}} = \theta_{\text{old}} + (2\xi - 1)\delta\theta_{\text{max}} \quad (8.42)$$

$$\frac{\rho_{\text{new}}}{\rho_{\text{old}}} = \exp(-\Delta\mathcal{V}/k_{\text{B}}T) \frac{\sin \theta_{\text{new}}}{\sin \theta_{\text{old}}} \quad (8.43)$$

This second approach may give problems if θ_{old} equals zero.

A disadvantage of the Euler angle approach is that the rotation matrix contains a total of six trigonometric functions (sine and cosine for each of the three Euler angles). These trigonometric functions are computationally expensive to calculate. An alternative is to use *quaternions*. A quaternion is a four-dimensional vector such that its components sum to 1: $q_0^2 + q_1^2 + q_2^2 + q_3^2 = 1$. The quaternion components are related to the Euler angles as follows:

$$q_0 = \cos \frac{1}{2}\theta \cos \frac{1}{2}(\phi + \psi) \quad (8.44)$$

$$q_1 = \sin \frac{1}{2}\theta \cos \frac{1}{2}(\phi + \psi) \quad (8.45)$$

$$q_2 = \sin \frac{1}{2}\theta \sin \frac{1}{2}(\phi + \psi) \quad (8.46)$$

$$q_3 = \cos \frac{1}{2}\theta \sin \frac{1}{2}(\phi + \psi) \quad (8.47)$$

The Euler angle rotation matrix can then be written

$$\mathbf{A} = \begin{pmatrix} q_0^2 + q_1^2 - q_2^2 - q_3^2 & 2(q_1q_2 + q_0q_3) & 2(q_1q_3 - q_0q_2) \\ 2(q_1q_2 - q_0q_3) & q_0^2 - q_1^2 + q_2^2 - q_3^2 & 2(q_2q_3 + q_0q_1) \\ 2(q_1q_3 + q_0q_2) & 2(q_2q_3 - q_0q_1) & q_0^2 - q_1^2 - q_2^2 + q_3^2 \end{pmatrix} \quad (8.48)$$

To generate a new orientation, it is necessary to rotate the quaternion vector to a new (random) orientation. As it is a four-dimensional vector, the orientation must be performed in four-dimensional space. This can be achieved as follows [Vesely 1982]:

1. Generate pairs of random numbers (ξ_1, ξ_2) between -1 and 1 until $S_1 = \xi_1^2 + \xi_2^2 < 1$.
2. Do the same for pairs ξ_3 and ξ_4 until $S_2 = \xi_3^2 + \xi_4^2 < 1$.
3. Form the random unit four-dimensional vector $(\xi_1, \xi_2, \xi_3\sqrt{(1-S_1)/S_2}, \xi_4\sqrt{(1-S_1)/S_2})$.

To achieve an appropriate acceptance rate the angle between the two vectors that describe the new and old orientations should be less than some value; this corresponds to sampling randomly and uniformly from a region on the surface of a sphere.

The introduction of an orientational component as well as a translational component increases the number of maximum displacement parameters that determine the acceptance ratio. It is important to check that the desired acceptance ratio is achieved, and also that an appropriate proportion of orientational and translational moves are made. Trial and error is often the most effective way to find the best combination of parameters.

8.5.2 Monte Carlo Simulations of Flexible Molecules

Monte Carlo simulations of flexible molecules are often difficult to perform successfully unless the system is small, or some of the internal degrees of freedom are frozen out, or special models or methods are employed. The simplest way to generate a new configuration of a flexible molecule is to perform random changes to the Cartesian coordinates of individual atoms, in addition to translations and rotations of the entire molecule. Unfortunately, it is often found that very small atomic displacements are required to achieve an acceptable acceptance ratio, which means that the phase space is covered very slowly. For example, even small movements away from an equilibrium bond length will cause a large increase in the energy. One obvious tactic is to freeze out some of the internal degrees of freedom, usually the 'hard' degrees of freedom such as the bond lengths and the bond angles. Such algorithms have been extensively used to investigate small molecules such as butane. However, for large molecules, even relatively small bond rotations may cause large movements of atoms down the chain. This invariably leads to high-energy configurations as illustrated in Figure 8.6. The rigid bond and rigid angle approximation must be used with care, for freezing out some of the internal degrees of freedom can affect the distributions of other internal degrees of freedom.

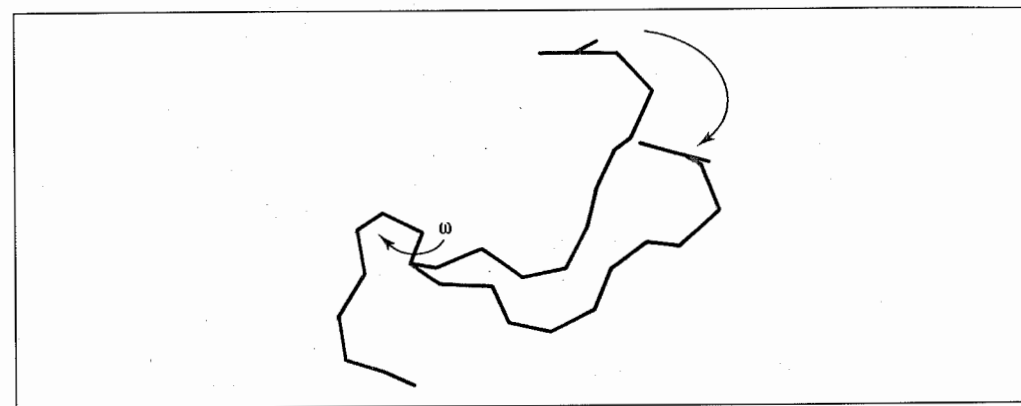


Fig. 8.6: A bond rotation in the middle of a molecule may lead to a large movement at the end.

8.6 Models Used in Monte Carlo Simulations of Polymers

A polymer is a macromolecule that is constructed by chemically linking together a sequence of molecular fragments. In simple synthetic polymers such as polyethylene or polystyrene all of the molecular fragments comprise the same basic unit (or monomer). Other polymers contain mixtures of monomers. Proteins, for example, are polypeptide chains in which each unit is one of the twenty amino acids. Cross-linking between different chains gives rise to yet further variations in the constitution and structure of a polymer. All of these features may affect the overall properties of the molecule, sometimes in a dramatic way. Moreover, one

may be interested in the properties of the polymer under different conditions, such as in solution, in a polymer melt or in the crystalline state. Molecular modelling can help to develop theories for understanding the properties of polymers and can also be used to predict their properties.

A wide range of time and length scales are needed to completely describe a polymer's behaviour. The timescale ranges from approximately 10^{-14} s (i.e. the period of a bond vibration) through to seconds, hours or even longer for collective phenomena. The size scale ranges from the 1–2 Å of chemical bonds to the diameter of a coiled polymer, which can be several hundreds of ångströms. Many kinds of model have been used to represent and simulate polymeric systems and predict their properties. Some of these models are based upon very simple ideas about the nature of the intra- and intermolecular interactions within the system but have nevertheless proved to be extremely useful. One famous example is Flory's rotational isomeric state model [Flory 1969]. Increasing computer performance now makes it possible to use techniques such as molecular dynamics and Monte Carlo simulations to study polymer systems.

Most simulations on polymers are performed using empirical energy models (though with faster computers and new methods it is becoming possible to apply quantum mechanics to larger and larger systems). Moreover, there are various ways in which the configurational and conformational degrees of freedom may be restricted so as to produce a computationally more efficient model. The simplest models use a lattice representation in which the polymer is constructed from connected interaction centres, which are required to occupy the vertices of a lattice. At the next level of complexity are the bead models, where the polymer is composed of a sequence of connected 'beads'. Each bead represents an 'effective monomer' and interacts with the other beads to which it is bonded and also with other nearby beads. The ultimate level of detail is achieved with the atomistic models, in which each non-hydrogen atom is explicitly represented (and sometimes all of the hydrogens as well). Our aim here is to give a flavour of the way in which Monte Carlo methods can be used to investigate polymeric systems. We divide the discussion into lattice and continuum models but recognise that there is a spectrum of models from the simplest to the most complex.

8.6.1 Lattice Models of Polymers

Lattice models have provided many insights into the behaviour of polymers despite the obvious approximations involved. The simplicity of a lattice model means that many states can be generated and examined very rapidly. Both two-dimensional and three-dimensional lattices are used. The simplest models use cubic or tetrahedral lattices in which successive monomers occupy adjacent lattice points (Figure 8.7). The energy models are usually very simple, in part to reflect the simplicity of the representation but also to permit the rapid calculation of the energy.

More complex models have been developed in which the lattice representation is closer to the 'true' geometry of the molecule. For example, in Figure 8.8 we show the bond fluctuation model of polyethylene, in which the 'bond' between successive monomers on the lattice

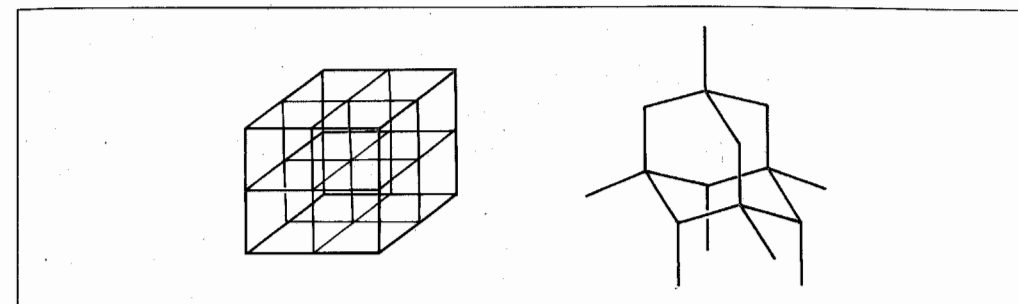


Fig. 8.7: Cubic and tetrahedral (diamond) lattices, which are commonly used for lattice simulations of polymers.

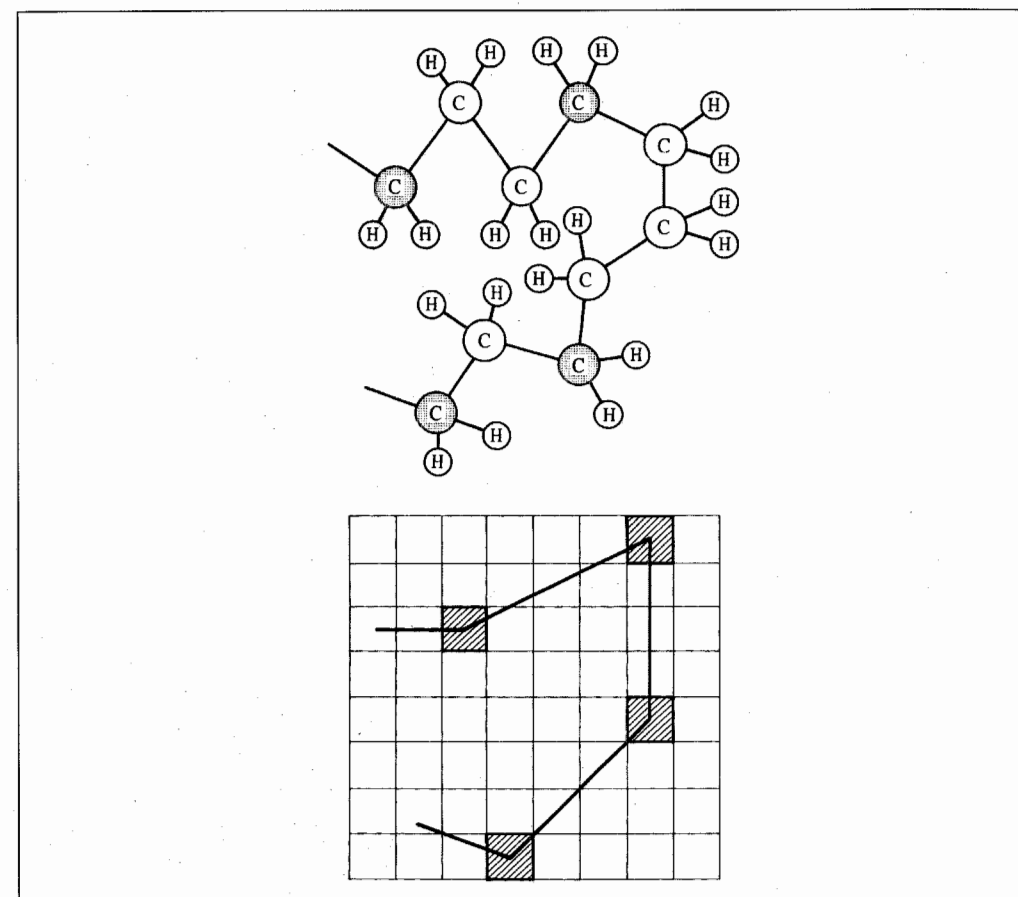


Fig. 8.8: The bond fluctuation model. In this example three bonds in the polymer are incorporated into a single 'effective bond' between 'effective monomers'. (Figure adapted from Baschnagel J, K Binder, W Paul, M Laso, U Suter, I Batoulis, W Jilge and T Bürger 1991. On the Construction of Coarse-Grained Models for Linear Flexible Polymer-Chains - Distribution-Functions for Groups of Consecutive Monomers. Journal of Chemical Physics 95:6014–6025.)

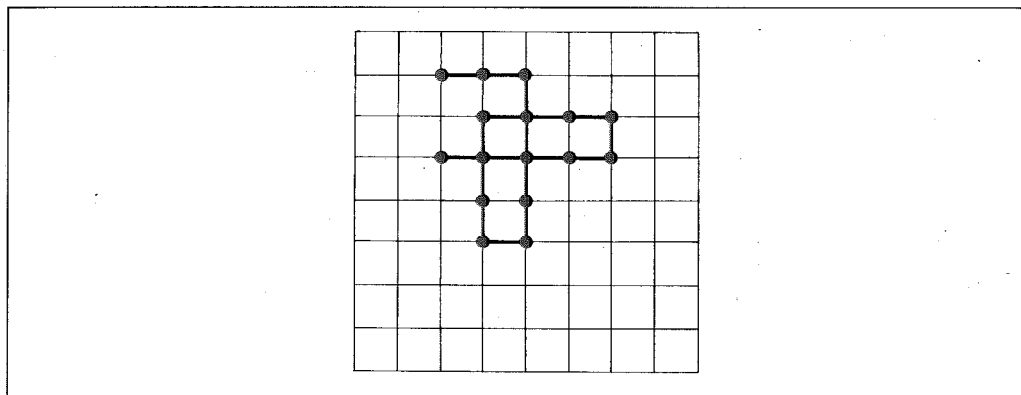


Fig. 8.9: In a random walk on a square lattice the chain can cross itself.

represent three bonds in the actual molecule [Baschnagel *et al.* 1991]. In this model each monomer is positioned at the centre of a cube within the lattice and five different distances are possible for the monomer–monomer bond lengths.

Lattices can be used to study a wide variety of polymeric systems, from single polymer chains to dense mixtures. The simplest type of simulation is a 'random walk', in which the chain is randomly grown in the lattice until it contains the desired number of bonds (Figure 8.9). In this model the chain is free to cross itself (i.e. excluded volume effects are ignored). Various properties can be calculated from such simulations, by averaging the results over a large number of trials. For example, a simple measure of the size of a polymer is the mean square end-to-end distance, $\langle R_n^2 \rangle$. For the random walk model $\langle R_n^2 \rangle$ is related to the number of bonds (n) and the length of each bond (l) by:

$$\langle R_n^2 \rangle = nl^2 \quad (8.49)$$

The radius of gyration is another commonly calculated property; this is the root mean square distance of each atom (or monomer) from the centre of mass. For the random walk model the radius of gyration $\langle s^2 \rangle$ is given in the asymptotic limit by:

$$\langle s^2 \rangle = \langle R_n^2 \rangle / 6 \quad (8.50)$$

The ability of the chain to cross itself in the random walk may seem to be a serious limitation, but it is found to be valid under some circumstances. When excluded volume effects are not important (also known as 'theta' conditions) then a subscript '0' is often added to properties such as the mean square end-to-end distance, $\langle R_n^2 \rangle_0$. Excluded volume effects can be taken into account by generating a 'self-avoiding walk' of the chain in the lattice (Figure 8.10). In this model only one monomer can occupy each lattice site. Self-avoiding walks have been used to exhaustively enumerate all possible conformations for a chain of a given length on the lattice. If all states are known then the partition function can be determined and thermodynamic quantities calculated. The 'energy' of each state may be calculated using an appropriate interaction model. For example, the energy may be proportional to the number of adjacent pairs of occupied lattice sites. A variation on this is to use polymers

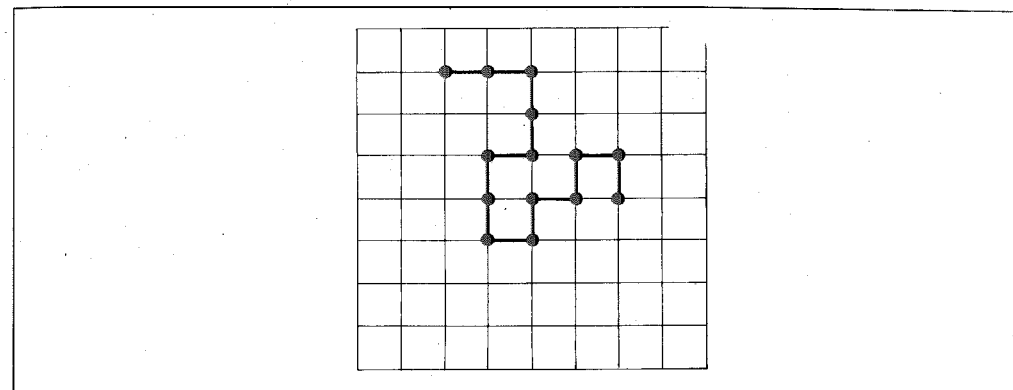


Fig. 8.10: Self-avoiding walk: only one monomer can occupy each lattice site.

consisting of two types of monomer (A and B), which have up to three different energy values: A–A, B–B and A–B. Again, the energy is determined by counting the number of occupied adjacent lattice sites. The relationship between the mean square end-to-end distance and the length of the chain (n) has been investigated intensively; with the self-avoiding walk the result obtained is different from the random walk, with $\langle R_n^2 \rangle$ being proportional to $n^{1.18}$ in the asymptotic limit.

Having grown a polymer onto the lattice, we now have to consider the generation of alternative configurations. Motion of the entire polymer chain or large-scale conformational changes are often difficult, especially for densely packed polymers. In variants of the Verdier–Stockmayer algorithm [Verdier and Stockmayer 1962] new configurations are generated using combinations of 'crankshaft', 'kink jump' and 'end rotation' moves (Figure 8.11). Another widely used algorithm in Monte Carlo simulations of polymers (not just in lattice models) is the 'slithering snake' model. Motion of the entire polymer chain is very difficult, especially for densely packed polymers, and one way in which the polymer can move is by wriggling around obstacles, a process known as *reptation*. To implement a slithering snake algorithm, one end of the polymer chain is randomly chosen as the 'head' and an attempt is made to grow a new bead at one of the available adjacent lattice positions. Each of the remaining beads is then advanced to that of its predecessor in the chain as illustrated in Figure 8.12. The procedure is then repeated. Even if it is impossible to move the chosen 'head', the configuration must still be included when ensemble averages are calculated.

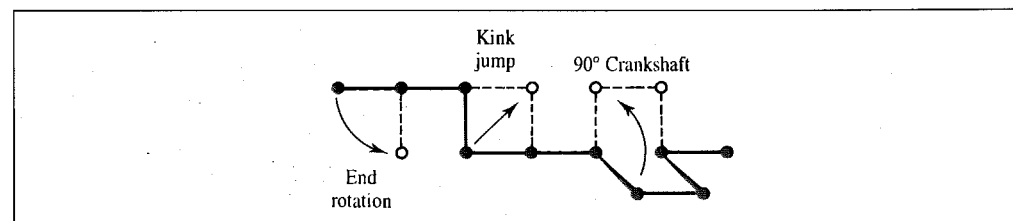


Fig. 8.11: The 'crankshaft', 'kink jump' and 'end rotation' moves used in Monte Carlo simulations of polymers.

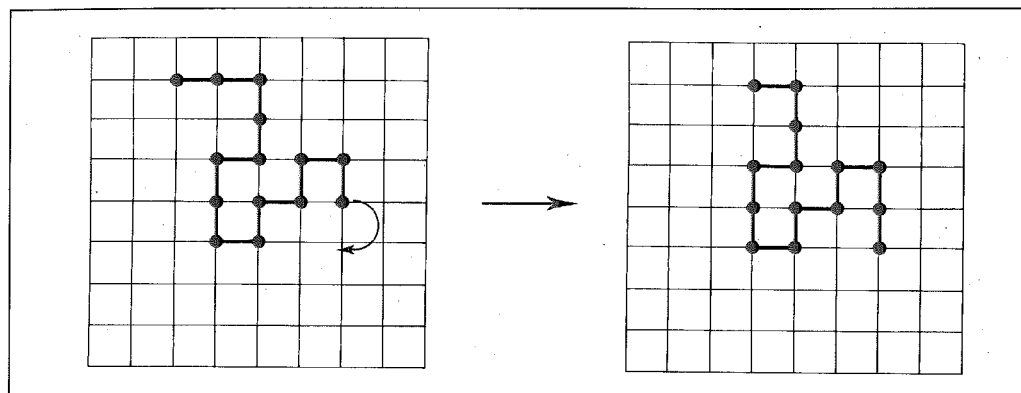


Fig. 8.12: The 'slithering snake' algorithm.

8.6.2 'Continuous' Polymer Models

The simplest of the continuous polymer models consists of a string of connected beads (Figure 8.13). The beads are freely jointed and interact with the other beads via a spherically symmetric potential such as the Lennard-Jones potential. The beads should not be thought of as being identical to the monomers in the polymer, though they are often referred to as such ('effective monomers' is a more appropriate term). Similarly, the links between the beads should not be thought of as bonds. The links may be modelled as rods of a fixed and invariant length or may be permitted to vary using a harmonic potential function.

In Monte Carlo studies with this freely jointed chain model the beads can sample from a continuum of positions. The *pivot algorithm* is one way that new configurations can be generated. Here, a segment of the polymer is randomly selected and rotated by a random amount, as illustrated in Figure 8.13. For isolated polymer chains the pivot algorithm can give a good sampling of the configurational/conformational space. However, for polymers in solution or in the melt, the proportion of accepted moves is often very small due to high-energy steric interactions.

The most unrealistic feature of the freely jointed chain model is the assumption that the bond angles can vary continuously. In the *freely rotating chain model* the bond angles are held fixed but free rotation is possible about the bonds, such that any torsion angle value between 0°

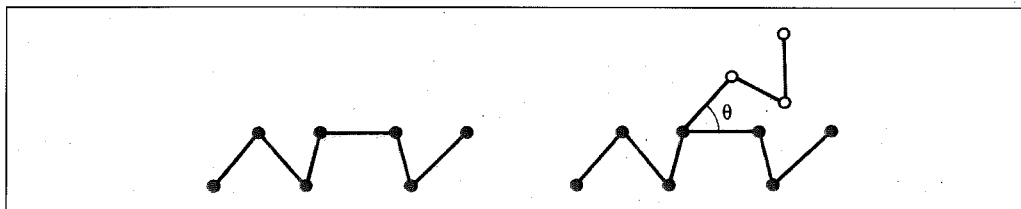


Fig. 8.13: The bead model for polymer simulations. The beads may be connected by stiff rods or by harmonic springs.

and 360° is equally likely. Fixing the bond angles in this way obviously affects the properties of the chain when compared to the freely jointed chain; one way to quantify this is via the characteristic ratio C_n , which is defined as:

$$C_n = \frac{\langle R_n^2 \rangle_0}{nl^2} \quad (8.51)$$

The characteristic ratio approximately indicates how extended the chain is. For the freely rotating chain the characteristic ratio is given by:

$$C_n = \frac{1 + \cos \theta'}{1 - \cos \theta'} - \frac{2 \cos \theta'}{n} \frac{1 - \cos^n \theta'}{(1 - \cos \theta')^2} \quad (8.52)$$

where θ' is the supplement of the normal bond angle (i.e. $\theta' = 180^\circ - \theta$). For an infinitely long chain the characteristic ratio becomes:

$$C_\infty = \frac{1 + \cos \theta'}{1 - \cos \theta'} \quad (8.53)$$

To move up the scale of complexity one now needs to consider the energetics of rotation about each bond. The simplest approach is to assume that each bond can be treated independently and that the total energy of the chain is the sum of the individual torsional energies for each bond. However, this particular model has some serious shortcomings arising from the assumption of independence.

The *rotational isomeric state model* (RIS) developed by Flory [Flory 1969] is probably the best known of the 'approximate' approaches to modelling polymer chains. Each bond is assumed to adopt one of a small number of discrete rotational states, which usually correspond to minima in the potential energy. For example, one might use three rotational states for a typical polyalkane, corresponding to the *trans*, *gauche(+)* and *gauche(-)* conformations. A key part of the RIS approach is its elegant use of various matrices to simplify the calculation. *Generator matrices* are used to establish certain conformation-dependent properties. Thus for a property A one would write:

$$A(\tau_1 \dots \tau_n) = \prod_{i=1}^n \mathbf{F}_i \quad (8.54)$$

where \mathbf{F}_i is the generator matrix for the particular property for bond i (with torsion angle τ_i). An example is the generator matrix for the square end-to-end distance R^2 , which takes the following form:

$$\mathbf{G}_i = \begin{bmatrix} 1 & 2\mathbf{l}^T \mathbf{T} & l^2 \\ 0 & \mathbf{T} & \mathbf{l} \\ 0 & 0 & 1 \end{bmatrix} \quad (8.55)$$

The vector \mathbf{l} is the bond vector for bond i and \mathbf{T} is the 3×3 matrix that transforms coordinates in the reference frame for bond $(i+1)$ to those in the frame of bond i . In this case the square end-to-end distance can be calculated from:

$$R^2 = \mathbf{G}_{[1} \mathbf{G}_2^{n-2} \mathbf{G}_n \quad (8.56)$$

The nomenclature is such that $\mathbf{G}_{[1}$ represents the first row of the matrix \mathbf{G}_1 and \mathbf{G}_n represents the last column of \mathbf{G}_n .

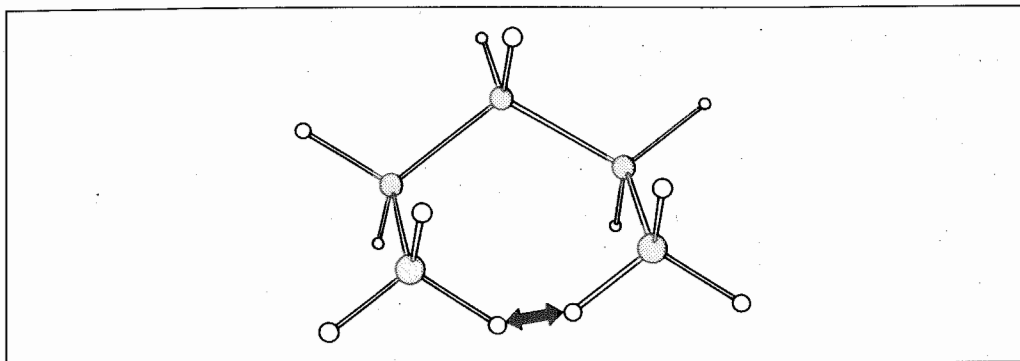


Fig. 8.14: Pentane violation.

In order to calculate average properties of the polymer chain one uses a standard statistical mechanical approach involving a summation over all possible conformations, with each term multiplied by the appropriate Boltzmann factor. This involves the use of a *statistical weights matrix*, which Flory introduced to deal with the influence a bond has on its neighbours. The pentane violation is the most important of these, wherein a sequence of gauche(+) and gauche(-) bonds gives rise to an unfavourable high-energy interaction (Figure 8.14). The statistical weights matrix associated with bond i has v_{i-1} rows and v_i columns, which correspond to the v_{i-1} rotational states of bond $(i-1)$ and the v_i rotational states of bond i . For example, for a typical polymer with the trans, gauche(+) and gauche(-) rotational states the statistical weights matrix is:

$$\mathbf{U}_i = \begin{bmatrix} u_{tt} & u_{tg^+} & u_{tg^-} \\ u_{g^+t} & u_{g^+g^+} & u_{g^+g^-} \\ u_{g^-t} & u_{g^-g^+} & u_{g^-g^-} \end{bmatrix} \quad (8.57)$$

Some typical values for the elements of this matrix would be:

$$\mathbf{U}_i = \begin{bmatrix} 1.0 & 0.54 & 0.54 \\ 1.0 & 0.54 & 0.05 \\ 1.0 & 0.05 & 0.54 \end{bmatrix} \quad (8.58)$$

The key point to note is the small weight (0.05) for adjacent gauche(+)-gauche(-) bonds. By combining the generator and statistical weights matrices it is possible to reduce the problem of calculating the average value of a property from a set of complex integrals to a series of straightforward matrix multiplications. Some of the properties that can be determined from the RIS model include the mean square end-to-end distance, the mean square radius of gyration and the mean square dipole moment.

The RIS model can be combined with the Monte Carlo simulation approach to calculate a wider range of properties than is available from the simple matrix multiplication method. In the RIS Monte Carlo method the statistical weight matrices are used to generate chain conformations with a probability distribution that is implied in their statistical weights.

Each conformation is generated by starting at one end of the chain and setting the backbone torsion angles one at a time until the entire chain has been constructed. The probability that a particular torsional state is selected for a given bond depends upon the *a priori* probabilities of each state and also upon the torsional state selected for the previous bond in the chain. These probabilities are used at each step by the Monte Carlo procedure to generate the whole chain. A large number of such chains is grown, calculating for each the properties of interest, which are then averaged. Properties which can be determined by the RIS-MC approach include the pair correlation function (which gives the relative probability of finding two atoms within the same chain separated by a distance r), the scattering function (which indicates how the polymer may scatter neutrons or X-rays) and the force-elongation relation (which gives the mean end-to-end distance of a chain subjected to an external force).

The ultimate level of detail in polymer modelling is achieved with the atomistic models, which as the name implies explicitly represent the atoms in the system. An atomistic model is clearly the closest to 'reality' and is necessary if one wishes to calculate accurately certain properties. One of the major problems with simulations of polymers that is particularly pertinent to the atomistic models is how to generate an initial configuration of the system. Amorphous polymers by definition do not adopt a characteristic and reproducible three-dimensional structure. It is important that the properties of the initial configuration are similar to the state one wishes to simulate else the computer time needed to move to the required state can be prohibitive. For short chains containing approximately 20–30 backbone bonds it is feasible to start from a regular crystalline structure, which can then be melted, but to 'melt' a long chain may require a prohibitive amount of computer time. For longer chains an initial configuration may be generated using a random walk and periodic boundary conditions (Figure 8.15). Such an arrangement will inevitably contain high-energy overlaps. These unfavourable interactions may be removed by relaxing the system using minimisation and/or computer simulation, during which the force field potentials are gradually turned on.

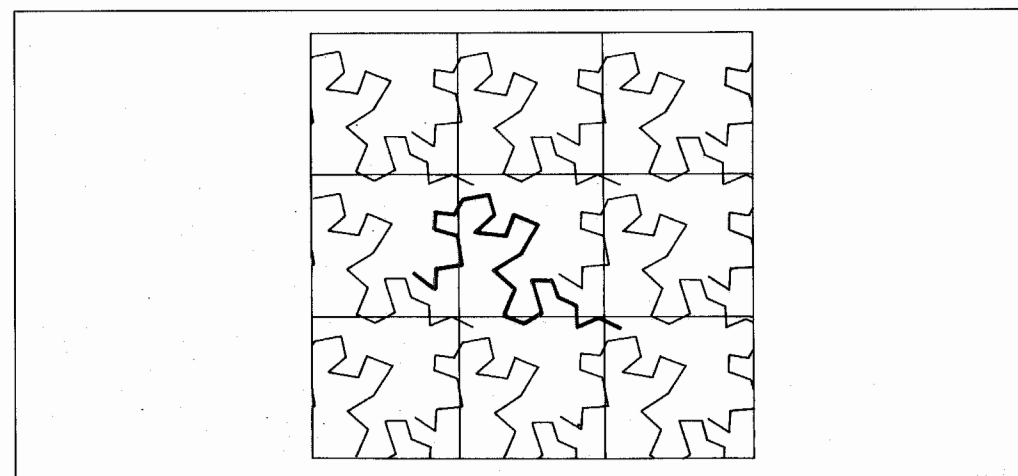


Fig. 8.15: Generation of an initial configuration of a polymer using periodic boundary conditions.

8.7 'Biased' Monte Carlo Methods

In some situations one is particularly interested in the behaviour of just a part of the system. For example, if we were simulating a solute-solvent system that contained a single solute molecule surrounded by a large number of solvent molecules then the behaviour of the solute and its interactions with the solvent would be of most interest. Solvent molecules far from the solute would be expected to behave almost like bulk solvent. A variety of techniques have been developed which can enhance the ability of the Monte Carlo method to explore the most important regions of phase space in such cases. One relatively simple procedure is *preferential sampling*, where the molecules in the vicinity of the solute are moved more frequently than those further away. This can be implemented by defining a cutoff region around the solute; molecules outside the cutoff region are moved less frequently than those inside the region as determined by a probability parameter p . At each Monte Carlo iteration a molecule is randomly chosen. If the molecule is inside the cutoff region then it is moved; if it is outside the region then a random number is generated between 0 and 1 and compared to the probability p . If p is greater than the random number a trial move is attempted; otherwise no move is made, no averages are accumulated, and a new molecule is randomly selected. The closer p is to zero, the more often 'closer' molecules are moved than 'further' molecules.

An alternative to the use of a fixed cutoff region is to relate the probability of choosing a solvent molecule to its distance from the solute, usually by some inverse power of the distance:

$$p \propto r^{-n} \quad (8.59)$$

In preferential sampling it is necessary to ensure that the correct procedures are followed when deciding whether to accept or reject a move in a manner that is consistent with the principle of microscopic reversibility. Suppose a molecule inside the cutoff region is moved outside the cutoff region. In the preferential sampling scheme the chances of the molecule now being selected for an out \rightarrow in move are less than for the original in \rightarrow out move and this must be taken into account when determining the acceptance criteria.

The *force-bias* Monte Carlo method [Pangali *et al.* 1978; Rao and Berne 1979] biases the movement according to the direction of the forces on it. Having chosen an atom or a molecule to move, the force on it is calculated. The force corresponds to the direction in which a 'real' atom or molecule would move. In the force-bias Monte Carlo method the random displacement is chosen from a probability distribution function that peaks in the direction of this force. The *smart Monte Carlo method* [Rossky *et al.* 1978] also requires the forces on the moving atom to be calculated. The displacement of an atom or molecule in this method has two components; one component is the force, and the other is a random vector $\delta\mathbf{r}_i^G$:

$$\delta\mathbf{r}_i = \frac{A\mathbf{f}_i}{k_B T} + \delta\mathbf{r}_i^G \quad (8.60)$$

where \mathbf{f}_i is the force on the atom and A is a parameter. The random displacement $\delta\mathbf{r}_i^G$ is chosen from a normal distribution with zero mean and variance equal to $2A$.

The main difference between the force-bias and the smart Monte Carlo methods is that the latter does not impose any limit on the displacement that an atom may undergo. The displacement in the force-bias method is limited to a cube of the appropriate size centred on the atom. However, in practice the two methods are very similar and there is often little to choose between them. In suitable cases they can be much more efficient at covering phase space and are better able to avoid bottlenecks in phase space than the conventional Metropolis Monte Carlo algorithm. The methods significantly enhance the acceptance rate of trial moves, thereby enabling larger moves to be made as well as simultaneous moves of more than one particle. However, the need to calculate the forces makes the methods much more elaborate, and comparable in complexity to molecular dynamics.

8.8 Tackling the Problem of Quasi-ergodicity: J-walking and Multicanonical Monte Carlo

If there are high-energy barriers between the potential energy minima in a system then a normal Metropolis Monte Carlo simulation may become trapped in just a few of the low-energy regions and fail to properly sample large regions of the thermally accessible space. Such a simulation may appear to possess all the qualities of a good simulation, in terms of its convergence, yet may give results that are completely incorrect. Such a simulation is often referred to as *quasi-ergodic*. This problem arises when studying systems as diverse as rare gas clusters near their melting temperature or protein folding, but it can also be demonstrated in even the simplest of model systems, such as the one-dimensional double-well potential (Figure 8.16). At low temperatures the simulation is unable to cross the high-energy barriers because of the favourable Boltzmann factor. A variety of methods have been suggested for tackling this problem, of which we shall consider two: J-walking and the multicanonical Monte Carlo method.

8.8.1 J-walking

In the J-walking (or Jump-walking) method [Frantz *et al.* 1990] a low-temperature Monte Carlo simulation is permitted occasionally to attempt jumps to regions of space that are accessible to a simulation run at a high temperature. The simplest way to implement this method is to perform the two simulations (at the high and low temperatures) in tandem. The low-temperature simulation is periodically permitted to attempt a jump to the configuration of the high-temperature simulation (the J-walker). The same Metropolis criteria are applied when deciding whether or not to accept the move. The high-temperature simulation will still in principle tend to be biased towards the low-energy regions so there will still be a reasonable chance that one of these special attempted jumps will be accepted.

In practice, it is found that this simple implementation is not the most effective approach. There are two particular problems. First, when the two simulations are run in tandem then significant correlations can arise, which results in large systematic errors. There are a number of ways to avoid these correlations, such as moving the J-walker an extra number

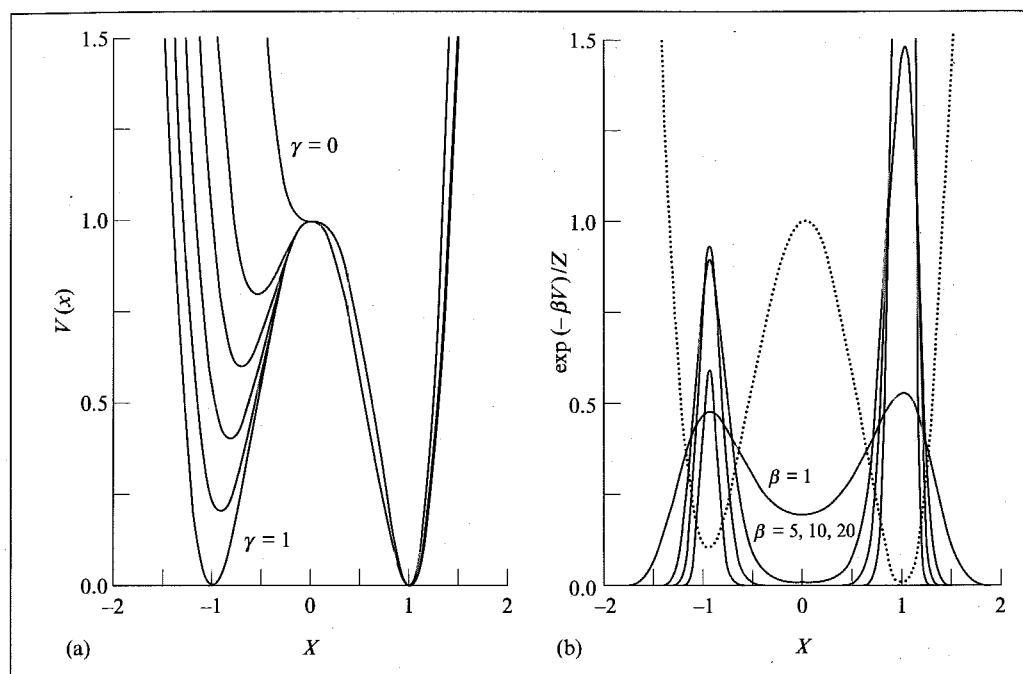


Fig. 8.16: Double-well potential which shows quasi-ergodicity. The potential is described by the quartic $\mathcal{V}(x) = 3\delta x^4 + 4\delta(\alpha - 1)x^3 - 6\delta\alpha x^2 + 1$, where $\delta = 1/(2\alpha + 1)$, and the potential is characterised by a parameter $\gamma (= (\alpha^4 + 2\alpha^3)/(2\alpha + 1))$. When $\gamma = 0$ there is just a single well and when $\gamma = 1$ the potential is symmetrical. On the right are Boltzmann distribution functions for the $\gamma = 0.9$ potential for various temperatures, expressed in terms of $\beta (= 1/k_B T)$. (Figure redrawn from Frantz, D D, D L Freeman and J D Doll 1990. Reducing quasi-ergodic behavior in Monte Carlo simulations by J-walking: applications to atomic clusters. *Journal of Chemical Physics* 93:2769–2784.)

of steps whenever a jump is attempted or running several high-temperature J-walkers and selecting jumps from them at random. However, it is found that a more effective approach is to perform the high-temperature simulation first. The low-temperature simulation first reads in and stores the configurations from the high-temperature simulation, which are then selected at random for each special jump. This approach obviously requires that sufficient storage is available to save the high-temperature configurations; however, it is not necessary to store every configuration (as there will be a high degree of correlation between successive configurations) but rather a representative sample.

One of the first applications of J-walking was to argon clusters [Frantz *et al.* 1990]. These clusters (containing up to 30 atoms) are experimentally known to exhibit strange melting behaviour, wherein the melting temperature is very heavily dependent upon the number of atoms in the cluster. Thus argon clusters Ar_n with seven, thirteen or nineteen atoms have particularly high melting temperatures, whereas those clusters with eight, fourteen or twenty atoms have particularly low melting temperatures. Moreover, certain clusters have different melting and freezing temperatures, implying that there is a range of

temperatures where both solid-like and liquid-like forms coexist. Simulation of these clusters in this transition region is difficult because of the problems of quasi-ergodicity and in obtaining satisfactory convergence. When applying the J-walking method to these clusters it was found necessary to generate the J-walker distribution in stages. The objective was to study an Ar_{13} cluster over a temperature range of 24–41 K using a J-walker at 50 K. However, the distribution of potential energies for the 50 K J-walker did not overlap with the distribution for a 20 K walker, which would mean that at 20 K hardly any of the attempted jumps to the J-walker distribution would be accepted. A series of simulations was thus performed. First, a 50 K J-walker was used to generate a distribution at 40 K, which in turn was used to produce a distribution at 30 K. Finally, the 30 K distribution acted as the J-walker for the 20 K simulation. Significantly better convergence for the cluster configurational energy and the heat capacity when compared with an equivalent standard Metropolis Monte Carlo procedure were noted, even when the simulations were started from a random configuration rather than the icosohedral low-energy geometry.

A related approach designed to be used for conformationally flexible molecules is the ‘jumping between wells’ (JBW) approach [Senderowitz *et al.* 1995]. Here a conformational analysis (see Chapter 9) is first performed on the molecule in order to identify its thermally accessible minimum-energy conformations (e.g. all within 5 kcal/mol of the global minimum-energy conformation). These minimum-energy conformations are stored in a list. The changes in internal coordinates that would be required to interconvert each pair of minimum-energy conformations are then determined. At each stage of the iterative cycle the minimum-energy conformation that is closest to the current structure is identified. A minimum-energy conformation is then selected at random from the conformation list and the appropriate transformation applied to the current structure. Small random changes in this new structure are then made to give a new trial structure, which is then accepted or rejected using the Metropolis criterion with reference to the initial starting structure. The process is then repeated. It is important with this method to avoid oversampling some of the potential energy wells which can occur if the combination of a conformational jump and the subsequent randomisation enters the space of a different minimum-energy conformation. This problem can be overcome by only including minimum-energy conformations that are significantly different from each other and by making the randomisation step small relative to the distance between the conformations.

8.8.2 The Multicanonical Monte Carlo Method

In a canonical ensemble the probability $P_{\text{canon}}(T, E)$ of visiting a point in phase space with an energy E is proportional to the Boltzmann factor, $w_B = \exp(-E/k_B T)$, multiplied by the density of states, $n(E)$, where the number of states between E and $E + dE$ is given by $n(E)\delta E$. Thus:

$$P_{\text{canon}}(T, E) \propto n(E)w_B(E) \quad (8.61)$$

The density of states increases rapidly with energy but the Boltzmann factor decreases exponentially, meaning that $P_{\text{canon}}(T, E)$ is bell-shaped, with values that can vary by many orders of magnitude as the energy changes. In the multicanonical method the simulation

is performed in an artificial multicanonical ensemble in which the probability of visiting a state is independent of its energy over a certain energy range. This is equivalent to replacing the Boltzmann factor by a multicanonical weight factor, $w_{\text{mu}}(E)$:

$$P_{\text{mu}}(E) \propto n(E)w_{\text{mu}}(E) = \text{constant} \quad (8.62)$$

This implies that the multicanonical weight factor is proportional to $n(E)^{-1}$. A simulation performed in the multicanonical ensemble is able to overcome any energy barrier, in contrast to the situation in a normal, canonical simulation. The main task in performing a multicanonical simulation is to determine the weight factor, which is not known *a priori* (unlike the case for the canonical ensemble, where it is equal to the Boltzmann factor). The multicanonical weight factor is usually determined in an iterative fashion from a series of short simulations. One approach is as follows [Okamoto and Hansmann 1995]. First, a simulation is carried out at a temperature (T_0) which is high enough (e.g. 1000 K) to ensure that all energy barriers can be overcome. An array $S(E)$ is established with all its elements initially set to zero. Each element of this array refers to a particular energy range E to $E + \delta E$, where δE might be (say) 1 kcal/mol. From this initial simulation a histogram is constructed which gives the number of times a state with an energy in the range E to $E + \delta E$ is determined. These histogram values are stored in an array $H(E)$. Each of the values in this array ($H(E)$) should thus initially approximate the energy distribution at the temperature T_0 :

$$H(E) \propto n(E) \exp(-E/k_B T_0) \quad (8.63)$$

During this simulation the minimum and maximum energies visited during the simulation are recorded (E_{min} and E_{max} , respectively). The array $S(E)$ is now updated according to the following formula:

$$S(E) = S(E) + \ln H(E) \quad (8.64)$$

for all energy values E where the energy is between E_{min} and E_{max} and where the value of the appropriate array element $H(E)$ is greater than some minimum value (e.g. 20). In other words, that particular energy level must have been visited at least twenty times during the simulation. The following two parameters are now calculated:

$$\beta(E) = \begin{cases} 1/k_B T_0 & E \geq E_{\text{max}} \\ 1/k_B T_0 + \frac{S(E') - S(E)}{E' - E} & E_{\text{min}} \leq E \leq E_{\text{max}} \\ \beta(E_{\text{min}}) & E < E_{\text{min}} \end{cases} \quad (8.65)$$

$$\alpha(E) = \begin{cases} 0 & E \geq E_{\text{max}} \\ \alpha(E') + (\beta(E') - \beta(E))E' & E < E_{\text{max}} \end{cases} \quad (8.66)$$

E' refers to that element in the array $S(E)$ which succeeds E (i.e. E and E' are adjacent elements). Having determined the parameters $\alpha(E)$ and $\beta(E)$, a multicanonical weight factor is calculated as:

$$w_{\text{mu}}(E) = \exp[-\beta(E)E - \alpha(E)] \quad (8.67)$$

A new simulation is now initiated using this multicanonical weight factor (rather than the Boltzmann factor), from which new values of $S(E)$ and thus $H(E)$ can be determined. This cycle is continued until the distribution in $H(E)$ is reasonably flat in the energy range being considered.

Once the final multicanonical weight factor has been derived it provides the distribution for the production simulation in which high-energy configurations will be sampled adequately and high-energy barriers can be crossed with ease. Moreover, from this single simulation it is possible to derive the canonical distribution $P_{\text{canon}}(T, E)$ at any temperature (hence the name 'multicanonical'):

$$P_{\text{canon}}(T, E) \propto P_{\text{mu}}(E)w_{\text{mu}}^{-1} e^{-E/k_B T} \quad (8.68)$$

The average value of any property A at a temperature T can be determined from the multicanonical simulation using the following formula:

$$\langle A \rangle_T = \frac{\int A(E)P_{\text{canon}}(T, E) dE}{\int P_{\text{canon}}(T, E) dE} \quad (8.69)$$

In practice, one is restricted to energies between E_{min} and E_{max} , and a range of temperatures $T_{\text{min}} \leq T \leq T_{\text{max}}$ where the range of permitted temperatures is determined by calculating the expectation value of the energy at temperature T :

$$E_{\text{min}} \leq \langle E \rangle_T \leq E_{\text{max}} \quad (8.70)$$

The array $S(E)$ is intimately related to the entropy of the system as can be demonstrated by expanding the logarithm of $H(E)$ as:

$$\ln H(E) = \ln n(E) - E/k_B T_0 + \text{constant} \quad (8.71)$$

and recalling that the entropy is related to the density of states by $S(E) = \ln n(E)$.

The multicanonical Monte Carlo method can be used to study a wide range of systems and is particularly useful where traditional Metropolis Monte Carlo methods encounter difficulties. In addition to the simulation of systems such as clusters of rare gas atoms, the multicanonical method has been used to study the properties of macromolecular systems. There has been particular interest in using the approach to study the properties of amino acid polymers, and in particular their ability to form certain types of regular structure such as alpha helices. These structures will be discussed in more detail in Chapter 10. It suffices for now to note that certain types of amino acid confer a greater propensity to adopt a helical structure. Traditional simulation techniques can be used to study the equilibrium between the helical and 'random' (or coil) structures but typically have to start from the regular (i.e. helical) structure, which is then 'unfolded' using molecular dynamics or Monte Carlo simulations. The large number of minima on the potential energy surface means that it is not practical to start from the unfolded structure and observe helix formation. However, the multicanonical Monte Carlo method does provide a mechanism for more completely exploring the energy surface (including the helical conformations), starting from a random structure. In one such study, Okamoto and Hansmann were able to compare the thermodynamics of the equilibrium between the helix and coil structures for three amino acids (alanine, valine and glycine) which have different observed propensities to form helix structures [Okamoto and Hansmann 1995].

One of the drawbacks of the multicanonical method is that, during the simulations to derive the weight factor, the energy distribution in $H(E)$ can oscillate rather than steadily approaching a limiting distribution. Another drawback is that it can fail to properly

sample the low-energy regions adequately. The multicanonical method samples energies within a certain range with approximately equal probability, but at the ends of this range the probability drops dramatically. Thus little sampling is done from the low-energy regions. These low-energy regions are proportionally more important at low temperatures, leading to poor statistics. To some extent the problems with J-walking are complementary to the limitations of the multicanonical method, and so attempts have been made to combine the two [Xu and Berne 1999]. In the combined method (termed 'multicanonical jump walking') the multicanonical weight factor is first derived and then used to perform a long multicanonical simulation. The configurations generated by this multicanonical simulation are saved and used as the 'high-temperature' component in the subsequent J-walking simulation. The key modification is that the standard acceptance criterion for the jump steps during this last phase must be multiplied by the ratio of the weight factors for the two energies (i.e. rather than comparing the usual Boltzmann factor, $\exp(-\Delta\mathcal{V}/k_B T)$, to the random number a modified factor $\exp(-\Delta\mathcal{V}/k_B T)[w(\mathcal{V}_{\text{old}})/w(\mathcal{V}_{\text{new}})]$ is used). This combined approach appears to provide more efficient sampling of phase space for a given number of Monte Carlo steps when compared to the regular J-walking or multicanonical sampling method, on both low-dimensional trial potentials and for systems such as rare gas clusters.

8.9 Monte Carlo Sampling from Different Ensembles

A Monte Carlo simulation traditionally samples from the constant NVT (canonical) ensemble, but the technique can also be used to sample from different ensembles. A common alternative is the isothermal-isobaric, or constant NPT , ensemble. To simulate from this ensemble, it is necessary to have a scheme for changing the volume of the simulation cell in order to keep the pressure constant. This is done by combining random displacements of the particles with random changes in the volume of the simulation cell. The size of each volume change is governed by the maximum volume change, δV_{max} . Thus a new volume is generated from the old volume as follows:

$$V_{\text{new}} = V_{\text{old}} + \delta V_{\text{max}}(2\xi - 1) \quad (8.72)$$

As usual, ξ is a random number between 0 and 1. When the volume is changed, it is in principle necessary to recalculate the interaction energy of the entire system, not just the interactions involving the one atom or molecule that has been displaced. However, for simple interatomic potentials the change in energy associated with a volume change can be calculated very rapidly by using *scaled coordinates*. For a set of particles that are modelled by a Lennard-Jones potential in a cubic box of length L_{old} , the potential energy can be written:

$$\mathcal{V}_{\text{old}}(\mathbf{r}^N) = 4\epsilon \sum_{i=1}^N \sum_{j=i+1}^N \left(\frac{\sigma}{L_{\text{old}} s_{ij}} \right)^{12} - 4\epsilon \sum_{i=1}^N \sum_{j=i+1}^N \left(\frac{\sigma}{L_{\text{old}} s_{ij}} \right)^6 \quad (8.73)$$

where s_{ij} is a scaled coordinate which is related to the actual interatomic distance by $s_{ij} = L_{\text{old}}^{-1} r_{ij}$. It is necessary to write the energy as the sum of two components, one from

the repulsive part of the Lennard-Jones potential and the other from the attractive part:

$$\mathcal{V}_{\text{old}}(\mathbf{r}^N) = \mathcal{V}_{\text{old}}(12) + \mathcal{V}_{\text{old}}(6) \quad (8.74)$$

The advantage of using scaled coordinates is that they are independent of the size of the simulation box. Thus the energy of the configuration in a different-sized box (with side L_{new}) is:

$$\mathcal{V}_{\text{new}}(\mathbf{r}^N) = 4\epsilon \sum_{i=1}^N \sum_{j=i+1}^N \left(\frac{\sigma}{L_{\text{new}} s_{ij}} \right)^{12} - 4\epsilon \sum_{i=1}^N \sum_{j=i+1}^N \left(\frac{\sigma}{L_{\text{new}} s_{ij}} \right)^6 \quad (8.75)$$

The energy $\mathcal{V}_{\text{new}}(\mathbf{r}^N)$ is related to the energy $\mathcal{V}_{\text{old}}(\mathbf{r}^N)$ as follows:

$$\mathcal{V}_{\text{new}}(\mathbf{r}^N) = \mathcal{V}_{\text{old}}(12) \left\{ \frac{L_{\text{old}}}{L_{\text{new}}} \right\}^{12} + \mathcal{V}_{\text{old}}(6) \left\{ \frac{L_{\text{old}}}{L_{\text{new}}} \right\}^6 \quad (8.76)$$

The change in energy from the old to the new system is thus:

$$\Delta\mathcal{V}(\mathbf{r}^N) = \mathcal{V}_{\text{old}}(12) \left\{ \frac{L_{\text{old}}}{L_{\text{new}}} - 1 \right\}^{12} + \mathcal{V}_{\text{old}}(6) \left\{ \frac{L_{\text{old}}}{L_{\text{new}}} - 1 \right\}^6 \quad (8.77)$$

Any long-range corrections to the potential must also be taken into account when the volume changes. One way to deal with these is to assume that the non-bonded cutoff scales with the box length. Under such circumstances, the long-range corrections to both the repulsive and attractive parts of the potential scale in exactly the same manner as the short-range interactions. However, the use of this assumption can give rise to serious problems, particularly for techniques such as the Gibbs ensemble Monte Carlo simulation (see Section 8.12) where two coupled simulation boxes of different dimensions are involved. The boxes contain identical particles, but this would be compromised by the use of different non-bonded cutoffs and long-range corrections.

This simple scaling method cannot be used when simulating molecules, for a change in the scaled coordinates would have the effect of introducing large and energetically unfavourable changes in the internal coordinates, such as the bond lengths. It is therefore necessary to recalculate the total interaction energy of the system each time a volume change is made. This is computationally expensive to do, but it is in any case advisable to change the volume relatively infrequently compared to the rate at which the particles are moved. One way to speed up the energy calculation associated with a volume change is to write the potential energy change as a Taylor series expansion of the box size.

The criterion used to accept or reject a new configuration is slightly different for the isothermal-isobaric simulation than for a simulation in the canonical ensemble. The following quantity is used:

$$\Delta H(\mathbf{r}^N) = \mathcal{V}_{\text{new}}(\mathbf{r}^N) - \mathcal{V}_{\text{old}}(\mathbf{r}^N) + P(V_{\text{new}} - V_{\text{old}}) - Nk_B T \ln \left(\frac{V_{\text{new}}}{V_{\text{old}}} \right) \quad (8.78)$$

If ΔH is negative then the move is accepted; otherwise, $\exp(-\Delta H/k_B T)$ is compared to a random number between 0 and 1 and the move accepted according to:

$$\text{rand}(0, 1) \leq \exp(-\Delta H/k_B T) \quad (8.79)$$

To check that an isothermal-isobaric simulation is working properly, the pressure can be calculated from the virial as outlined in Section 6.2.3, including the appropriate long-range correction (which will not, of course, be constant as the volume of the box changes). Its value should be equal to the input pressure that appears in Equation (8.78).

8.9.1 Grand Canonical Monte Carlo Simulations

In the grand canonical ensemble the conserved properties are the chemical potential, the volume and the temperature. It can sometimes be more convenient to perform a grand canonical simulation at constant activity, z , which is related to the chemical potential μ by:

$$\mu = k_B T \ln \Lambda^3 z \quad (8.80)$$

where Λ is the de Broglie wavelength given by $\Lambda = \sqrt{h^2/2\pi m k_B T}$.

The key feature about the grand canonical Monte Carlo method is that the number of particles may change during the simulation. There are three basic moves in a grand canonical Monte Carlo simulation:

1. A particle is displaced, using the usual Metropolis method.
2. A particle is destroyed.
3. A particle is created at a random position.

The probability of creating a particle should be equal to the probability of destroying a particle. To determine whether to accept a destruction move the following quantity is calculated:

$$\Delta D = \frac{[\mathcal{V}_{\text{new}}(\mathbf{r}^N) - \mathcal{V}_{\text{old}}(\mathbf{r}^N)]}{k_B T} - \ln \left(\frac{N}{zV} \right) \quad (8.81)$$

For a creation step the equivalent quantity is:

$$\Delta C = \frac{[\mathcal{V}_{\text{new}}(\mathbf{r}^N) - \mathcal{V}_{\text{old}}(\mathbf{r}^N)]}{k_B T} - \ln \left(\frac{zV}{N+1} \right) \quad (8.82)$$

If ΔD or ΔC is negative then the move is accepted; if positive, then the exponential $\exp(-\Delta D/k_B T)$ or $\exp(-\Delta C/k_B T)$ as appropriate is calculated and compared with a random number between 0 and 1 in the usual way.

It is important that the possibility of creating a new particle is the same as the probability of destroying an old one. The ratio of particle creation/destruction moves to translation/rotation moves can vary, but the most rapid convergence is often achieved if all types of move occur with approximately equal frequency.

In grand canonical Monte Carlo simulations of liquids there can be some practical problems in achieving statistically accurate results. This is because the probability of achieving a successful creation or destruction step is often very small. Creation steps fail because the fluid is so dense that it is difficult to insert a new particle without causing significant

overlaps with neighbouring particles. Destruction steps fail because the particles in a fluid often experience significant attractive interactions, which are lost when the particle is removed. These problems are particularly acute for long-chain molecules. However, some of the newer Monte Carlo techniques such as the configurational bias Monte Carlo method do enable such systems to be simulated and accurate results obtained. These techniques will be discussed in Section 8.11.

8.9.2 Grand Canonical Monte Carlo Simulations of Adsorption Processes

One application of the grand canonical Monte Carlo simulation method is in the study of the adsorption and transport of fluids through porous solids. Mixtures of gases or liquids can be separated by the selective adsorption of one component in an appropriate porous material. The efficacy of the separation depends to a large extent upon the ability of the material to adsorb one component in the mixture much more strongly than the other component. The separation may be performed over a range of temperatures and so it is useful to be able to predict the adsorption isotherms of the mixtures.

A typical example of such a calculation is the simulation of Cracknell, Nicholson and Quirke [Cracknell *et al.* 1994] who studied the adsorption of a mixture of methane and ethane onto a microporous graphite surface. Four types of move were employed in their simulations: particle moves, particle deletions, particle creations and attempts to exchange particles. Methane was modelled as a single Lennard-Jones particle and ethane as two Lennard-Jones particles separated by a fixed bond length. The graphite surfaces were modelled as Lennard-Jones atoms with a slit-shaped pore being constructed from two layers of graphite separated by an appropriate distance. Triangle-shaped pores can also be used. The simulations were used to calculate the selectivity of the solid for the two components as the ratio of the mole fractions in the pore to the ratio of the mole fractions in the bulk. The selectivity was determined as a function of the pressure for different pore sizes to give some indication of the effect of changing the physical nature of the solid. The pressure can be calculated directly from the input chemical potential using the following standard relationship (for an ideal gas):

$$P = \{\exp(\mu/k_B T) k_B T\} / \Lambda^3 \quad (8.83)$$

The selectivity showed a complicated dependence upon the pore size (Figure 8.17).

The selectivity is best interpreted by considering the interactions between ethane molecules and the walls of the pore. For the smallest pore sizes, the molecules are restricted to the centre of the pore and the ethane molecules are forced to lie flat. As the pore size increases, it becomes possible for ethane to adopt a particularly favourable orientation perpendicular to the walls, with each methyl group being in a potential energy minimum for interaction with the pore atoms. This particular pore size ($2.5\sigma_{\text{CH}_4}$) thus has the greatest selectivity for ethane over methane. As the pore size increases further the distribution of ethane becomes more complex, with some layers of ethane lying flat on the pore wall and some in the centre of the pore, with ethane molecules spanning the space between. These arrangements are shown in Figure 8.18.

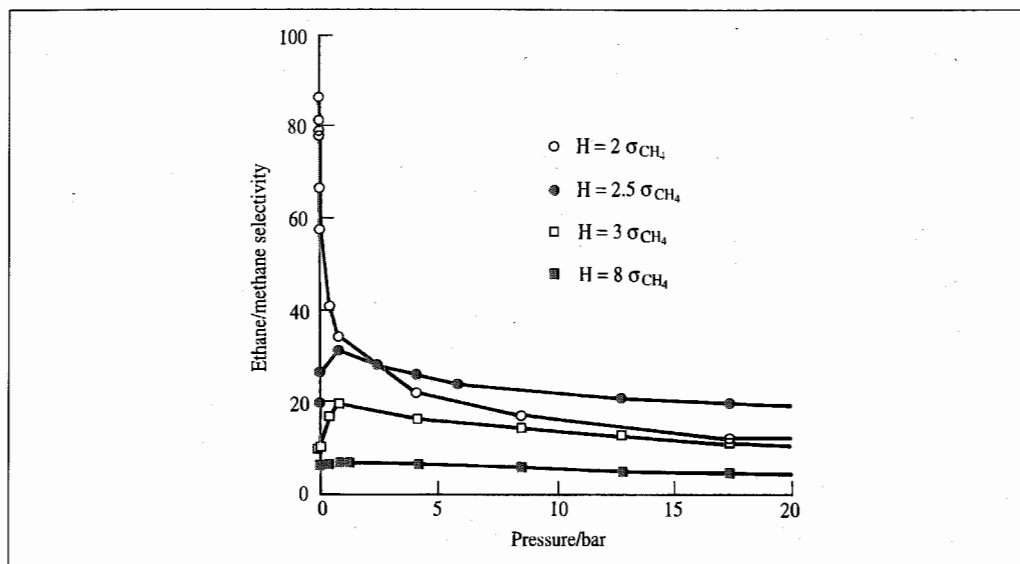


Fig. 8.17: Ethane/methane selectivity calculated from grand canonical Monte Carlo simulations of mixtures in slit pores at a temperature of 296 K. The selectivity is defined as the ratio of the mole fractions in the pore to the ratio of the mole fractions in the bulk. H is the slit width defined in terms of the methane collision diameter σ_{CH_4} . (Figure redrawn from Cracknell R F, D Nicholson and N Quirke 1994. A Grand Canonical Monte Carlo Study of Lennard-Jones Mixtures in Slit Pores; 2: Mixtures of Two-Centre Ethane with Methane. *Molecular Simulation* 13:161-175.)

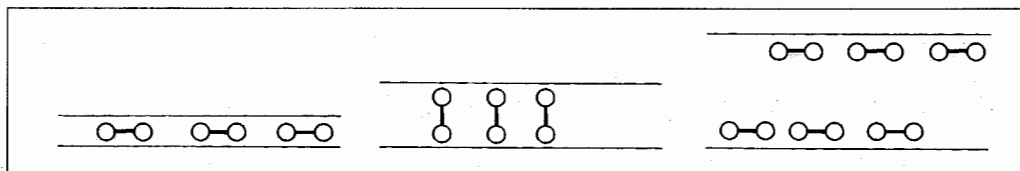


Fig. 8.18: Schematic illustration of the arrangements of ethane molecules in slits of varying sizes. In the slit of width $(2.5\sigma_{\text{CH}_4})$ each methyl group is able to occupy a potential minimum from the pore (middle).

8.10 Calculating the Chemical Potential

In a grand canonical simulation the chemical potential is constant. One may also wish to determine how the chemical potential varies during a simulation. The chemical potential is usually determined using an approach due to Widom [Widom 1963], in which a 'test' particle is inserted into the system and the resulting change in potential energy is calculated. The Widom approach is applicable to both molecular dynamics and Monte Carlo simulations. Consider a system containing $N - 1$ particles, into which we insert another particle at a random position. The inserted particle causes the internal potential energy to change by an amount $\mathcal{V}(\mathbf{r}^{\text{test}})$, i.e. $\mathcal{V}(\mathbf{r}^N) = \mathcal{V}(\mathbf{r}^{N-1}) + \mathcal{V}(\mathbf{r}^{\text{test}})$. Then the configurational

integral for the $N\mathcal{V}$ particle system is given by:

$$Z_N = \int d\mathbf{r}^N \exp[-\mathcal{V}(\mathbf{r}^N)/k_B T] \quad (8.84)$$

or

$$Z_N = \int d\mathbf{r}^N \exp[-\mathcal{V}(\mathbf{r}^{\text{test}})/k_B T] \exp[-\mathcal{V}(\mathbf{r}^{N-1})/k_B T] \quad (8.85)$$

By substituting unity in the form Z_{N-1}/Z_{N-1} it is possible to show that $Z_N = Z_{N-1} V \langle \exp[-\mathcal{V}(\mathbf{r}^{\text{test}})/k_B T] \rangle$.

The excess chemical potential, that is the difference between the actual value and that of the equivalent ideal gas system, is given by:

$$\mu_{\text{excess}} = -k_B T \ln \langle \exp[-\mathcal{V}(\mathbf{r}^{\text{test}})/k_B T] \rangle \quad (8.86)$$

The excess chemical potential is thus determined from the average of $\exp[-\mathcal{V}(\mathbf{r}^{\text{test}})/k_B T]$. In ensembles other than the canonical ensemble the expressions for the excess chemical potential are slightly different. The ghost particle does not remain in the system and so the system is unaffected by the procedure. To achieve statistically significant results many Widom insertion moves may be required. However, practical difficulties are encountered when applying the Widom insertion method to dense fluids and/or to systems containing molecules, because the proportion of insertions that give rise to low values of $\mathcal{V}(\mathbf{r}^{\text{test}})$ falls dramatically. This is because it is difficult to find a 'hole' of the appropriate size and shape.

8.11 The Configurational Bias Monte Carlo Method

Various techniques have been developed to tackle the problem of calculating the chemical potential in cases where the routine Widom method does not give converged results. Of these methods, the configurational bias Monte Carlo (CBMC) method, which was originally introduced by Siepmann [Siepmann 1990], is particularly exciting as it can be applied to assemblies of chain molecules. The configurational bias Monte Carlo method also provides a way to overcome the problems associated with Monte Carlo simulations of assemblies of chain molecules, where many proposed moves are rejected because of high-energy overlaps. The problem of calculating the chemical potential in such cases is clear from the following example. The probability of successfully inserting a single monomer into a fluid of typical liquid density is of the order of 0.5%, or 1 in 200. If one wishes to insert a molecule consisting of n such monomers, the probability is thus approximately 1 in 200^n . For an eight-segment molecule, this probability is less than 1 in 10^{18} , making such calculations impractical. The configurational bias Monte Carlo simulation technique can dramatically improve the chances of making a successful insertion.

The essence of the configurational bias Monte Carlo method is that a growing molecule is preferentially directed (i.e. biased) towards acceptable structures. The effects of these biases can then be removed by modifying the acceptance rules. The configurational bias methods are based upon work published in 1955 by Rosenbluth and Rosenbluth

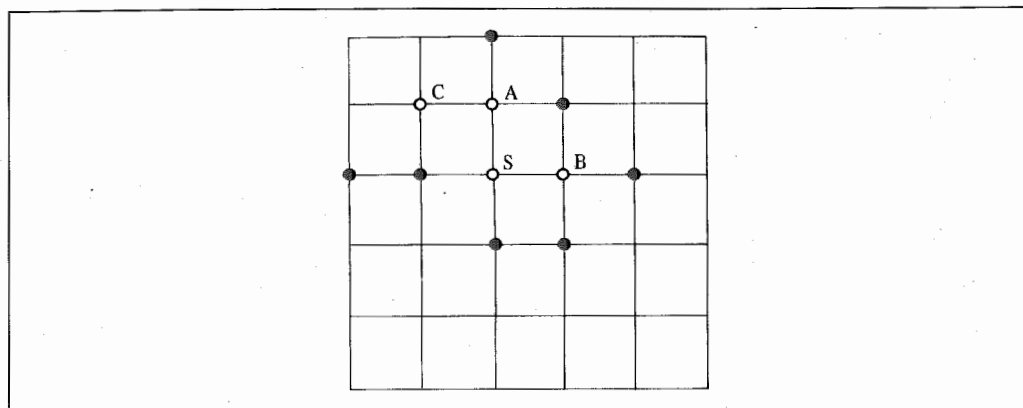


Fig. 8.19: The insertion of a three-unit molecule onto the lattice shown, starting at point S, can be achieved in only one way (see text). (Figure adapted from Siepmann J I 1990. *A Method for the Direct Calculation of Chemical Potentials for Dense Chain Systems*. *Molecular Physics* 70:1145–1158.)

[Rosenbluth and Rosenbluth 1955] and can be applied to both lattice models and to systems with arbitrary molecular potentials and conformations. The method is most easily explained using a two-dimensional lattice model. Suppose we wish to insert a three-unit molecule onto the lattice shown in Figure 8.19. First we consider how the conventional approach would tackle this problem. The initial step is to select a lattice point at random. Suppose we select the lattice point labelled S in Figure 8.19. We then choose one of the four neighbours of S at random. Of the four neighbouring sites, two are occupied and two are free (A and B in Figure 8.19). There is thus a 50% probability that the move will be rejected at this stage. If we select site B then the molecule can be grown no further as all of the adjacent sites are filled. Were we to grow onto A then we would select one of the three neighbouring sites at random. Of these only site C is available. On average, only one trial in twelve will successfully grow a molecule from S using a conventional Monte Carlo algorithm.

Let us now consider how the configurational bias Monte Carlo method would deal with this problem. Again, the first site (S) is chosen at random. We next consider where to place the second unit. The sites adjacent to S are examined to see which are free. In this case, only two of the four sites are free. One of these free sites is chosen at random. Note that the conventional Monte Carlo procedure selected from all four adjoining sites at random, irrespective of whether it is occupied or not. A *Rosenbluth weight* for the move is then calculated. The Rosenbluth weight for each step i is given by:

$$W_i = \frac{n'}{n} W_{i-1} \quad (8.87)$$

where W_{i-1} is the weight for the previous step ($W_0 = 1$), n' is the number of available sites and n is the total number of neighbouring sites (not including the one occupied by the previous unit). In the case of our lattice, $W_1 = 2/4 = 1/2$. If site B is chosen then there is no site available for the third unit, and so the attempt has to be abandoned. If site A is chosen, its adjacent sites are examined to see which are free. In this case there is only one

free site where the third and final unit can be placed. The Rosenbluth weight for this step is $1/3 \times 1/2 = 1/6$. The overall statistical weight for the move is obtained by multiplying the number of successful trials by the Rosenbluth weight of each trial; as half the trials succeed, the statistical weight is therefore $1/2 \times 1/6 = 1/12$. This is exactly the same result that would be obtained with a conventional sampling scheme, though recall that in a conventional scheme only one trial in twelve results in a successful insertion. By contrast, with the configurational bias method the proportion of successful trials is one in two.

The configurational bias algorithm can be extended to take account of intermolecular interactions between the growing chain and its lattice neighbours. If the energy of segment i when occupying a particular site Γ is $v_{\Gamma}(i)$ then that site is chosen with a probability given by:

$$p_{\Gamma}(i) = \frac{\exp[-v_{\Gamma}(i)/k_B T]}{Z_i} \quad (8.88)$$

Z_i is the sum of the Boltzmann factors for all of the b positions considered:

$$Z_i = \sum_{\Gamma=1}^b \exp[-v_{\Gamma}(i)/k_B T] \quad (8.89)$$

The site can be chosen using a biased roulette wheel algorithm, in which the interval between 0 and 1 is divided into b adjacent segments each with a size proportional to the probabilities $p_1(i), p_2(i), \dots, p_b(i)$ (Figure 8.20). The site within whose interval a random number between 0 and 1 lies is the one chosen. The chain is thus biased towards those sites with a higher Boltzmann weighting; the sum of the Boltzmann factors plays the role of n' in Equation (8.87). The Rosenbluth weight for the entire chain (of length l) can be calculated as:

$$W_l = \exp[-\mathcal{V}_{\text{tot}}(l)/k_B T] \prod_{i=2}^l \frac{Z_i}{b} \quad (8.90)$$

where $\mathcal{V}_{\text{tot}}(l)$ is the total energy of the chain, equal to the sum of the individual segment energies $v_{\Gamma}(i)$. The average Rosenbluth weight is directly related to the excess chemical potential:

$$\mu_{\text{ex}} = k_B T \ln \langle W_l \rangle \quad (8.91)$$

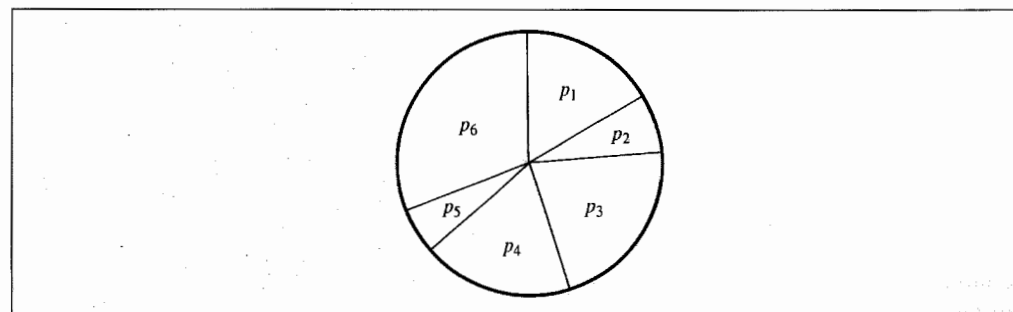


Fig. 8.20: A biased roulette wheel chooses states according to their probabilities.

If a segment has a zero Rosenbluth weight then growth of the chain is terminated. However, such chains must still be included in the averaging used to determine the excess chemical potential.

So far, we have only considered a fixed number of neighbouring sites for each segment. The method can be extended to cover fully flexible chains, where the set of possible neighbouring positions is infinite [De Pablo *et al.* 1992, 1993]. When growing each segment, a subset containing k random directions is chosen. These trial directions need not be uniformly distributed in space. For each of these orientations the energy $v_T(i)$ is calculated and so is the Boltzmann factor. An orientation is then chosen with probability:

$$p_T(i) = \frac{\exp[-v_T(i)/k_B T]}{\sum_{\Gamma=1}^k \exp[-v_T(i)/k_B T]} \quad (8.92)$$

The Rosenbluth factor is accumulated as follows:

$$W_i = W_{i-1} \frac{1}{k} \sum_{\Gamma=1}^k \exp[-v_T(i)/k_B T] \quad (8.93)$$

To implement this method it is necessary to determine the appropriate number of trial directions, k . If $k = 1$ then the method is equivalent to the original Widom particle insertion method. If k is too large then too much time is taken calculating the Rosenbluth factors for trial positions that are very close in phase space. Frenkel and colleagues have investigated how the choice of k influences the accuracy of the results and the efficiency with which those results were obtained [Frenkel *et al.* 1991]. The system they examined was of a flexible chain containing up to 20 segments in a moderately dense atomic fluid. The conventional particle insertion method failed completely for this system. Not surprisingly, the results showed that as the length of the chain increases so the number of random orientations that need to be considered also increases. At least four trial orientations were used at each step and k was chosen to increase logarithmically with the number of segments to be grown. The limiting value of k was considered to be reached when so many trials were required that the configurational bias method was no more efficient than alternative methods of regrowing chains, such as reptation algorithms. For example, for a 6-segment chain the proportion of accepted configurations (once the initial monomer had been inserted successfully) was 0.00001% for $k = 1$, 3.2% for $k = 10$ and 35% for $k = 50$. For a 20-segment chain the proportion of accepted configurations was 0.0001% for $k = 20$, 0.66% for $k = 50$ and 2.0% for $k = 100$.

The Rosenbluth algorithm can also be used as the basis for a more efficient way to perform Monte Carlo sampling for fully flexible chain molecules [Siepmann and Frenkel 1992], which, as we have seen, is difficult to do as bond rotations often give rise to high energy overlaps with the rest of the system.

The configurational bias Monte Carlo method involves three types of move. Two of these are translational or rotational moves of the entire molecule, which are performed in the conventional way. The third type of move is a conformational change. A chain is selected at random and one of the segments within it is also randomly chosen. That part of the chain that lies above or below the segment (chosen with equal probability) is discarded and an

attempt is made to regrow the discarded portion. Let us consider first the case where each segment is restricted to a given number of discrete orientations, either because the chain is restricted to a lattice or because the model discretely samples the conformational space (e.g. it only permits *gauche* and *trans* conformations to a hydrocarbon chain). At each stage, the Boltzmann weights of the b discrete conformations are determined and one of the sites is chosen with a probability given by Equation (8.92). The Rosenbluth weight is determined for the growing chain using Equation (8.93).

Having generated a trial conformation it must be decided whether to accept it or not. To do this a random number is generated in the range 0–1 and compared with the ratio of the Rosenbluth weights for the trial conformation ($W_{i,\text{trial}}$) and the old conformation ($W_{i,\text{old}}$). The new chain is then accepted using the following criterion:

$$\text{rand}(0, 1) \leq \frac{W_{i,\text{trial}}}{W_{i,\text{old}}} \quad (8.94)$$

A similar approach can be adopted with continuous chains. Here it is also possible to enhance the sampling by guiding the choice of trial sites towards those with a particularly favourable intramolecular energy. This can be achieved by generating random vectors on the surface of a sphere of unit radius for each segment. The potential energy (angle-bending and torsional) for a bond directed along this vector is calculated. The vector is then accepted or rejected using the Metropolis criterion. If it is accepted, the vector is scaled to the appropriate bond length. This procedure continues until the desired number of trial sites have been generated. A trial site is then selecting using Boltzmann factors, which only consider the intra- and intermolecular non-bonded interactions of the sites with the chain and with the rest of the system. The Rosenbluth weights are similarly calculated and the move is accepted according to the ratio of the new and old Rosenbluth weights. Again, the correct choice of the number of trial sites is crucial to the efficiency of the method.

For branched molecules some modifications are required to the configurational bias method as described so far. This is because there may be bond angles which share the same central atom and torsion angles which have the same central two atoms in common. Thus in 2-methylalkanes the bond angles to the two terminal methyl groups share the 2-carbon atom. In 3,4-dimethylhexane there is a potential torsion problem. In the 'standard' configurational bias method one of the methyl groups would be grown, followed by the second. What is sometimes observed, however, is that the distributions of these bond angles is not equal (as it should be, as they are equivalent). Two possible ways to deal with this problem are to grow both atoms simultaneously [Dijkstra 1997] or to use a small Monte Carlo simulation to generate the trial positions [Vlugt *et al.* 1999]. When there are multiple torsion angles these two methods are not suitable; indeed, for a molecule such as 2,3-dimethylbutane the entire molecule must be generated in a single step. Martin and Siepmann suggested that it was possible to decouple the selection of the different energy terms [Martin and Siepmann 1999]. Suppose that the Lennard-Jones, torsional and bond-angle terms are decoupled. Then the probability of generating a particular configuration is given by:

$$P = \prod_{n=1}^{n_{\text{step}}} \left[\frac{\exp(-v_{\text{LJ}}(i)/k_B T)}{W_{\text{L}}(n)} \right] \left[\frac{\exp(-v_{\text{tor}}(j)/k_B T)}{W_{\text{T}}(n)} \right] \left[\frac{\exp(-v_{\text{bend}}(k)/k_B T)}{W_{\text{B}}(n)} \right] \quad (8.95)$$

The relevant Rosenbluth weights are:

$$W_L(n) = \sum_{i=1}^{n_{LJ}} \exp(-v_{LJ}(i)/k_B T) \quad (8.96)$$

$$W_T(n) = \sum_{j=1}^{n_{tor}} \exp(-v_{tor}(j)/k_B T) \quad (8.97)$$

$$W_B(n) = \sum_{k=1}^{n_{bend}} \exp(-v_{bend}(k)/k_B T) \quad (8.98)$$

where n_{LJ} , n_{tor} and n_{bend} are the number of trial sites for the Lennard-Jones, torsional and angle-bending interactions, respectively. Under these conditions, the move is accepted with a probability:

$$P_{acc} = \min \left[1, \frac{\prod_{n=1}^{n_{step}} W_L(n)_{new} W_T(n)_{new} W_B(n)_{new}}{\prod_{n=1}^{n_{step}} W_L(n)_{old} W_T(n)_{old} W_B(n)_{old}} \right] \quad (8.99)$$

The advantage of this decoupling method is that a large number of trial sites can be chosen for the computationally less expensive bond angle selection without increasing the cost of performing the other selections. Once the bond angle distribution has been chosen by this biased method, it is used as input to a biased selection of the torsional and Lennard-Jones interactions. An extension to this decoupling procedure involves grouping the torsional and Lennard-Jones together and having each biased selection of bond angles send many possible conformations forward to the next step. This coupling and decoupling of terms is claimed to provide a great deal of flexibility when designing a configurational bias scheme for any particular molecule and would also be applicable to force field models that included additional terms such as bond stretching or cross terms.

8.11.1 Applications of the Configurational Bias Monte Carlo Method

The CBMC method has been used to investigate a number of systems involving long-chain alkanes. Siepman and McDonald examined a monolayer of 90 $\text{CH}_3(\text{CH}_2)_{15}\text{SH}$ molecules chemisorbed onto a gold surface [Siepman and McDonald 1993a, b]. The thiol group forms a bond with the gold surface atoms, thus producing a high degree of surface ordering of the adsorbed molecules. Spectroscopic experiments indicated that the chains were tilted relative to the surface and adopted a predominantly *trans* conformation for the alkyl links. Both discrete and continuous versions of the configurational bias Monte Carlo method were employed; in the discrete model, each CH_2CH_2 segment was restricted to the *trans* and two *gauche* conformations. In the continuous simulation, six trial sites were used for each segment. The molecules were initially placed on a triangular lattice in an extended conformation perpendicular to the surface.

In the structure obtained at the end of the simulation the chains were ordered in an approximately hexagonal pattern. During equilibration *gauche* conformations were introduced into the alkyl chains, causing the system to tilt. However, once the molecules had all tilted the

gauche defects were gradually squeezed out to give chains with predominantly *trans* links. The final configuration is shown in Figure 8.21 (colour plate section).

The configurational bias Monte Carlo method has also been used to investigate the adsorption of alkanes in zeolites. Such systems are of especial interest in the petrochemical industry. One interesting experimental result obtained for the zeolite silicalite was that short-chain alkanes (C_1 to C_5) and long-chain alkanes (C_{10}) have simple adsorption isotherms but hexane and heptane show kinked isotherms. Such systems are obvious candidates for theoretical investigations because the experimental data is difficult and time-consuming to obtain. Moreover, the simulation can often provide a detailed molecular explanation for the observed behaviour. The simulation of such systems is difficult using conventional methods; the Monte Carlo method suffers from the problems of low acceptance ratios or a very slow exploration of phase space, and long simulation times would be required with molecular dynamics as the diffusion of long-chain alkanes is very slow. The configurational bias Monte Carlo method enabled effective and efficient simulations to be performed, providing both thermodynamic properties and the spatial distribution of the molecules within the zeolite [Smit and Siepman 1994; Smit and Maesen 1995]. The adsorption isotherms (i.e. the number of molecules adsorbed as a function of the applied pressure) were calculated using grand canonical simulations in which the zeolite was coupled to a reservoir at constant temperature and chemical potential.

Silicalite has both straight and zig-zag channels, which are connected via intersections (Figure 8.22). An analysis of the configurations showed that the distribution of a short alkane, such as butane, was approximately equal between the two types of channel. However, as the length of the alkane chain was increased, so there was a greater probability of finding it in a straight channel than in a zig-zag channel. Hexane is an interesting case, for its length is almost equal to the period of the zig-zag channels. At low pressure the hexane molecules move freely in the zig-zag channels and occupy the intersections for part of the time. To fill the zeolite with hexane it is first necessary for the alkane molecules to occupy just the zig-zag channels and not the intersections. This is accompanied by a loss of entropy, which must be compensated for by a higher chemical potential and so gives the kinked isotherm. The straight channels can then be filled with hexane. Different behaviour

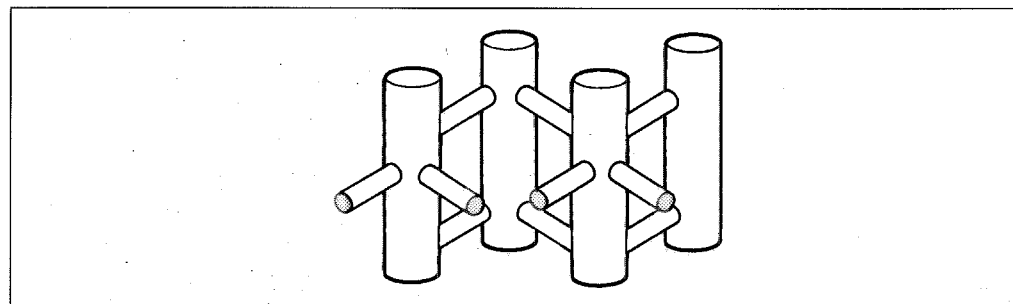


Fig. 8.22: Schematic structure of the zeolite silicalite showing the straight and zig-zag channels. (Figure adapted from Smit B and J I Siepman 1994. *Simulating the Adsorption of Alkanes in Zeolites*. Science 264:1118-1120.)

is observed for smaller alkanes because more than one molecule can occupy the zig-zag channels. Longer alkanes always partially occupy the intersection and so there is no benefit from freezing the molecules in the zig-zag channels. It is also possible to simulate the behaviour of branched alkanes [Vlugt *et al.* 1998] and compare these with their linear equivalents. Thus, whereas *n*-butane has an equal probability of being in either channel, isobutane has a preference for the intersection. Moreover, once all the intersections are full it requires considerable energy to place isobutane elsewhere in the zeolite. This requires a much higher pressure, giving rise to an inflection in the adsorption isotherms.

8.12 Simulating Phase Equilibria by the Gibbs Ensemble Monte Carlo Method

The most 'obvious' way to investigate phase equilibria is to set up an appropriate system with a conventional simulation technique. Unfortunately, simulations of systems with more than one phase usually require inordinate amounts of computer time. There are several reasons why the use of conventional simulation methods to investigate phase equilibria is difficult. First, it would take a very long time to equilibrate such a system, which would need to separate into its two phases (e.g. liquid and vapour). The properties of the fluid in the interfacial region differ substantially from the properties in the bulk and so to obtain a 'bulk' measurement all of the interfacial atoms must be ignored. Smit has calculated the percentage of the number of particles in the interfacial region for systems of varying sizes; these percentages range from 10% in the interfacial region for a system of 50 000 particles to 95% for a system of 100 particles [Smit 1993]. To simulate phase equilibria directly would thus require long simulations to be performed on systems containing many particles.

The Gibbs ensemble Monte Carlo simulation method, invented by Panagiotopoulos [Panagiotopoulos 1987], enables phase equilibria to be studied directly using small numbers of particles. Rather than trying to form an interface within a single simulation, two simulation boxes are used, each representing one of the two phases. There is no physical interface between the two boxes, which are subject to the usual periodic boundary conditions (Figure 8.23). Three types of move are possible. The first type of move comprises particle displacements within each box, as in a conventional Monte Carlo simulation. The second type of move involves volume changes of the two boxes by equal and opposite amounts so that the total volume of the system remains constant. The third type of move involves the removal of a particle from one box and its attempted placement in the other box. This is identical to the Widom insertion method for calculating the chemical potential. Indeed, as the energy of the inserted particle must be calculated, it is possible to determine the chemical potential in the Gibbs ensemble without any additional computational cost. These three types of move are often performed in strict order, but it may be better to choose each type of move at random, ensuring that, on average, the appropriate numbers of each type of move are made.

The properties of the Gibbs ensemble Monte Carlo simulation method have been examined in great detail using simple systems such as the Lennard-Jones fluid and simple gases. A

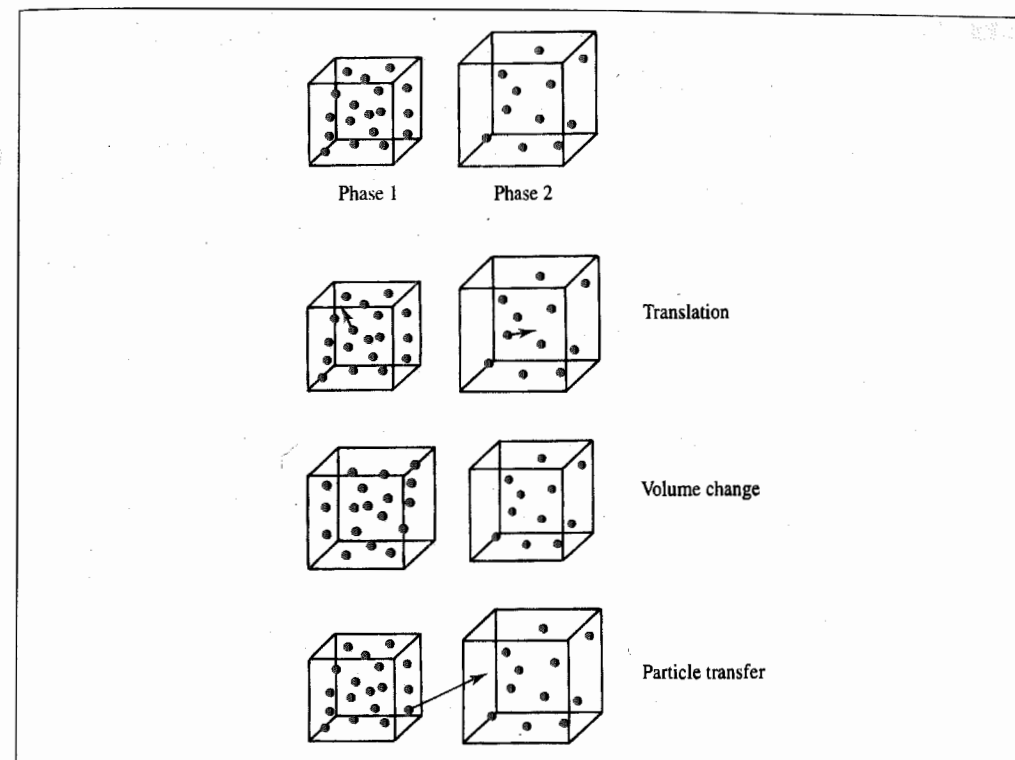


Fig. 8.23: The Gibbs ensemble Monte Carlo simulation method uses one box for each of the two phases. Three types of move are permitted: translations within either box; volume changes (keeping the total volume constant) and transfer of a particle from one box to the other.

particularly exciting development is the use of the configurational bias Monte Carlo method in conjunction with the Gibbs ensemble method to construct the phase diagrams of complex, long-chain molecules. For example, the vapour-liquid phase equilibria of *n*-pentane and *n*-octane have been investigated by Siepmann, Karaborni and Smit using this combined approach on systems containing 200 pentane or 160 octane molecules [Siepmann *et al.* 1993a]. The calculated properties of these two systems agreed very well with the available experimental data, particularly for the shorter alkane. Their studies were subsequently extended to much longer alkanes (up to C_{48}) [Siepmann *et al.* 1993b]. One particularly noteworthy result was that the density at the critical point increased with the length of the carbon chain up to *n*-octane but then *decreased* as the chain increased in length. Until shortly before the simulations were performed it had been assumed that the critical density for longer chains could be extrapolated from the experimental data obtained with short chains under the assumption that the critical density increased with the length of the chain. Later experiments were able to examine longer chains and did indeed demonstrate that the critical point density passed through a maximum at octane and then decreased for shorter chain lengths.

8.13 Monte Carlo or Molecular Dynamics?

In principle, the modeller has the choice of using either the Monte Carlo or molecular dynamics technique for a given simulation. In practice one technique must be chosen over the other. Sometimes the decision is a trivial one, for example because a suitable program is readily available. In other cases there are clear reasons for choosing one method instead of the other. For example, molecular dynamics is required if one wishes to calculate time-dependent quantities such as transport coefficients. Conversely, Monte Carlo is often the most appropriate method to investigate systems in certain ensembles; for example, it is much easier to perform simulations at exact temperatures and pressures with the Monte Carlo method than using the sometimes awkward and ill-defined constant temperature and constant pressure molecular dynamics simulation methods. The Monte Carlo method is also well suited to certain types of models such as the lattice models.

The two methods can differ in their ability to explore phase space. A Monte Carlo simulation often gives much more rapid convergence of the calculated thermodynamic properties of a simple molecular liquid (modelled as a rigid molecule), but it may explore the phase space of large molecules very slowly due to the need for small steps unless special techniques such as the configurational bias Monte Carlo method are employed. However, the ability of the Monte Carlo method to make non-physical moves can significantly enhance its capacity to explore phase space in appropriate cases. This may arise for simulations of isolated molecules, where there are a number of minimum energy states separated by high barriers. Molecular dynamics may not be able to cross the barriers between the conformations sufficiently often to ensure that each conformation is sampled according to the correct statistical weight. Molecular dynamics advances the positions and velocities of all the particles simultaneously and it can be very useful for exploration of the local phase space whereas the Monte Carlo method may be more effective for conformational changes, which jump to a completely different area of phase space.

Given that the two techniques in some ways complement each other in their ability to explore phase space, it is not surprising that there has been some effort to combine the two methods. Some of the techniques that we have considered in this chapter and in Chapter 7 incorporate elements of the Monte Carlo and molecular dynamics techniques. Two examples are the stochastic collisions method for performing constant temperature molecular dynamics, and the force bias Monte Carlo method. More radical combinations of the two techniques are also possible.

An obvious way to combine Monte Carlo with molecular dynamics is to use each technique for the most appropriate part of a simulation. For example, when simulating a solvated macromolecule, the equilibration phase is usually performed in a series of stages. In the first stage, the solute is kept fixed while the solvent molecules (and any ions, if present) are allowed to move under the influence of the solute's electrostatic field. This solvent equilibration may often be performed more effectively using a Monte Carlo simulation as the solvent and ions do not have any appreciable conformational flexibility. To simulate the whole system, molecular dynamics is then the most appropriate method. Such a protocol has been used to perform long simulations of DNA molecules [Swaminathan *et al.* 1991].

A variety of hybrid molecular dynamics/Monte Carlo methods have been devised, in which the simulation algorithm alternates between molecular dynamics and Monte Carlo. The aim of such methods is to achieve better sampling, and thereby more rapid convergence of thermodynamic properties. *In extremis*, each molecular dynamics (or stochastic dynamics) step is followed by a Monte Carlo step, the velocities being unaffected by acceptance or rejection of the Monte Carlo step. Such a method has been devised by Guarnieri and Still [Guarnieri and Still 1994]. An alternative is to perform a block of molecular dynamics steps to generate a new state, which is then accepted or rejected on the basis of the total energy (potential plus kinetic) using the usual Metropolis criterion. If the new coordinates are rejected then the original coordinates from the start of the block are restored and molecular dynamics is run again, but with an entirely new set of velocities that is chosen from a Gaussian distribution. This approach is very similar to the stochastic collisions method for temperature control discussed in Section 7.7.1 but with the addition of a Monte Carlo acceptance or rejection step [Duane *et al.* 1987]. A simulation using this hybrid algorithm samples from the canonical ensemble (constant temperature) and was shown by Clamp and colleagues to be more effective than conventional molecular dynamics or Monte Carlo methods for exploring the phase space of both simple model systems and proteins [Clamp *et al.* 1994].

Appendix 8.1 The Marsaglia Random Number Generator

The Marsaglia random number generator [Marsaglia *et al.* 1990] is known as a *combination generator* because it is constructed from two different generators. It has a period of about 2^{144} . The first generator is a lagged Fibonacci generator that performs the following binary operation on two real numbers x and y :

$$x \bullet y = x - y \text{ if } x \geq y; \quad x \bullet y = x - y - 1 \text{ if } x < y \quad (8.100)$$

The values x and y are chosen from numbers earlier in the sequence, so that the n th value in the sequence is calculated by:

$$x_n = x_{n-r} \bullet x_{n-s} \quad (8.101)$$

r and s are the *lags*, which are chosen to give numbers that are satisfactorily random and have a long period. Marsaglia chose $r = 97$ and $s = 33$. The algorithm does therefore require the last 97 numbers to be stored at all stages.

The second generator is an arithmetic sequence method that generates random numbers using the following mathematical operation:

$$c \circ d = c - d \text{ if } c \geq d; \quad c \circ d = c - d + 16777213/16777216 \text{ if } c < d \quad (8.102)$$

The n th value in this sequence is given by:

$$c_n = c_{n-1} \circ (7654321/16777216) \quad (8.103)$$

The n th number, U_n , in the combined sequence is then obtained as:

$$U_n = x_n \circ c_n \quad (8.104)$$

The c sequence requires one initial seed value and the x sequence requires 97 initial seeds (which should themselves be reasonably random). These can be supplied by the user but in the published algorithm these 97 values were obtained from another combination generator comprising a lagged Fibonacci generator and a congruential algorithm.

Further Reading

- Adams D J 1983. Introduction to Monte Carlo Simulation Techniques. In Perran J W (Editor) *Physics of Superionic Conductors and Electrode Materials*. New York, Plenum, pp. 177–195.
- Allen M P and D J Tildesley 1987. *Computer Simulation of Liquids*. Oxford, Oxford University Press.
- Colbourn E A (Editor) 1994. *Computer Simulation of Polymers*. Harlow, Longman.
- Frenkel D. Monte Carlo Simulations: A Primer. In van Gunsteren W F, P K Weiner and A J Wilkinson (Editors). *Computer Simulation of Biomolecular Systems Volume 2*. Leiden, ESCOM, pp. 37–66.
- Galaiatsatos V 1995. Computational Methods for Modelling Polymers: An Introduction. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry Volume 6*. New York, VCH Publishers, pp. 149–208.
- Kaols M H and P A Whitlock 1986. *Monte Carlo Methods, Volume 1: Basics*. New York, John Wiley & Sons.
- Kermer K 1993. Computer Simulation of Polymers. In Allen M P and D J Tildesley (Editors). *Computer Simulation in Chemical Physics*. Dordrecht, Kluwer, NATO ASI Series 397:397–459.
- Rubinstein R Y 1981. *Simulation and Monte Carlo Methods*. New York, John Wiley & Sons.

References

- Barker J A and R O Watts 1969. Structure of Water; A Monte Carlo Calculation. *Chemical Physics Letters* 3:144–145.
- Baschnagel J, K Binder, W Paul, M Laso, U Suter, I Batoulis, W Jilge and T Bürger 1991. On the Construction of Coarse-Grained Models for Linear Flexible Polymer Chains – Distribution Functions for Groups of Consecutive Monomers. *Journal of Chemical Physics* 95:6014–6025.
- Clamp M E, P G Baker, C J Stirling and A Brass 1994. Hybrid Monte Carlo: An Efficient Algorithm for Condensed Matter Simulation. *Journal of Computational Chemistry* 15:838–846.
- Cracknell R F, D Nicholson and N Quirke 1994. A Grand Canonical Monte Carlo Study of Lennard-Jones Mixtures in Slit Pores; 2: Mixtures of Two-Centre Ethane with Methane. *Molecular Simulation* 13:161–175.
- De Pablo J J, M Laso, J I Siepmann and U W Suter 1993. Continuum–Configurational Bias Monte Carlo Simulations of Long-chain Alkanes. *Molecular Physics* 80:55–63.
- De Pablo J J, M Laso, and U W Suter 1992. Estimation of the Chemical Potential of Chain Molecules by Simulation. *Journal of Chemical Physics* 96:6157–6162.
- Dijkstra M. 1997. Confined Thin Films of Linear and Branched Alkanes. *Journal of Chemical Physics* 107:3277–3288.
- Duane S, A D Kennedy and B J Pendleton 1987. Hybrid Monte Carlo. *Physics Letters* B195:216–222.
- Flory P J 1969. *Statistical Mechanics of Chain Molecules*. New York, Interscience.
- Frantz D D, D L Freeman and J D Doll 1990. Reducing Quasi-ergodic Behavior in Monte Carlo Simulations by J-walking: Applications to Atomic Clusters. *Journal of Chemical Physics* 93:2769–2784.
- Frenkel D D, C A M Mooij and B Smit 1991. Novel Scheme to Study Structural and Thermal Properties of Continuously Deformable Materials. *Journal of Physics Condensed Matter* 3:3053–3076.

- Guarnieri F and W C Still 1994. A Rapidly Convergent Simulation Method: Mixed Monte Carlo/Stochastic Dynamics. *Journal of Computational Chemistry* 15:1302–1310.
- Marsaglia G, A Zaman and W W Tsang 1990. Towards a Universal Random Number Generator. *Statistics and Probability Letters* 8:35–39.
- Martin M G and J I Siepmann 1999. Novel Configurational-bias Monte Carlo Method for Branched Molecules. Transferable Potentials for Phase Equilibria. 2. United-atom Description of Branched Alkanes. *Journal of Physical Chemistry* 103:4508–4517.
- Metropolis N, A W Rosenbluth, M N Rosenbluth, A H Teller and E Teller 1953. Equation of State Calculations by Fast Computing Machines. *Journal of Chemical Physics* 21:1087–1092.
- Okamoto Y and U H E Hansmann 1995. Thermodynamics of Helix-coil Transitions Studied by Multicanonical Algorithms. *Journal of Physical Chemistry* 99:11276–11287.
- Panagiotopoulos A Z 1987. Direct Determination of Phase Coexistence Properties of Fluids by Monte Carlo Simulation in a New Ensemble. *Molecular Physics* 61:813–826.
- Pangali C, M Rao and B J Berne 1978. On a Novel Monte Carlo Scheme for Simulating Water and Aqueous Solutions. *Chemical Physics Letters* 55:413–417.
- Rao M and B J Berne 1979. On the Force Bias Monte Carlo Simulation of Simple Liquids. *Journal of Chemical Physics* 71:129–132.
- Rosenbluth M N and A W Rosenbluth 1955. Monte Carlo Calculation of the Average Extension of Molecular Chains. *Journal of Chemical Physics* 23:356–359.
- Rosky P J, J D Doll and H L Friedman 1978. Brownian Dynamics as Smart Monte Carlo Simulation. *Journal of Chemical Physics* 69:4628–4633.
- Senderowitz H, F Guarnieri and W C Still 1995. A Smart Monte Carlo Technique for Free Energy Simulations of Multicanonical Molecules. Direct Calculations of the Conformational Populations of Organic Molecules. *Journal of the American Chemical Society* 117:8211–8219.
- Sharp W E and C Bays 1992. A Review of Portable Random Number Generators. *Computers and Geosciences* 18:79–87.
- Siepmann J I 1990. A Method for the Direct Calculation of Chemical Potentials for Dense Chain Systems. *Molecular Physics* 70:1145–1158.
- Siepmann J I and D Frenkel 1992. Configurational Bias Monte Carlo: A New Sampling Scheme for Flexible Chains. *Molecular Physics* 75:59–70.
- Siepmann J I, S Karaborni and B Smit 1993a. Vapor–Liquid Equilibria of Model Alkanes. *Journal of the American Chemical Society* 115:6454–6455.
- Siepmann J I, S. Karaborni and B Smit 1993b. Simulating the Crucial Behaviour of Complex Fluids. *Nature* 365:330–332.
- Siepmann J I and I R McDonald 1993a. Domain Formation and System-size Dependence in Simulations of Self-assembled Monolayers. *Langmuir* 9:2351–2355.
- Siepmann J I and I R McDonald 1993b. Monte Carlo Study of the Properties of Self-assembled Monolayers Formed by Adsorption of $\text{CH}_3(\text{CH}_2)_{15}\text{SH}$ on the (111) Surface of Gold. *Molecular Physics* 79:457–473.
- Smit B 1993. Computer Simulation in the Gibbs Ensemble. In Allen M P and D J Tildesley (Editors). *Computer Simulation in Chemical Physics*. Dordrecht, Kluwer. NATO ASI Series 397, pp. 173–210.
- Smit B and T L M Maesen 1995. Commensurate ‘Freezing’ of Alkanes in the Channels of a Zeolite. *Nature* 374:42–44.
- Smit B and J I Siepmann 1994. Simulating the Adsorption of Alkanes in Zeolites. *Science* 264:1118–1120.
- Stephenson G 1973. *Mathematical Methods for Science Students*. London, Longman.
- Swaminathan S, G Ravishanker and D L Beveridge 1991. Molecular Dynamics of B-DNA Including Water and Counterions – A 140-ps Trajectory for d(CGCGAATTCGCG) Based on the Gromos Force Field. *Journal of the American Chemical Society* 113:5027–5040.

- Verdier P H and W H Stockmayer 1962. Monte Carlo Calculations on the Dynamics of Polymers in Dilute Solution. *Journal of Chemical Physics* **36**:227-235.
- Vesely F J 1982. Angular Monte Carlo Integration Using Quaternion Parameters: A Spherical Reference Potential for CCl_4 . *Journal of Computational Physics* **47**:291-296.
- Vlugt T J H, R Krishna and B Smit 1999. Molecular Simulations of Adsorption Isotherms for Linear and Branched Alkanes and Their Mixtures in Silicalite. *Journal of Physical Chemistry* **103**:1102-1118.
- Vlugt T J H, W Zhu, F Kapteijn, J A Moulijn, B Smit and R Krishna 1998. Adsorption of Linear and Branched Alkanes in the Zeolite Silicalite-1. *Journal of the American Chemical Society* **120**:5599-5600.
- Widom B 1963. Topics in the Theory of Fluids. *Journal of Chemical Physics* **39**:2808-2812.
- Xu H and B J Berne 1999. Multicanonical Jump Walking: A Method for Efficiently Sampling Rough Energy Landscapes. *Journal of Chemical Physics* **110**:10299-10306.

CHAPTER NINE

Conformational Analysis

9.1 Introduction

The physical, chemical and biological properties of a molecule often depend critically upon the three-dimensional structures, or *conformations*, that it can adopt. Conformational analysis is the study of the conformations of a molecule and their influence on its properties. The development of modern conformational analysis is often attributed to D H R Barton, who showed in 1950 that the reactivity of substituted cyclohexanes was influenced by the equatorial or axial nature of the substituents [Barton 1950]. An equally important reason for the development of conformational analysis at that time was the introduction of analytical techniques such as infrared spectroscopy, NMR and X-ray crystallography, which actually enabled the conformation to be determined.

The conformations of a molecule are traditionally defined as those arrangements of its atoms in space that can be interconverted purely by rotation about single bonds. This definition is usually relaxed in recognition of the fact that small distortions in bond angles and bond lengths often accompany conformational changes, and that rotations can occur about bonds in conjugated systems that have an order between one and two.

A key component of a conformational analysis is the *conformational search*, the objective of which is to identify the 'preferred' conformations of a molecule: those conformations that determine its behaviour. This usually requires us to locate conformations that are at minimum points on the energy surface. Energy minimisation methods therefore play a crucial role in conformational analysis. An important feature of methods for performing energy minimisation is that they move to the minimum point that is closest to the starting structure. For this reason, it is necessary to have a separate algorithm which generates the initial starting structures for subsequent minimisation. It is these algorithms for generating initial structures that will be a major focus of this chapter. It is important to recognise the difference between a conformational search and a molecular dynamics or Monte Carlo simulation; the conformational search is concerned solely with locating minimum energy structures, whereas the simulation generates an ensemble of states that includes structures not at energy minima. However, as we shall see, both the Monte Carlo and molecular dynamics methods can be used as part of a conformational search strategy.

If possible, it is desirable to identify all minimum energy conformations on the energy surface. However, the number of minima may be so large that it is impractical to contemplate finding them all. Under such circumstances it is usual to try to find all the accessible minima. The relative populations of a molecule's conformations can be calculated using statistical mechanics via the Boltzmann distribution, though it is important to remember that the statistical weights involve contributions from all the degrees of freedom, including the vibrations as

well as the energies. Solvation effects may also be important, and various schemes are now available for calculating the solvation free energy of a conformation, which can be added to the intramolecular energy. These solvation schemes (which will be discussed in more detail in Section 11.9) provide computationally efficient ways to include the effects of the solvent on conformational equilibria. For some molecules such as proteins there are so many minima on the energy surface that it is impractical to try to find them all. Under such circumstances, it is often assumed that the native (i.e. naturally occurring) conformation is the one with the very lowest value of the energy function. This conformation is usually referred to as the *global minimum energy conformation*. One should usually be wary of algorithms which find only a single conformation. For example, even though the global minimum energy conformation has the lowest energy, it may not be the most highly populated because of the contribution of the vibrational energy levels to the statistical weight of each structure. Moreover, the global minimum energy conformation may not be the active (i.e. functional) structure. Indeed, in some cases it is possible that the active conformation does not correspond to any minimum on the energy surface of the isolated molecule. It may even be necessary for a molecule to adopt more than one conformation. For example, a substrate might bind in one conformation to an enzyme and then adopt a different conformation prior to reaction.

Conformational search methods can be conveniently divided into the following categories: systematic search algorithms, model-building methods, random approaches, distance geometry and molecular dynamics. Before discussing these methods, we should note that conformational analysis can sometimes be performed quite effectively using Dreiding or CPK mechanical models. The invention of these models should be regarded as an important development in conformational analysis and molecular modelling. Mechanical models do, however, have some shortcomings. For example, they provide no quantitative information about the relative energies of the various conformations. It is often quite difficult to make accurate measurements of a molecule's internal coordinates such as the distance between two atoms that are on opposite sides of the structure. They are subject to the forces of gravity (which makes them unwieldy for large molecules) and the hands of marauding colleagues! It is also difficult to construct models that have significant deviations from standard bond lengths and angles. Nevertheless, manual models can be very useful, particularly as they are portable and because they can be easily manipulated in a way that is often not possible with computer images (although this may change with the development of 'virtual reality' molecular modelling systems).

We next introduce the basic algorithms and then describe some of the many variants upon them. We then discuss two methods called evolutionary algorithms and simulated annealing, which are generic methods for locating the globally optimal solution. Finally, we discuss some of the ways in which one might analyse the data from a conformational analysis in order to identify a 'representative' set of conformations.

9.2 Systematic Methods for Exploring Conformational Space

As the name suggests, a systematic search explores the conformational space by making regular and predictable changes to the conformation. The simplest type of systematic search

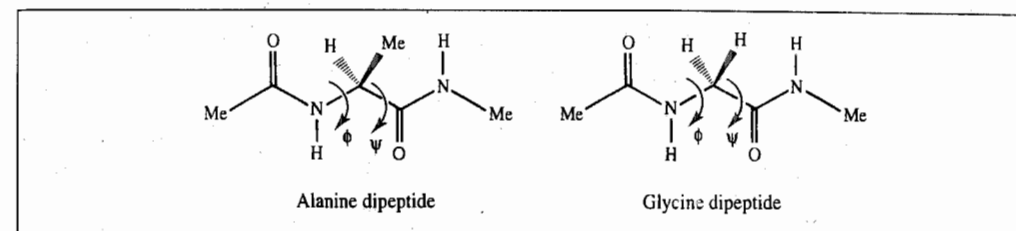


Fig. 9.1: The alanine dipeptide and the glycine dipeptide.

(often called a *grid search*) is as follows. First, all rotatable bonds in the molecule are identified. The bond lengths and angles remain fixed throughout the calculation. Each of these bonds is then systematically rotated through 360° using a fixed increment. Every conformation so generated is subjected to energy minimisation to derive the associated minimum energy conformation. The search stops when all possible combinations of torsion angles have been generated and minimised. To illustrate the grid search algorithm, let us consider the conformational energy surface for the 'alanine dipeptide' $\text{CH}_3\text{CONHCHMeCONHCH}_3$ (Figure 9.1), which is used as a model for the conformational behaviour of amino acids in proteins. If we assume that the bond lengths and bond angles are fixed and that the amide bonds adopt *trans* conformations then only the two torsion angles labelled ϕ and ψ in Figure 9.1 can vary. The energy is then a function of just these two variables, and as such it can be represented as a contour diagram as shown in Figure 9.2. This contour plot is known as a *Ramachandran map*, after G N Ramachandran who showed that the amino acids were restricted to a limited range of conformations [Ramachandran *et al.* 1963]. The accessible areas on the contour maps calculated by Ramachandran do indeed correspond to those conformations that are observed in X-ray structures of proteins (Figure 9.3). Two regions are particularly important; these correspond to the α -helix and β -strand structures,

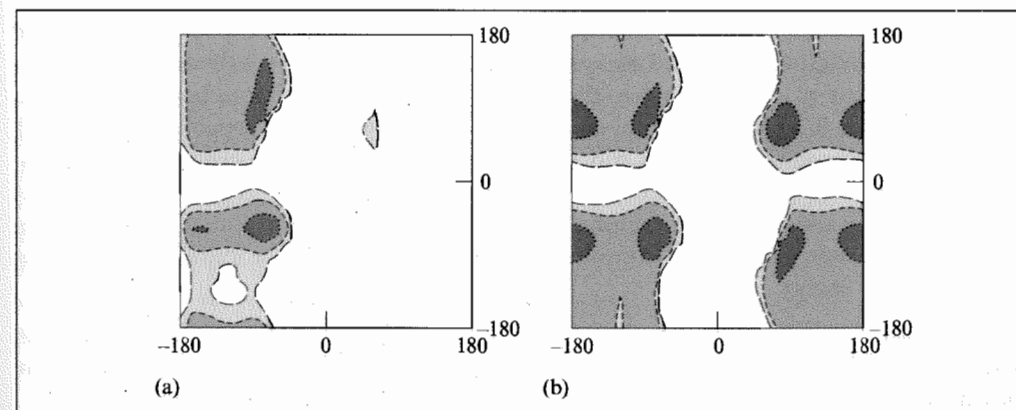


Fig. 9.2: Ramachandran map for the alanine dipeptide (a) and glycine dipeptide (b), calculated using the AMBER force field [Weiner *et al.* 1984]. In both cases contours are drawn at 1.0, 2.0 and 3.0 kcal/mol above the lowest-energy conformation found.

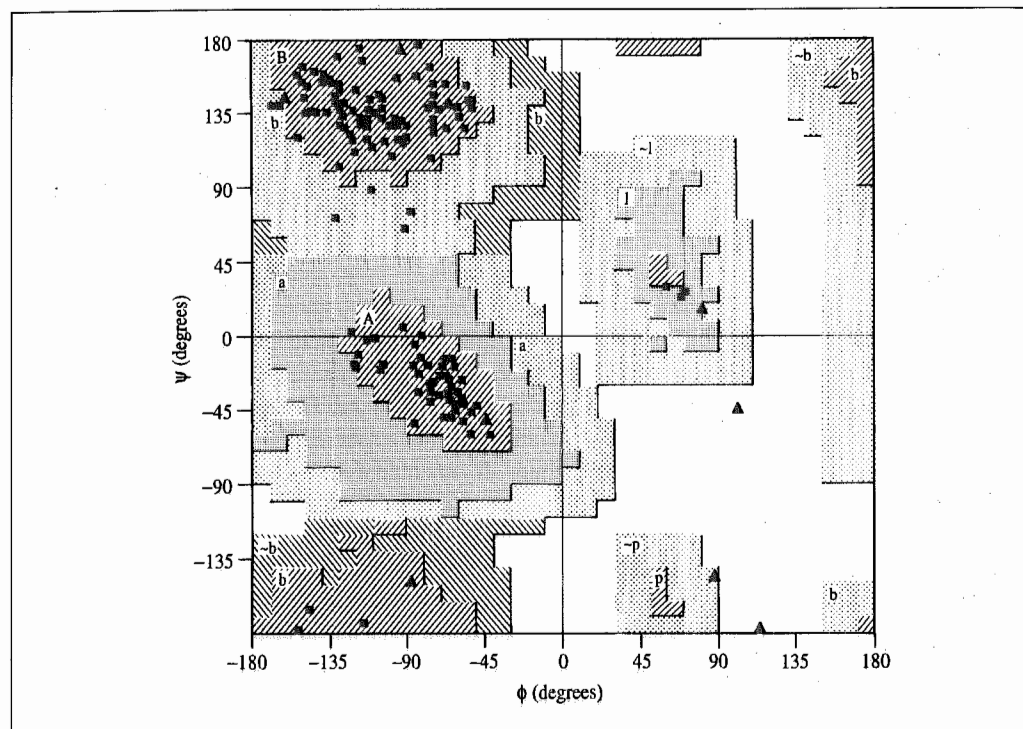


Fig. 9.3: Experimentally observed distribution of (ϕ, ψ) angles in the enzyme dihydrofolate reductase. The symbols are the actual values; the shaded areas correspond to the (ϕ, ψ) distribution averaged over many protein structures.

which will be discussed in more detail in Section 10.2. The amino acid glycine, which has no side chain (Figure 9.1), has a wider range of accessible conformations than the other amino acids, as can be seen from the Ramachandran map in Figure 9.2.

To perform a grid search of the conformational space of the alanine dipeptide, a series of conformations would be generated by systematically varying ϕ and ψ between 0° and 360° . This is equivalent to drawing a two-dimensional grid over the Ramachandran contour diagram in Figure 9.2; each grid point corresponds to a conformation generated by the grid search with some combination of ϕ and ψ . It can readily be seen that, even for a relatively large torsional increment, the number of conformations generated by the grid search is much larger than the number of minima on the surface; many of the initial conformations minimise to the same minimum energy structure. Moreover many of the initial conformations are very high in energy.

A major drawback of the grid search is that the number of structures to be generated and minimised increases in an exceptional fashion with the number of rotatable bonds, a phenomenon known as the *combinatorial explosion*. The number of structures generated is given by:

$$\text{Number of conformations} = \prod_{i=1}^N \frac{360}{\theta_i} \quad (9.1)$$

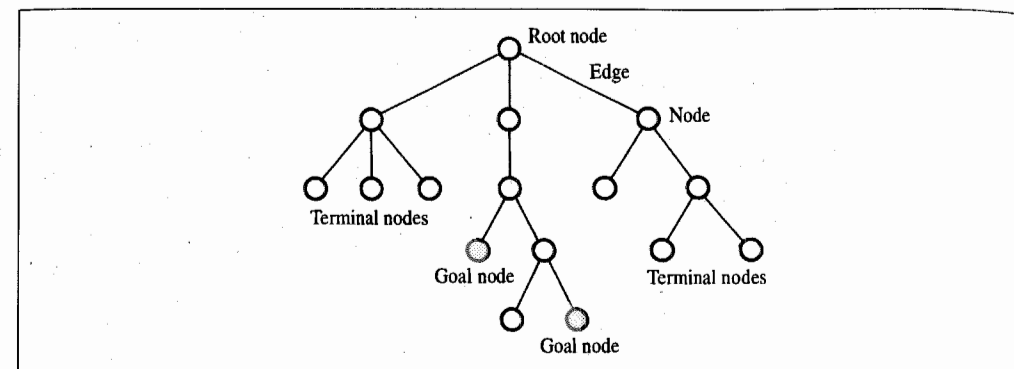


Fig. 9.4: Schematic illustration of a search tree.

where θ_i is the dihedral increment chosen for bond i . For example, if there are five bonds and an increment of 30° is used for each bond, then 248 832 structures will be generated. If the number of bonds is increased to seven, then the number of structures increases to almost 36 million. To put these figures into context, suppose each structure takes just one second to minimise. The five-bond problem will then require 69 hours to complete and the seven-bond problem will require 415 days. Despite this apparent limitation, systematic search algorithms are routinely employed to consider problems involving 10–15 bonds. This is achieved by eliminating from the time-consuming energy minimisation stage structures that have a very high energy or some other problem. A good way to understand these enhanced systematic search methods is to use a *tree representation* of the problem.

Search trees are widely used to represent the different states that a problem can adopt. An example is shown in Figure 9.4 from which it should be clear where the name derives, especially if the page is turned upside down. A tree contains *nodes* that are connected by *edges*. The presence of an edge indicates that the two nodes it connects are related in some way. Each node represents a state that the system may adopt. The *root node* represents the initial state of the system. *Terminal nodes* have no child nodes. A *goal node* is a special kind of terminal node that corresponds to an acceptable solution to the problem.

Suppose we wish to use a grid search to explore the conformational space of a simple alkane, *n*-hexane. We will assume that rotation of the terminal methyl groups can be ignored and so just three bonds can vary. If we permit each of the variable bonds to adopt just three values, corresponding to the *trans*, *gauche*(+) and *gauche*(-) conformations*, then the search tree for this problem contains 27 terminal nodes ($\equiv 3 \times 3 \times 3$) and is shown in Figure 9.5. The root node represents the starting point, where none of the rotatable bonds have been assigned a torsion angle. When the first rotatable bond is set to its first value (i.e. the *trans* conformation with a torsion angle of 180°), this corresponds to moving from the root node to the node numbered 1 in the tree. The second bond is now set to a *trans* conformation; this corresponds

* The *trans* conformation corresponds to a torsion angle of 180° , the *gauche*(+) conformation to one of $+60^\circ$ and the *gauche*(-) conformation to -60° . These approximately correspond to the torsion angles of the three minimum energy conformations of butane.

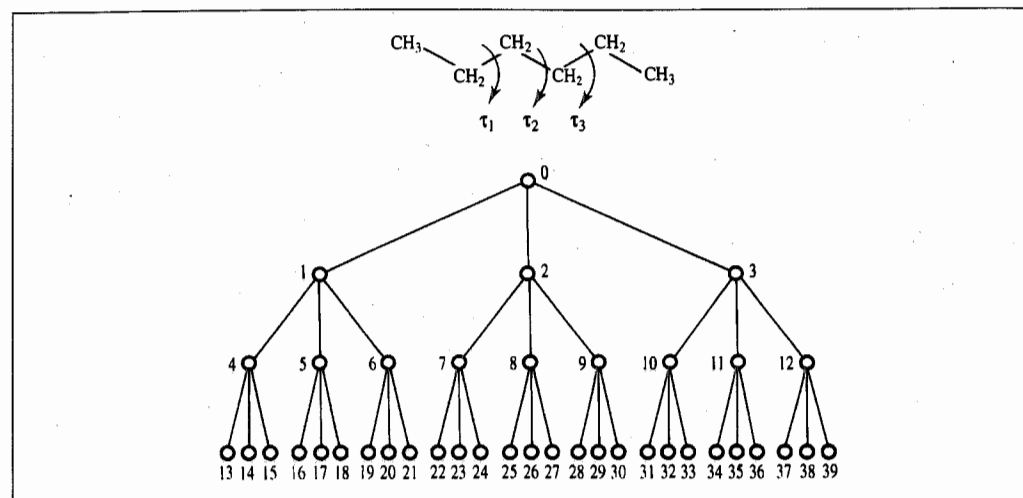


Fig. 9.5: Tree representation of the conformational search problem for hexane. Unlike the tree in Figure 9.4 the path length from the root node to any of the terminal nodes is constant.

to a move to node 4. When the third bond is assigned *trans* we reach node 13, which is a terminal node and corresponds to a conformation ready for minimisation. To generate further conformations, it is necessary to change one of the three torsion angles. The most convenient way to do this is to assign a new value to the last bond to be set (i.e. bond 3). Setting bond 3 to a *gauche*(+) conformation is equivalent to moving back up the tree from node 13 to node 4 and then down to the terminal node 14. This gives a second completed conformation. By proceeding in this fashion through the search tree (a process called *backtracking*) all conformations of the molecule can be generated. The search algorithm we have described is known as the *depth-first search*.

The efficiency of a depth-first search can be enhanced by discarding structures that violate some energetic or geometric criterion. Structures with high-energy steric interactions are

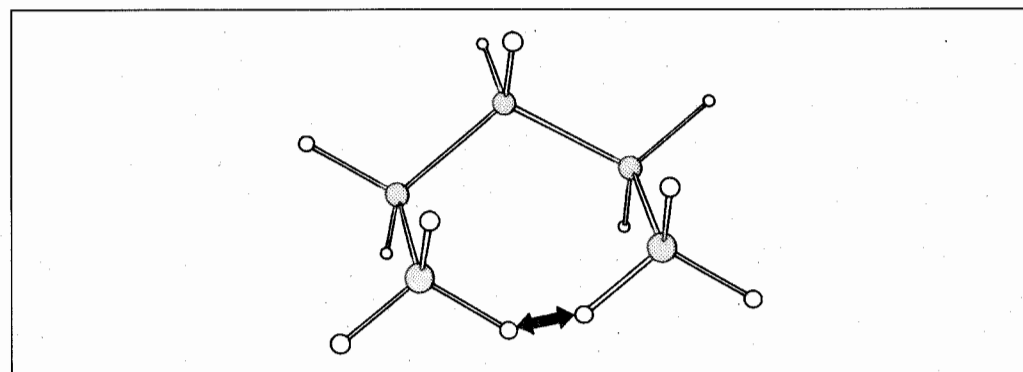


Fig. 9.6: A pentane violation arises when there are successive *gauche*(+) and *gauche*(-) torsion angles in an alkane chain.

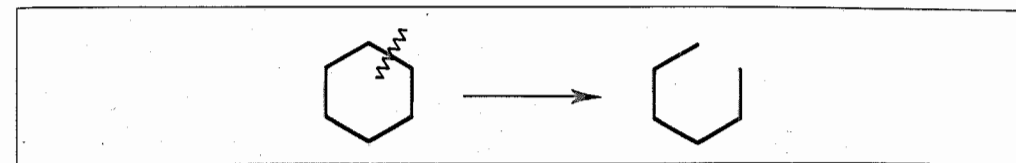


Fig. 9.7: A 'pseudo-acyclic' molecule is generated by breaking the ring.

then rejected before the energy-minimisation stage. We can further enhance the efficiency of the systematic search by checking partially constructed conformations before all the torsion angles have been assigned. Suppose we generate a partial structure containing two non-bonded atoms that are very close in space. In hexane, such a high-energy structure is generated if the first rotatable bond is set to a *gauche*(+) conformation and the second rotatable bond is assigned *gauche*(-) (a pentane violation, Figure 9.6). Whatever value is assigned to the third torsion angle, this high-energy steric problem will remain. All structures that lie below that node in the search tree (number 9 in Figure 9.5) can thus be eliminated, or *pruned*. It is important to stress that this is only possible if those parts of the molecule that are in violation will not be moved relative to each other by a subsequent torsional assignment.

Cyclic molecules are often quite difficult to analyse using a systematic search. The usual strategy is to break the ring, giving a 'pseudo-acyclic' molecule that can then be treated as a normal acyclic molecule. This process is illustrated for cyclohexane in Figure 9.7. When searching the conformational space of cyclic molecules additional checks must be included to ensure that the rings are properly formed. For example, an all-*trans* structure is a perfectly acceptable conformation of *n*-hexane, but is not an acceptable conformation of cyclohexane due to the unreasonable bond length between the ring-closure atoms. It is therefore common practice to check several intramolecular parameters when using a systematic search to explore the conformational space of a cyclic system; these parameters usually include the bond length between the ring-closure bonds, together with the bond angles at these atoms (Figure 9.8). In some programs other internal parameters (e.g. the torsion angles adjacent to the ring closure bond) are also checked. The main reason why rings are problematic

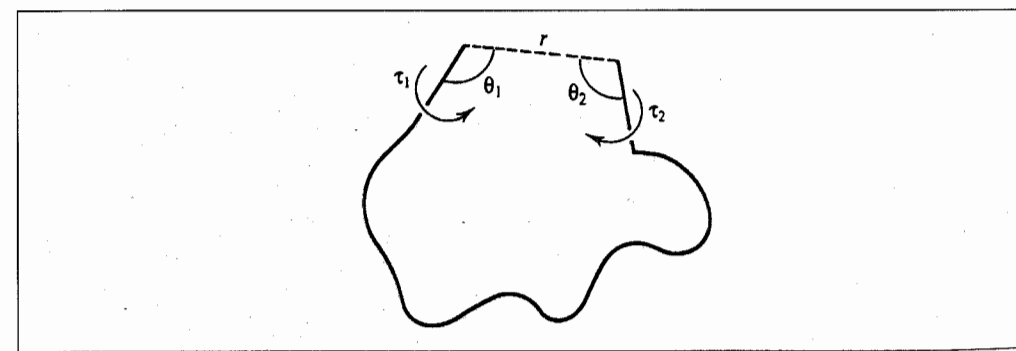


Fig. 9.8: The intramolecular parameters that may be checked when exploring the conformational space of a ring system.

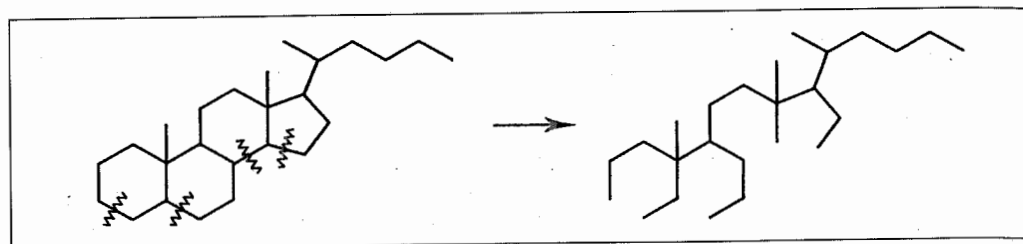


Fig. 9.9: Creation of pseudo-acyclic molecule for a system with many rings.

for the systematic search is that these checks can often be applied only very late in the analysis; it is often necessary almost to complete the ring before the structure is rejected or accepted. One simple check that can be used when constructing cyclic molecules is to ensure that at all stages the distance of the growing chain from the start atom is short enough to enable the remaining bonds to close the ring.

The systematic search is most efficient when the rotatable bonds are processed in a unidirectional fashion. This ensures that once an atom has been fixed in space relative to the atoms already considered then it will not be moved again. For an acyclic molecule this means that the search starts at one end of the molecule and moves down the chain. Molecules containing rings are treated by opening the cycles to give the pseudo-acyclic molecule as described above and then processed (Figure 9.9).

Any systematic search ultimately requires a balance to be made between the resolution of the grid and the available computer resources. Too fine a grid and the search may take too long; too coarse and important minima may be missed. The non-bonded criteria, which determine whether a structure is to be rejected, must also be assigned. The non-bonded criterion is often referred to as a 'bump check' and is usually set to a modest value (say, 2.0 Å) as the energy-minimisation step will be able to remove minor problems in the structure. For cyclic molecules the ring-closure criteria may also affect the results. It must also be remembered that the various cutoffs are interdependent, so that changing one may require others to be reassigned.

9.3 Model-building Approaches

One way to at least partially alleviate the combinatorial explosion that inevitably accompanies a systematic search is to use larger 'building blocks', or *molecular fragments*, to construct the conformations [Gibson and Scheraga 1987; Leach *et al.* 1988, 1990]. Fragment- or model-building approaches to conformational analysis construct conformations of a molecule by joining together three-dimensional structures of molecular fragments. This approach would be expected to be more efficient than the normal systematic search because there are usually many fewer combinations of fragment conformations than combinations of torsion angle values. This is particularly so for cyclic fragments, which are in any case problematic for the systematic search method. For example, the molecule in Figure 9.10 could be constructed by joining together the fragments indicated. Many molecular modelling systems offer a facility

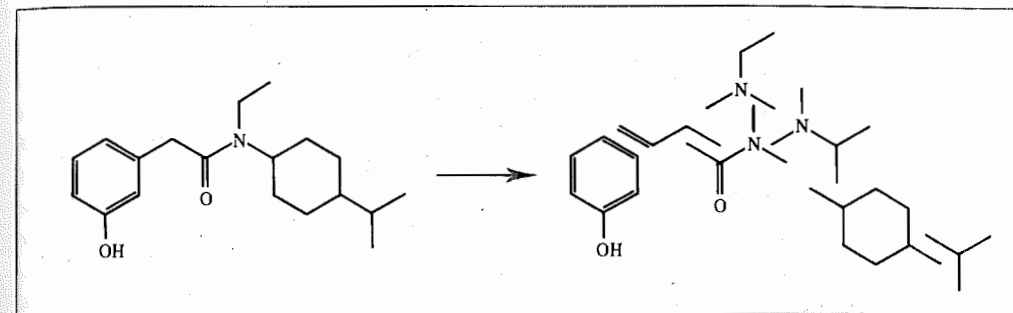


Fig. 9.10: A conformation may be obtained by joining appropriate fragments.

for constructing structures from molecular fragments, though the user usually has to specify manually which fragments are to be joined and how this is to be achieved. Clearly, if each fragment can adopt a number of conformations, then it is impractical to tackle the problem manually and some means of automating the method is required.

A program to explore conformational space automatically using the fragment-building approach must first decide which fragments are needed to construct the molecule [Leach *et al.* 1990]. This is done using a *substructure search algorithm*, which determines whether each of the fragments that the program 'knows about' is present in the molecule and how the atoms in the fragment match onto the atoms in the molecule. Having identified the fragments that are required, conformations can be generated. The conformations available to each fragment (often called templates) should span the range of conformations the fragment can adopt. For example, cyclohexane rings adopt the chair, twist-boat and boat conformations in molecules and so templates corresponding to these structures should be available. A conformation of the molecule is constructed by assigning a template to each fragment and then attempting to join the templates together. The search problem can be represented as a tree, as for a systematic search, and so all of the usual tree-searching algorithms are applicable. The search can be significantly enhanced by tree pruning.

The fragment-based approach to conformational analysis relies upon two assumptions. The first assumption is that each fragment must be conformationally independent of the other fragments in the molecule. The second assumption is that the conformations stored for each fragment must cover the range of structures that are observed in fully constructed molecules. The fragment conformations can be obtained from a variety of sources; two common approaches are by analysing a structural database (see Section 9.11) or from some other conformational search method. A third limitation is that one can obviously only analyse molecules for which there are fragments available.

9.4 Random Search Methods

A random search is, in many ways, the antithesis of a systematic search. A systematic search explores the energy surface of the molecule in a predictable fashion, whereas it is not

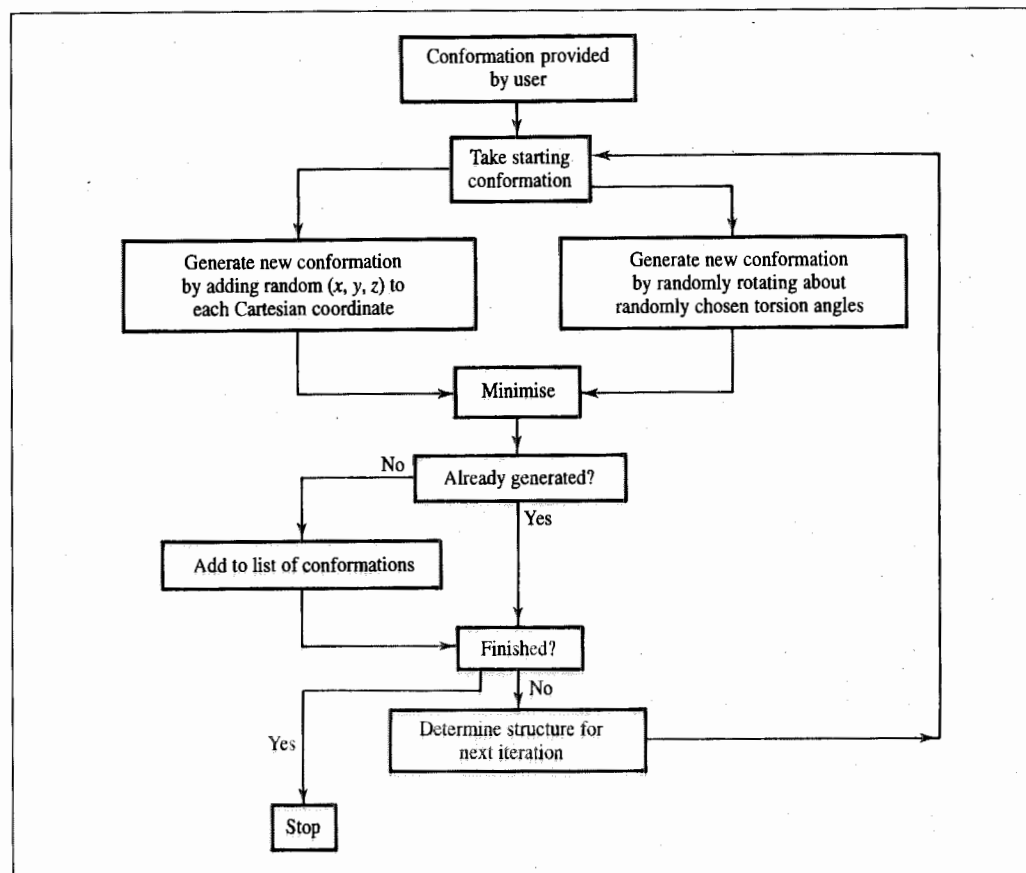


Fig. 9.11: Flow chart steps followed by a random conformational search.

possible to predict the order in which conformations will be generated by a random method. A random search can move from one region of the energy surface to a completely unconnected region in a single step. A random search can explore conformational space by changing either the atomic Cartesian coordinates or the torsion angles of rotatable bonds. Both types of algorithm use a similar approach, which is outlined in flow chart form in Figure 9.11. At each iteration, a random change is made to the 'current' conformation. The new structure is then refined using energy minimisation. If the minimised conformation has not been found previously, it is stored. The conformation to be used as the starting point for the next iteration is then chosen and the cycle starts again. The procedure continues until a given number of iterations have been performed or until it is decided that no new conformations can be found.

The Cartesian and dihedral versions of the random search differ in the way in which each new structure is generated. The Cartesian method adds a random amount to the x , y and z coordinates of all the atoms in the molecule [Saunders 1987; Ferguson and Raber 1989]

whereas the dihedral method generates new conformations by making random changes to the rotatable bonds, with the bond lengths and bond angles being kept fixed [Li and Scheraga 1987; Chang *et al.* 1989]. The Cartesian method is extremely simple to implement but does rely heavily upon the minimisation step as the initial structures that are generated by the randomisation procedure can be very distorted and extremely high in energy; in some implementations the coordinates can change by 3 Å or more. The advantage of the dihedral search method is that many fewer degrees of freedom need to be considered. However, special procedures are required when applying the dihedral method to molecules that contain rings; typically these are broken in a manner similar to the systematic search described above to give a pseudo-acyclic molecule (Figure 9.9). After randomisation each ring is checked to ensure that any ring-closure constraints are satisfied. In the random dihedral method it is possible to change all dihedrals or just a randomly chosen subset of them.

There are many ways in which the structure for input to the next iteration of the search can be selected. A simple approach is to take the structure obtained from the previous step. An alternative is to select randomly a structure from those generated previously, weighting the choice towards those structures that have been selected the least (a *uniform usage* protocol). A third method is to use the lowest-energy structure found so far, or to bias the selection towards the lowest-energy structures. The Metropolis Monte Carlo scheme is often used to make the choice. Each newly generated structure (after energy minimisation) is accepted as the starting point for the next iteration if it is lower in energy than the previous structure or if the Boltzmann factor of the energy difference, $\exp[-(\mathcal{V}_{\text{new}}(\mathbf{r}^N) - \mathcal{V}_{\text{old}}(\mathbf{r}^N))/k_B T]$ is larger than a random number between 0 and 1. If not, the previous structure is retained for the next iteration. There is no fundamental reason why any of these methods should be preferred over another, but some are reported to be more efficient at exploring the conformational space or finding the global minimum energy conformation.

In a systematic search there is a defined endpoint to the procedure, which is reached when all possible combinations of bond rotations have been considered. In a random search, there is no natural endpoint; one can never be absolutely sure that all of the minimum energy conformations have been found. The usual strategy is to generate conformations until no new structures can be obtained. This usually requires each structure to be generated many times and so the random methods inevitably explore each region of the conformational space a large number of times.

9.5 Distance Geometry

One way to describe the conformation of a molecule other than by Cartesian or internal coordinates is in terms of the distances between all pairs of atoms. There are $N(N-1)/2$ interatomic distances in a molecule, which are most conveniently represented using an $N \times N$ symmetric matrix. In such a matrix, the elements (i, j) and (j, i) contain the distance between atoms i and j and the diagonal elements are all zero. Distance geometry explores conformational space by randomly generating many distance matrices, which are then converted into conformations in Cartesian space. The crucial feature about distance geometry (and the reason why it works) is that it is not possible to arbitrarily assign values to the

interatomic distances in a molecule and always obtain a low-energy conformation. Rather, the interatomic distances are closely interrelated and indeed many combinations of distances are geometrically impossible. This can be illustrated using a simple three-atom molecule (ABC). Simple trigonometry requires that the sum of the distances AB and AC must be greater than or equal to the distance BC. Thus, a conformation in which the distances are $AB = 1.5 \text{ \AA}$, $AC = 1.4 \text{ \AA}$ and $BC = 3.5 \text{ \AA}$ is not geometrically possible.

Distance geometry uses a four-stage process to derive a conformation of a molecule [Crippen 1981; Crippen and Havel 1988]. First, a matrix of upper and lower interatomic distance bounds is calculated. This matrix contains the maximum and minimum values permitted to each interatomic distance in the molecule. Values are then randomly assigned to each interatomic distance between its upper and lower bounds. In the third step, the distance matrix is converted into a trial set of Cartesian coordinates, which in the fourth step are then refined.

Some of the interatomic distance bounds can be determined from simple chemical principles. For example, X-ray crystallographic studies have shown that bond lengths are restricted to a small range of values that are determined primarily by the atomic number and hybridisation of the two atoms. The distance between two atoms which are both bonded to a third atom (i.e. are in 1,3 relationship) is also severely restricted and can be calculated from the angle at the central atom and the lengths of the two bonds. The distance between two atoms that are separated by three bonds (i.e. are in a 1,4 relationship) can vary with the torsion angle of the central bond, the minimum distance corresponding to a torsion angle of 0° and the maximum distance to a torsion angle of 180° . These three cases are shown in Figure 9.12. It is not so easy to determine limits on the other interatomic distances (i.e. between atoms in a 1, n relationship where $n > 4$) but it is usual to require that such atom pairs do not approach closer than the sum of the van der Waals radii of the two atoms. The upper bound is then usually assigned an arbitrarily large value.

A procedure called *triangle smoothing* is then used to refine the initial set of distance bounds. Triangle smoothing uses two simple trigonometrical restrictions on groups of three atoms, which are illustrated in Figure 9.13. The first restriction is that the distance between two atoms A and C can be no greater than the sum of the maximum values of the distances AB and BC. This can be written:

$$u_{AC} \leq u_{AB} + u_{BC} \quad (9.2)$$

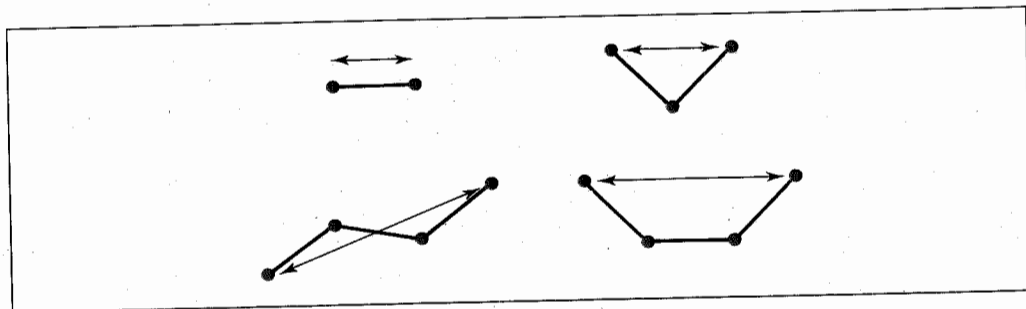


Fig. 9.12: The upper and lower distance bounds for atoms in 1,2, in 1,3 and 1,4 relationships can be derived from simple chemical principles.

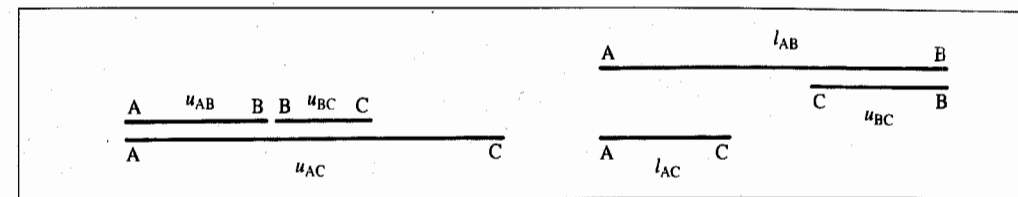


Fig. 9.13: The two triangle inequalities used in distance geometry.

where u_{AB} indicates the upper bound on the AB distance. The second restriction is that the minimum value of the AC distance can be no less than the difference between the lower bound on AB and the upper bound on BC:

$$l_{AC} \geq l_{AB} - u_{BC} \quad (9.3)$$

where l_{AB} is used to indicate the lower bound distance. These two inequalities are repeatedly applied to the set of distances bounds until the entire set of distance bounds is self-consistent and all possible interatomic distance triplets satisfy both inequalities. Triangle smoothing need be performed only once for each molecule.

We can now proceed to the generation of conformations. First, random values are assigned to all the interatomic distances between the upper and lower bounds to give a trial distance matrix. This distance matrix is now subjected to a process called *embedding*, in which the 'distance space' representation of the conformation is converted to a set of atomic Cartesian coordinates by performing a series of matrix operations. We calculate the *metric matrix*, G , each of whose elements (i, j) is equal to the scalar product of the vectors from the origin to atoms i and j :

$$G_{ij} = \mathbf{i} \cdot \mathbf{j} \quad (9.4)$$

The elements G_{ij} can be calculated from the distance matrix using the cosine rule:

$$G_{ij} = (d_{i0}^2 + d_{j0}^2 - d_{ij}^2)/2 \quad (9.5)$$

where d_{i0} is the distance from the origin to atom i and d_{ij} is the distance between atoms i and j .

It is usual to take the centre of the molecule as the origin of the coordinate system. The distance of each atom from the centre can be calculated directly from the interatomic distances using the following expression:

$$d_{i0}^2 = \frac{1}{N} \sum_{j=1}^N d_{ij}^2 - \frac{1}{N^2} \sum_{j=2}^N \sum_{k=1}^{j-1} d_{jk}^2 \quad (9.6)$$

The metric matrix G is a square symmetric matrix. A general property of such matrices is that they can be decomposed as follows:

$$G = \mathbf{V} \mathbf{L}^2 \mathbf{V}^T \quad (9.7)$$

The diagonal elements of \mathbf{L}^2 are the eigenvalues of G and the columns of \mathbf{V} are its eigenvectors. The atomic coordinates can be derived from the metric matrix by rewriting

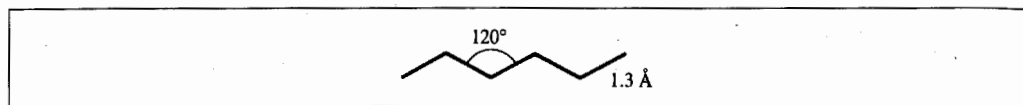


Fig. 9.14: Five-carbon fragment to illustrate distance geometry algorithm.

Equation (9.4) as

$$\mathbf{G} = \mathbf{X}\mathbf{X}^T \quad (9.8)$$

where \mathbf{X} is a matrix containing the atomic coordinates. Equating Equations (9.7) and (9.8) gives:

$$\mathbf{X} = \mathbf{V}\mathbf{L} \quad (9.9)$$

As \mathbf{L} has only diagonal entries, the matrix \mathbf{L} is identical to its transpose: $\mathbf{L} = \mathbf{L}^T$. The atomic coordinates are thus obtained by multiplying the square roots of the eigenvalues by the eigenvectors.

The triangle-smoothing and embedding steps of distance geometry are best understood using a specific example. Let us consider a five-atom, all-carbon fragment (Figure 9.14), in which all of the carbon-carbon bonds are assumed to have an optimal length of 1.3 Å and all internal angles are 120°. If we further assume that the carbon van der Waals radius is 1.4 Å, then the initial bounds matrix is as follows:

$$\begin{pmatrix} 0.0 & 1.3 & 2.2517 & 3.4395 & 99.0 \\ 1.3 & 0.0 & 1.3 & 2.2517 & 3.4395 \\ 2.2517 & 1.3 & 0.0 & 1.3 & 2.2517 \\ 2.6 & 2.2517 & 1.3 & 0.0 & 1.3 \\ 2.8 & 2.6 & 2.2517 & 1.3 & 0.0 \end{pmatrix} \quad (9.10)$$

Note that the lower bound for the distance between atoms 1 and 5 equals the sum of their van der Waals radii and that the upper bound has been set to an arbitrarily large value of 99 Å. All other distances have been allocated on the basis of geometric arguments. In a 'real' example the upper and lower bounds for bonded atoms would usually be slightly different (by approximately 0.1 Å) to reflect the fact that bond lengths in real molecules do vary a little. The 1,3 bounds would also be made slightly different. After triangle smoothing only one distance bound is changed, the upper bound of the distance between atoms 1 and 5. This distance is changed to a value that is equal to the sum of the upper bounds of the distances between atoms 1 and 3 and 3 and 5. The smoothed bounds matrix that results is:

$$\begin{pmatrix} 0.0 & 1.3 & 2.2517 & 3.4395 & 4.5033 \\ 1.3 & 0.0 & 1.3 & 2.2517 & 3.4395 \\ 2.2517 & 1.3 & 0.0 & 1.3 & 2.2517 \\ 2.6 & 2.2517 & 1.3 & 0.0 & 1.3 \\ 2.8 & 2.6 & 2.2517 & 1.3 & 0.0 \end{pmatrix} \quad (9.11)$$

Suppose interatomic distances are now randomly assigned between the lower and upper bounds to give the following distance matrix:

$$\begin{pmatrix} 0.0 & 1.3 & 2.25 & 3.11 & 3.42 \\ & 0.0 & 1.3 & 2.25 & 2.85 \\ & & 0.0 & 1.3 & 2.25 \\ & & & 0.0 & 1.3 \\ & & & & 0.0 \end{pmatrix} \quad (9.12)$$

The corresponding metric matrix is:

$$\begin{pmatrix} 3.571 & 1.569 & -0.427 & -2.276 & -2.436 \\ 1.569 & 1.256 & 0.105 & -1.122 & -1.808 \\ -0.427 & 0.105 & 0.644 & 0.261 & -0.583 \\ -2.276 & -1.122 & 0.261 & 1.569 & 1.569 \\ -2.436 & -1.808 & -0.583 & 1.569 & 3.259 \end{pmatrix} \quad (9.13)$$

The eigenvalues of this matrix are 8.18, 1.74, 0.26, 0.10 and 0.0, with the matrix of eigenvectors being:

$$\mathbf{W} = \begin{pmatrix} 0.621 & 0.455 & -0.425 & 0.164 \\ 0.355 & -0.184 & 0.800 & 0.020 \\ 0.0 & -0.573 & -0.368 & -0.580 \\ -0.408 & -0.287 & -0.153 & 0.727 \\ -0.567 & 0.590 & 0.145 & -0.330 \end{pmatrix} \quad (9.14)$$

The 'best' three-dimensional structure is obtained by taking the eigenvectors that correspond to the three largest eigenvalues, providing they are all positive. If these eigenvalues are λ_1 , λ_2 and λ_3 and \mathbf{W} is the matrix containing the associated eigenvectors, then the Cartesian coordinates (x_i , y_i , z_i) of each atom i are calculated as follows:

$$x_i = \sqrt{\lambda_1} W_{i1} \quad (9.15)$$

$$y_i = \sqrt{\lambda_2} W_{i2} \quad (9.16)$$

$$z_i = \sqrt{\lambda_3} W_{i3} \quad (9.17)$$

For our five-carbon example, the coordinates obtained using the three highest eigenvalues are:

Atom	x coordinate	y coordinate	z coordinate
1	1.777	0.601	-0.218
2	1.014	-0.244	0.410
3	-0.001	-0.757	-0.188
4	-1.166	-0.379	-0.079
5	-1.623	0.799	0.075

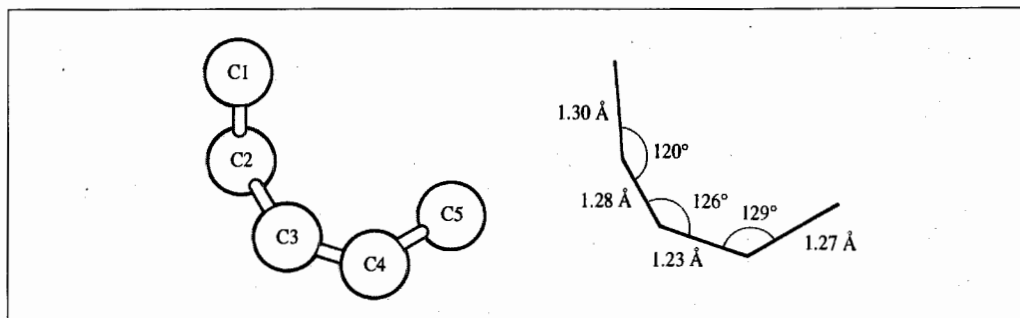


Fig. 9.15: Conformation of the five-carbon fragment generated by distance geometry.

This conformation is illustrated schematically in Figure 9.15. The interatomic distance matrix for this conformation is:

$$\begin{pmatrix} 0.0 & 1.299 & 2.24 & 3.10 & 3.42 \\ & 0.0 & 1.29 & 2.24 & 2.85 \\ & & 0.0 & 1.23 & 2.25 \\ & & & 0.0 & 1.25 \\ & & & & 0.0 \end{pmatrix} \quad (9.18)$$

Note that the distances in this conformation do not equal the original randomly chosen distances, nor do they all lie between the upper and lower bound values in the bounds matrix. For example, the distance between atoms 4 and 5 is 1.25 Å rather than 1.3 Å. This is because it may be necessary to use more than just three dimensions to find a conformation that satisfies the distances in the initial distance matrix. The number of non-zero eigenvectors of the metric matrix equals the dimensionality of the space in which a solution can be found. In general, if there are N interatomic distances then a solution can be found in $N - 1$ dimensions. This is in part a consequence of the fact that three-dimensional objects must not only satisfy triangle inequalities but also tetrangle, pentangle and hexangle relationships. Moreover, triangle smoothing is often only applied to the bounds matrix; the distance matrix that is used as input to the embedding stage may, in fact, contain combinations of distances that violate the triangle inequalities. Improved sampling of conformational space can be achieved if the trial distances are selected so that they do satisfy the triangle inequalities (a process known as *metrisation*), but for reasons of computational cost it is often not used in its full form.

If we add the coordinates corresponding to the fourth eigenvalue, then the original distance matrix is reproduced exactly. These fourth-dimensional coordinates are as follows:

Atom	4th coordinate
1	0.053
2	0.006
3	-0.188
4	0.235
5	-0.107

The distance between atoms 4 and 5 in this four-dimensional space is exactly 1.3 Å.

In the final step of the distance geometry algorithm the coordinates are refined so that the conformation better satisfies the initial distance bounds. A conjugate gradients minimisation algorithm is often employed for this step. The function to be minimised has a positive value for distances that are outside the permitted range but is zero otherwise. The penalty functions most commonly used are:

$$E = \sum_i \sum_{j>i} \begin{cases} (d_{ij}^2 - u_{ij}^2)^2 & d_{ij} > u_{ij} \\ 0 & l_{ij} \leq d_{ij} \leq u_{ij} \\ (l_{ij}^2 - d_{ij}^2)^2 & d_{ij} < l_{ij} \end{cases} \quad (9.19)$$

$$E = \sum_i \sum_{j>i} \begin{cases} [(d_{ij}^2 - u_{ij}^2)/u_{ij}^2]^2 & d_{ij} > u_{ij} \\ 0 & l_{ij} \leq d_{ij} \leq u_{ij} \\ [(l_{ij}^2 - d_{ij}^2)/d_{ij}^2]^2 & d_{ij} < l_{ij} \end{cases} \quad (9.20)$$

where u_{ij} is the upper bound distance between atoms i and j and l_{ij} is the lower bound distance. The first function weights long distances more than short distances whereas the second error function weights all distances equally. Both functions are zero when all the distances are between the upper and lower bounds. A conformation in which all distance bounds are satisfied is not necessarily at an energy minimum, and so the final structure may subsequently be subjected to force field energy minimisation to derive the associated minimum energy structure.

During the optimisation of the structure against the distance constraints it is usual to incorporate *chiral constraints*. These are used to ensure that the final conformation is the desired stereoisomer. Chiral constraints are necessary because the interatomic distances in two enantiomeric conformations are identical and as a consequence the 'wrong' isomer may quite legitimately be generated. Chiral constraints are usually incorporated into the error function as a chiral volume, calculated as a scalar triple product. For example, to maintain the correct stereochemistry about the tetrahedral atom number 4 in Figure 9.16, the following scalar triple product must be positive:

$$(\mathbf{v}_1 - \mathbf{v}_4) \cdot [(\mathbf{v}_2 - \mathbf{v}_4) \times (\mathbf{v}_3 - \mathbf{v}_4)] \quad (9.21)$$

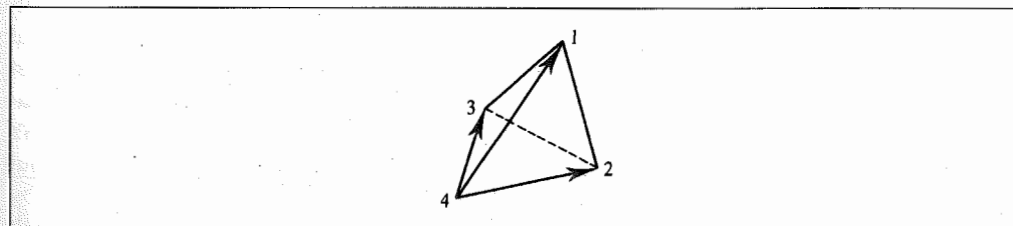


Fig. 9.16: The stereochemistry about tetrahedral atoms can be maintained with an appropriate chiral constraint.

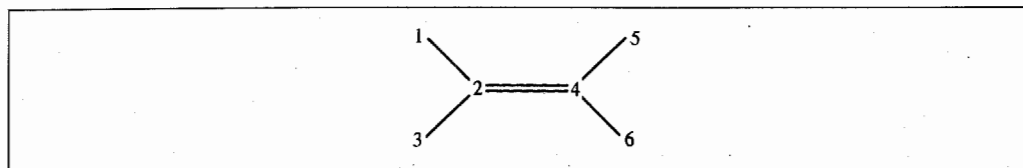


Fig. 9.17: A double bond can be forced to adopt a planar conformation through the use of appropriate chiral constraints.

The other stereoisomer corresponds to a negative chiral volume. Chiral constraints are included in the penalty function by adding terms of the following form:

$$(V_{\text{ch}} - V_{\text{ch}}^*)^2 \quad (9.22)$$

where V_{ch}^* is the desired value of the chiral constraint. Chiral constraints can also be used to force groups of atoms to lie in the same plane by requiring the chiral volume to have a value of zero. More than one such constraint may be required for each planar group. For example, to force all six atoms about the double bond in Figure 9.17 to lie in the same plane, the three sets of chiral volumes defined by the atoms 1, 2, 3, 4; 2, 4, 5, 6; 1, 2, 4, 5 must be zero. A commonly used strategy in many distance geometry programs is to perform the first few steps of refinement using a conformation defined in four dimensions, as this can help to invert any incorrect chiral centres. The minimisation then switches to three-dimensional conformations for the final stages of the distance geometry refinement. A force field energy minimisation may also be used. When the conformation has been refined, the next structure is generated, starting with the assignment of random distances.

Many enhancements have been made to the basic distance geometry method. Some of the most useful enhancements result from the incorporation of chemical information. For example, if the lower bound for the 1,4 distances is set to a value equivalent to a torsion angle of 60° rather than one of 0° then eclipsed conformations can be avoided. Similarly, amide bonds can be forced to adopt a nearly planar structure by an appropriate choice of distance bounds and chiral constraints.

9.5.1 The Use of Distance Geometry in NMR

One of the most important uses of distance geometry is for deriving conformations that are consistent with experimental distance information, especially distances obtained from NMR experiments. The NMR spectroscopist has at his or her disposal a range of experiments that can provide a wealth of information about the conformation of a molecule. Two of the most commonly used NMR experiments that provide such conformationally dependent information are the 2D-NOESY (nuclear Overhauser enhancement spectroscopy) and the 2D-COSY (correlated spectroscopy) experiments [Derome 1987]. NOESY provides information about the distances between atoms which are close together in space but may be separated by many bonds. The strength of the NOESY signal is inversely proportional to the sixth power of the distance and so by analysing the nuclear Overhauser spectrum it is possible

to calculate approximate values for the distance between relevant pairs of atoms. COSY experiments are often used to provide information about atoms which are covalently separated by three bonds (i.e. torsion angles). Both types of experiment provide information about interatomic distances, and so distance geometry is a natural technique to use for generating conformations that are consistent with the experimental data. Distance geometry has been particularly useful for solving the structures of proteins and nucleic acids, where the amount of data is so large that it is impossible to perform the task manually. The distance information provided by NMR experiments does of course supplement the geometrical constraints on the interatomic distances that are derived from the internal geometry (i.e. bond lengths and angles).

Distance geometry is, at heart, a random technique. It is therefore usual to generate more than one conformation in order to try to explore the conformational space that is consistent with the experimentally derived distances. The resulting set of structures is often displayed as a superimposed set; this enables the similarities and differences between the structures to be easily identified. For example, in Figure 9.18 (colour plate section) we show an ensemble of conformations of RANTES, a small protein called a chemokine that is implicated in inflammation [Chung *et al.* 1995]. It is often found that some parts of the molecule adopt very similar conformations in all the structures whereas other regions show considerable variation. This is often interpreted as an indication of conformational flexibility, but it is important to remember that it may also indicate a lack of experimental data for those atoms.

9.6 Exploring Conformational Space Using Simulation Methods

The Monte Carlo and molecular dynamics simulation methods can be used to explore the conformational space of molecules. During such a simulation the system is able to

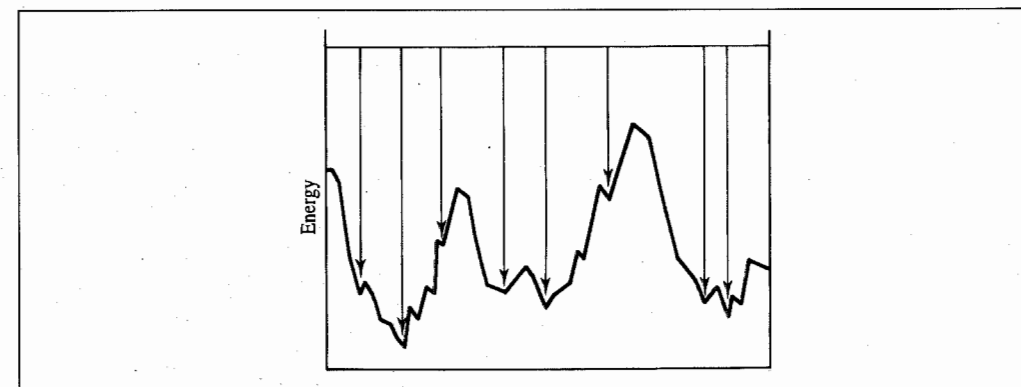


Fig. 9.19: Schematic illustration of an energy surface. A high-temperature molecular dynamics simulation may be able to overcome very high energy barriers and so explore conformational space. On minimisation, the appropriate minimum energy conformation is obtained (arrows).

overcome energy barriers and so explore different regions of the conformational space. A distinction must be made between a 'true' Monte Carlo simulation (Chapter 8) and the minimisation-based random search methods discussed in Section 9.4. A true Monte Carlo simulation does not include any energy minimisation, and each randomly generated conformation would be accepted or rejected using the Metropolis criterion. There are difficulties in applying the Monte Carlo technique in simulations of flexible molecules, as we have discussed in Chapter 8.

Molecular dynamics is widely used for exploring conformational space. A common strategy is to perform the simulation at a very high, physically unrealistic temperature. The additional kinetic energy enhances the ability of the system to explore the energy surface and can prevent the molecule getting stuck in a localised region of conformational space. This is schematically illustrated in Figure 9.19. Structures are then selected at regular intervals from the trajectory for subsequent energy minimisation.

9.7 Which Conformational Search Method Should I Use? A Comparison of Different Approaches

With such an array of methods for exploring conformational space, it can be difficult to decide which to choose. Each method has its own strengths and weaknesses. Systematic searches are subject to the effects of combinatorial explosion, and they are not naturally suited to molecules with rings. However, they do have a definite endpoint; when the search has finished, one can be guaranteed to have found all conformations for a given dihedral increment. Random search methods can require long runs to ensure that the conformational space has been covered, and they can generate the same structure many times. Distance geometry is particularly useful when experimental information can be incorporated, as is restrained molecular dynamics.

A comparison of various methods for searching conformational space has been performed for cycloheptadecane ($C_{17}H_{34}$) [Saunders *et al.* 1990]. The methods compared were the systematic search, random search (both Cartesian and torsional), distance geometry and molecular dynamics. The number of unique minimum energy conformations found with each method within 3 kcal/mol of the global minimum after 30 days of computer processing were determined (the study was performed in 1990 on what would now be considered a very slow computer). The results are shown in Table 9.1.

Method	Total unique conformers found after 30 days processing
Systematic search	211
Random Cartesian search	222
Random dihedral search	249
Distance geometry	176
Molecular dynamics	169

Table 9.1: A comparison of five different conformational searching algorithms. (Data from [Saunders *et al.* 1990].)

Combining the results from all the different methods revealed a grand total of 262 conformations within 3 kcal/mol of the global energy minimum. No one method found all of them but the random dihedral search did give the best performance in this case. The global energy minimum would be expected to constitute only about 8% of the total population of conformational states for this molecule if it is assumed that the entropies of all the conformations are the same. Largely as a result of this study cycloheptadecane is now often considered as the proving ground for new conformational search methods (despite the fact that it is not a very good representative for a typical 'organic' molecule).

9.8 Variations on the Standard Methods

Each year usually sees the publication of a handful of new methods for exploring conformational space. Many can be considered as variants on one of the approaches discussed thus far but which may provide some advantage in terms of the efficiency and effectiveness with which they explore conformational space. Some of these alternative methods are designed for quite specific types of molecule (such as ring systems) and as may not be particularly 'general' approaches to conformational analysis. Here we describe in more detail two of these newer methods, one which extends the systematic search and one which uses an alternative approach to generating the initial structure prior to energy minimisation.

9.8.1 The Systematic Unbounded Multiple Minimum Method

One of the possible limitations of a regular systematic conformational search is that successive conformations are often very similar to each other. In a typical depth-first search strategy successive conformations often differ by just one or two torsion angles. An additional potential limitation is that the torsional increment is usually specified at the start of the search, so that if one wants to perform a search at a higher resolution one typically has to rediscover all the conformations generated at the lower resolution. The systematic unbounded multiple minimum (SUMM) method [Goodman and Still 1991] is designed to address these two concerns. SUMM generates conformations by first selecting a structure from those generated previously (using the uniform usage protocol) and changing one or more of its torsion angles. This gives a new structure, which is then energy minimised, checked to determine whether it has already been generated, and if not it is added to the list of conformational minima. Central to the approach is that the changes in torsion angles are determined in a preordained, systematic manner. This is achieved by setting up a sequence of torsional modifications such that the first components in the sequence correspond to changes of a single torsion angle through 120° . Later components correspond to changes of two torsion angles through 120° , and so on. For a given number of torsion angle changes and a given torsional increment the actual changes are mixed up so that successive components represent modifications to very different parts of the molecule. For example, the first two modifications in a normal systematic search of hexane could correspond to changes of 0° , 0° , 120° and 0° , 0° , 240° for torsions τ_1 , τ_2 and

τ_3 , respectively, whereas in SUMM the second modification might correspond to a rather different sequence (e.g. 0° , 240° , 0°). If one were to run the search to completion then all possible torsional variations would be considered both with and without this mixing protocol. However, if the search was terminated after a fixed number of steps then the mixing protocol increases the chances of achieving greater coverage of conformational space. The algorithm maintains a record of the torsional modifications that have been made to each distinct minimum energy conformation so, when that structure is next selected to act as the starting point for a conformational change, it knows which torsional changes should be performed.

The SUMM approach can be applied to both cyclic and acyclic molecules, though for cyclic systems it is still necessary to check for ring-closure violations. When rings are present then SUMM can use a preoptimisation procedure to reduce the length of the often abnormally long ring-closure bond. If a structure with such a long bond is subjected to normal energy minimisation then the bond length will be rapidly corrected but this often leads to significant distortions to the rest of the molecule, with torsion angles changing significantly from their initial values. This could be construed as undermining the rationale of a systematic search. The preoptimisation procedure makes small sequential changes to those torsion angles that affect each of the ring-closure bonds in an iterative fashion, in order to gradually bring the ring-closure bond(s) closer to their ideal value. SUMM is considered to be particularly efficient for locating all low energy conformations of a molecule. A random search method spends more and more time generating structures that have already been identified earlier in the search, in contrast to systematic methods such as SUMM.

9.8.2 Low-mode Search

The low-mode search method [Kolossvary and Guida 1996] is closely related to the methods described in Section 5.9 for locating saddle points on energy surfaces. As we discussed there, one way to locate a transition state is to follow the 'path of shallowest ascent' from a minimum. In favourable cases, this path will largely correspond to one of the low-frequency normal modes of vibration. By continuing to move along the path one might expect to locate the second minimum that is connected via the saddle point with the starting structure. Locating saddle points can be a very difficult and time-consuming process and so some modifications are required to make such an approach practical for conformational searching. In the low-mode search an initial minimum energy conformation is subjected to normal mode analysis. Those low-frequency modes which are below a user-specified frequency threshold (e.g. 250 cm^{-1}) are identified and searched by changing the atomic coordinates in a manner given by the relevant eigenvector. This perturbation of the initial structure is performed in discrete steps until either the energy of the structure increases beyond a specified threshold or until the energy first increases but then starts to fall. This latter case may correspond to a movement over the saddle point and into the locality of a nearby minimum. In such cases (which are relatively rare), the structure is then fully energy minimised to give a new minimum energy conformation.

A particular advantage of the low-mode search is that it can be applied to both cyclic and acyclic molecules without any need for special ring closure treatments. As the low-mode search proceeds a series of conformations is generated which themselves can act as starting points for normal mode analysis and deformation. In a sense, the approach is a systematic one, bounded by the number of low-frequency modes that are selected. An extension of the technique involves searching random mixtures of the low-frequency eigenvectors using a Monte Carlo procedure.

9.9 Finding the Global Energy Minimum: Evolutionary Algorithms and Simulated Annealing

Evolutionary algorithms and simulated annealing are two methods that have found widespread use in molecular modelling. Their use is by no means restricted to the problem of finding the global minimum energy conformation of a molecule, but they have been applied to problems as diverse as protein-ligand docking, molecular design, QSAR and pharmacophore mapping [Clark and Westhead 1996; Jones 1998; Judson 1997]. We will consider some of these alternative applications in Chapter 12. Nevertheless, conformational analysis is a very good problem with which to introduce and describe these two methods.

9.9.1 Genetic and Evolutionary Algorithms

Evolutionary algorithms (EA) are a group of methods based on ideas of biological evolution that are designed to find optimal solutions to problems. There are currently three basic classes of evolutionary algorithm: genetic algorithms (GA), evolutionary programming (EP) and evolution strategies (ES). There are many similarities between these three classes but some key differences. Common to all three is the idea of creating a 'population' of possible solutions to the problem. The members of the population are scored using a 'fitness function' that measures how 'good' they are. The population changes over time and (hopefully) evolves towards better solutions. This process of generating new solutions is often referred to as 'breeding', with the new solutions being the 'children' that are generated from the 'parents' of the previous 'generation'. We will give an outline of these methods using the problem of finding the global minimum energy conformation as an example.

Probably the best-known of the three classes is the genetic algorithm [Goldberg 1989]. The following is a description of the basic method (or *canonical* genetic algorithm). The first step is to create a population of μ possible solutions. In conformational analysis, this initial population would correspond to a set of randomly generated conformations of the molecule. Each member of the population is coded by a 'chromosome'. This is usually stored as a linear string of bits (i.e. 0s and 1s). The chromosome codes for the values of the torsion angles of the rotatable bonds in the molecule, as illustrated in Figure 9.20. The initial population is most easily obtained by randomly setting bits to 0 or 1 in the chromosomes. After decoding each chromosome and assigning the torsion angles to the

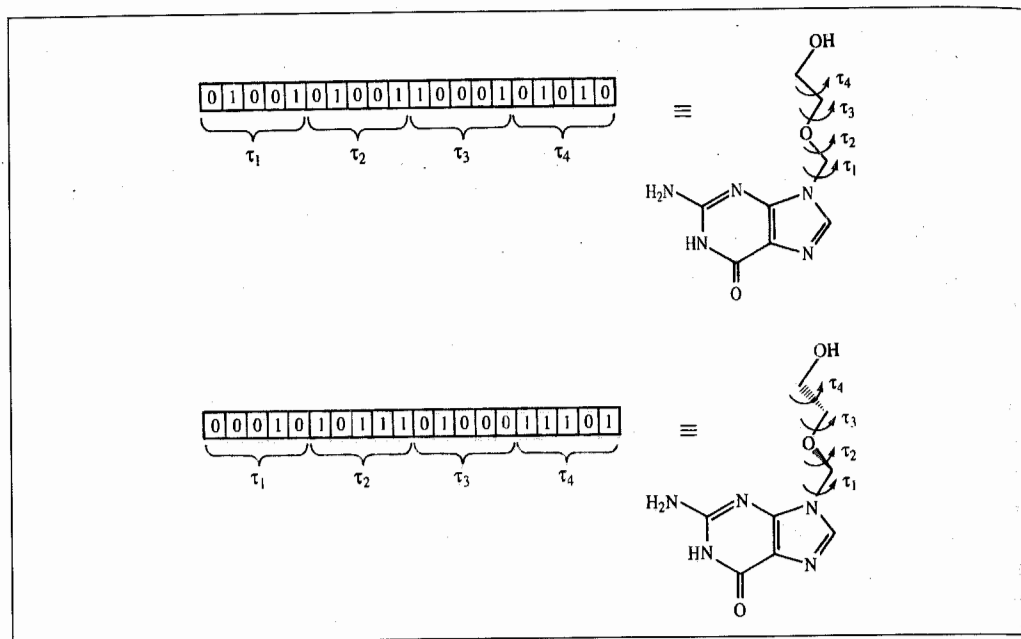


Fig. 9.20: The chromosome in a genetic algorithm codes for the torsion angles of the rotatable bonds.

appropriate values in the molecule, the fitness of each member of the population can be calculated. In conformational analysis, an appropriate fitness function would be the internal energy, as might be calculated using molecular mechanics. A new population is then generated. In the canonical genetic algorithm, $\mu/2$ pairs of parents are selected from the current population. These pairs are chosen at random, but with a bias towards the most fit individuals. A technique called *roulette wheel selection* is often used to achieve this bias by using slot sizes in the roulette wheel that are proportional to the values of the fitness function. A simple example is shown in Figure 9.21. The use of roulette wheel selection means that particularly fit members of the population may be able to produce many offspring. The new population is then subjected to genetic operators, the two most commonly used of which are *crossover* (or *recombination*) and *mutation*. In crossover, a cross

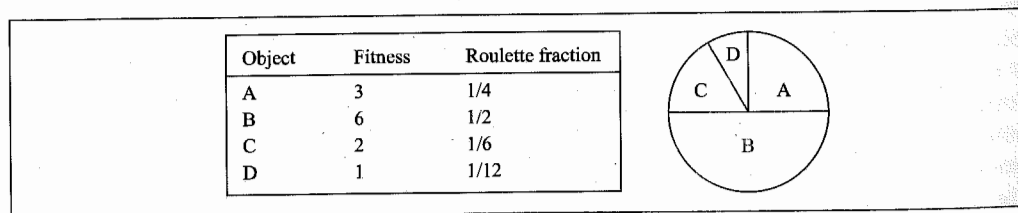


Fig. 9.21: The basis of roulette wheel selection showing how the more fit members of the population are selected in proportion to their fitness values.

position i is randomly selected ($1 \leq i \leq l-1$, where l is the length of the chromosome). Two new strings are then created by swapping the bits between positions $i+1$ and l . For example, suppose we have the following two chromosomes:

```
00100011110001
11000011001100
```

and the crossover point is chosen to be 6. Then the two new strings are:

```
00100011001100
11000011110001
```

The crossover operator is applied to the selected pairs of parents with a probability P_c , a typical value being 0.8 (i.e. there is an 80% chance that any of the $\mu/2$ pairs will actually undergo this type of recombination). Following the crossover phase mutation is applied to all individuals in the population. Here, each bit may be inverted (0 to 1 and vice versa) with a probability P_m . The mutation operator is usually assigned a low probability (e.g. 0.01).

This completes one complete cycle of the genetic algorithm. The new population then becomes the current population ready for a new cycle. The algorithm repeatedly applies this sequence for a predetermined number of iterations and/or until it converges.

Many variants on the canonical genetic algorithm have been suggested. For example, it is common practice to carry forward the highest-ranking individuals unchanged: this is often referred to as an 'elitist' strategy and it ensures that the best individuals are not lost. In a *steady-state* genetic algorithm each iteration involves just one operator (mutation or crossover) with the resulting one or two individuals replacing the worst members of the population. A variety of different crossover schemes have been suggested, such as the two-point crossover. Real-valued chromosomes are an alternative to the binary representation of the canonical genetic algorithm. A real-valued chromosome consists of a string of real numbers, which correspond to the parameters of the problem (e.g. the torsion angles in the molecule). The main difference is that mutation is typically effected by adding a random increment (often chosen from a Gaussian distribution) to a randomly chosen parameter.

One of the main issues when implementing a genetic algorithm is the need to prevent premature convergence. The 'selection pressure' can be an important factor in determining whether such premature convergence occurs (or, conversely, that the program takes too long to converge to a solution). The selection pressure is defined as the relative probability that the fittest individual in a population will be chosen as a parent relative to an individual with average fitness. The selection pressure can be controlled by rescaling the fitness values when applying roulette wheel selection. Another way to deal with premature convergence is to use the *island model* in which one maintains a number of separate sub-populations within the whole and to introduce another operator that corresponds to the movement of an individual from one island to another. *Niching* is a related technique for achieving the same goal; here one tries to force the individuals in the population away from the most heavily populated regions of the space. This can be achieved by calculating a 'distance' between all pairs of members in the population, which is then used to reduce the chances of selecting pairs of individuals that are very similar.

The main difference between the genetic algorithm and evolutionary programming is that the latter does not use a crossover operator. Rather, new individuals are generated from their parents using mutation alone. In addition, individuals in evolutionary programming are typically represented using a sequence of real numbers, rather than a binary representation (though it should be noted that chromosomes containing integers and real numbers have been successfully incorporated into genetic algorithms). At the start of each iteration of the evolutionary programming method one child is bred from each of the members of the current population using a mutation operator. During mutation each of the real variables in the chromosome is modified by adding a randomly generated real number, usually taken from a Gaussian distribution. These μ children are scored using the fitness function. The 2μ parents and children then compete for survival into the next generation. This is achieved by performing a series of 'tournaments'. In the tournament stage, each individual is compared with a number M of opponents selected at random from the 2μ population of parents and children. The individuals are then ranked according to the number of 'wins' they achieve and the appropriate number selected from the top of the set to give the next population. A win is recorded for each opponent with a worse fitness score. As the number of opponents, M , increases so the selection pressure increases, and it is necessary to choose an appropriate value for M to find the appropriate compromise between premature convergence and taking too long to find a solution.

Evolutionary strategies are very similar to evolutionary programming but differ in two key respects. First, crossover operators are permitted and, second, the probabilistic tournament is replaced with a straightforward ranking. At each iteration, λ children are generated using crossover and mutation from the current population. Typically, λ would be about seven times larger than μ . The children are scored and then the set of $(\mu + \lambda)$ parents and offspring are ranked according to fitness, with the top μ individuals being selected to form the next generation. A slight alternative to this approach is to select the μ individuals for the next population from just the λ new offspring. This is referred to as (μ, λ) selection.

Genetic and evolutionary algorithms are primarily intended for performing global optimisation. However, they do involve a significant random element and so they are not guaranteed to produce the same solution (e.g. the global minimum energy conformation) from each run, except for rather simple problems. What they are particularly useful for is producing solutions very close to the global optimum in a reasonable amount of time. An additional advantage of genetic and evolutionary algorithms is that as they maintain a population of possible solutions one may obtain several 'reasonable' solutions from a single run. Nevertheless, it is common practice to perform several runs in order to obtain a variety of different solutions and to investigate the nature of the energy surface.

Judson and co-workers were among the first to investigate the use of genetic algorithms in conformational analysis [Judson *et al.* 1993; McGarrah and Judson 1993]. Their implementation was tested on a variety of types of molecule, including a cyclic hexapeptide and a selection of more 'drug-like' molecules extracted from the Cambridge Structural Database. A key conclusion from these studies was the need to avoid premature convergence and to maintain a diverse population. When compared to a straightforward systematic search procedure the genetic algorithm was found to be particularly effective for the more flexible

molecules, especially those molecules with more than eight rotatable bonds [Meza *et al.* 1996].

9.9.2 Simulated Annealing

Annealing is the process in which the temperature of a molten substance is slowly reduced until the material crystallises to give a large single crystal. It is a technique that is widely used in many areas of manufacture, such as the production of silicon crystals for computer chips. A key feature of annealing is the use of very careful temperature control at the liquid–solid phase transition. The perfect crystal that is eventually obtained corresponds to the global minimum of the free energy. Simulated annealing is a computational method that mimics this process in order to find the 'optimal' or 'best' solutions to problems which have a large number of possible solutions [Kirkpatrick *et al.* 1983].

In simulated annealing, a cost function takes the role of the free energy in physical annealing and a control parameter corresponds to the temperature. To use simulated annealing in conformational analysis the cost function would be the internal energy. At a given temperature the system is allowed to reach 'thermal equilibrium' using a molecular dynamics or Monte Carlo simulation. At high temperatures, the system is able to occupy high-energy regions of conformational space and to pass over high energy barriers. As the temperature falls, the lower energy states become more probable in accordance with the Boltzmann distribution. At absolute zero, the system should occupy the lowest-energy state (i.e. the global minimum energy conformation). To guarantee that the globally optimal solution is actually reached would require an infinite number of temperature steps, at each of which the system would have to come to thermal equilibrium. Careful temperature control is required when the energy of the system is comparable with the height of the barriers that separate one region of conformational space from another. This is often difficult to achieve in practice and so simulated annealing cannot *guarantee* to find the global minimum, much as a genetic algorithm cannot guarantee to identify the globally optimal solution. However, if the same answer is obtained from several different runs then there is a high probability that it corresponds to the true global minimum. Several simulated annealing runs may enable a series of low-energy conformations of a molecule to be obtained.

9.10 Solving Protein Structures Using Restrained Molecular Dynamics and Simulated Annealing

A particularly important application of molecular dynamics, often in conjunction with the simulated annealing method, is in the refinement of X-ray and NMR data to determine the three-dimensional structures of large biological molecules such as proteins. The aim of such refinement is to determine the conformation (or conformations) that best explain the experimental data. A modified form of molecular dynamics called *restrained molecular dynamics* is usually used in which additional terms, called *penalty functions*, are added to the potential energy function. These extra terms have the effect of penalising conformations

that do not agree with the experimental data. Molecular dynamics is used to explore the conformational space in order to find a conformation (or conformations) that not only has a low intrinsic energy but is also consistent with the experimental data. Simulated annealing can often be a convenient way to ensure that the conformational space is explored effectively.

9.10.1 X-ray Crystallographic Refinement

X-ray crystallography is a powerful technique for elucidating the structures of molecules. An X-ray diffraction pattern arises because of constructive and destructive interference between X-rays scattered from different parts of the crystal. An X-ray beam scattered by an electron at a point r travels a different distance to the detector than a beam scattered by an electron at the origin (Figure 9.22). As a consequence, the two scattered X-ray beams will have different phases and will interfere. As the detector is moved through different scattering angles, θ (Figure 9.22), the intensity of the scattered radiation will fluctuate between zero (destructive interference) and twice the amplitude of the original beam (constructive interference). In a real sample the amplitude of the scattered radiation from a point is proportional to the electron density at that point. The total signal reaching the detector is obtained by integrating the electron density over the whole crystal and is expressed as the structure factor, F . The structure factor is a complex number that can be written $F = |F|e^{i\phi}$, where $|F|$ is the amplitude and $e^{i\phi}$ is the phase. If the electron distribution is known (i.e. if we know the three-dimensional structure) it is possible to determine the structure factor for all scattering angles, and so we can calculate the X-ray diffraction pattern. The X-ray crystallographer is faced with the reverse problem: to determine the electron distribution (and thereby the three-dimensional structure) from the diffraction pattern. The difficulty is that it is only possible to measure the intensities of the spots (which are equal to the amplitudes $|F|^2$), but not the phases; this is the famous *phase problem*, which is one of the major obstacles in solving an X-ray structure.

To obtain the electron density distribution it is necessary to guess, calculate or indirectly estimate the phases. Various methods have been developed to tackle the phase problem. For proteins the most common strategy is multiple isomorphous replacement in which the protein crystals are soaked in solutions containing salts of heavy metals such as mercury,

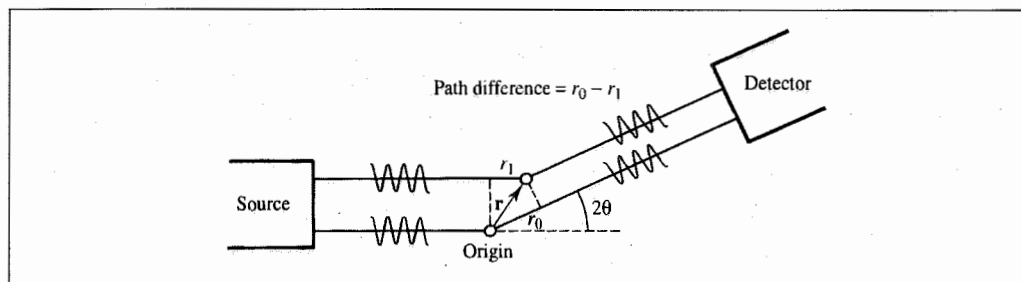


Fig. 9.22: Schematic illustration of an X-ray scattering experiment. The X-ray beam travels a different distance when scattered by an electron at the origin compared to an electron situated at r .

platinum or silver. These heavy atoms may bind to specific parts of the protein (e.g. mercury ions may react with exposed SH groups). By comparing the diffraction patterns of the native crystals and the crystals of the heavy-atom derivatives it can be possible to estimate some of the phases (under the assumption that no structural change occurs). Once some of the phases are known, others can be determined, and eventually an initial electron density map can be obtained. The electron density map is often represented as a three-dimensional surface by contouring at a constant value (Figure 9.23, colour plate section). An initial model of the molecule is then fitted to the electron density. When the diffraction experiment is performed at high resolution then the locations of individual atoms are often easy to identify. However, at lower resolution it can be difficult to find the optimal fit of the atoms in the model to the electron density as individual features are not so well defined. This is often the case in proteins.

The objective of the refinement is to obtain a structure that gives the best possible agreement with the experimental data. This is done by gradually changing the structure to give better and better agreement between the calculated and observed structure factor amplitudes. This degree of agreement is quantified by the value of the crystallographic R factor, which is defined as the difference between the observed ($|F_{\text{obs}}|$) and calculated ($|F_{\text{calc}}|$) structure factor amplitudes:

$$R = \frac{\sum (|F_{\text{obs}}| - |F_{\text{calc}}|)}{\sum |F_{\text{obs}}|} \quad (9.23)$$

Traditionally, least-squares methods have been used to refine protein crystal structures. In this method, a set of simultaneous equations is set up whose solutions correspond to a minimum of the R factor with respect to each of the atomic coordinates. Least-squares refinement requires an $N \times N$ matrix to be inverted, where N is the number of parameters. It is usually necessary to examine an evolving model visually every few cycles of the refinement to check that the structure looks reasonable. During visual examination it may be necessary to alter a model to give a better fit to the electron density and prevent the refinement falling into an incorrect local minimum. X-ray refinement is time consuming, requires substantial human involvement and is a skill which usually takes several years to acquire.

Jack and Levitt introduced molecular modelling techniques into the refinement in the form of an energy minimisation step (using a force field function) that was performed alternately with the least-squares refinement [Jack and Levitt 1978]. This approach was shown to give convergence to better structures. More recently, restrained molecular dynamics methods were introduced by Brunger, Kuriyan and Karplus [Brunger *et al.* 1987]. These methods have had a dramatic impact on the refinement of X-ray and NMR structure of proteins.

In the restrained molecular dynamics approach the total 'potential energy' is written as the sum of the usual potential energy and the penalty term, as usual:

$$E_{\text{tot}} = \mathcal{V}(\mathbf{r}^N) + E_{\text{sf}} \quad (9.24)$$

The additional penalty function that is added to the empirical potential energy function in restrained dynamics X-ray refinement has the form:

$$E_{\text{sf}} = S \sum (|F_{\text{obs}}| - |F_{\text{calc}}|)^2 \quad (9.25)$$

where E_{sf} describes the differences between the observed structure factor amplitudes and those calculated from the atomic model. S is a scale factor which is chosen so that the gradient of E_{sf} is comparable to the gradient of the potential energy part of the function. The conformational space is explored using molecular dynamics with simulated annealing; very high temperatures are used in the initial stages to permit the system to range widely over the energy surface. The temperature is then gradually reduced as the structure settles into a conformation which not only has a low energy but also a low R factor.

9.10.2 Molecular Dynamics Refinement of NMR Data

We have already discussed in Section 9.5.1 the type of information that NMR experiments can provide about the conformation of a molecule and the use of distance geometry for determining structures that are consistent with the experimental data. In the simplest molecular dynamics approach, we could incorporate harmonic restraint terms of the form $k(d - d_0)^2$ where d is the distance between the atoms in the current conformation and d_0 is the desired distance dynamics approach derived from the NMR spectrum. k is a force constant, the value of which determines how tightly the restraint should be applied. The information provided by the COSY experiment can also be expressed as a torsion angle via the Karplus equation; torsional restraints may be incorporated into the molecular dynamics energy function as an alternative to the use of distances. There are many other ways in which the restraints can be incorporated; for example, some practitioners prefer to penalise a structure only if the distance exceeds the target:

$$v(d) = k(d - d_0)^2 \quad d > d_0 \quad (9.26)$$

$$v(d) = 0 \quad d \leq d_0 \quad (9.27)$$

The atoms are prevented from coming too close by the van der Waals terms in the force field. More sophisticated functional forms have also been used which try to take into account the imprecise nature of the experimental values. A simple Hooke's law relationship implies that an exact value is known for the distance, whereas there can be significant uncertainty about its value. A more appropriate functional form has the following form:

$$v(d) = k_l(d - d_l)^2 \quad d < d_l \quad (9.28)$$

$$v(d) = 0 \quad d_l \leq d \leq d_u \quad (9.29)$$

$$v(d) = k_u(d - d_u)^2 \quad d_u < d \quad (9.30)$$

This potential is shown schematically in Figure 9.24. d_l and d_u are the lower and upper distances that are considered to be consistent with the experimental data. $(d_l + d_u)/2$ is thus the assigned target distance obtained from a measurement of the NOESY intensity and the error associated with that measurement is $\pm(d_u - d_l)/2$. A distance between d_l and d_u incurs no penalty. Outside this region the restraint is applied using two harmonic potentials. These restraining potentials may have different force constants and so be of different steepness. In some functional forms, the harmonic potential is eventually replaced by a linear function, as illustrated in Figure 9.24.

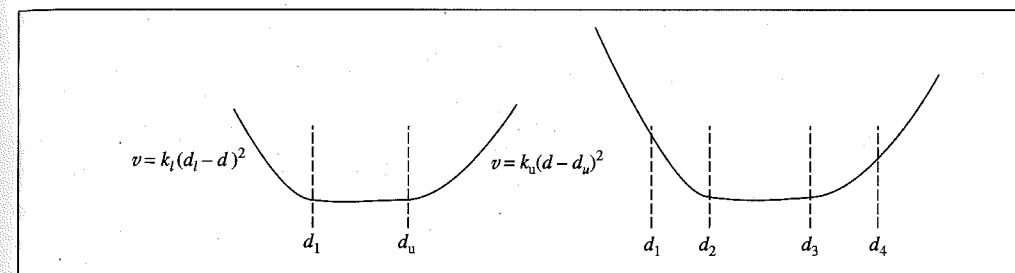


Fig. 9.24: A restraining potential that does not penalise structures in which the distance lies between the lower and upper distances d_l and d_u and uses harmonic functions outside this range (left). The harmonic potentials may also be replaced by linear restraints further from this region (right).

9.10.3 Time-averaged NMR Refinement

If the molecule interconverts between two or more conformations on a timescale that is rapid compared to the chemical shift timescale then the NMR spectrum only shows the average of the signals from the individual conformations. This behaviour is illustrated schematically in Figure 9.25, where an atom or group (such as a leucine side chain) interconverts between two energy minima within the protein. The NMR spectrum comprises a single peak that is a weighted average of the resonances from the two individual conformations. If the two conformations make distinct interactions then two sets of distance restraints can be derived, and a standard refinement procedure would attempt to satisfy both sets of restraints simultaneously. This would lead to a conformation in which the group is positioned at the top of the barrier between the two minima. This incorrect result is a consequence of assuming that one single structure is consistent with all of the experimental data, rather than recognising that the experimental data may arise from more than one conformation.

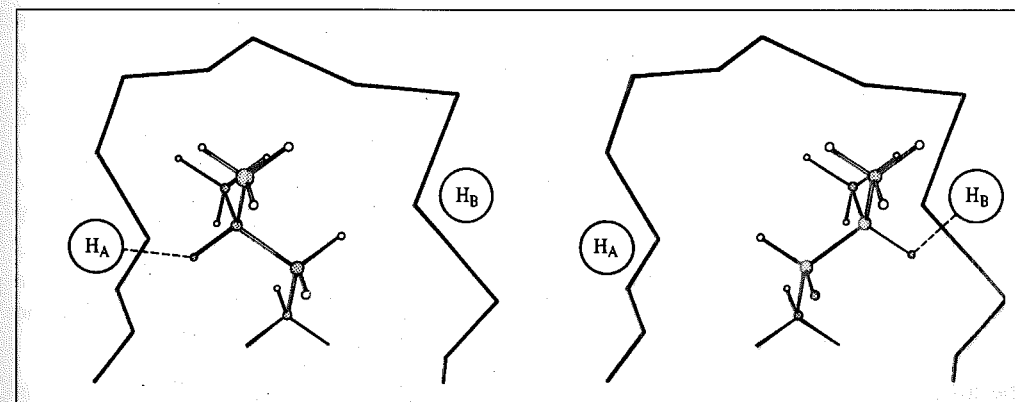


Fig. 9.25: If the leucine side chain interconverts rapidly between two conformations then the NMR spectrum will be an average of them. With a traditional refinement this leads to a structure that simultaneously tries to satisfy all restraints and is at the top of the energy barrier between the two minima.

Time-averaged restraints may be able to overcome this problem [Torda *et al.* 1990]. Rather than using the instantaneous value of a distance in the restraint function, the time-averaged restraint method uses a value that is averaged over time. The simple harmonic error function then becomes:

$$v(d) = k(\langle d(t) \rangle - d_0)^2 \quad (9.31)$$

where $\langle d(t) \rangle$ is the time-averaged value of the distance, as obtained from the molecular dynamics simulation. At a time t' , $\langle d(t') \rangle$ is given by:

$$\langle d(t') \rangle = \frac{1}{t'} \int_0^{t'} d(t) dt \quad (9.32)$$

As the intensity of the NOESY signal is proportional to the inverse sixth power of the distance, the 'distance' to use in this case is actually given by:

$$d_{\text{NOESY}} = \langle d(t')^{-6} \rangle^{-1/6} \quad (9.33)$$

The time-averaged value of the distance that should be used in the error function is thus

$$\langle d(t') \rangle = \left[\frac{1}{t'} \int_0^{t'} d(t)^{-6} dt \right]^{-1/6} \quad (9.34)$$

One way to implement time-averaged restraints is to evaluate $\langle d(t') \rangle$ using Equation (9.34) as the simulation proceeds and incorporate the value in the error function (Equation (9.31)). If the simulation is run for long enough, all the accessible conformational states should be visited and be included in the calculation of the average distances. However, it is rarely possible to achieve the length of simulation needed to ensure that the conformation space has been adequately covered. We require a method that can supply an accurate picture of the dynamics of the molecule with minimal computational effort. The normalisation factor $1/t'$ in Equation (9.34) becomes progressively larger as the simulation proceeds, thus making $\langle d(t') \rangle$ increasingly less sensitive to the current value of the distance. What we require is a means to bias the instantaneous value of $\langle d(t') \rangle$ towards the values from the most recent part of the simulation. In this way, if the 'current' value of $\langle d(t') \rangle$ is incompatible with the restraint, then the penalty function should be proportionately increased. This can be achieved using an exponential 'memory function' which has the effect of weighting the recent history more heavily. Various memory functions are possible; one functional form is:

$$\langle d(t') \rangle = \left(\frac{\int_0^{t'} e^{-(t-t')/\tau} d(t)^{-6} dt}{\int_0^{t'} e^{-(t-t')/\tau} dt} \right)^{-1/6} \quad (9.35)$$

where τ is the time constant for the exponential damping factor. Small values of τ give a higher weighting to recent values of the distance. If τ is infinite, all the past history of the simulation is given equal weight.

The time-averaged restraint method is quite complicated to implement, and some skill is required when choosing the most appropriate functional form and the damping constant. The data produced by the simulation must also be interpreted with care. The technique is only truly applicable where the conformations are relatively close together, so that

interconversion between the different conformations can be achieved relatively easily. Nevertheless, the technique does provide a more accurate representation of the dynamics of the real system, and it does enable the conformation to fluctuate more. One drawback of any restraint method is that the additional penalty terms represent an unnatural perturbation of the forces within the molecule. When using 'static' restraints the size of the force constants for the restraint terms can be quite large, which can often cause the conformations to have rather high energies. Smaller force constants can often be used with time-averaged restraints, which means that the conformations generally have lower energies.

9.11 Structural Databases

Experimental information about the structures of molecules can often be extremely useful for forming theories of conformational analysis and helping to predict the structures of molecules for which no experimental information is available. The most important technique currently available for determining the three-dimensional structure of molecules is X-ray crystallography. The international crystallographic community has established centres where crystallographic data is collected and then distributed in electronic form. Two particularly important databases for the molecular modeller are the Cambridge Structural Database (CSD) [Allen *et al.* 1979], which contains crystal structures of organic and organometallic molecules; and the Protein Databank (PDB) [Bernstein *et al.* 1977; Berman *et al.* 2000], which contains structures of proteins and some DNA fragments. Other databases are also available, such as the Inorganic Structural Database of inorganic compounds and complexes [Bergerhoff *et al.* 1983].

A database is of little use without software tools to search, extract and manipulate the data. A simple use of a database is for extracting information about a particular molecule or group of molecules. For example, one may wish to retrieve the crystal structure of ranitidine (Figure 9.26). The molecule(s) may be specified in a variety of ways, such as by name, molecular formula or literature citation. The data may also be identified by creating a two-dimensional representation of the molecule (as in Figure 9.26) and using a substructure search program (see Section 12.2) to search the database. In fact, the CSD contains two entries for ranitidine: one is the crystal structure of the hydrochloride salt and the other is the structure of the oxalate salt. Crystallographic databases have also been used to develop an understanding of the factors that influence the conformations of molecules, and of the ways in which molecules interact with each other. For example, the CSD has been comprehensively analysed to characterise how the lengths of chemical

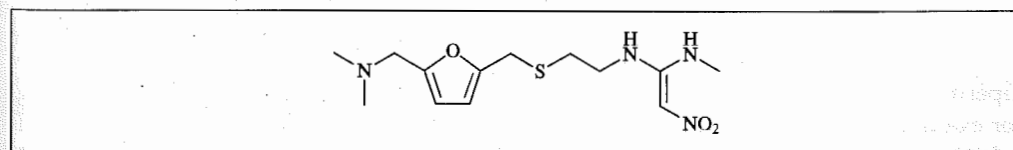


Fig. 9.26: Ranitidine.

bonds depend upon the atomic numbers, hybridisations and environment of the atoms involved [Allen *et al.* 1987]. A major use of the CSD is searching for molecules which contain a particular fragment, in order to investigate the conformation(s) that the fragment adopts. Intermolecular interactions can also be investigated. For example, analyses of intermolecular hydrogen bonding have revealed distinct distance and angular preferences [Murray-Rust and Glusker 1984; Glusker 1995]. This type of analysis can now be applied to a wide range of functional groups and molecular fragments, as illustrated in Figure 9.27 (colour plate section), which shows the distribution of OH groups around thiazole rings [Bruno *et al.* 1997]. This shows that the nitrogen in thiazole is a much stronger hydrogen bond acceptor than the sulphur atom.

The protein database has provided much useful information about the structures that proteins adopt and the PDB has been extensively analysed to try to understand the principles that determine why a given amino acid sequence folds into one specific conformation. We shall discuss the use of molecular modelling methods to predict protein structure in more detail in Chapter 10. Here we shall just mention one interesting way in which the information contained in the protein database has been put to a practical purpose. One of the steps in determining the structure of a protein by X-ray crystallography involves fitting the polypeptide chain to the electron density. This can be a complex and time-consuming task, even with today's sophisticated molecular graphics. Computer programs have been developed which extract the conformations of short polypeptide fragments (up to four amino acids long) from known X-ray structures [Jones and Thirup 1986]. These fragments are then used to generate a chain that fits the electron density. This method is feasible because a given segment of polypeptide chain often adopts a limited selection of conformations in protein structures, and a significant proportion of an unknown protein structure can often be constructed using such 'spare parts' taken from other proteins.

It should be remembered that a crystallographic database can only provide information about the crystalline state of matter, and that the possible influence of crystal packing forces should always be taken into account. This is less of a concern for proteins than for 'small' molecules as protein crystals contain a large amount of water and indeed NMR studies have established that proteins have approximately the same structure in solution as in the crystal. A second, more subtle, bias is that crystallographic databases contain only molecules that can be crystallised and indeed only those molecules whose X-ray structures were considered important enough to be published. The structures in a crystallographic database may therefore not necessarily be a wholly representative set.

9.12 Molecular Fitting

Fitting is the procedure whereby two or more conformations of the same or different molecules are oriented in space so that particular atoms or functional groups are optimally superimposed upon each other. Fitting methods are widely used in molecular modelling. For example, fitting is an integral part of many conformational search algorithms, particularly those that require each conformation to be compared with those generated previously in order to check for duplicates.

A molecular fitting algorithm requires a numerical measure of the 'difference' between two structures when they are positioned in space. The objective of the fitting procedure is to find the relative orientations of the molecules in which this function is minimised. The most common measure of the fit between two structures is the root mean square distance between pairs of atoms, or RMSD:

$$\text{RMSD} = \sqrt{\frac{\sum_{i=1}^{N_{\text{atoms}}} d_i^2}{N_{\text{atoms}}}} \quad (9.36)$$

where N_{atoms} is the number of atoms over which the RMSD is measured and d_i is the distance between the coordinates of atom i in the two structures, when they are overlaid.

When fitting two structures, the aim is to find the relative orientations of the two molecules in which the RMSD is a minimum. Many methods have been devised to perform this seemingly innocuous calculation. Some algorithms, such as that described by Ferro and Hermans [Ferro and Hermans 1977] use an iterative procedure in which the one molecule is moved relative to the other, gradually reducing the RMSD. Other methods locate the best fit directly, such as Kabsch's algorithm [Kabsch 1978].

If the molecules are flexible then a better fit might be achieved if one or both of the molecules can change their conformation (for example, by rotating about single bonds). This is often referred to as *flexible fitting* or *template forcing*. In its simplest form flexible fitting is achieved by minimising the RMSD using a special minimisation algorithm that permits rotation about single bonds as well as translation and rotation in space. An alternative approach is to use restrained molecular dynamics, which may enable a more thorough exploration of the conformational space in order to find the best fit. Here, restraints are placed on the distance between pairs of matched atoms, which are incorporated into the energy function as additional penalty terms.

9.13 Clustering Algorithms and Pattern Recognition Techniques

Molecular modelling programs can generate a large amount of data, which must often be processed and analysed. Many of the conformational search algorithms that we have considered can generate conformations that are very similar, if not identical. Under such circumstances it is desirable to be able to select from the data set a smaller, 'representative' set of conformations for subsequent analysis. This can be done using cluster analysis, which groups together 'similar' objects, from which the representatives can be extracted (Figure 9.28).

There is no 'correct' method of performing cluster analysis and a large number of algorithms have been devised from which one must choose the most appropriate approach. There can also be a wide variation in the efficiency of the various cluster algorithms, which may be an important consideration if the data set is large.

A cluster analysis requires a measure of the 'similarity' (or dissimilarity) between pairs of objects. When comparing conformations, the RMSD would be an obvious measure to use.

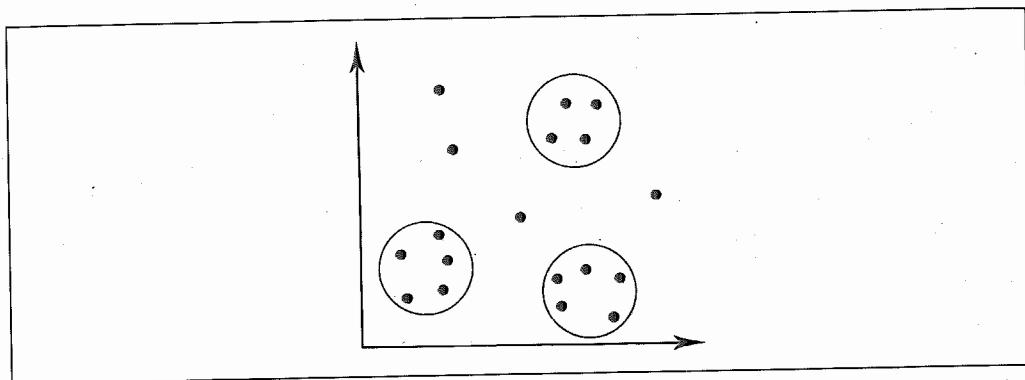


Fig. 9.28: The aim of cluster analysis is to group together 'similar' objects.

Alternatively, the 'distance' between two conformations can be measured in terms of their torsion angles. Here, there may be more than one way in which the 'distance' can be calculated. The Euclidean distance between two conformations would be calculated using:

$$d_{ij} = \sqrt{\sum_{m=1}^{N_{\text{tor}}} (\omega_{m,i} - \omega_{m,j})^2} \quad (9.37)$$

where $\omega_{m,i}$ is the value of torsion angle m in conformation i . N_{tor} is the total number of torsion angles. An alternative is the Hamming distance (also known as the Manhattan or city-block distance, Figure 9.29):

$$d_{ij} = \sum_{m=1}^{N_{\text{tor}}} |\omega_{m,i} - \omega_{m,j}| \quad (9.38)$$

When using torsion angles to calculate 'distances' between conformations it is important to remember that a torsion angle is a cyclic measure and that the difference should be measured along the shortest path, in either a clockwise or an anticlockwise direction. The clusters produced using the RMSD and the torsion angle measures may be very different. This is due to a 'leverage' effect when using torsion angles, which arises because small

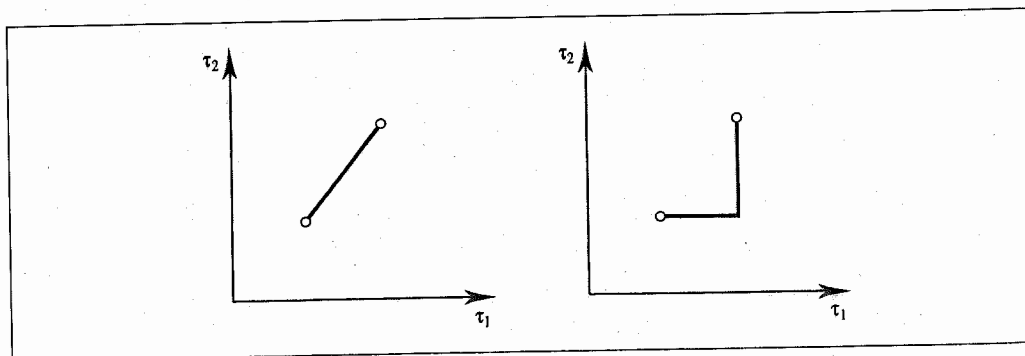


Fig. 9.29: Euclidean and Hamming distance measures of torsional similarity.

changes in the torsion angles in the middle of a molecule can give rise to large movements near the ends. The RMSD produces clusters in which the molecules have a similar shape.

One family of relatively straightforward clustering algorithms is the *linkage methods*. These algorithms first require the distance between each pair of conformations to be calculated. At the start of the cluster analysis the data set contains as many clusters as there are conformations; each cluster contains just a single conformation. At each step the total number of clusters is reduced by one by merging the 'closest' or 'most similar' pair of clusters into a single cluster. Thus in the first step the closest two conformations are merged into a single cluster. In the next step, the closest two clusters are merged, and so on. Clustering continues until the distance between the closest pair of clusters exceeds a predetermined value, until the number of clusters falls below a specified maximum number, or until all the conformations have been merged into a single cluster. Such algorithms are referred to as *agglomerative methods*, in contrast to *divisive clustering algorithms*, which start with a single cluster containing all of the data that is then partitioned into clusters. A representative conformation may then be chosen from each cluster, for example the conformation that is closest to the average structure of the cluster. The linkage methods differ in the way in which they calculate the distance between two clusters.

In the *single-linkage* or *nearest-neighbour* method the 'distance' between a pair of clusters is equal to the shortest distance between any two members, one from each cluster. The *complete-linkage* or *furthest-neighbour* method is the logical opposite of the single-linkage method in that it considers the furthest pair of objects in a pair of clusters. The *group average* method computes the average of the similarities between all pairs of objects in the two clusters.

We can contrast these methods using the data shown in Figure 9.30, which were obtained by searching the Cambridge Structural Database for the ribose phosphate fragment also shown

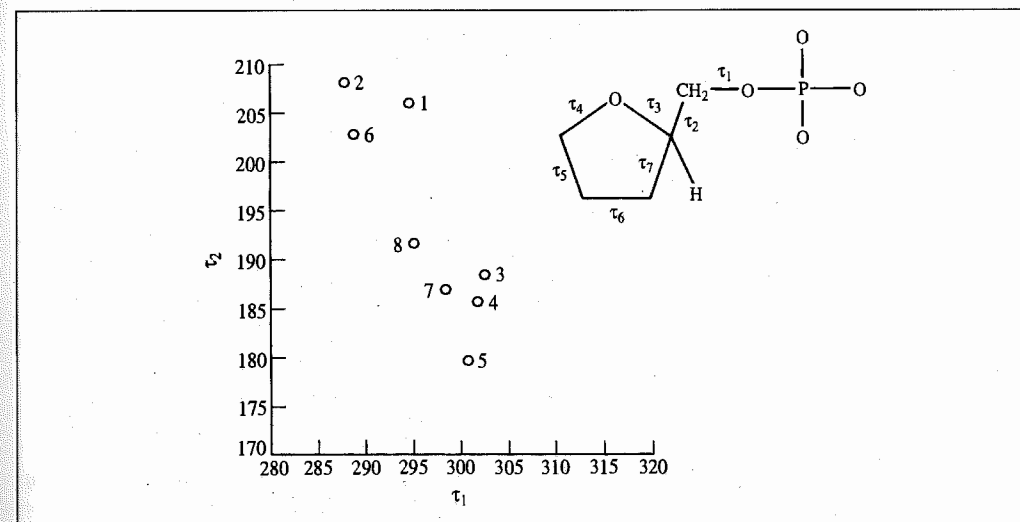


Fig. 9.30: Ribose phosphate fragment used to extract data from Cambridge Structural Database and eight sets of torsion angle values for τ_1 and τ_2 .

	1	2	3	4	5	6	7	8
1	0.0	7.2	19.3	21.4	26.9	6.6	19.4	14.2
2	7.2	0.0	24.6	26.3	31.1	5.3	23.6	17.9
3	19.3	24.6	0.0	2.7	8.8	19.9	4.5	8.1
4	21.4	26.3	2.7	0.0	6.1	21.4	3.7	8.9
5	26.9	31.1	8.8	6.1	0.0	26.0	7.6	13.2
6	6.6	5.3	19.9	21.4	26.0	0.0	18.5	12.8
7	19.4	23.6	4.5	3.7	7.6	18.5	0.0	5.7
8	14.2	17.9	8.1	8.9	13.2	12.8	5.7	0.0

Table 9.2: Distance matrix for eight ribose phosphate fragments.

in Figure 9.30. A total of 44 molecules were found to contain this fragment. The values of the two torsion angles τ_1 and τ_2 indicated in Figure 9.28 were determined for each occurrence of the fragment (some molecules contained more than one representative of the fragment). For simplicity, the results for just eight fragments are plotted in Figure 9.28. The distance matrix for these eight sets of two torsion angles, calculated using a Euclidean measure, is given in Table 9.2.

All of the clustering methods first join the two structures that are 'closest' (conformations 3 and 4), to which is then added conformation 7. In the third step, conformations 2 and 6 are connected. In the fourth step, the single-linkage and complete-linkage methods differ; the single-linkage algorithm joins conformation 8 to the cluster 3-4-7, whereas the complete-linkage method joins conformation 1 to the cluster 2-6. The order in which the points are clustered together for the three linkage methods is given in Table 9.3. As can be seen, the three linkage methods form the clusters in a similar but not identical order.

These three linkage methods are all *hierarchical agglomerative* clustering methods, because there is a specific order in which the clusters are formed and amalgamated. The same basic approach underlies all such methods, involving a series of iterations at each of which the two closest clusters are identified and combined into a larger cluster. The process continues until just a single cluster remains. These methods have the advantage of being simple to program, and they also produce a clustering that is independent of the order in

Step number	Single linkage	Complete linkage	Group average
1	3-4 (2.7)	3-4 (2.7)	3-4 (2.7)
2	3-4-7 (3.7)	3-4-7 (4.5)	3-4-7 (4.1)
3	2-6 (5.3)	2-6 (5.3)	2-6 (5.3)
4	3-4-7-8 (5.7)	2-6-1 (7.2)	2-6-1 (6.9)
5	3-4-7-8-5 (6.1)	3-4-7-5 (8.8)	3-4-7-5 (7.5)
6	2-6-1 (6.6)	3-4-7-5-8 (13.2)	3-4-7-5-8 (9.0)
7	2-6-1-3-4-7-8-5 (12.8)	2-6-1-3-4-7-5-8 (31.1)	2-6-1-3-4-7-5-8 (21.3)

Table 9.3: A comparison of the single-linkage, complete-linkage and average-linkage cluster methods using the data in Table 9.2. The figures in parentheses indicate the 'distance' between the clusters as they are formed. In this particular case Ward's clustering follows the same order of cluster formation as the group average method.

which the objects are stored. However, they do suffer from some drawbacks. For example, the commonly used single-linkage method tends to produce long, elongated clusters. In addition, simple implementations require an $M \times M$ similarity matrix to be calculated, which can severely limit their applicability when clustering large data sets.

A fourth hierarchical method that is quite popular is Ward's method [Ward 1963]. This method merges those two clusters whose fusion minimises the 'information loss' due to the fusion. Information loss is defined in terms of a function which for each cluster i corresponds to the total sum of squared deviations from the mean of the cluster:

$$E_i = \sum_{j=1}^{N_i} (|r_j - \bar{r}_i|)^2 \quad (9.39)$$

The summation runs over the N_i objects in cluster i , each located at r_j and where the mean of the cluster is \bar{r}_i . The total information loss is calculated by adding together the values for each cluster. At each iteration that pair of clusters which gives rise to the smallest increase in the total error function are merged. Two more hierarchical clustering algorithms are the centroid method, which determines the distance between two clusters as the distance between their centroids, and the median method, which represents each cluster by the coordinates of the median value. Fortunately, all six hierarchical agglomerative methods can be represented by a single equation, first proposed by Lance and Williams [Lance and Williams 1967], with the different algorithms having different coefficients.

A hierarchical clustering can be represented visually by constructing a *dendrogram*, which indicates the relationship between the items in the data set. A sample dendrogram is shown in Figure 9.31 for the single-linkage clustering described above. Along the x axis are represented the M individual objects. The y axis indicates the intercluster distance.

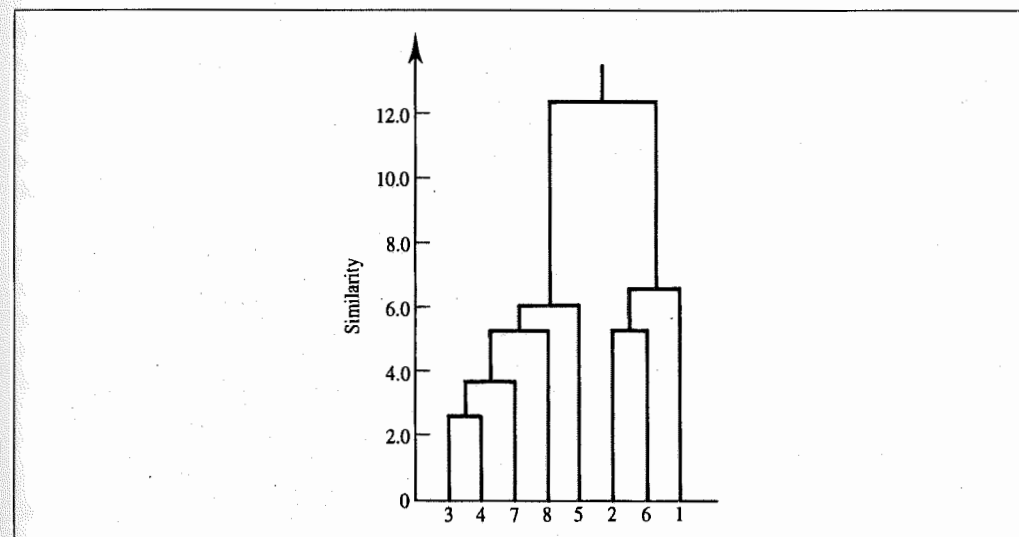


Fig. 9.31: Dendrogram of the single-linkage data in Table 9.3.

The dendrogram enables us to identify how many clusters there are at any stage, and what the members of those clusters are. A dendrogram can thus be a very useful way to show the underlying structure of the data and for suggesting the appropriate number of clusters to choose. A line drawn across the dendrogram enables one to read off how many clusters there are at that particular distance measure. For example, there are four clusters at a value of 6.0. For the data in Figure 9.30 it would probably be decided that the data fall into two clusters, one containing conformations 1, 2 and 6 and the other containing 3, 4, 5, 7 and 8. As this example illustrates, deciding how many clusters there are can be somewhat subjective; a small threshold can lead to a large number of 'tight' clusters (and frequently many clusters with just one member) whereas a larger threshold can produce clusters that are spread out.

An example of a non-hierarchical clustering method is the Jarvis–Patrick algorithm [Jarvis and Patrick 1973]. The Jarvis–Patrick method uses a 'nearest-neighbours' approach. The nearest neighbours of each conformation are the conformations that are the shortest distance away. Two conformations are considered to be in the same cluster in the Jarvis–Patrick method if they satisfy the following criteria:

1. They are in each other's list of m nearest neighbours.
2. They have p (where $p < m$) nearest neighbours in common.

Conformations can thus be placed in clusters and clusters fused together (because any two individual elements satisfy the two criteria) without any hierarchical relationships. The Jarvis–Patrick method can also be extended to take account not only of the number of nearest neighbours but also the position of each conformation within the neighbour list. In addition, it is possible to require that a molecule's nearest neighbours must be within some defined distance. This ensures that the nearest neighbours of each conformation are not too dissimilar.

To illustrate the use of the Jarvis–Patrick method, let us consider the data in Figure 9.30 once more. The three nearest neighbours of each fragment are given in Table 9.4. Suppose we require that two out of the three nearest neighbours should be common. If we examine the pair 1, 2 we find that each neighbour list contains the other fragment and the remaining two nearest neighbours are the same (i.e. 6, 8). These objects would therefore be placed in the same cluster. However, fragments 2 and 6 would not be considered in the same cluster

Fragment	Nearest neighbours
1	2, 6, 8
2	1, 6, 8
3	4, 7, 8
4	3, 5, 7
5	3, 4, 7
6	1, 2, 8
7	3, 4, 8
8	3, 4, 7

Table 9.4: The three nearest neighbours of each fragment in Figure 9.30.

according to these criteria; although they are in each other's list they do not have two out of three nearest neighbours in common. One advantage of the Jarvis–Patrick algorithm is that it can be used to cluster very large data sets which may be too large for any of the hierarchical methods to handle, due to their typically larger computational requirements. The K-means method, which is another non-hierarchical approach, is also applicable to larger sets. The K-means algorithm first chooses a set of c 'seed' objects, usually at random. The remaining objects are assigned to the nearest seed to give an initial set of c clusters. The centroids of each of these clusters are then determined and the objects are reassigned to the nearest of these new cluster centroids. New centroids are then determined, and the process continues until no objects change clusters. The K-means method is obviously dependent upon the initial set of (random) cluster centroids and different results will usually result from different initial seeds.

A common use of cluster analysis is in selecting a set of representative molecules from a large chemical database; the advent of robotic methods for high-throughput screening has made this of particular interest and some studies have been published comparing various approaches [Downs *et al.* 1994]. It is always crucial to bear in mind that, in addition to the differences between clustering algorithms, the performance of a cluster analysis also depends critically upon the methods used to calculate the distances between the objects in the data set.

9.14 Reducing the Dimensionality of a Data Set

The *dimensionality* of a data set is the number of variables that are used to describe each object. For example, a conformation of a cyclohexane ring might be described in terms of the six torsion angles in the ring. However, it is often found that there are significant correlations between these variables. Under such circumstances, a cluster analysis is often facilitated by reducing the dimensionality of a data set to eliminate these correlations. *Principal components analysis* (PCA) is a commonly used method for reducing the dimensionality of a data set.

9.14.1 Principal Components Analysis

Consider the data shown in Figure 9.32. It is easy to see that there is a high degree of correlation between the x and the y values. If we were to define a new variable, $z = x + y$, then we could express most of the variation in the data as the values of this new variable z . The new variable is called a *principal component*. In general, a principal component is a linear combination of the variables:

$$p_i = \sum_{j=1}^v c_{i,j} x_j \quad (9.40)$$

where p_i is the i th principal component and $c_{i,j}$ is the coefficient of the variable x_j . There are v such variables. The first principal component of a data set corresponds to that linear combination of the variables which gives the 'best fit' straight line through the data when

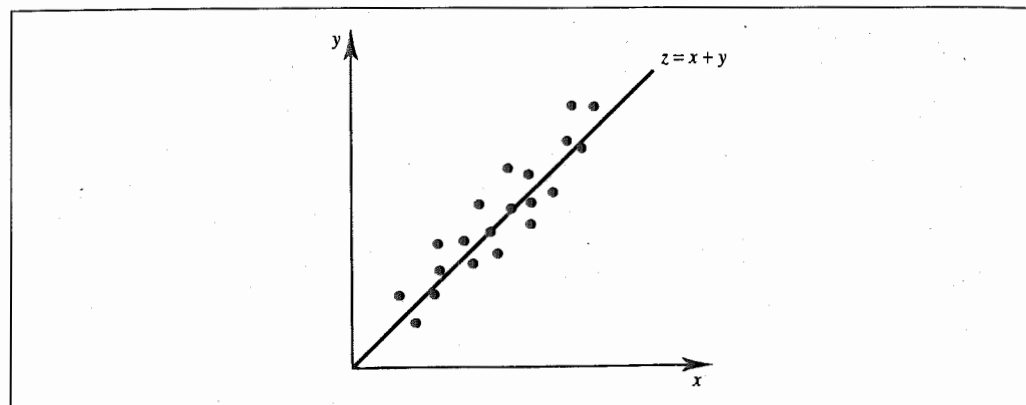


Fig. 9.32: Most of the variance in this set of highly correlated data values can be explained in terms of a new variable, $z = x + y$.

they are plotted in the v -dimensional space. More specifically, the first principal component maximises the *variance* in the data so that the data have their greatest 'spread' of values along the first principal component. This is clear in the two-dimensional example shown in Figure 9.32. The second and subsequent principal components account for the maximum variance in the data not already accounted for by previous principal components. Each principal component corresponds to an axis in a v -dimensional space, and each principal component is orthogonal to all the other principal components. There can clearly be as many principal components as there are dimensions in the original data, and indeed in order to explain all of the variation in the data it is usually necessary to include all the principal components. However, in many cases only a few principal components may be required to explain a significant proportion of the variation in the data. If only one or two principal components can explain most of the data then a graphical representation is possible.

The principal components are calculated using standard matrix techniques [Chatfield and Collins 1980]. The first step is to calculate the variance-covariance matrix. If there are s observations, each of which contains v values, then the data set can be represented as a matrix D with v rows and s columns. The variance-covariance matrix Z is:

$$Z = D^T D \quad (9.41)$$

The eigenvectors of Z are the coefficients of the principal components. As Z is a square symmetric matrix its eigenvectors will be orthogonal (provided there are no degenerate eigenvalues). The eigenvalues and their associated eigenvectors can be obtained by solving the secular equation $|Z - \lambda I| = 0$ or by matrix diagonalisation. The first principal component corresponds to the largest eigenvalue, the second principal component to the second largest eigenvalue, and so on. The i th principal component accounts for a proportion $\lambda_i / \sum_{j=1}^v \lambda_j$ of the total variance in the data. The first m principal components therefore account for $\sum_{j=1}^m \lambda_j / \sum_{j=1}^v \lambda_j$ of the total variation in the data.

As an example of the application of principal components analysis, we shall consider the conformations adopted by the five-membered ribose ring in our set of conformations

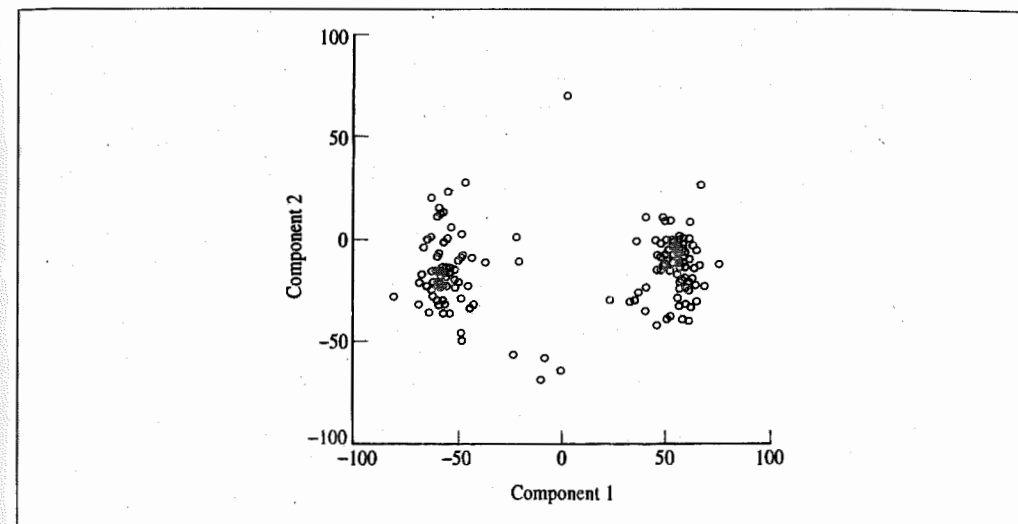


Fig. 9.33: Scatterplot of the first two principal components for the ring torsion angles τ_3 – τ_7 .

extracted from the Cambridge Structural Database. The conformation of a five-membered ring can be described in terms of five torsion angles ($\tau_3 \dots \tau_7$ in Figure 9.30). As we cannot visualise points in a five-dimensional space it would clearly be useful to reduce the dimensionality of the data set. When a principal components analysis is performed on the data, the following results are obtained:

Principal component	Proportion of variance explained (%)	$c(\tau_3)$	$c(\tau_4)$	$c(\tau_5)$	$c(\tau_6)$	$c(\tau_7)$
1	85.9	-0.14	-0.26	0.55	-0.61	0.48
2	14.0	-0.63	0.59	-0.31	-0.06	0.41
3	0.0002	-0.19	0.50	0.65	-0.004	-0.53
4	0.0001	-0.47	-0.38	0.12	0.71	0.19
5	0.0001	0.58	0.43	0.28	0.35	0.53

It can thus be seen that most of the variation in the data (85.9%) is explained by the first principal component, with all but a fraction being explained by the first two components. These two principal components can be plotted as a scatter graph, as shown in Figure 9.33, suggesting that there does indeed seem to be some clustering of the conformations of the five-membered ring in this particular data set.

9.15 Covering Conformational Space: Poling

As we have discussed, a common strategy in conformational analysis is to perform a two-stage process involving, first, the generation of a large number of minimum energy

conformations followed by the selection of a subset using a technique such as cluster analysis. This subset may then be considered 'representative' of the conformational space in subsequent calculations. There can be both practical and scientific objections to this approach. One of the practical problems is that it may take some considerable computational effort to first explore the conformational space and then to cluster the resulting conformations. One of the scientific objections is that, by restricting the initial search to minimum energy conformations, one may not adequately cover the conformational space. Consider the case of a broad, shallow minimum. It might be better to describe this region of the conformational space using an ensemble of structures rather than just a single minimum energy structure. A technique termed 'poling' has been described which is intended to promote the generation of diversity in the conformational coverage [Smellie *et al.* 1995a, b]. The poling approach introduces a penalty function into the geometry optimisation step that is a common component of most conformational searching procedures. This function is designed to penalise a conformation that is too close to any of the conformations already generated.

The effect of poling on the conformational space is shown using a one-dimensional energy surface in Figure 9.34, which contains two energy minima. Suppose we first generate conformation 1. The poling function is now introduced to modify the energy surface in the region of this conformation. This can have the effect of introducing new minima (labelled 2 in Figure 9.34) into the energy surface. In the third iteration poling functions are introduced around conformations 1 and 2, enabling conformation 3 to be produced. Note that in this example the unperturbed energy surface contains only two minima,

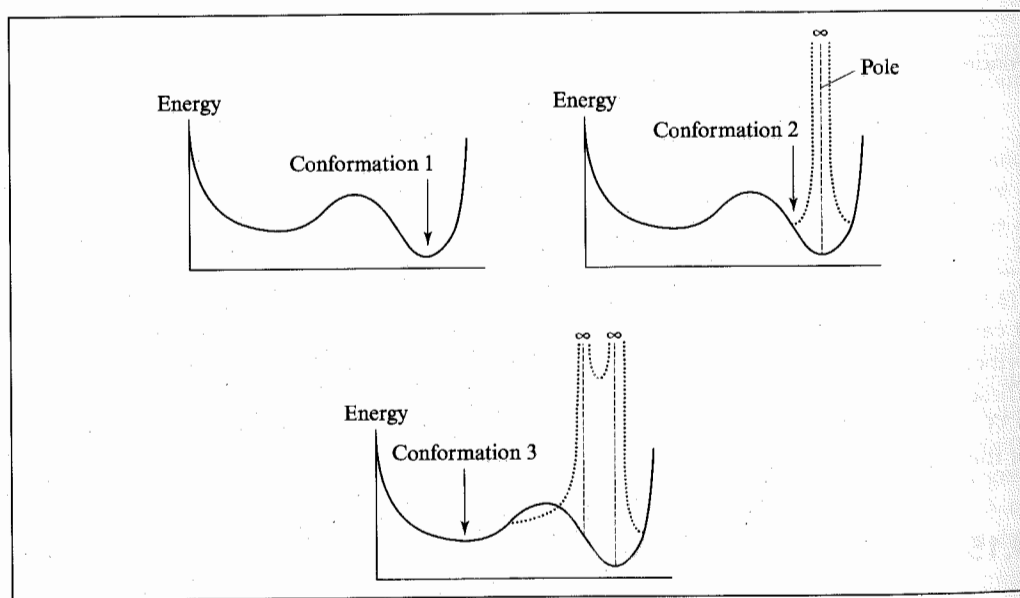


Fig. 9.34: Modification of the energy surface using poling. (Figure adapted from Smellie A S, S L Teig and P Towbin 1995b. Poling: Promoting Conformational Variation. *Journal of Computational Chemistry* 16:171-187.)

which would be the only structures identified by a traditional conformational search algorithm. The use of the poling function not only ensures that previously explored regions of conformational space are avoided but can also lead to a wider range of conformations being generated, that might be considered to better represent the accessible conformational space.

The poling function typically adopts the following general functional form:

$$F_{\text{pole}} = W_{\text{pole}} \sum_i \frac{1}{(D_i)^N} \quad (9.42)$$

$$D_i = \left(\frac{\sum_{j=1}^{N_d} (d_{j,\text{curr}} - d_{j,i})^2}{N_d} \right)^{1/2} \quad (9.43)$$

Here, N_d poling distances are being compared between the current conformation and a previously generated conformation, i . Thus $d_{j,\text{curr}}$ represents one of these N_d poling distances in the current conformation and $d_{j,i}$ represents the equivalent distance in the i th conformation. D_i is thus the root-mean-square difference between the current conformation and conformation i over these poling distances. The steepness of these poling functions can be changed by modifying the power N . One way that the poling distance could be implemented would be to sum over all interatomic distances in the molecule. However, this would be rather inefficient as the number of such distances increases with the square of the number of atoms in the molecule. Rather, in the published work a set of 'chemically significant features', such as hydrogen bond donors and hydrogen bond acceptors, was identified and the poling distance was defined from each such feature to the centroid of all the features.

9.16 A 'Classic' Optimisation Problem: Predicting Crystal Structures

Many molecules are obtained and used in a crystalline form, the nature of which can have a significant impact on their properties and behaviour. Moreover, it is sometimes possible for a given material to exist in more than one crystalline form, depending upon the conditions under which it was prepared. This is the phenomenon of *polymorphism*. This can be important because the various polymorphs may themselves have different properties. It is therefore of interest to be able to predict the three-dimensional atomic structure(s) that a given molecule may adopt, for those cases where it is difficult to obtain experimental data and also where one might wish to prioritise molecules not yet synthesised.

Many different approaches have been suggested as possible approaches to this problem, from the 1960s onwards [Verwer and Leusen 1998]. What is obvious from all of these efforts is that this is an extremely difficult problem. Both thermodynamics and kinetics can be important in determining which crystalline form is obtained under a certain set of experimental conditions. Kinetic effects are particularly difficult to take into account and so are usually ignored. A proper treatment of the thermodynamic factors would require one to deal with the relative free energies of the different possible polymorphs.

These relative free energies are dependent upon internal energies, crystal densities and entropies:

$$\Delta G = \Delta U + P\Delta V - T\Delta S \quad (9.44)$$

Of these three contributions, that due to differences in density can safely be ignored (at least at normal pressures). The entropy differences are also invariably neglected, due to difficulties in calculating these contributions. This leaves the internal energy (at 0K) as the metric by which the relative stabilities of polymorphs are predicted. In practice, the aim of most crystal structure prediction methods is to suggest a (hopefully small) number of solutions, according to their relative energies. These may subsequently be distinguished experimentally. For example, whilst it might prove impossible to obtain high-quality single-crystal diffraction data it may be feasible to acquire powder diffraction data, which, when combined with the computational results, can lead to a plausible solution.

When viewed purely as a search problem, one can readily appreciate its complexity. Not only should one consider the conformational flexibility of the molecule but one also has to suggest how these conformations might be able to pack into a low-energy structure. In addition, real crystals not only contain molecules of the compound of interest (sometimes in more than one conformation) but also solvent molecules and counterions, with the stoichiometry often being unknown prior to the experimental determination. One redeeming feature of the problem is that there are some well-established constraints on the way in which the molecules pack together, namely that the final structure must fall into one of the 230 space groups. In addition, in order to reduce the scale of the problem an algorithm may be limited to the more commonly occurring space groups and/or restricted to consider just one molecule in the asymmetric unit (see Section 3.8.1 for a brief description of some of these crystallographic terms).

Here we will consider just two of the more recent methods, which have much in common but also some significant differences. Gavezzotti's PROMET method [Gavezzotti 1991, 1994] starts by constructing clusters (called 'crystal nuclei') containing two molecules. The molecules are provided to the program in a predefined conformation that remains constant throughout the calculation. The relative locations of the molecules in these clusters are generated by applying common symmetry operators. Each symmetry operator can give rise to a number of possible cluster geometries, each of which is assessed by calculating its intermolecular energy. The most favourable clusters are selected for the subsequent steps, which may involve the application of additional symmetry operators or an attempt to construct a full crystal structure by translating the cluster to give a lattice in three dimensions. The intermolecular energy is again used to guide this process, which proceeds by building a sequence of clusters in one dimension, then two and finally three. Only if an improvement in the energy is achieved does the algorithm proceed to the next stage. In essence, the process is somewhat akin to a systematic search, but one which uses a variety of criteria to prune the search tree. In addition to the energetic criteria the growing lattice must also meet the condition that it has to belong to one of the known space groups. Further, analysis of known crystal structures reveals additional criteria that can be applied to restrict the ranges of some of the unit cell dimensions.

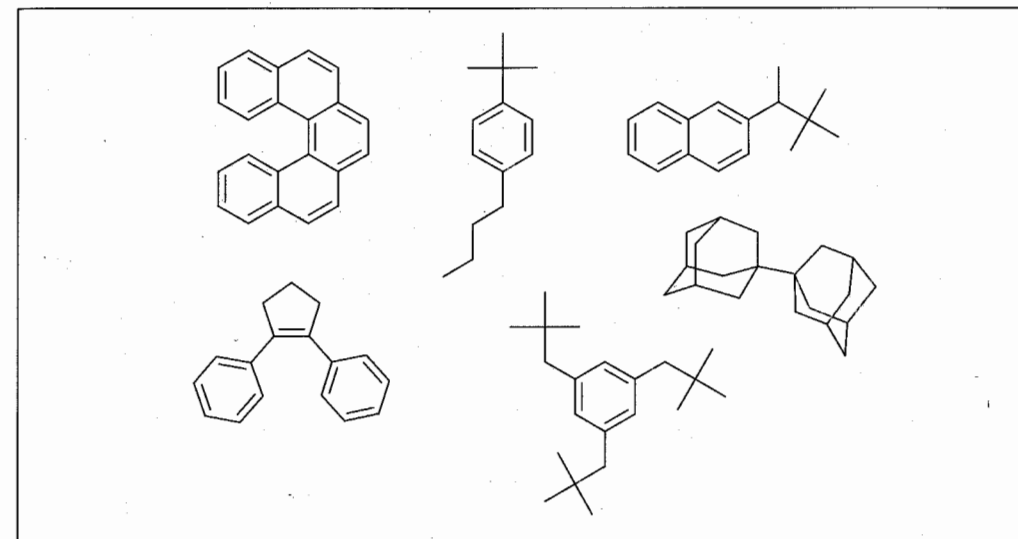


Fig. 9.35: Some typical molecules tested using the PROMET approach to predicting crystal structures.

The PROMET method has been applied to a variety of systems, though at the time of writing it is limited to rigid molecules containing C, H, N, O, S and Cl atoms in the eleven space groups which account for more than 90% of organic crystal structures. Some of these molecules are shown in Figure 9.35. To try to simulate polymorph prediction several cases were identified from the literature where the structure of one polymorph was completely characterised and where mention was made of another polymorph for which only the cell parameters and space group could be determined [Gavezzotti and Filippini 1996]. In these cases the computational method was used to try to predict the complete structure of the unknown polymorph. In the small number of test cases examined the approach did prove quite successful, in that it was possible to identify low-energy packing arrangements that agreed with the experimental data. Of perhaps most interest was the fact that the energies of the various polymorphs were not in particularly good agreement with their relative stability but that the geometries were often rather well predicted. This points to a synergy between experiment and theory, in that when some experimental data were available (e.g. unit cell dimensions and space group) then one might be quite confident of being able to predict its structure.

The second type of approach we shall consider can be thought of as more 'ab initio' in nature, in that there is no assumption that the lattice must be constructed from energetically favourable arrangements of molecules in the crystal nuclei. Rather, all packing arrangements in all possible space groups are generated. The following procedure is typically employed. Initially, molecules are placed into an oversized unit cell such that the symmetry relationships between the molecules are consistent with a particular space group. The molecules are then permitted to move, using either a random or systematic algorithm. This first phase can lead to a large number of trial structures (often several thousands), which are then clustered to identify duplicates. The lowest-energy structure from each cluster is then minimised, followed by a final clustering. It can sometimes be appropriate to employ this final clustering before the

structures are properly and completely minimised, as it is usually considered important to use as complete a force field model as possible for this final step (for example, using Ewald summations and very tight convergence criteria).

Probably the main difference between the different variants lies in the way in which the molecules move in the first step. One approach is to use Monte Carlo simulated annealing

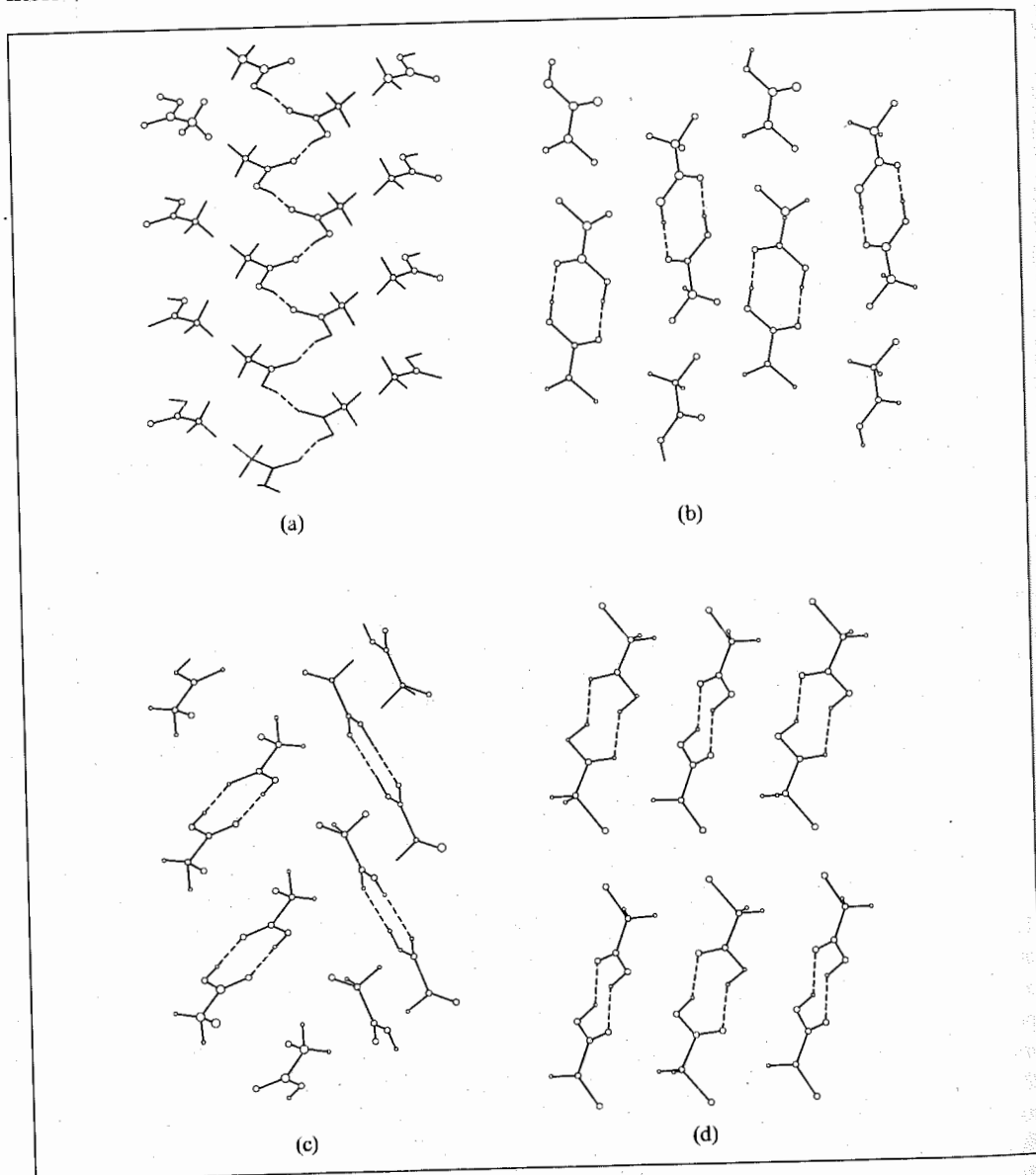


Fig. 9.36: Hydrogen bonding patterns in the crystal structures of acetic acid and its halo derivatives: (a) acetic acid, (b) fluoro acetic acid, (c) chloro acetic acid, (d) bromo acetic acid.

[Gdanitz 1992; Karfunkel and Gdanitz 1992; Karfunkel *et al.* 1993; Leusen 1996]. In this particular implementation it is the angular degrees of freedom that are varied during the Monte Carlo search. These angular variables comprise the cell angles, the Euler angles which describe the rigid-body rotations of the molecules in the cell, and the Euler angles which describe the actual location of the molecules in the cell. Having chosen a new set of angular parameters (or a subset of them) the translational parameters (which are the cell lengths and the distances between the molecules in the cell) are adjusted to relieve any close contacts that resulted. The new configuration is then accepted or rejected according to the Metropolis Monte Carlo criterion. A standard simulated annealing procedure is used, involving a slow decrease in temperature from several thousand kelvin to 300 K. This typically provides about 2000 accepted structures per space group for the next stage (clustering), from 4000–5000 Monte Carlo steps. An alternative to the Monte Carlo simulated annealing is to use a systematic search method. For example, one approach [van Eijck *et al.* 1995] involves a brute-force grid search in which the relevant parameters are varied systematically, with each trial structure being subjected to a few cycles of minimisation to locate approximately the nearest energy minimum.

One rather simple molecule which has been the subject of detailed comparisons of a number of methods is acetic acid. The reason for the interest in this molecule is that acetic acid (together with formic acid) adopts a chain-like structure, with each molecule forming one hydrogen bond with each of the two neighbouring molecules in the chain (Figure 9.36). Almost all other monocarboxylic acids form dimer-based structures, including fluoro, chloro and bromo acetic acid. Moreover, polymorphism has been detected for both chloro and bromo acetic acid. It is clearly of interest to be able to understand the reasons for this behaviour and also to determine whether it can be predicted by current computational methods. Acetic acid itself has been considered using both the grid search and the Monte Carlo simulated annealing approach, though it should also be noted that other aspects of the procedure differ as well (clustering algorithm, force field, minimisation method, etc.) [Mooij *et al.* 1998; Payne *et al.* 1998]. In both cases a number of low-energy structures were identified. These included the known crystal structure but also alternatives that were very similar in energy (i.e. dimers were found for acetic acid and chain structures for the halogenated derivatives). It thus appears that these methods are currently able to explore the search space quite effectively but the force fields currently used for the final assessment do not provide sufficient discrimination between the low-energy alternatives.

Further Reading

- Aldenderfer M S and R K Blahfield 1984. *Cluster Analysis*. Newbury Park, CA. Sage; New York, Garland Publishing.
- Blaney J M and J. S Dixon 1994. Distance Geometry in Molecular Modeling, In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 5. New York, VCH Publishers, pp. 299–335.
- Chatfield C and A J Collins 1980. *Introduction to Multivariate Analysis*. London, Chapman & Hall.
- Desiraju G R 1997. Crystal Gazing: Structure Prediction and Polymorphism. *Science* 278:404–405.
- Everitt B.S. 1993 *Cluster Analysis*. Chichester, John Wiley & Sons.

- Gavezzotti A (Editor) 1997. *Theoretical Aspects and Computer Modeling of the Molecular Solid State*. Chichester, John Wiley & Sons.
- Gavezzotti A 1998. The Crystal Packing of Organic Molecules: Challenge and Fascination below 1000Da. *Crystallography Reviews* 7:5-121.
- Leach A R 1991. A Survey of Methods for Searching the Conformational Space of Small and Medium-Sized Molecules. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 2. New York, VCH Publishers, pp. 1-55.
- Perutz M 1992. *Protein Structure. New Approaches to Disease And Therapy*. New York, W H Freeman.
- Scheraga H A 1993. Searching Conformational Space. In van Gunsteren W F, P K Weiner and A J Wilkinson (Editors) *Computer Simulation of Biomolecular Systems* Volume 2. Leiden, ESCOM.
- Schulz G E and R H Schirmer 1979. *Principles of Protein Structure*. New York, Springer-Verlag.
- Torda A E and W F van Gunsteren 1992. Molecular Modeling Using NMR Data. In Lipkowitz K B and D B Boyd (Editors). *Reviews in Computational Chemistry* Volume 2. New York, VCH Publishers, pp. 143-172.
- Verwer P and F J J Leusen 1998. Computer Simulation to Predict Possible Crystal Polymorphs. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 12. New York, VCH Publishers, pp. 327-365.

References

- Allen F H, S A Bellard, M D Brice, B A Cartwright, A Doubleday, H Higgs, T Hummelink, B G Hummelink-Peters, O Kennard, W D S Motherwell, J R Rodgers and D G Watson 1979. The Cambridge Crystallographic Data Centre: Computer-based Search, Retrieval, Analysis and Display of Information. *Acta Crystallographica* B35:2331-2339.
- Allen F H, O Kennard, D G Watson, L Brammer, A G Orpen and R Taylor 1987. Tables of Bond Lengths Determined by X-ray and Neutron Diffraction. 1. Bond Lengths in Organic Compounds. *Journal of the Chemical Society Perkin Transactions II*:S1-S19.
- Barton D H R 1950. The Conformation of the Steroid Nucleus. *Experientia* 6:316-320.
- Bergerhoff G, R Hundt, R Sievers and I S Brown 1983. The Inorganic Crystal Structure Database. *Journal of Chemical Information and Computer Sciences* 23:66-69.
- Berman H M, J Westbrook, Z Feng, G Gilliland, T N Bhat, H Weissig, I N Shindyalor and P E Bourne 2000. The Protein Data Bank. *Nucleic Acids Research* 28:235-242.
- Bernstein F C, T F Koetzle, G J B Williams, E Meyer, M D Bryce, J R Rogers, O Kennard, T Shikanouchi and M Tasumi 1977. The Protein Data Bank: A Computer-based Archival File for Macromolecular Structures. *Journal of Molecular Biology* 112:535-542.
- Brunger A T, J Kuriyan and M Karplus 1987. Crystallographic R-factor Refinement by Molecular Dynamics. *Science* 235:458-460.
- Bruno I J, J C Cole, J P M Lommerse, R S Rowland, R Taylor and M L Verdonk 1997. Isostar: A Library of Information about Nonbonded Interactions. *Journal of Computer-Aided Molecular Design* 11:525-537.
- Chang G, W C Guida and W C Still 1989. An Internal Coordinate Monte Carlo Method for Searching Conformational Space. *Journal of the American Chemical Society* 111:4379-4386.
- Chatfield C and A J Collins 1980. *Introduction to Multivariate Analysis*. London, Chapman & Hall.
- Chung C-W, R M Cooke, A E I Proudfoot and T N C Wells 1995. The Three-dimensional Structure of RANTES. *Biochemistry* 34:9307-9314.
- Clark D E and D R Westhead 1996. Evolutionary Algorithms in Computer-aided Molecular Design. *Journal of Computer-Aided Molecular Design* 10:337-358.
- Crippen G M 1981. *Distance Geometry and Conformational Calculations*. Chemometrics Research Studies Series 1. New York, John Wiley & Sons.

- Crippen G M and T F Havel 1988. *Distance Geometry and Molecular Conformation*. Chemometrics Research Studies Series 15. New York, John Wiley & Sons.
- Derome A E 1987. *Modern NMR Techniques for Chemistry Research*. Oxford, Pergamon.
- Downs G M, P Willett and W Fisanick 1994. Similarity Searching and Clustering of Chemical Structure Databases using Molecular Property Data. *Journal of Chemical Information and Computer Sciences* 34:1094-1102.
- Ferguson D M and D J Raber 1989. A New Approach to Probing Conformational Space with Molecular Mechanics: Random Incremental Pulse Search. *Journal of the American Chemical Society* 111:4371-4378.
- Ferro D R and J Hermans 1977. A Different Best Rigid-body Molecular Fit Routine. *Acta Crystallographica* A33:345-347.
- Gavezzotti A 1991. Generation of Possible Crystal Structures from the Molecular Structure for Low-polarity Organic Compounds. *Journal of the American Chemical Society* 113:4622-4629.
- Gavezzotti A 1994. Are Crystal Structures Predictable? *Accounts of Chemical Research* 27:309-314.
- Gavezzotti A and G Filippini 1996. Computer Prediction of Organic Crystal Structures Using Partial X-ray Diffraction Data. *Journal of the American Chemical Society* 118:7153-7157.
- Gdanitz, R J 1992. Prediction of Molecular Crystal Structures by Monte Carlo Simulated Annealing Without Reference to Diffraction Data. *Chemical Physics Letters* 190:391-396.
- Gibson K D and H A Scheraga 1987. Revised Algorithms for the Build-up Procedure for Predicting Protein Conformations by Energy Minimization. *Journal of Computational Chemistry* 8:826-834.
- Glusker J P 1995. Intermolecular Interactions Around Functional Groups in Crystals: Data for Modeling the Binding of Drugs to Biological Macromolecules. *Acta Crystallographica* D51:418-427.
- Goldberg D E 1989. *Genetic Algorithms in Search, Optimization and Machine Learning*. Reading, MA., Addison-Wesley.
- Goodman J M and W C Still 1991. An Unbounded Systematic Search of Conformational Space. *Journal of Computational Chemistry* 12:1110-1117.
- Jack A and M Levitt 1978. Refinement of Large Structures by Simultaneous Minimization of Energy and R-factor. *Acta Crystallographica* A34:931-929.
- Jarvis R A and E A Patrick 1973. Clustering Using a Similarity Measure Based on Shared Near Neighbours. *IEEE Transactions in Computers* C-22:1025-1034.
- Jones G 1998. Genetic and Evolutionary Algorithms. In Schleyer, P v R, N L Allinger, T Clark, J Gasteiger, P A Kollman, H F Schaefer III and P R Schreiner (Editors) *The Encyclopedia of Computational Chemistry*. Chichester, John Wiley & Sons.
- Jones T A and S Thirup 1986. Using Known Substructures in Protein Model Building and Crystallography. *EMBO Journal* 5:819-822.
- Judson R 1997. Genetic Algorithms and Their Use in Chemistry. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 10. New York, VCH Publishers, pp. 1-73.
- Judson R S, W P Jaeger, A M Treasurywala and M L Peterson 1993. Conformational Searching Methods for Small Molecules. 2. Genetic Algorithm Approach. *Journal of Computational Chemistry* 14:1407-1414.
- Kabsch W 1978. A Discussion of the Solution for the Best Rotation to Relate Two Sets of Vectors. *Acta Crystallographica* A34:827-828.
- Karfunkel H R and R J Gdanitz 1992. *Ab initio* Prediction of Possible Crystal Structures for General Organic Molecules. *Journal of Computational Chemistry* 13:1171-1183.
- Karfunkel H R, B Rohde, F J J Leusen, R J Gdanitz, and G Rihs 1993. Continuous Similarity Measure Between Nonoverlapping X-ray Powder Diagrams of Different Crystal Modifications. *Journal of Computational Chemistry* 14:1125-1135.
- Kirkpatrick S, C D Gelatt and M P Vecchi 1983. Optimization by Simulated Annealing. *Science* 220:671-680.

- Kolossvary I and W C Guida 1996. Low Mode Search. An Efficient, Automated Computational Method for Conformational Analysis: Application to Cyclic and Acyclic Alkanes and Cyclic Peptides. *Journal of the American Chemical Society* **118**:5011–5019.
- Lance G N and W T Williams 1967. A General Theory of Classificatory Sorting Strategies 1. Hierarchical Systems. *Computer Journal* **9**:373–380.
- Leach A R, D P Dolata and K Prout 1990. Automated Conformational Analysis and Structure Generation: Algorithms for Molecular Perception. *Journal of Chemical Information and Computer Science* **30**:316–324.
- Leach A R, K Prout and D P Dolata. 1988. An Investigation into the Construction of Molecular Models using the Template Joining Method. *Journal of Computer-Aided Molecular Design* **2**:107–123.
- Leusen F J L 1996. *Ab initio* Prediction of Polymorphs. *Journal of Crystal Growth* **166**:900–903.
- Li Z Q and H A Scheraga 1987. Monte-Carlo-minimization Approach to the Multiple-minima Problem in Protein Folding. *Proceedings of the National Academy Of Sciences USA* **84**:6611–6615.
- McGarrah D B and R S Judson 1993. Analysis of the Genetic Algorithm Method of Molecular-conformation Determination. *Journal of Computational Chemistry* **14**:1385–1395.
- Meza J C, R S Judson, T R Faulkner and A M Treasurywala 1996. A Comparison of a Direct Search Method and a Genetic Algorithm for Conformational Searching. *Journal of Computational Chemistry* **17**:1142–1151.
- Mooij W T M, B P van Eijck, S L Price, P Verwer and J Kroon 1998. Crystal Structure Predictions for Acetic Acid. *Journal of Computational Chemistry* **19**:459–474.
- Murray-Rust P M and J P Glusker 1984. Directional Hydrogen Bonding to sp^2 and sp^3 -hybridized Oxygen Atoms and Its Relevance to Ligand-Macromolecule Interactions. *Journal of the American Chemical Society* **106**:1018–1025.
- Payne R S, R J Roberts, R C Crowe and R Docherty 1998. Generation of Crystal Structures of Acetic Acid and Its Halogenated Analogs. *Journal of Computational Chemistry* **19**:1–20.
- Ramachandran G N, C Ramakrishnan and V Sasiekhara 1963. Stereochemistry of Polypeptide Chain Configurations. *Journal of Molecular Biology* **7**:95–99.
- Saunders M 1987. Stochastic Exploration of Molecular Mechanics Energy Surface: Hunting for the Global Minimum. *Journal of the American Chemical Society* **109**:3150–3152.
- Saunders M, K N Houk, Y-D Wu, W C Still, M Lipton, G Chang and W C Guida 1990. Conformations of Cycloheptadecane. A Comparison of Methods for Conformational Searching. *Journal of the American Chemical Society* **112**:1419–1427.
- Smellie A S, S D Kahn and S L Teig 1995a. Analysis of Conformational Coverage. 1. Validation and Estimation of Coverage. *Journal of Chemical Information and Computer Science* **35**:285–294.
- Smellie A S, S L Teig and P Towbin 1995b. Poling: Promoting Conformational Variation. *Journal of Computational Chemistry* **16**:171–187.
- Torda A E, R M Scheek and W F van Gunsteren 1990 Time-averaged Nuclear Overhauser Effect Distance Restraints Applied to Tendamistat. *Journal of Molecular Biology* **214**:223–235.
- Van Eijck B P, W T M Mooij and J Kroon 1995. Attempted Prediction of the Crystal Structures of Six Monosaccharides. *Acta Crystallographica* **B51**:99–103.
- Verwer P and F J L Leusen 1998. Computer Simulation to Predict Possible Crystal Polymorphs. In Lipkowitz K B and Boyd D B (Editors) *Reviews in Computational Chemistry*. New York, Wiley-VCH, pp. 327–365.
- Ward J H 1963. Hierarchical Grouping to Optimise an Objective Function. *American Statistical Association Journal*: 236–244.
- Weiner S J, P A Kollman, D A Case, U C Singh, C Ghio, G Alagona, S Profeta and P Weiner 1984. A New Force Field for Molecular Mechanical Simulation of Nucleic Acids and Proteins. *Journal of the American Chemical Society* **106**:765–784.

CHAPTER TEN

Protein Structure Prediction, Sequence Analysis and Protein Folding

10.1 Introduction

Peptides and proteins are polymers constructed from sequences of amino acids. They perform many functions essential to life. There are twenty common naturally occurring amino acids, shown in Figure 10.1. The amino acids are linked together via amide bonds to give a polypeptide chain. All the naturally occurring amino acids have the same relative stereochemistry at the alpha-carbon (referred to as 'L'). The side chains have different sizes, shapes, hydrogen-bonding capabilities and charge distributions, which enable proteins to display the vast array of biological functions required by living systems.

Protein biosynthesis is a very complex process. The amino acid sequence of a protein is determined by the DNA sequence of the corresponding gene. Each amino acid is coded by three adjacent DNA bases. However, DNA is not used directly in protein biosynthesis. Rather, an RNA copy is made from the DNA template; this *messenger RNA* (mRNA) in turn acts as the template for the protein synthesis. This process is known as *transcription*. In the subsequent *translation* step, the mRNA template is read by transfer RNA (tRNA), which also brings the actual amino acids to the site of synthesis. This two-stage, unidirectional flow of genetic information was proposed by Francis Crick and is known as the 'Central Dogma'. It is often represented by the diagram shown in Figure 10.2(a). Some modifications were required to the theory following the discovery of retroviruses, which can transfer genetic information from RNA to DNA (the dotted lines in Figure 10.2(b)) but the Central Dogma as originally proposed still holds true for most organisms.

The biological function of a protein or peptide is often intimately dependent upon the conformation(s) that the molecule can adopt. In contrast to most synthetic polymers where the individual molecules can adopt very different conformations, a protein usually exists in a single native state. These native states are found under conditions typically found in living cells (aqueous solvents near neutral pH at 20–40°C). Proteins can be unfolded (or *denatured*) using high-temperature, acidic or basic pH or certain non-aqueous solvents. However, this unfolding is often reversible and so proteins can be folded back to their native structure in the laboratory.

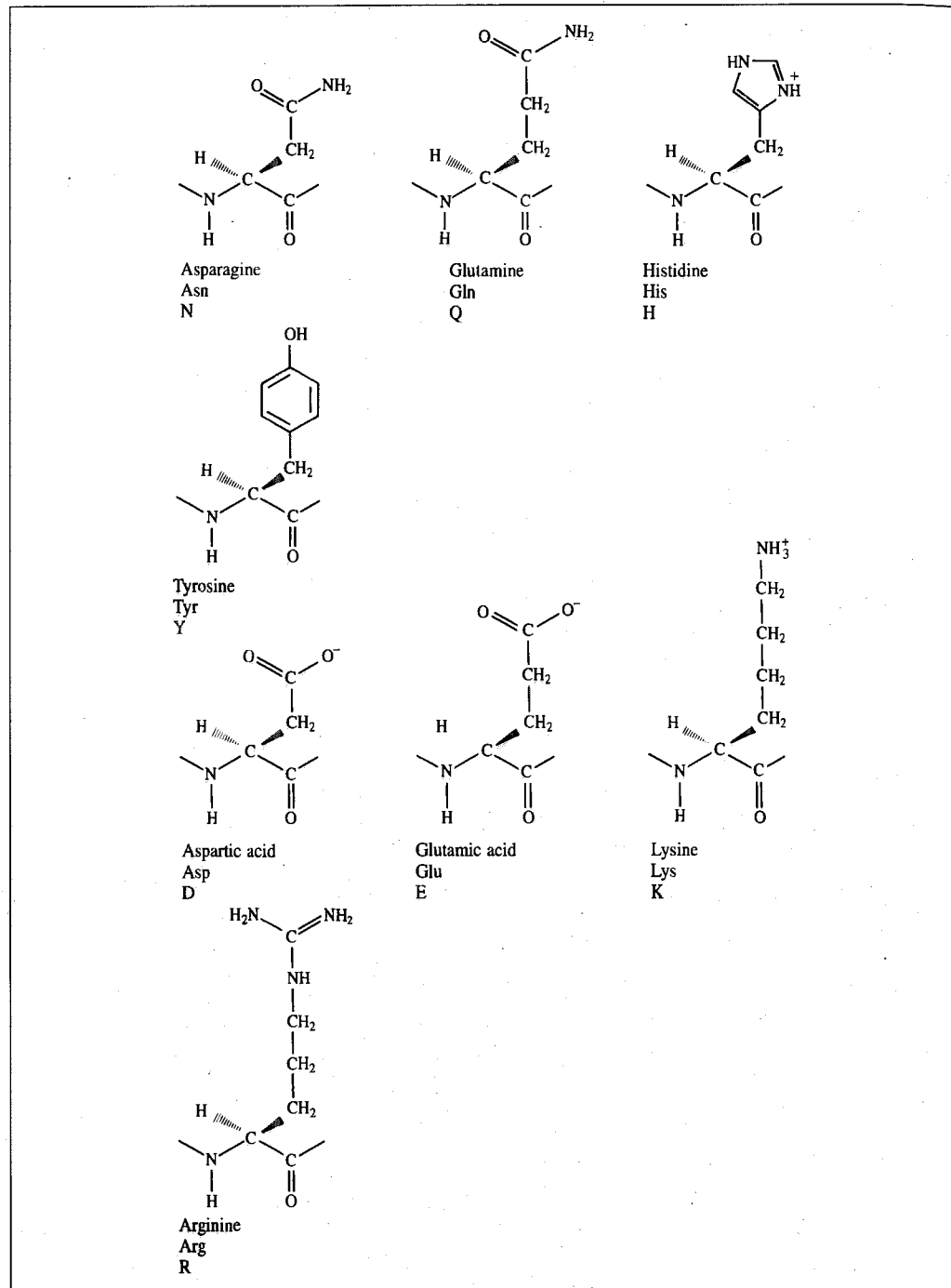


Fig. 10.1: The twenty naturally occurring amino acids and their codes.

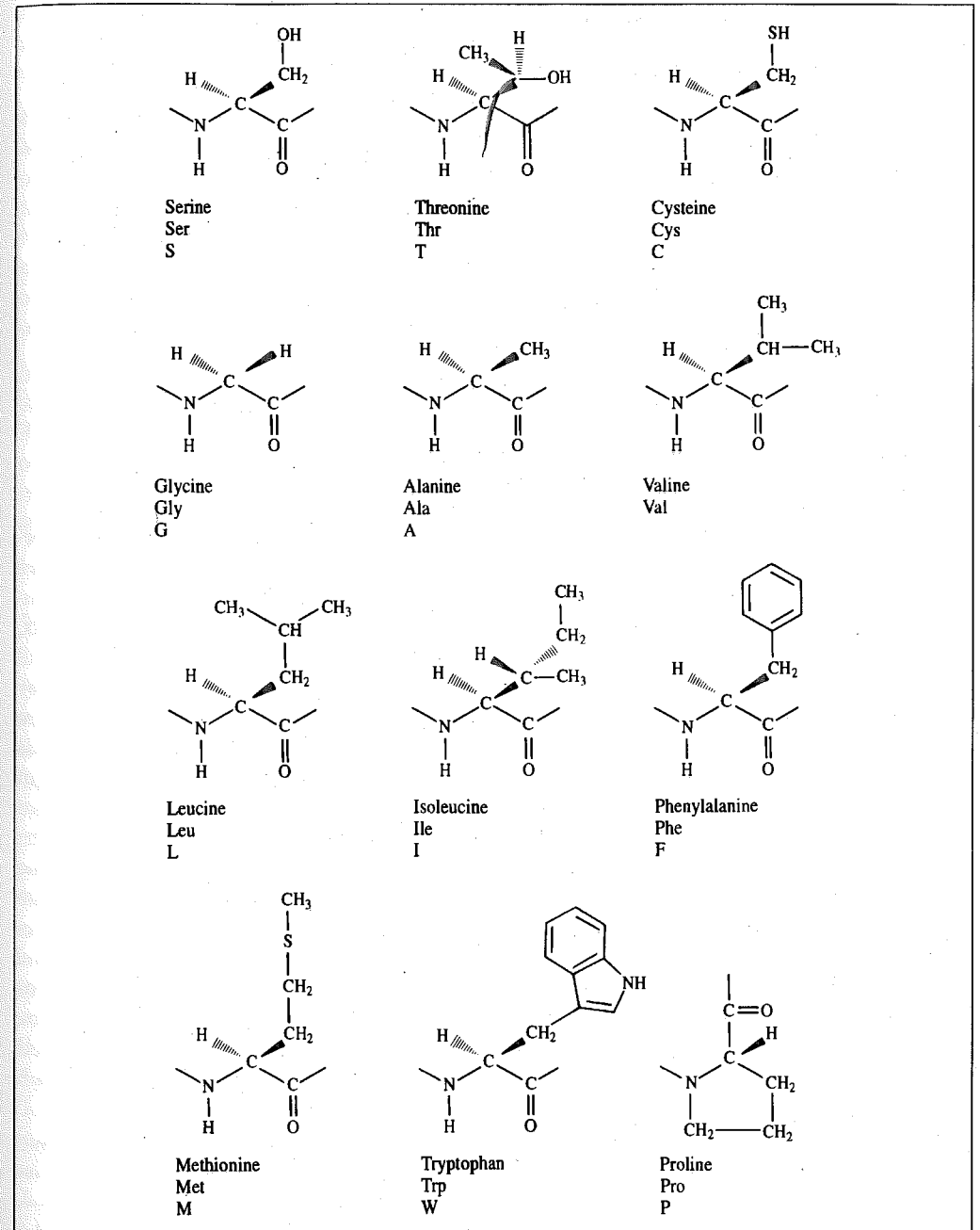


Fig. 10.1: Continued.

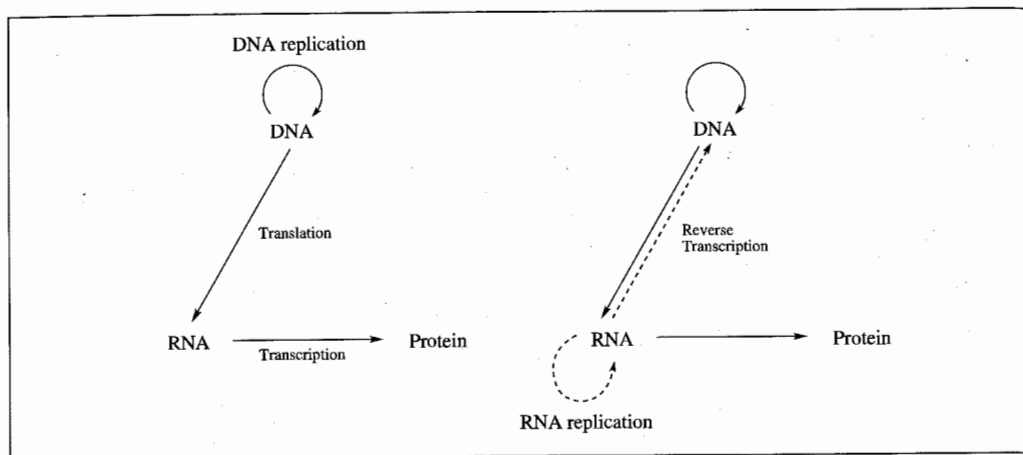


Fig. 10.2: The original Central Dogma of molecular biology (left) and its modification in light of the discovery of retroviruses (right).

X-ray crystallography and NMR are the methods most widely used to provide detailed information about protein structures. Unfortunately, the rate at which new protein sequences are being determined far exceeds the rate at which protein structures are determined experimentally. This is a particularly pertinent problem due to the efforts of the Human Genome Project, which is expected to have sequenced the entire human genome by 2003, if not earlier. It will then be necessary to determine the amino acid sequences of the proteins that are encoded by the DNA. This is not quite so straightforward as might be imagined due to the complex nature of the transcription/translation process and also because experimental sequencing methods do not deliver single 'genes'. Moreover, not all DNA actually codes for protein and in many genes the biological information is contained in distinct units called *exons* (the intervening non-coding regions being *introns*). *Functional genomics* is concerned with characterising the proteins expressed by the genome and assigning a biological function. It is possible to some extent to assign function just from an analysis of the sequence alone. However, the intimate relationship between the three-dimensional structure and function of a protein makes functional assignment based upon the structure more appealing (sometimes known as *structural genomics**). The general difficulties in obtaining protein structures using experimental techniques means that there is considerable interest in theoretical methods for predicting the three-dimensional structure of proteins from the amino acid sequence: this is often referred to as the *protein folding problem*. The results of a sequence analysis or structure prediction

* Functional genomics and structural genomics are widely used (and abused) terms. They are sometimes considered to apply only to large-scale, high-throughput technologies. Moreover, although computational methods have an important role to play, input from experimental techniques can also be critical. An example of this was the use of X-ray crystallography to predict the function of a protein from a hyperthermophile bacterium, *Methanococcus jannaschii* [Zarembinski *et al.* 1998]. The crystal structure had ATP bound, suggesting that the function of the protein was either an ATPase or an ATP-binding molecular switch, subsequently confirmed by experiments. The protein did have structural similarity to other proteins but in this particular case the information was not functionally useful.

are often referred to as *annotations*. This is a generic term used to describe additional information that is attached to the sequence or structure, such as key sequence or structural features, the location of catalytic residues or a proposed function.

Bioinformatics is a relatively new discipline that is concerned with the collection, organisation and analysis of biological data. It is beyond our scope to provide a comprehensive overview of this discipline; a few textbooks and reviews that serve this purpose are now available (see the suggestions for further reading). However, we will discuss some of the main methods that are particularly useful when trying to predict the three-dimensional structure and function of a protein. To help with this, Appendix 10.1 contains a limited selection of some of the common abbreviations and acronyms used in bioinformatics and Appendix 10.2 lists some of the most widely used databases and other resources.

In the rest of this chapter we will first introduce some of the key principles of protein structure and then discuss a number of approaches to tackling the protein folding problem. Another area that we will consider is protein folding: how a protein manages to fold into its own unique three-dimensional structure. A number of experimental and theoretical techniques can be used to investigate protein folding, which has led to a greater understanding of this phenomenon.

10.2 Some Basic Principles of Protein Structure

The first X-ray structures revealed that proteins did not adopt regular or symmetrical structures but were much more complex. However, certain structural motifs were observed to occur frequently. The most common motifs are the α -helix and the β -strand, shown in Figure 10.3. These constitute the *secondary structure* of a protein (the primary structure being the amino acid sequence and the tertiary structure the detailed three-dimensional conformation). Linus Pauling predicted that the α -helix would be a stable element of polypeptide structure well before the first protein structure was solved [Pauling *et al.* 1951]. His prediction was based upon mechanical models constructed after a careful analysis of the geometry of the peptide unit in crystal structures of small molecules and can be considered a classic example of the predictive power of molecular modelling. Another type of helix, the 3_{10} helix, is also found infrequently. The β -strands often form extended structures called β -sheets in which the strands are hydrogen-bonded to each other. In a β -sheet the strands can run in either parallel or anti-parallel directions, as shown in Figure 10.4. Secondary structural elements are connected by regions often referred to as 'loops', which adopt less regular structures. Nevertheless, common conformations can be identified in certain types of loop structure, such as conformations that are commonly adopted by the ' β -turn' regions between β -strands [Wilmot and Thornton 1988].

If we ignore the small variations in bond angles and bond lengths then the conformation of an amino acid residue in a protein or peptide can be classified according to the torsion angles of its rotatable bonds. There are three backbone torsion angles, labelled ϕ , ψ and ω (Figure 10.5). The conformations of the side chains are characterised by the torsion angles χ_1 , χ_2 , etc. The amide bond has a relatively high energy barrier for rotation away from

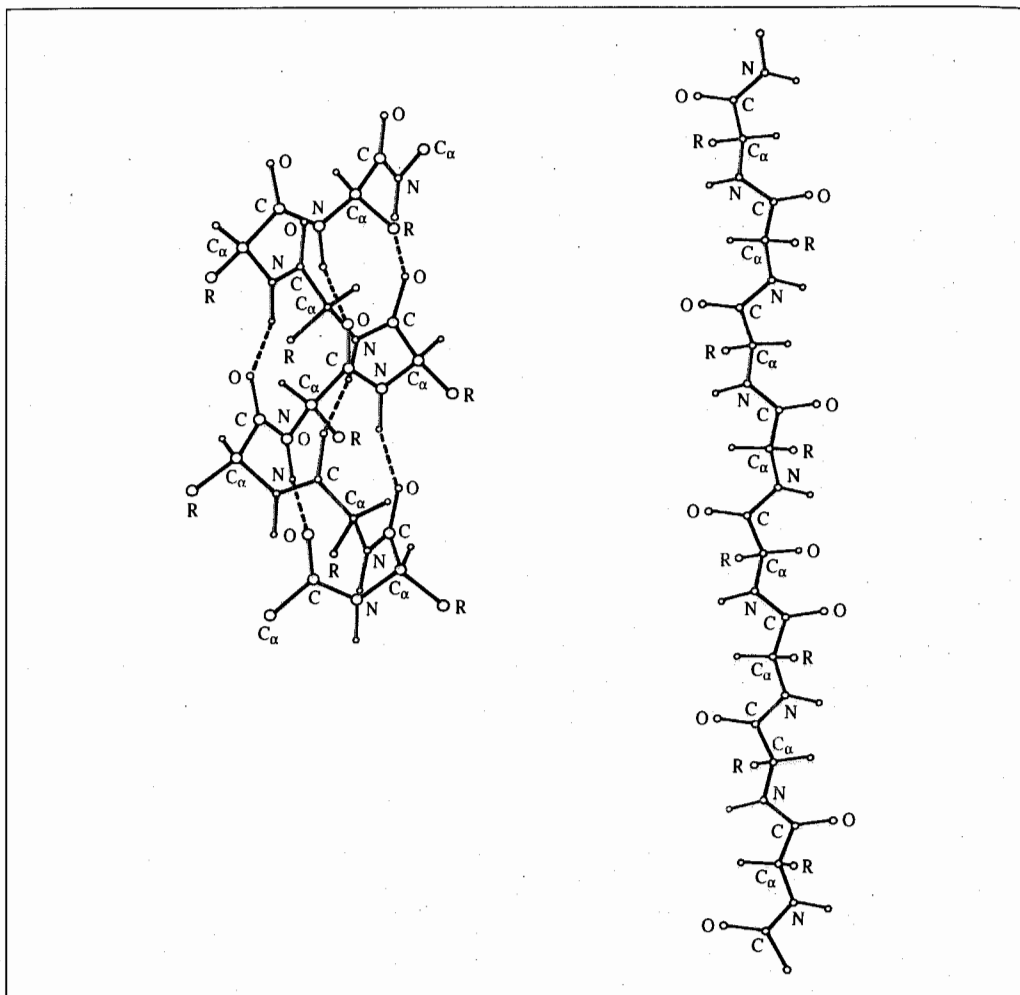


Fig. 10.3: The α -helix and β -strand structures.

planarity and so ω rarely deviates significantly from 0° or 180° . Moreover, there is a significant preference for the *trans* ($\omega = 180^\circ$) conformation (except for proline, which shows a relatively high proportion of *cis*-peptide linkages). We have already noted the contribution of Ramachandran to our understanding of protein structure in our discussion of conformational analysis (Section 9.2 and Figure 9.3). Examination of the X-ray structures of proteins shows that most amino acids occupy one of the low-energy regions in the Ramachandran contour map. Indeed, it is now common practice, when assessing an X-ray or NMR structure determination, to construct its Ramachandran map and to examine closely any residues which adopt a conformation outside the preferred regions. The side chains also tend to adopt preferred conformations, though there are many examples of unusual or higher-energy structures [Ponder and Richards 1987]. Subsequent investigations have revealed

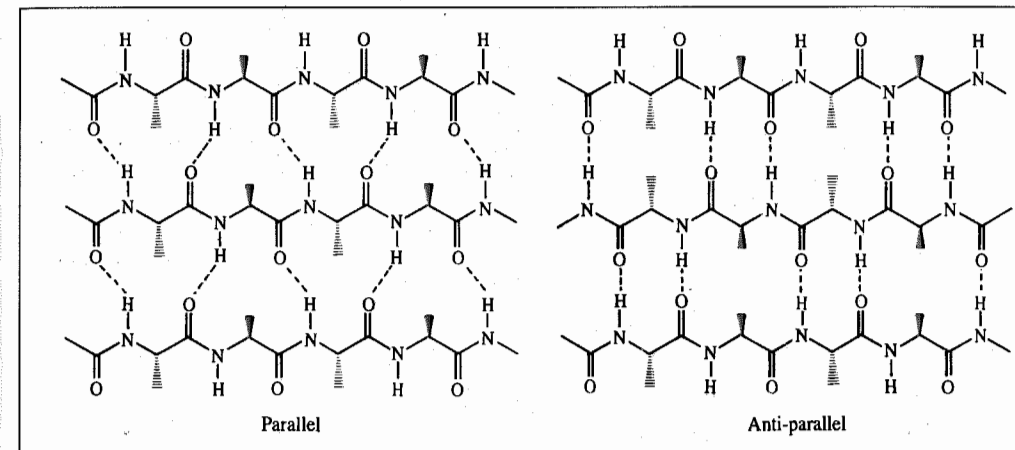


Fig. 10.4: The formation of parallel and anti-parallel β -sheets.

that the side-chain conformations are often correlated with the backbone structure; for certain conformations of the backbone only particular side-chain structures are possible [Summers *et al.* 1987; Dunbrack and Karplus 1993].

As more protein structures became available it was observed that some contained more than one distinct region, with each region often having a separate function. Each of these regions is usually known as a *domain*, a domain being defined as a polypeptide chain that can fold independently into a stable three-dimensional structure.

10.2.1 The Hydrophobic Effect

Water-soluble globular proteins usually have an interior composed almost entirely of non-polar, hydrophobic amino acids such as phenylalanine, tryptophan, valine and leucine with polar and charged amino acids such as lysine and arginine located on the surface of the molecule. This packing of hydrophobic residues is a consequence of the *hydrophobic effect*, which is the most important factor that contributes to protein stability. The molecular basis for the hydrophobic effect continues to be the subject of some debate but is generally considered to be entropic in origin. Moreover, it is the entropy change of the solvent that is

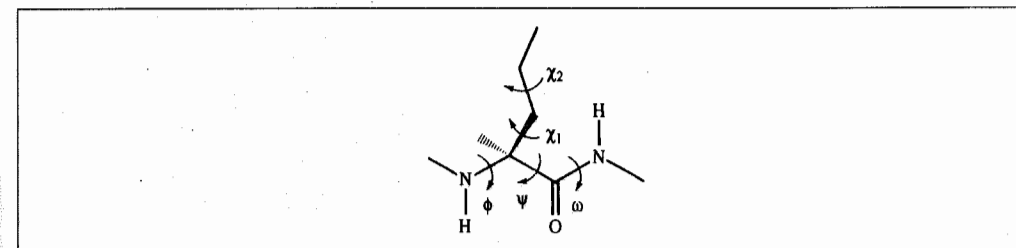


Fig. 10.5: The torsion angles that define the conformation of an amino acid.

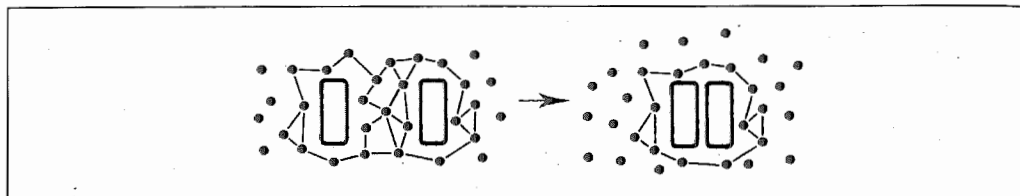


Fig. 10.6: The hydrophobic effect. Water molecules around a non-polar solute form a cage-like structure, which reduces the entropy. When two non-polar groups associate, water molecules are liberated, increasing the entropy.

important. The contribution to the overall free energy of folding due to the packing of the non-polar amino acids is positive (i.e. unfavourable) on both enthalpic and entropic grounds. The enthalpy change to remove the non-polar amino acids from water is positive due to dipole/induced-dipole electrostatic interactions between the polar water molecules and the hydrocarbon side chains. When packed, there are only (weaker) dispersion interactions between the side chains. The entropy change associated with packing the amino acids is negative because the unfolded state is less ordered than the packed state (for example, more conformational degrees of freedom are accessible). Water molecules are believed to form a cage-like structure around a non-polar solute, which has been likened to a local 'iceberg'. The water molecules in this region are locally ordered with most of the hydrogen bonding network of pure water intact. The area of the non-polar interface is much larger for the unfolded protein (Figure 10.6) and so the entropy change of the water when the protein folds is large. The enthalpy change of the water is negative as the disruption to the hydrogen bonding network is less for the folded protein. Of these four contributions, the two enthalpy terms are believed to be small, with the entropy change associated with the ordering of the solvent water molecules being the dominant term. This is just one possible model for the hydrophobic effect; unfortunately, experimental data (e.g. from X-ray crystallography or NMR) are scarce. It is also worth noting that preliminary molecular dynamics simulations were unable to find any evidence for the enhancement of water structure at a hydrophobic protein interface [Kovacs *et al.* 1997].

Not all proteins are water soluble; a very important class is the membrane-bound proteins, which include receptors and ion channels. The arrangement of the amino acids in these proteins is very different in the membrane-spanning regions. The membrane provides a very hydrophobic environment and so hydrophobic residues are often located on the outside, towards the membrane. It is very difficult to obtain X-ray crystal structures of membrane-bound proteins due to the problems of obtaining satisfactory crystals. The crystal structure of the photosynthetic reaction centre, which earned Michel, Deisenhofer and Huber the Nobel Prize in 1988, was obtained after much painstaking work in which the protein was crystallised from a detergent solution. Electron microscopy has been used to determine the structures of membrane-bound proteins; in favourable cases, the resolution of this technique approaches that of X-ray crystallography but is usually much lower. Henderson and Unwin have pioneered the application of this method to membrane proteins with their determination of the structures of bacteriorhodopsin and rhodopsin [Henderson *et al.* 1990; Havelka *et al.* 1995]. Both of these proteins contain seven *trans*-membrane helices, which are connected by loops in the extracellular and intracellular regions.

No universal solution has yet been found to the protein folding problem, but a variety of promising approaches have been developed. In the next sections we shall consider some of these methods for predicting the structures of proteins and peptides. Our discussion will first consider methods that attempt to predict the structures of proteins from first principles. We will then discuss methods that use a stepwise approach, in which elements of secondary structure are first identified and then these elements are packed together. Finally, we will consider the prediction of protein structures by homology modelling (sometimes referred to as comparative modelling), where the structure of the unknown protein is based upon the known structure(s) of related (i.e. homologous) protein(s). As part of this we will also describe some of the methods that can be used for sequence analysis.

10.3 First-principles Methods for Predicting Protein Structure

The most ambitious approaches to the protein folding problem attempt to solve it from first principles (*ab initio*). As such, the problem is to explore the conformational space of the molecule in order to identify the most appropriate structure. The total number of possible conformations is invariably very large and so it is usual to try to find only the very lowest energy structure(s). Some form of empirical force field is usually used, often augmented with a solvation term (see Section 11.12). The global minimum in the energy function is assumed to correspond to the naturally occurring structure of the molecule.

All of the conformational search methods that were described in Sections 9.2–9.7 have been used at some stage to explore the conformational space of small peptides. Here we will describe some of the methods designed specifically for tackling the problem for peptides and proteins.

H A Scheraga has devised many novel methods with his colleagues for exploring the conformational space of peptides and proteins [Scheraga 1993]. Each new method is rigorously evaluated using a standard test molecule, met-enkephalin (H–Tyr–Gly–Phe–Met–OH). One method is the 'build-up' approach, in which the peptide is constructed from three-dimensional amino acid templates [Gibson and Scheraga 1987]. Each template corresponds to a low-energy region of the Ramachandran map. To explore the conformational space of a peptide, a dipeptide fragment is first constructed by joining together all possible pairs of templates available to the first two amino acids. Each dipeptide fragment is minimised and the lowest-energy structures are retained for the next step, in which the third amino acid is connected. The peptide is gradually built up in this way, with energy minimisation and selection of the lowest-energy structures at each stage.

The simplest of Scheraga's random search methods is a random dihedral search, in which a single dihedral is selected at each iteration and randomly rotated [Li and Scheraga 1987]. The resulting structure is minimised and then accepted or rejected according to the Metropolis criterion. Such an algorithm is equally applicable to organic molecules as to proteins. The electrostatically driven Monte Carlo method [Ripoll and Scheraga 1988, 1989] is a more complex random search method which recognises the importance of long-range electrostatic interactions in polypeptides and proteins. It is based upon the observation that the local

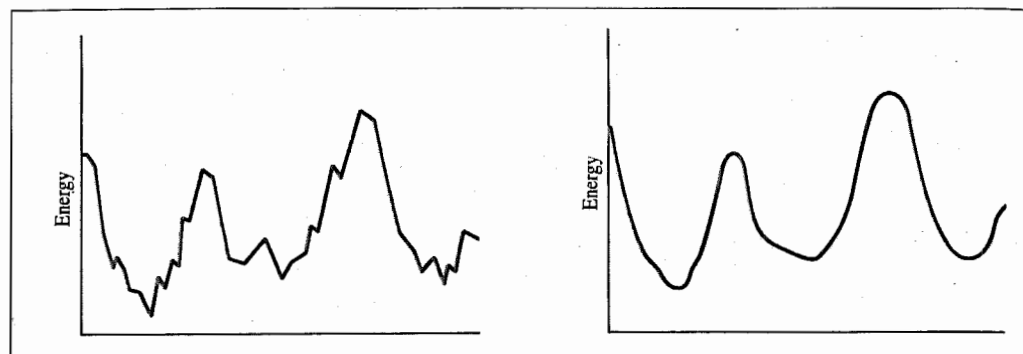


Fig. 10.7: Schematic energy surfaces for 'all-atom' (left) and simplified (right) models.

dipoles of amide units often adopt a favourable alignment in the electrostatic field of the protein. Two different types of move are thus used in the scheme. In the first type of move, an amide unit is randomly selected and the backbone (ϕ , ψ) torsion angles are changed to enable its dipole to be optimally aligned in its local electrostatic field. The resulting conformation is minimised and then accepted or rejected according to the Metropolis criterion. The second type of move involves a random change to a randomly selected dihedral followed by minimisation and acceptance or rejection in the usual way. This approach thus combines moves designed to optimise the long-range electrostatic interactions with moves that have a more local influence on the conformation.

Proteins have many more rotatable bonds than peptides, and so it is common to use some form of simplified energy model to make the problem tractable. The energy surface of a model with fewer degrees of freedom should have a smaller number of minima than the energy surface of a more detailed model; it must be assumed that the energy surface of the simplified model reproduces the general features of the more detailed representation, but without the fine structure (Figure 10.7). Various simplified models have been developed for investigating the conformational space of proteins. Many of these models are analogous to the models used to perform Monte Carlo simulations of polymers, such as the lattice and 'bead' models. An optimisation procedure based on molecular dynamics/simulated annealing or a genetic algorithm is often used with such simplified models to first identify families of low-energy structures, which may then be converted into a more detailed representation for subsequent refinement.

10.3.1 Lattice Models for Investigating Protein Structure

One reason for the interest in lattice models is that they can be used to try to answer some of the fundamental questions about protein structure. For example, it may be feasible to enumerate all possible conformations for a chain of a given length on the lattice. From this set of states statistical mechanics can be used to derive thermodynamic properties and to investigate the relationship between the structure and the sequence. In the 'HP' model [Chan and Dill 1993], a protein is modelled as a sequence of hydrophobic (H) and

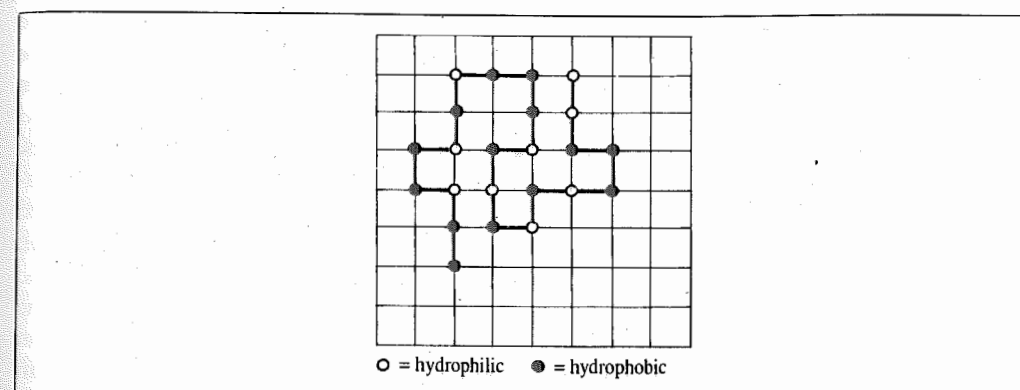


Fig. 10.8: The HP model of Chan and Dill.

hydrophilic (P) monomers. The sequence is grown onto a two-dimensional lattice using a self-avoiding walk, and the energy of the resulting conformation is calculated by summing interactions between pairs of monomers that occupy adjacent lattice sites but are not covalently bonded (Figure 10.8). Such interactions between pairs of hydrophobic monomers are favoured by a constant energy increment with all other interaction energies set to zero. Exhaustive enumeration of all the conformations was possible for chains of 30 or so monomers, and the global energy minimum for each chain was determined. Several interesting features arose from studies of this model. When the hydrophobic-hydrophobic interaction energy is small a large number of conformations are accessible. As the hydrophobic interaction energy increases there is a sharp decrease in the number of compact conformations containing hydrophobic cores. Another interesting feature of this and similar models is that α -helices and β -sheets naturally arise in the compact cores of such models. This suggests that the formation of secondary structure in a protein is not driven by specific hydrogen-bonding interactions between amino acids but rather by the compact nature of the core; conformations other than helices or sheets are not viable.

The simple lattice models are intended to address general questions about protein folding and structure. More sophisticated lattice models have been designed which are used to predict the actual structures of specific proteins. With such methods it is usually not possible to exhaustively explore the conformational space even on the lattice and so methods such as Monte Carlo simulated annealing are used to generate low-energy structures. Skolnick has developed several lattice models, which are used in a three-stage procedure to construct a model of the protein [Godzik *et al.* 1993; Skolnick *et al.* 1997]. In the first stage, a 'coarse' lattice is used in which five different types of move are permitted, excluded volume effects due to side-chain packing are taken into account and the interaction energy model contains a total of seven terms. A set of low-energy structures is obtained using Monte Carlo simulated annealing; these are then refined using a finer lattice model. This second model is closer to the actual structures of proteins and uses a more accurate representation of the side chains. The conformations obtained from the finer lattice model may then be converted to a full atomic model for refinement using a technique such as energy minimisation with a standard force field.

Some simplified models of proteins represent each amino acid residue as one or more 'pseudo-atoms', according to the size and chemical nature of the amino acid. These models are analogous to the 'bead' polymer models and full conformational freedom is possible. All of the standard types of calculation can be performed with these simplified models, encompassing techniques such as energy minimisation and molecular dynamics. An empirical model is used to calculate the residue-residue interaction energies. The parameters for these empirical models can be derived in a variety of ways. One option is to parametrise the simple model to reproduce the results of a more detailed, all-atom model. An early attempt to develop such a representation was made by Levitt [Levitt 1976] who used energy minimisation to predict the structures of small proteins. In this model the interaction between each pair of residues is equal to the average of the calculated interaction over all spatial orientations of the two residues. Minimisation of a polypeptide chain from an initial open structure resulted in compact conformations with the same size and shape as experimentally determined protein structures, together with features such as secondary structure and β -turns. Some of Levitt's observations are still very pertinent. In particular, he notes that the 'wrong' structure may still have a lower energy (as predicted by the energy function) than the 'correct' structure; this is also found to be the case with more complex molecular mechanics functions [Novotny *et al.* 1988].

To summarise, first-principles methods have been successfully used to predict the naturally occurring conformations of small peptides but are not yet sufficiently reliable to predict accurately the structures of proteins, though in some cases the general fold of the molecule is quite similar to the native structure. However, as we shall see later in Section 10.8 lattice models have been very useful in helping to understand protein folding.

10.3.2 Rule-based Approaches Using Secondary Structure Prediction

Most protein structures contain a significant amount of secondary structure (α -helices and β -strands). An obvious way to tackle the problem of predicting a protein's three-dimensional structure is first to determine which stretches of amino acids should adopt each type of secondary structure, and then pack these secondary structural elements together.

The first step in this procedure requires the secondary structural elements to be predicted. In other words, each amino acid must be assigned to one of three classes: α -helix, β -strand or coil (i.e. neither helix nor strand). Some approaches also predict whether an amino acid is present in a turn structure. One of the first methods for secondary structure prediction was devised by Chou and Fasman [Chou and Fasman 1978]. Theirs is a statistical method, based upon the observed propensity of each of the 20 amino acids to exist as α -helix, β -strand and coil. These propensities were originally determined by analysing 15 protein X-ray structures. The fractional occurrence of each amino acid in each of these three states was calculated, as was the fractional occurrence of the amino acid over all 15 structures. The propensity of that residue for a given type of secondary structure then equals the ratio of these two values. Each residue was also classified according to its propensity to act as an 'initiator' or as a 'breaker' of α -helices and β -strands. To predict the secondary structure, the amino acid sequence is searched for potential α -helix or β -strand initiating

residues. The helix or strand is then extended so long as the average propensity value for a window of five or six residues exceeds a threshold value. β -turns are also predicted using a statistical measure of the propensity of an amino acid to exist in such structures.

Many other methods have been proposed for predicting the secondary structure of a protein from its sequence, including approaches based upon information theory [Garnier *et al.* 1978] and neural networks [Ning and Sejnowski 1988]. However, the performance of even the best methods was often barely more than 65–70% (a success rate of 33% would be achieved purely by chance, if helix, sheet and coil structures were present in equal amounts). Moreover, some of these prediction rates are probably higher than they should be due to the use of the same protein structures to develop and evaluate the models. More recent methods for secondary structure prediction do not use just the sequence of interest but also other related sequences, in the form of a multiple sequence alignment. These related sequences can be found using a standard sequence-searching program such as BLAST (described below). The use of multiple sequences improves the performance of secondary structure prediction, because the algorithm is then able to search for a consensus over the aligned sequences rather than being misled by some chance effect if only a single sequence is analysed. Two methods that use multiple sequences are PHD [Rost and Sander 1993] (a neural network approach) and DSC [King *et al.* 1997]. Methods such as these are truly able to achieve 70% prediction accuracy using unrelated proteins for development and testing. Moreover, by combining the results from more than one method it is possible to make small improvements over any single individual method [Cuff and Barton 1999]. Nevertheless, it may be that there is an inherent upper limit to the performance of secondary structure prediction, because it only considers local interactions, neglecting interactions between amino acids that are far apart in the sequence but close in three-dimensional space.

Having predicted the secondary structural elements, it is then necessary to determine how they could pack together in order to achieve a low-energy structure [Cohen and Presnell 1996]. Cohen, Sternberg and Taylor analysed the packing of α -helices and β -sheets in a number of proteins and deduced a series of rules that could be used to derive favourable packing arrangements [Cohen *et al.* 1982]. For example, from an analysis of 18 protein structures they observed that an α -helix usually packs against a β -sheet in a parallel arrangement involving two rows of non-polar residues on the helix. These rules were then used to pack the α -helices and β -sheets into a stable core structure [Sternberg *et al.* 1982]. The number of possible packing arrangements was usually very large, but this number could be drastically reduced using two simple filters. First, there had to be a sufficient number of residues between sequential elements of secondary structure to span the distance between them, and second there should be no unfavourable interactions between the helices and sheets in the packed structure. Having generated one or more approximate structures the model was submitted to an energy refinement. The results from a secondary structure prediction can also be used as a restraint in lattice models [Ortiz *et al.* 1998].

The rule-based approach to protein structure prediction is obviously very reliant on the quality of the initial secondary structure prediction, which may not be particularly accurate. The method tends to work best if it is known to which structural class the protein belongs; this can sometimes be deduced from experimental techniques such as circular dichroism.

For example, some proteins contain only α -helices, which obviously makes the problem of predicting the secondary structure considerably easier. Overall, such approaches have had variable success at predicting protein structure. Nevertheless, secondary structure prediction is increasingly used as part of more general approaches to predicting protein structure.

10.4 Introduction to Comparative Modelling

There are striking similarities between the three-dimensional structures of some proteins. For example, the three-dimensional structures of trypsin, chymotrypsin and thrombin are shown in Figure 10.9 (colour plate section), from which it is obvious that they adopt very similar conformations. These proteins are all members of the trypsin-like serine protease family of enzymes but it is also possible for biologically unrelated proteins to show significant structural similarity. For example, many proteins have a structure consisting of eight twisted parallel β -strands arranged in a barrel-like structure with the β -strands connected by α -helices (Figure 10.10). This structure is often referred to as a 'TIM barrel' after triosephosphate isomerase, which was the first protein with this framework to have its structure determined by X-ray crystallography.

Comparative modelling* exploits the structural similarities between proteins by constructing a three-dimensional structure based upon the known structure(s) of one or more related proteins. To do this, it is necessary to decide which protein structure(s) to use as the 3D templates, and then to decide how to match the amino acids in the unknown structure with the amino acids in the known structure(s).

If the biological function of the protein is known it may be relatively straightforward to decide which protein(s) one might wish to consider to build the model. In other cases, the function of the protein may not be known, but it may be possible to deduce to which family it belongs by searching a sequence database for the presence of particular combinations of amino acids (called *motifs*) that often imply a particular function or structural feature. Sometimes the template is the protein whose sequence is the closest match for the unknown protein. Identifying and quantifying such matches is the role of sequence alignment methods.

10.5 Sequence Alignment

We have already seen that trypsin, chymotrypsin and thrombin have very similar 3D structures. If we now overlay the three-dimensional structures of these enzymes we find that identical amino acids are found at many positions in space, including the active site serine, histidine and aspartic acid residues, as shown in Figure 10.11 (colour plate section). The amino acid sequences of these proteins can be arranged into a *sequence alignment* as

* The term 'comparative' modelling is now preferred to the older name 'homology' modelling; the latter implies some similarity of function between the unknown protein and the template, but this may not necessarily be the case.

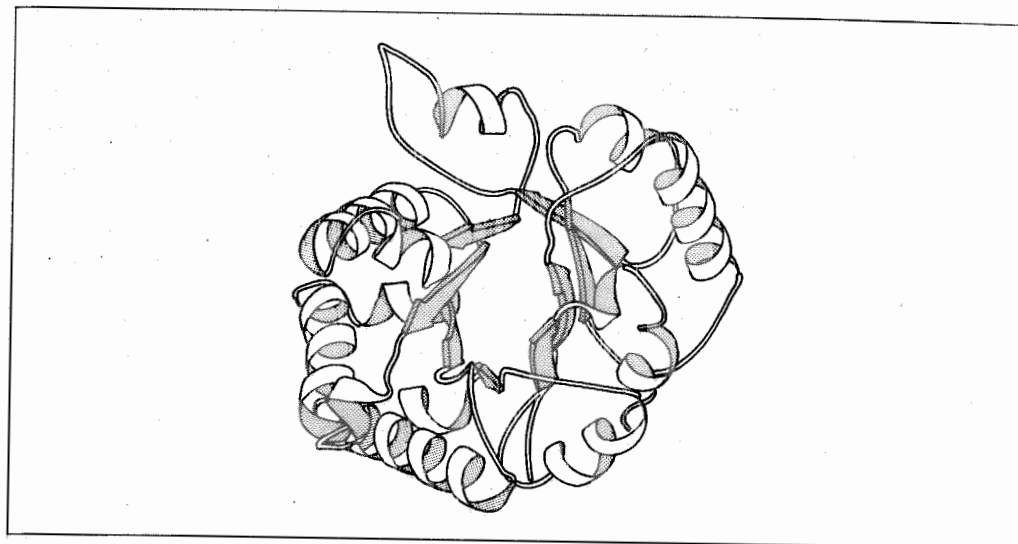


Fig. 10.10: The 'TIM barrel' [Noble et al. 1991].

shown in Figure 10.12, written using the one-letter codes for the amino acids. The relationship between sequence and structure was examined by Chothia and Lesk in 1986 [Chothia and Lesk 1986] who showed that proteins with similar sequences tend to have similar three-dimensional structures. The objective of a sequence-alignment algorithm is to position the amino acid sequences so that the matched stretches of amino acids correspond to common structural or functional features (such as the secondary structure or catalytic residues). Gaps in the aligned sequences correspond to regions where polypeptide loops are deleted or inserted. As such, sequence alignment is a key component of many procedures for predicting the structure of a new protein whose sequence has just been determined. In

Trypsin	SQWVSAAHC	YKSGIQVRLG	EDNINVVEGN	E.QFISASKS
Chymotrypsin	EDWVVTAAHC	GVTTSDVVVA	GEFDQGLETE	DTQVLKIGKV
Thrombin	DRWVLTAAHC	LLYPPWDKNF	TVDDLLVRIG	KHSRTRYERK	VEKISMLDKI
Trypsin	IVHPSYN.SN	TLNNDIMLIK	LKSAASLNSR	VASISLP...	TSCA..SAGT
Chymotrypsin	FKNPKFS.IL	TVRNDITLLK	LATPAQFSET	VSAVCLP...	SADEDFPAGM
Thrombin	YIHPRYNWKE	NLDRDIALLK	LKRPIELSDY	IHPVCLPDKQ	TAAKLLHAGF
Trypsin	QCLISGWGN.	...TKSSGT	SYPDLVKCLK	APILSDSSCK	SAYPGQITSN
Chymotrypsin	LCATFGWGK.	...TKYNAL	KTPDKLQOAT	LPIVSNNTDCR	KYWGSRVTDV
Thrombin	KGRVTGWGNR	RETWTTSVAE	VQPSVLQVVN	LPLVERPVCK	ASTRIRITDN
Trypsin	MFCAGYLEGG	...KDSCQGD	SGGPVV..CS	GK...LQGI	VSWGSGCAQK
Chymotrypsin	MICAG..ASG	...VSSCMGD	SGGPLV..CQ	KNGAWTLAGI	VSWGSSTCST
Thrombin	MFCAGYKPGE	GKRGDACEGD	SGGPFVMSKSP	YNNRWYQMG I	VSWGEGCDRD

Fig. 10.12: Sequence alignment of trypsin, chymotrypsin and thrombin (bovine). The active sites histidine, aspartic acid and serine are highlighted.

the following discussion, we will focus on the alignment of amino acid sequences but the algorithms can also be used (sometimes with small modifications) for DNA sequences.

Three general types of sequence-alignment method can be identified. Some algorithms attempt to match two sequences along their entire length. A modification of this approach is to search for local alignments involving sections (not necessarily continuous) from the sequences. The best-known examples of these first two methods were developed by Needleman and Wunsch [Needleman and Wunsch 1970] and Smith and Waterman [Smith and Waterman 1981], respectively. Both of these methods are fairly computationally intensive and not particularly suited to searching the large (and rapidly growing) sequence databases. For this purpose more approximate, heuristic methods are preferred, two major ones being BLAST and FASTA.

We will consider these algorithms in this chronological order, although in a typical comparative modelling exercise one would probably first use a heuristic algorithm to determine possible sequences of interest, then the Smith-Waterman method to identify appropriate sub-sequences, and finally the Needleman-Wunsch algorithm to derive the alignment to use in the actual construction of the model. It is always a good idea when possible to manually check any automatic alignment; the results of most automatic alignment programs can often be improved by some manual intervention. We will illustrate the various methods using the alignment of protein sequences but all of these algorithms can also be used to align nucleic acid sequences.

Any alignment algorithm requires a means for 'scoring' an arbitrary alignment of the two sequences. The objective is to find the alignment that gives the 'best' score. The simplest type of score is the *percentage sequence identity*, which gives the percentage of amino acids that are the same in the two sequences; thus identical pairs score 1 and all others score 0. An alternative approach recognises that topologically equivalent residues in two structurally homologous proteins may not be identical but often have very similar shape, electronic, hydrogen-bonding and hydrophobic properties. Such 'conservative' substitutions can frequently be made with little disruption to the three-dimensional structure of the protein and so it is desirable to take this into account in the scoring scheme. For example, in the alignment of the serine proteases position 54 is restricted to either threonine or serine due to the need to form a hydrogen bond to position 43. Dayhoff and co-workers analysed substitution frequencies in aligned sequences and have published a series of tables which give the probability of mutating one amino acid to another [Dayhoff 1978]. These probabilities are usually stored as 20×20 matrices known as PAM matrices (PAM stands for point-accepted mutation per 100 residues). One PAM corresponds to a change (on average) in 1% of all amino acid positions. The PAM concept can also be considered as a measure of 'evolutionary distance'. The best-known PAM matrix, and the one originally published by Dayhoff, is the PAM250, corresponding to 250 cycles of PAM evolution. The mutation probability matrices for evolutionary distances of 1 PAM and 250 PAM are given in Appendices 10.3 and 10.4, respectively. Each element M_{ij} of these matrices gives the probability that an amino acid in column i will have mutated to the amino acid in row j after the relevant amount of evolutionary time. It is important to realise that not every residue will necessarily have changed over this period; some will not have changed at all whereas others will have

mutated several times, possibly returning to the original state. Thus after 1 PAM there is a 0.23% probability that histidine will have mutated to glutamine whereas after 250 PAM the probability is 8%. The PAM250 matrix suggests that about 20% of the amino acids are the same after this period of evolution, with 55% of the tryptophan residues being unchanged but only 7% of the methionine residues. The matrix for any evolutionary period can be determined simply by multiplying the basic single PAM matrix an appropriate number of times. It is common to encounter a PAM matrix in a symmetrical 'log-odds' form, as shown in Figure 10.13. Each element S_{ij} of the log-odds matrix is obtained from the basic matrix by dividing M_{ij} by the relative frequency of occurrence of the amino acid i and then taking logarithms. Each element in the log-odds matrix thus represents the probability of amino acid replacement per occurrence of amino acid i per occurrence of amino acid j . An amino acid pair with S_{ij} greater than zero replaces each other more often (i.e. are likely mutations) than would be the case for random sequences of the same composition, whereas

	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val	
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
Ala	A	2																			
Arg	R	-2	6																		
Asn	N	0	0	2																	
Asp	D	0	-1	2	4																
Cys	C	-2	-4	-4	-5	4															
Gln	Q	0	1	1	2	-5	4														
Glu	E	0	-1	1	3	-5	2	4													
Gly	G	1	-3	0	1	-3	-1	0	5												
His	H	-1	2	2	1	-3	3	1	-2	6											
Ile	I	-1	-2	-2	-2	-2	-2	-3	-2	5											
Leu	L	-2	-3	-3	-4	-6	-2	-3	-4	2	6										
Lys	K	-1	3	1	0	-5	1	0	-2	0	-2	-3	5								
Met	M	-1	0	-2	-3	-5	-1	-2	-3	-2	2	4	0	6							
Phe	F	-4	-4	-4	-6	-4	-5	-5	-5	-2	1	2	-5	0	9						
Pro	P	1	0	-1	-1	-3	0	-1	-1	0	-2	-3	-1	-2	-5	6					
Ser	S	1	0	1	0	0	-1	0	1	-1	-1	-3	0	-2	-3	1	3				
Thr	T	1	-1	0	0	-2	-1	0	0	-1	0	-2	0	-1	-2	0	1	3			
Trp	W	-6	2	-4	-7	-8	-5	-7	-7	-3	-5	-2	-3	-4	0	-6	-2	-5	17		
Tyr	Y	-3	-4	-2	-4	0	-4	-4	-5	0	-1	-1	-4	-2	7	-5	-3	-3	0	10	
Val	V	0	-2	-2	-2	-2	-2	-2	-1	-2	4	2	-2	2	-1	-1	-1	0	-6	-2	4

Fig. 10.13: The PAM250 scoring matrix in the log-odds form [Dayhoff 1978]. Each element is given by $S_{ij} = 10(\log_{10} M_{ij}/f_i)$, where M_{ij} is the appropriate element of the mutation probability matrix (Appendix 10.4) and f_i is the frequency of occurrence of amino acid i (i.e. the probability that i will occur in a sequence by chance). The numbers are rounded to the nearest integer.

a pair with S_{ij} less than zero replaces each other less often (i.e. are less likely mutations). With the log-odds matrix the probabilities can be summed when comparing sequences without having to multiply them. The lower PAM matrices tend to score very similar sequences highly, whereas the higher PAM matrices can be used to find more distant relationships. Indeed, it is sometimes suggested that combinations of PAM matrices should be used to cover both possibilities. The values in the PAM matrices were obtained by considering a small number of closely related sequences and counting the observed amino acid substitutions. The BLOSUM matrices ('block substitution matrix') are obtained in a similar fashion but are often considered superior as they were derived from analyses on sequences less similar than for the PAM matrices [Henikoff and Henikoff 1992]. More recently, 'definitive' mutation matrices were obtained using a computationally elegant procedure that enabled an exhaustive match of an entire protein sequence database to be performed [Gonnet *et al.* 1992].

10.5.1 Dynamic Programming and the Needleman-Wunsch Algorithm

The Needleman-Wunsch algorithm is widely used for aligning pairs of sequences; this algorithm guarantees to find the optimal alignment based upon the scoring matrix used [Needleman and Wunsch 1970]. The algorithm uses *dynamic programming*, which forms the basis for a number of widely used methods in bioinformatics. Sequence alignment is a 'hard' problem, because there are an extremely large number of possible solutions (of the order of 10^{30} for two sequences of length 100). Here we describe the basic algorithm as it is commonly implemented today; this is equivalent to but not exactly the same as the original Needleman-Wunsch approach.

A matrix H is constructed with M rows and N columns to represent the M amino acids of protein A and the N amino acids of protein B. The elements of this matrix are filled in a sequential manner. Each element $H_{i,j}$ of this matrix corresponds to the optimal score for aligning two sub-sequences, $1 \dots i$ from the first sequence and $1 \dots j$ from the second ($1 \leq i \leq M$, $1 \leq j \leq N$). The algorithm works from the 'top left' to the 'bottom right' of the matrix (the original Needleman-Wunsch algorithm works in the reverse direction but gives the same result). The value assigned to each matrix element $H_{i,j}$ is determined from the three preceding elements $H_{i-1,j-1}$, $H_{i-1,j}$ and $H_{i,j-1}$ to the north-west, north and west, respectively, according to the following formula:

$$H_{i,j} = \max \begin{cases} H_{i-1,j-1} + w_{A_i,B_j} \\ H_{i-1,j} + w_{A_i,\Delta} \\ H_{i,j-1} + w_{\Delta,B_j} \end{cases} \quad (10.1)$$

These three moves to the point i, j are illustrated in Figure 10.14 and correspond to a match, a gap in sequence B and a gap in sequence A, respectively. The symbol Δ in Equation (10.1) is used to represent a gap. w_{A_i,B_j} represents the score associated with aligning residue A_i with residue B_j . In the simplest identity scoring scheme w_{A_i,B_j} would equal 1 if the residues were identical and zero otherwise. More typical would be the use of the PAM or BLOSUM scoring matrices. The remaining two scores, $w_{A_i,\Delta}$ and w_{Δ,B_j} , are *gap penalties*. The simplest scheme is

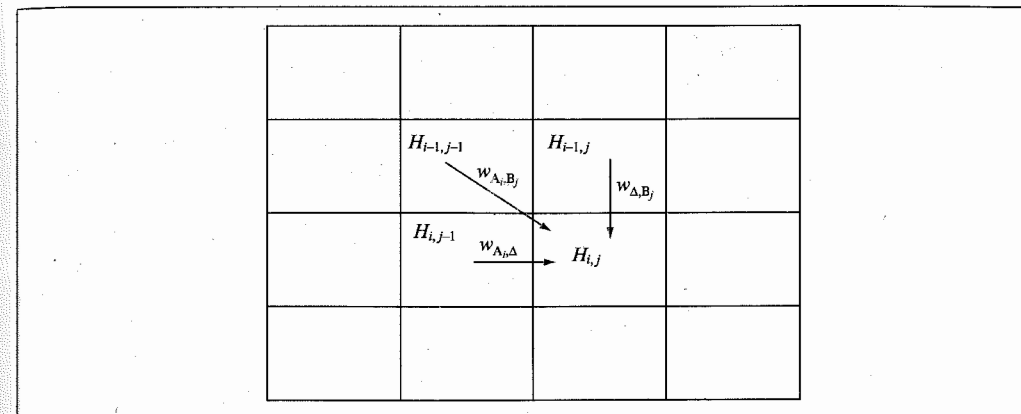


Fig. 10.14: The three moves used in dynamic programming to update the matrix element H_{ij} .

to use no gap penalty. This is illustrated for two stretches of polypeptide with sequences AECENRCKCRDP (A) and AVCNERCKLCKPM (B). Sequence A thus has 12 residues and sequence B has 13. The matrix H is shown in Figure 10.15; as can be seen it is usual to introduce a 'zeroth' row and column. For each of these outer elements of the matrix there is only one predecessor matrix element, which corresponds to matching each residue with a gap. The algorithm starts at $H_{0,0}$ and fills up the matrix one row at a time. For our simple scoring scheme, which uses sequence identity as the scoring scheme and no gap penalty, these outer elements are all zero. Let us consider the element $H_{6,6}$. This corresponds to matching an arginine from sequence A with another arginine from sequence B. The three preceding elements from the matrix are $H_{5,5}$, $H_{5,6}$ and $H_{6,5}$, which all have values of 3. The value of $w_{6,6}$ is 1 (as the two residues are identical) and so the scores corresponding to the three possible moves in Equation (10.1) are 4, 3 and 3. The value of $H_{6,6}$ is thus set to 4. Now consider $H_{10,10}$. $w_{10,10}$ is zero (arginine from A and cysteine from B). The values of the relevant elements $H_{10,9}$, $H_{9,9}$ and $H_{9,10}$ are 6, 6 and 7, respectively, and so in this case the maximum score (7) derives from a vertical move. It is sometimes possible for more than one move to give the same score. An example of this is $H_{5,5}$ (asparagine in A and glutamic acid in B). As the two residues are not identical w_{A_i,B_j} is zero. The scores for the three types of move are 2 (diagonal, from $H_{4,4}$), 3 (from $H_{4,5}$) and 3 (from $H_{5,4}$) and so the value assigned to $H_{5,5}$ is also 3.

Having completed the scoring matrix the overall score for matching the two sequences corresponds to the value of the final matrix element, $H_{M,N}$. In our example this overall score is 8. To determine the actual alignment it is necessary to trace back through the matrix. This can be achieved by storing for each matrix element $H_{i,j}$ which of the three possibilities gave the maximum value. As we have seen, in some cases a tie results and so alternative alignments may have the same score. This is the case for our sequences, due to the presence of a tie at element $H_{5,5}$. These are shown in Figure 10.15.

When no gap penalty is used then the Needleman-Wunsch alignment may contain a large, unrealistic number of gaps. The simplest type of gap penalty (other than not to have one) is to use a length-dependent scheme in which one assigns a fixed negative value for each insertion

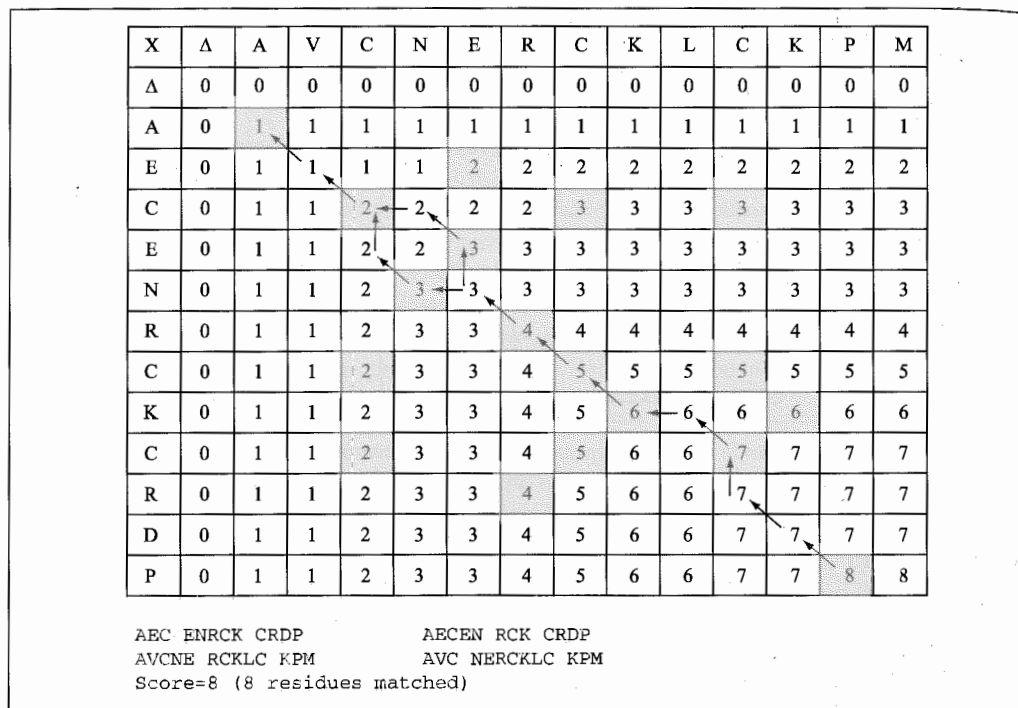


Fig. 10.15: Finding the optimal sequence alignment using dynamic programming with an identity scoring scheme and no gap penalty. Sequence A = AECENRCKCRDP; sequence B = AVCNERCKLCKPM. The value in each matrix element corresponds to the optimal score for the appropriate pair of sub-sequences. As can be seen, there are two alignments each with a score of 8.

or deletion (collectively known as *indels*). For example, if we introduce a gap penalty of -2 and make the score for a non-identical pair -1 then the dynamic programming matrix for our example changes to that shown in Figure 10.16. The alignment generated in this case happens to have no gaps (except at the end). Such an alignment is characterised by a straight diagonal.

The scoring scheme in this second example is somewhat artificial, designed primarily to illustrate the effect that a gap penalty can have on the alignment. More sophisticated types of gap penalty are possible. With the length-dependent scheme two isolated gaps contribute the same score as two consecutive gaps. Most dynamic programming methods permit a gap penalty of the form $v + uk$, where k is the gap length, v is the gap opening penalty and u is the gap extension penalty. It is most common to have a larger gap opening penalty and a smaller gap extension penalty. Yet more sophisticated are the position-specific penalty schemes. For example, if the 3D structure of one or both sequences is known then further improvements can be obtained by penalising even more severely gaps that occur in an α -helix or β -strand. The scoring matrix can also be modified to use position-specific weights. This would for example lead to an increase in the weight for aligning a residue known to be in the active site with a residue of the same type or reducing the penalty for gaps in solvent-exposed, peripheral regions.

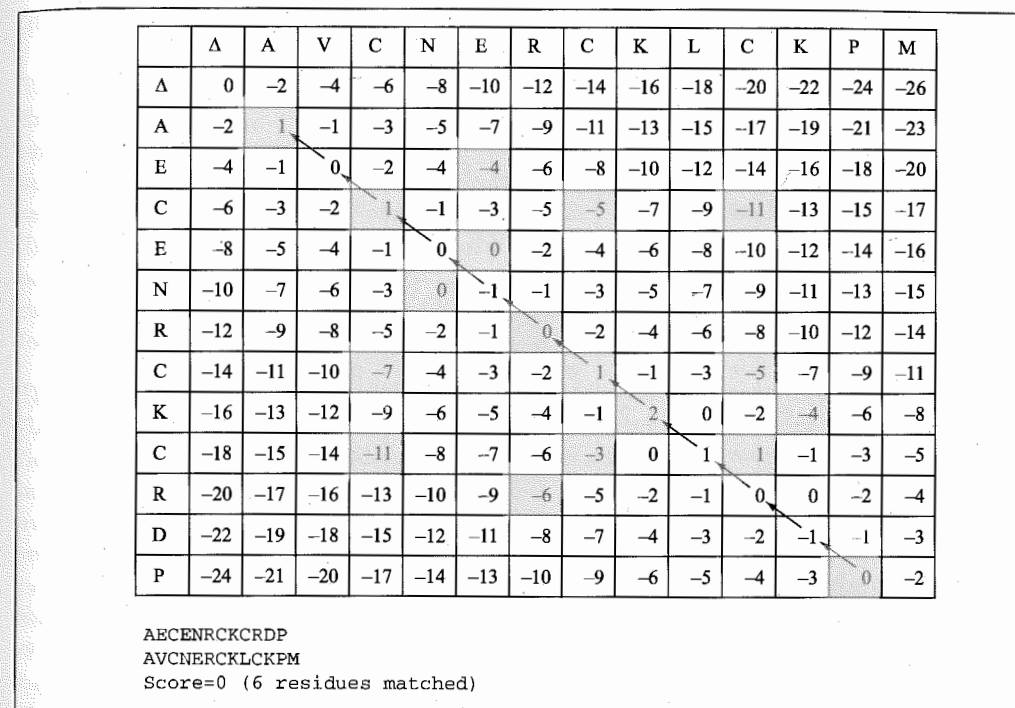


Fig. 10.16: Finding the optimal sequence alignment using dynamic programming with a scoring scheme in which a match scores 1, a mismatch scores -1 and the gap penalty is -2 .

10.5.2 The Smith–Waterman Algorithm

The Needleman–Wunsch algorithm finds a global alignment of the two sequences. This is appropriate for two sequences that are known to be similar over their whole lengths. However, it is quite common for sequences to show just local regions of similarity which would otherwise be missed in a global alignment. This can occur, for example, in multi-domain proteins, which contain a number of distinct folded sequences, each with a separate function. Even if our unknown sequence were homologous to one of the domains a global alignment against the multi-domain sequence might fail to correctly identify a match. The Smith–Waterman algorithm [Smith and Waterman 1981] is essentially the same as the method that we have described so far, except that a zero is added to the recurrence equation to give:

$$H_{i,j} = \max \left\{ \begin{array}{l} H_{i-1,j-1} + w_{A_i,B_j} \\ H_{i-1,j} + w_{A_i,\Delta} \\ H_{i,j-1} + w_{\Delta,B_j} \\ 0 \end{array} \right\} \quad (10.2)$$

The zero prevents negative similarity. The pair of segments with the maximum similarity is found by first locating the matrix element with the maximum value of $H_{i,j}$ and then tracing

	Δ	A	V	C	N	E	R	C	K	L	C	K	P	M
Δ	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	1	0	0	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	1	0	0	0	0	0	0	0	0
C	0	0	0	1	0	0	0	1	0	0	1	0	0	0
E	0	0	0	0	0	1	0	0	0	0	0	0	0	0
N	0	0	0	0	1	0	0	0	0	0	0	0	0	0
R	0	0	0	0	0	0	1	0	0	0	0	0	0	0
C	0	0	0	1	0	0	0	2	0	0	1	0	0	0
K	0	0	0	0	0	0	0	0	3	1	0	2	0	0
C	0	0	0	1	0	0	0	1	1	2	2	0	1	0
R	0	0	0	0	0	0	1	0	0	0	1	1	0	0
D	0	0	0	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0	0	0	0	0	0	0	0	0	0	1	0

Fig. 10.17: Finding the optimal local sequence alignment using the Smith–Waterman algorithm with a scoring scheme in which a match scores 1, a mismatch scores -1 and the gap penalty is -2 . The algorithm identifies the conserved RCK motif.

back in the same way as before, ending with an element equal to zero. The next-best pair of segments can be found by tracing back from the second-largest element of H not associated with the first traceback. The Smith–Waterman matrix for our two sequences (scoring 1 for a match, -1 for a mismatch and -2 for a gap) is shown in Figure 10.17, from which it can be seen that the algorithm identifies the conserved RCK motif in the middle of the two sequences.

Many variants on the basic dynamic programming method are possible, some of which we have already discussed. Other modifications of a more practical nature can increase the speed and decrease the memory requirements of the procedure.

It is important to consider the significance of an alignment. A common way to quantify this is to compare the score for a given global or local alignment to the distribution of scores obtained from aligning pairs of random sequences of the same length and amino acid composition. Such a distribution can be obtained by generating a suitably large number of random sequences, calculating their alignment scores and then determining the mean and standard deviation. The ‘true’ alignment score is then expressed as the number of standard deviation units it is above the mean of the random distribution. These scores are referred to as SD scores or Z scores; for proteins with 100–200 amino acids a value above 15 corresponds to an almost ideal alignment, whereas scores below 5 should be treated with caution. However, the distribution of scores is often skewed and it is usually possible to identify high-scoring sequences that actually have no structural similarity to the target. This skewed distribution has important consequences, which will be discussed further in the next section.

10.5.3 Heuristic Search Methods: FASTA and BLAST

The dynamic programming methods for global or local sequence alignment guarantee to find the optimal solution and can be efficiently implemented so that they find all alignments within some cutoff score. However, such methods can take a significant amount of time to search a large sequence database. This is increasingly important due to the growth of the sequence databases. Heuristic alignment methods were developed to tackle this problem. They do not guarantee always to find the globally optimal solution but in practice they rarely miss a particularly significant match. They generally work by rapidly identifying regions of potential interest using fast look-up methods and then expand these regions locally to identify the alignment.

The FASTA algorithm [Pearson 1990; Pearson and Lipman 1988] (and its predecessor, FASTP [Lipman and Pearson 1985]) uses a look-up table in the initial step, which involves the identification of all exact matches of length k (known as a k -tuple, abbreviated to $ktup$). For amino acid sequences there are 20^k possible k -tuples (4^k possibilities for DNA). The location of every k -tuple within both the query and target sequences is stored in an array of length 20^k , which thus immediately enables all matches of the length k between the two sequences to be determined. The next step is to merge pairs of k -tuples that are present on the same diagonal of the pairwise sequence matrix. In some cases, there may be a continuous section of matching k -tuples along a diagonal; in other cases, there may be gaps between the k -tuples. A simple formula is used to determine which of these diagonal regions has the most significant number of k -tuple matches. These diagonal regions are scored using a scoring matrix (e.g. PAM250) and the top-scoring regions retained. It may then be possible to join together some of these regions in order to give a longer alignment via the introduction of a joining penalty (similar to a gap penalty). Finally, an optimised alignment can be recalculated centred on this highest-scoring initial region. The result of this alignment is reported as the overall score for that particular database sequence. The highest-scoring sequences from the database are then taken and dynamic programming is used to optimise the alignment, restricting the range of the dynamic programming search to a narrow band centred on the top-scoring diagonals. This four-stage process is illustrated in Figure 10.18.

FASTA only reports the single best local alignment of the query against each of the database sequences. This can mean that a region of high similarity (but which is biologically irrelevant) may mask a lower-scoring but more significant region. The $ktup$ parameter is used to vary the speed and sensitivity of the search. For protein databases the standard value is $ktup = 2$, so only alignments where there is a pair of identical residues between the query and the database sequence are examined.

The BLAST program also searches for regions of local similarity but in its original form did not consider gaps (BLAST stands for Basic Local Alignment Search Tool [Altschul *et al.* 1990]). Given a query sequence and a database sequence the algorithm first finds all segment pairs of some length w (typically 3 for proteins) that have a score greater than some threshold (T) when using a suitable scoring matrix such as one of the PAM matrices. Each hit is then extended in both directions to check whether it lies within a longer alignment (called a maximal segment pair, or MSP) that has a ‘significant’ score. BLAST uses a parameter X

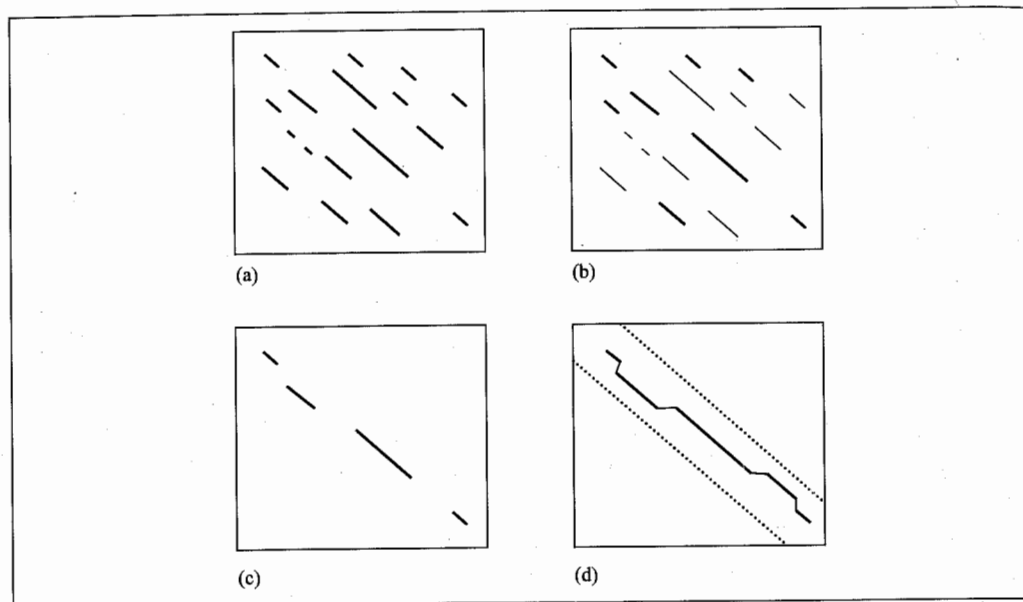


Fig. 10.18: Operation of the FASTA algorithm. (a) Locate regions of identity, (b) scan these regions using a scoring matrix and save the best ones, (c) optimally join initial regions to give a single alignment, (d) recalculate an optimised alignment centred around the highest scoring initial region. (Figure adapted from Pearson WR and D J Lipman 1988. *Improved tools for biological sequence comparison*. Proceedings of the National Academy of Sciences USA 85:2444–2448.)

to determine how long an extension will be attempted to raise the score above the required threshold score, S . The threshold score corresponds to the highest MSP score at which chance similarities are likely to appear (though in some cases the parameter X is used to reject segments that score more than X below the best found so far). The program can thus return one or more sets of local alignments that exceed the score S . The performance of the algorithm is dependent upon the values of the initial threshold, T , and the parameter X . A lower value of T reduces the possibility of missing MSPs at the expense of increasing the number of hits that proceed to the second, extension stage.

The threshold scores S used by BLAST are derived from the statistical analysis of a simple model in which the amino acids occur randomly with a probability P . The MSP scores obtained for pairs of random sequences do not follow a normal distribution (i.e. one which is symmetrical about the mean) but are skewed (formally known as an *extreme value distribution*, Figure 10.19). It turns out that the number of locally optimal segment pairs with a score of at least x is approximately distributed according to the Poisson distribution according to $KMN \exp[-\lambda x]$ where K and λ are constants that can be determined from the amino acid probabilities and the scoring matrix, and M and N are the lengths of the two sequences. This leads to the notion of a *p value*, which is the probability that a particular segment pair would occur by chance. A normalised score which enables all scoring schemes to be directly compared has also been defined. The normalised score S' is related to the basic

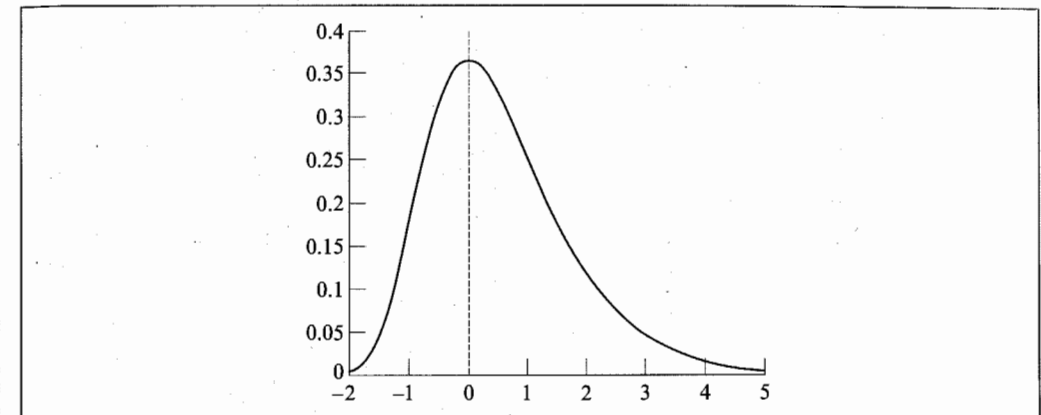


Fig. 10.19: The probability density of the extreme value distribution typical of the MSP scores for random sequences. The probability that a random variable with this distribution has a score of at least x is given by $1 - \exp[-e^{-\lambda(x-u)}]$, where u is the characteristic value and λ is the decay constant. The figure shows the probability density function (which corresponds to the function's first derivative) for $u = 0$ and $\lambda = 1$.

threshold score S by:

$$S' = \frac{\lambda S - \ln K}{\ln 2} \quad (10.3)$$

The number of distinct MSPs with a score of at least S' that are expected to occur by chance is represented by E , where E is:

$$E = \frac{MN}{2^{S'}} \quad (10.4)$$

Here, M is the length of the query sequence and N is the total length of the comparison sequence (which, for a database search, is obtained by adding the lengths of all the sequences in the database). The smaller the E value the more significant the match. Thus as the length of the query sequence or the size of the database increases so the normalised score S' must also increase to maintain a given significance level.

It is also possible to extend this concept to cover the presence of more than one distinct segment pair in a pair of sequences (for example, if there are three MSPs present with scores of 40, 45 and 50 then one can calculate the probability of finding three pairs with at least a score of 40 by chance). The ability of BLAST to provide a quantitative significance of any match found is a particularly useful feature of the program, which, with its continuing development and availability, has made it the most widely used method for sequence database searching.

Gapped-BLAST and PSI-BLAST are two significant extensions to the basic BLAST algorithm [Altschul *et al.* 1997]. As we indicated above, the original BLAST method does not permit gaps to be introduced into the MSPs. This could lead to statistically significant matches being missed where the introduction of a gap could have enabled several local alignments to be combined. The ability to introduce gaps means that only one of the alignments need be found as it can then be extended to include the others. A dynamic programming method,

able to extend a central segment in both directions, is used for this. Another modification, introduced at the same time as Gapped-BLAST, was to require two nearby overlapping hits to be present along the same diagonal before extensions were performed. PSI-BLAST stands for Position-Specific Iterated BLAST. This method uses a score matrix that is sensitive to the position of the amino acid in the query sequence and not just its identity. In this case, an iterative procedure is used whereby the significant alignments found from the first BLAST run are used to define a position-specific score matrix for the second run, and so on. PSI-BLAST is much more sensitive than BLAST. For example, it was able to detect the similarity between histidine triad proteins and galactose-1-phosphate uridylyl-transferase proteins with E values of less than 10^{-4} whereas a BLAST search could not determine these relationships with an E value threshold as high as 0.01 [Altschul *et al.* 1997]. These particular families of proteins had previously been identified as having a possible evolutionary relationship from a comparison of the three-dimensional structures; the BLAST family of programs (as with all the sequence-alignment methods we have considered) only work with the 'one-dimensional' sequence. When compared with the rigorous Smith-Waterman method, Gapped-BLAST missed eight of the 1739 significant similarities found by the dynamic programming method but ran 100 times faster. PSI-BLAST ran 40 times faster than the Smith-Waterman method and found all of the matches, together with many others. However, unsupervised use can sometimes be prone to introduce errors which may propagate over subsequent cycles.

10.5.4 Multiple Sequence Alignment

Simply put, a multiple sequence alignment is an alignment containing more than two sequences. If we know the sequences of other proteins that are suspected to be in the same family then it is usually preferable to create such an alignment. A multiple sequence alignment is often more reliable than a pairwise alignment as it is easier to detect any trends; with just two sequences it is easy to be misled by some chance correspondence. The alignment of the serine proteases in Figure 10.12 is an example of a multiple sequence alignment. Information from multiple sequences is also valuable even at the sequence-matching level, as the results with PSI-BLAST indicate (it uses the multiple hits to define the score matrix for the next iteration).

Perhaps the most obvious way to generate a multiple sequence alignment would be to extend the basic Needleman-Wunsch dynamic programming method to cover more than two sequences. Such a generalisation from two to N sequences is possible, though in practice the method is limited to comparisons of three sequences. With some approximations (such as the use of a 'window' centred on the diagonal to restrict which elements of the H matrix are considered) it is possible to increase the number of sequences that can be globally aligned. However, the most common approach to multiple sequence alignment is to use some form of hierarchical clustering approach. First, all pairwise sequence alignments are generated. A hierarchical cluster analysis (see Section 9.13) is then used to group the most similar pair of sequences, then the next most similar pair, and so on. Such an approach needs not only to align two sequences together but also to align one sequence with an alignment and one alignment with another alignment (an alignment contains two or more

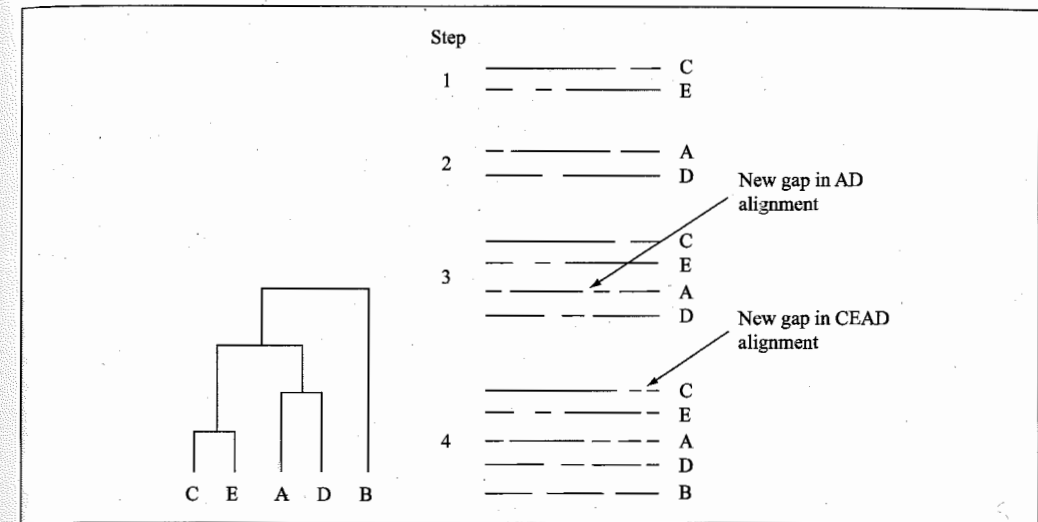


Fig. 10.20: Schematic illustration of the creation of a multiple sequence alignment for five sequences A-E. In the first step sequences C and E are aligned. In the second step sequences A and D are aligned. In the third step the pair CE is aligned with the pair AD. Finally, the quartet CEAD is aligned with B.

individual sequences). Position-specific scoring matrices (also known as *profiles*) are used during these steps; these can be obtained by averaging the substitution values for the amino acids at a particular location. For example, suppose two sequences A and B are aligned in the first step and it is then desired to align a third sequence C with the alignment AB. If at some position A contains serine and B contains threonine then the score for matching (say) an alanine residue from C at this position would be the average score for alanine/serine and alanine/threonine. Once formed, gaps are usually maintained. Moreover, it may be necessary to introduce additional gaps into an alignment in order to achieve an optimal match in later stages. Under such circumstances, the gap is introduced into all of the sequences that form that particular alignment. The process is illustrated schematically in Figure 10.20.

More sophisticated approaches apply weighting schemes so that very similar sequences have lower weights because they contain duplicated information. In addition, special procedures can be used to handle gaps – for example, the gap penalty can be made position-specific (e.g. lower at existing gaps, higher near existing gaps, residue-specific) [Thompson *et al.* 1994]. Nevertheless, despite these developments most automatic alignments benefit from some manual intervention. Computer graphics programs which can display all of the sequences, colour-coded by amino acid type or properties, can greatly facilitate this process.

A multiple sequence alignment can suggest whether certain residues are conserved more frequently than others, and regions where insertions and deletions are more common. This gives rise to the notion of a *profile* [Gribskov *et al.* 1987], which, as we indicated earlier, can be considered as a position-specific weighting scheme that indicates the score for matching an amino acid at a particular position, together with insertion and deletion penalties.

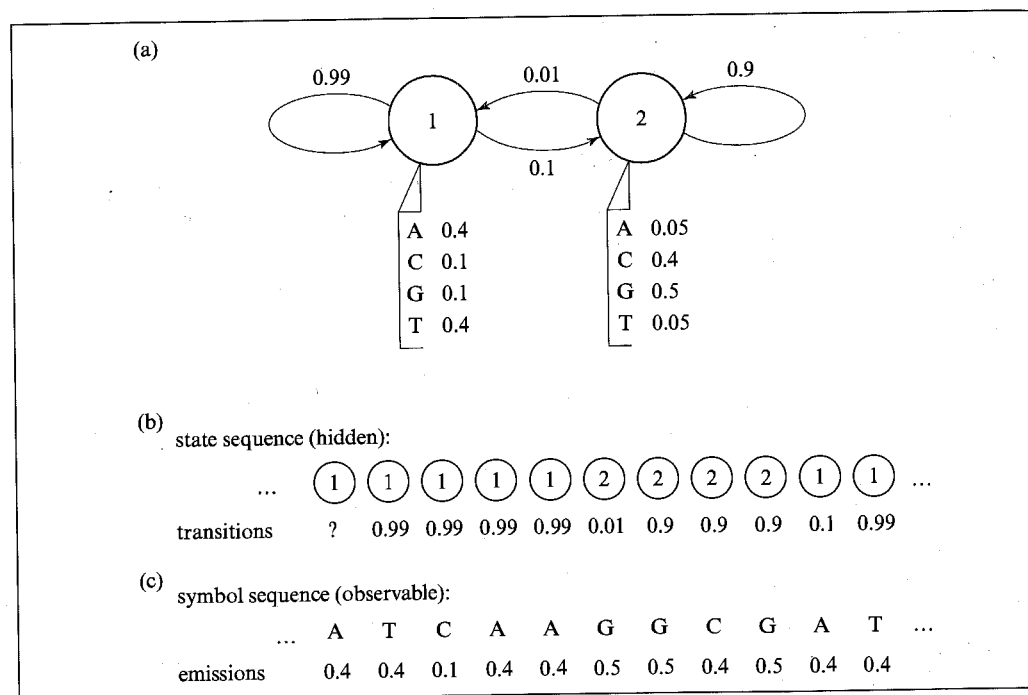


Fig. 10.21: Simple two-state hidden Markov model for the generation of DNA sequences. State 1 generates AT-rich sequences and state 2 generates GC-rich sequences according to the symbol emission probabilities. In addition, there are transition probabilities as indicated by the arrows. A sample hidden state sequence is shown together with an (observed) symbol sequence together with the associated probabilities. (Figure redrawn from Eddy S R 1996. *Hidden Markov Models*. *Current Opinion in Structural Biology* 6:361-365).

Hidden Markov models (HMMs) are a class of statistical modelling tools that can be used for multiple sequence alignment as well as other useful applications in bioinformatics. They have been extensively used in other areas of science and technology, particularly for speech recognition [Rabiner 1989]. The name derives from the fact that they are constructed from a series of 'states', each of which corresponds to the columns of a multiple alignment. These states are interconnected according to a series of transition probabilities; the choice of which state to occupy depends upon the current state and so the sequence of states is a Markov chain. They are 'hidden' because the sequence of states is not observed; rather, one observes the amino acid or nucleic acid sequence that it generates. One of the simplest HMMs with some biological relevance is shown in Figure 10.21 [Eddy 1996]; this has just two states, one of which preferentially generates AT-rich sequences and the other generates GC sequences. Having generated its symbol (according to the symbol emission probabilities) there is only a 1% probability of making a transition to the other state. This means that such a model will tend to generate sequences of either A and T or G and C with infrequent switches between the two.

For protein sequences a more complex model is used [Krogh *et al.* 1994]. First, there is a beginning and an end state with as many intervening states as there are columns in the

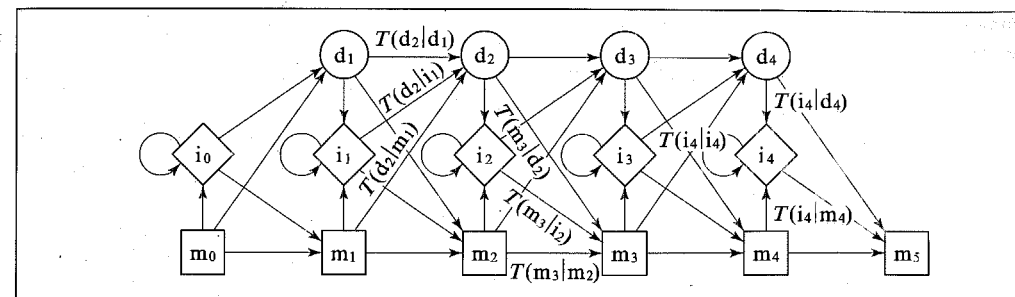


Fig. 10.22: Hidden Markov model used for protein sequence analysis. m_1 - m_4 are match states (corresponding in this case to a four-position alignment). m_0 and m_5 are the begin and end states, and i and d are the insert and delete states. There are three possible transitions from each state to other states. (Figure redrawn from Krogh A, M Brown, S Mian, K Sjölander and D Haussler 1994. *Hidden Markov Models in Computational Biology. Applications to Protein Modelling*. *Journal of Molecular Biology* 235:1501-1531).

multiple alignment. At each position there are three possibilities. Either an amino acid is generated according to the distribution of that state, or the state is skipped (a deletion) or an amino acid might be inserted. There are probabilities for moving between the various states (Figure 10.22). The symbol generation and state-transition probabilities are determined by an iterative training process. A key feature of this training phase is that it is not necessary to provide aligned data, unlike other approaches. The other key feature is that an HMM inherently contains position-specific gap penalties that are learned from the data. The HMM builds its profile during the process of actually performing the multiple alignment, rather than this being a separate task that is performed once the alignment has been generated.

Having built a hidden Markov model for a particular family of proteins, it can then be used to search a database. A score is computed for each sequence in the database and those sequences that score significantly more than other sequences of a similar length are identified. This was demonstrated for two key families of proteins, globins and kinases in the original paper [Krogh *et al.* 1994]. For the kinases, 296 sequences with a Z score above 6 were identified from the SWISSPROT database of protein sequences. Of these 296 sequences, 278 were already known to be kinases or were classified as such by a battery of other procedures and were thus considered to constitute 'certain' kinases. Of the remainder, some were considered 'false positives' (i.e. were not kinases) whilst for others no definite conclusion could be drawn. In addition, a handful of sequences below the $Z = 6$ cutoff were also identified as 'false negatives'.

10.5.5 Protein Structure Alignment and Structural Databases

Despite the great progress in sequence-alignment methods there are still some cases where the similarity can only be identified by considering the three-dimensional structures of the proteins. In such cases, it is necessary to have some means of aligning two proteins based upon structural criteria. Perhaps the most basic way to achieve such a structural comparison is using computer graphics and manual superposition. However, automated methods have also been developed, of which the so-called *double dynamic programming* approach is

particularly popular [Taylor and Orengo 1989]. The method is so named because it uses two dynamic programming steps. In the first step, it is necessary to determine the score for each pair of residues, one from each structure. These scores are used to fill a rectangular H matrix, to which dynamic programming is applied to determine the optimal alignment.

In order to determine the pairwise residue score between residue i from protein A and residue j from protein B a second rectangular matrix is used. Each element (l, m) of this matrix corresponds to a pair of residues, l from protein A and m from protein B. The matrix element (l, m) is set to the similarity value:

$$s = \frac{a}{[(^A\mathbf{V}_{il} - ^B\mathbf{V}_{jm})^2 + b]} \quad (10.5)$$

$^A\mathbf{V}_{il}$ is the 3D vector from residue i to residue l in protein A and $^B\mathbf{V}_{jm}$ is the vector from residue j to residue m in protein B (Figure 10.23). The more similar are the two vectors the greater the similarity score. a and b are constants. Having filled in the elements in the matrix corresponding to residues l and m , dynamic programming is used to find the optimal similarity S_{ij} between residues i and j . This is the value entered into the main matrix, which is used for the final dynamic programming step. Sequence information can be incorporated into the procedure by changing the numerator in Equation (10.5) from a to $(wD_{R_i R_j} + a)$, where D is one of the common scoring matrices used in normal sequence alignment for substitution of a residue R_i with residue R_j . w is a weighting factor that determines the relative contributions of structure and sequence. Other advances in the technique enable

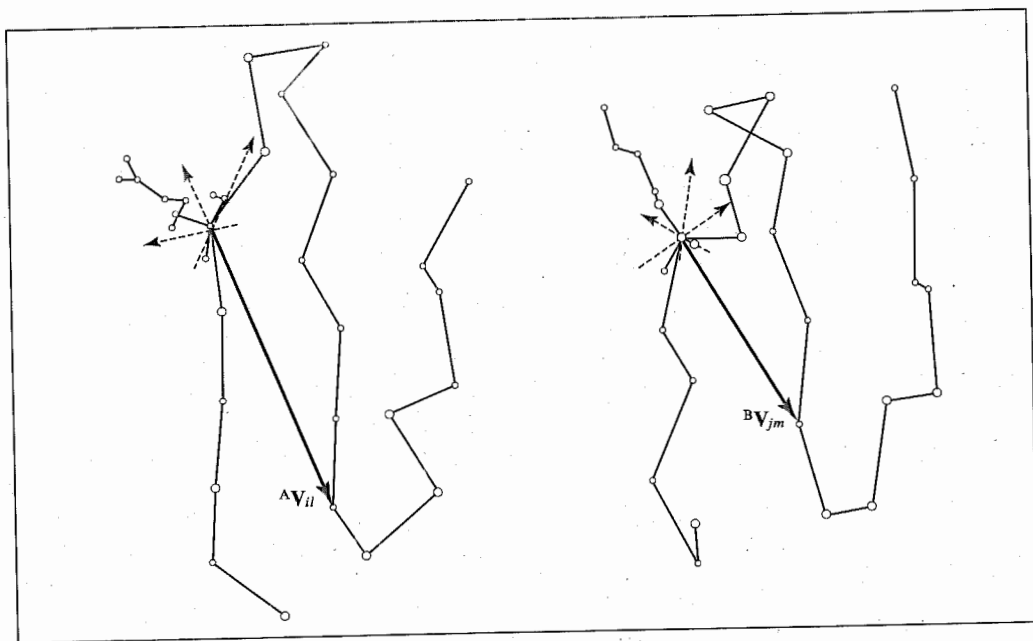


Fig. 10.23: Vectors used to calculate the 3D similarity for the second matrix used in the double dynamic programming method.

local regions of structural similarity to be identified [Orengo and Taylor 1990, 1993]. A particularly fast implementation uses an initial secondary structural filter [Orengo *et al.* 1992].

A number of structural databases have been developed to classify proteins according to their three-dimensional structures. Many of these are accessible via the World Wide Web. The protein databank (PDB [Bernstein *et al.* 1977]) is the primary source of data about the structures of biological macromolecules and contains a large number of structures, but many of these are of identical proteins (complexed with different ligands or determined at different resolutions) or are of close homologues.

The SCOP (Structural Classification of Proteins) database adopts a hierarchical approach to protein structure, with several different levels [Murzin *et al.* 1995]. Unusually, the SCOP database is constructed from a visual inspection and comparison of structures. Multi-domain proteins are split into their individual domains, which are then classified into different families, superfamilies, folds and fold classes. A *family* comprises proteins with 30% or greater sequence identity or where there is a very close match of function and structure. A *superfamily* comprises those proteins with low sequence identities but which have similar structures and functional features. A *fold* is defined as a particular set of secondary structural elements joined together in a specific topology. Finally, the *fold class* is usually one of five higher-level classifications: (a) all-alpha (contains only α -helices), (b) all-beta (only β -sheets), (c) alpha and beta (usually written α/β , where α -helices and β -strands are intermixed), (d) alpha plus beta (written $\alpha + \beta$, where the α -helices and β -strands are mostly segregated), (e) multi-domain. In 1997, there were more than 7600 entries in the PDB but after removing duplicates (the same protein from the same organism) this number falls to 1729. If just one structure from each homologous family is included (defined so that no two proteins have more than 25% sequence identity) then this gives 652 proteins with 463 superfamilies and 327 folds [Brenner *et al.* 1997]. Classifications such as these have led some to suggest that there is an upper limit to the number of different protein families, with about 1000 being a commonly suggested limit. Other databases which are based upon a structural classification include CATH [Orengo *et al.* 1993] and FSSP [Holm and Sander 1994]. The latter uses an algorithm called DALI [Holm and Sander 1993] to compare protein structures based upon a comparison of residue-residue distance matrices. FSSP includes a representative set of three-dimensional structures which are used for other applications, such as threading (see below). It is also possible to combine structural and sequence information as in the HSSP database, which uses alignment methods to identify sequences that are related to proteins of known three-dimensional structure and so by implication have the same secondary and tertiary structure [Holm and Sander 1999].

10.6 Constructing and Evaluating a Comparative Model

A sequence alignment establishes the correspondences between the amino acids in the unknown protein and the template protein (or proteins) from which it will be built. The three-dimensional structures of two or more related proteins are conveniently divided into *structurally conserved regions* (SCRs) and *structurally variable regions* (SVRs). The structurally conserved regions correspond to those stretches of maximum sequence identity or sequence

similarity where one expects the conformation of the unknown protein to be very similar to that of the template protein(s). The structurally conserved regions are often found in the core of the protein or the active site. The structurally variable regions usually correspond to polypeptide loops which connect secondary structure elements together. These loops can show significant differences in sequence and be of completely different lengths.

There are currently three different classes of method for constructing the three-dimensional model [Sali 1995]. The first method involves piecing together rigid bodies taken from the template protein(s). The second method assembles the target protein by joining together small segments or by reconstructing a set of coordinates. The third approach generates a series of spatial restraints from the templates, which are used in conjunction with an optimisation procedure to derive a structure of the target. Whilst each of these methods is in principle capable of producing a structure complete with loops and side chains, it is more common to consider the actual construction as a three-stage process. First, the amino acid backbone for the structurally conserved regions is generated. This gives the 'core' of the protein, to which the loops are then added. Finally, the side chains are placed. The model may then be subjected to some form of refinement, such as energy minimisation. Finally, the model should be validated to ensure that it conforms to a variety of rules about protein geometry derived from analyses of known protein structures.

The simplest type of rigid-body method involves simply transferring the backbone conformation of the core of the protein from a single template to the unknown protein. An alternative is to construct a framework by averaging the structures from a number of protein templates. Each template can be given a weight related to its sequence similarity to the unknown target [Srinivasan *et al.* 1996].

The segment-matching procedure starts with a basic framework, usually consisting of a set of alpha-carbon atoms. These coordinates are used to guide the fitting of the segments [Levitt 1992]. In the case of comparative modelling the initial framework would be derived from the structures of one or more homologous proteins. The conformations of the segments typically come from known protein structures, but an alternative is to use some form of geometrical algorithm to generate an energetically feasible set of atomic coordinates. The notion that the structure of a protein might be predicted by piecing together 'spare parts' from known structures was first demonstrated by Jones and Thirup [Jones and Thirup 1986]. This idea is at the heart of many methods used to construct frameworks, loops and side chains. More ambitious is the use of fragment assembly without using an underlying framework [Simons *et al.* 1997, 1999a, b]. In this case, the fragments are taken from proteins of known structure which show local sequence similarity to the unknown target. The initial structures resulting from this 'splicing' process are then subjected to simulated annealing using a scoring function that has sequence-dependent terms (representing factors such as the burial of hydrophobic residues and electrostatics) and sequence-independent terms (describing the packing of α -helices and β -strands). A number of runs are typically performed, from which the most promising are selected.

The third method, satisfaction of spatial restraints, adopts a rather different approach to the problem. One possible method that would fall into this category would be distance geometry, with the distance constraints being derived from related template structures.

An alternative is to use an optimisation procedure in Cartesian space; this is the basis for a program called Modeller [Sali and Blundell 1993]. In this method, a large number of restraints are derived. Some of these restraints come from an analysis of the sequence alignment of the target protein to homologous proteins of known structure; others are derived from a statistical analysis of the relationships between various features of protein structure. Typical features include the distribution of distances between alpha-carbon atoms, residue solvent accessibilities or side-chain torsion angles. Particularly relevant to comparative modelling are the associations between these features for two related proteins. Thus the backbone conformation of a particular residue may be restrained according to the residue type, the conformation of an equivalent residue in a related protein and the local sequence similarity between the two proteins. The restraints are expressed as *probability density functions* (pdf), each of which is a smooth function which gives the distribution of the feature as a function of the related variables. These individual probability density functions are combined to give a molecular function, which is then optimised. The optimisation uses a combination of conjugate gradients with molecular dynamics and simulated annealing. Local restraints are considered first, and then the global restraints.

For those methods which construct just a template for the structurally conserved regions the next task is to determine the conformations of the loop regions. These generally occur on the surface of the molecule. Each loop must obviously adopt a conformation that enables it to properly join together the appropriate parts of the core. The loop conformation should also have a low internal energy and not have any unfavourable interactions with the rest of the molecule. In certain cases, the loops may be restricted to a set of *canonical structures*. For example, it has been observed that the conformations of some antibody loops fall into a small number of classes [Chothia *et al.* 1989]. Similarly, the loops that connect certain types of secondary structure often show distinct conformations. The β -turns that connect strands of β -sheets have been classified into a small number of distinct families [Wilmot and Thornton 1988]. In other cases, we require an alternative method for predicting the loop conformations. Here we discuss just a few of the many methods that have been proposed for modelling polypeptide loops. These methods generally proceed by searching a database for suitable segments or by using some form of conformational search.

Loop conformations can be obtained by searching the protein databank for stretches of polypeptide chain that contain the appropriate number of amino acids and also have the correct spatial relationship between the two ends [Jones and Thirup 1986]. A test for amino acid homology may also be included in the criteria for loop selection. This procedure can be made very efficient by precalculating the necessary geometric information from loops in the protein databank and then using screening methods to identify the loops that can fit. This geometric screen uses information about interatomic distances between key atoms at the base of the loop. Loops that clash with the rest of the protein are rejected.

For loops that contain fewer than seven rotatable bonds, an algorithm devised by Go and Scheraga [Go and Scheraga 1970] can be used to calculate possible loop geometries directly. Go and Scheraga showed that it was possible to determine the torsion angles that would permit the end-to-end distance of the loop to achieve the desired value. The original Go and Scheraga method was developed for a model with fixed bond lengths and bond

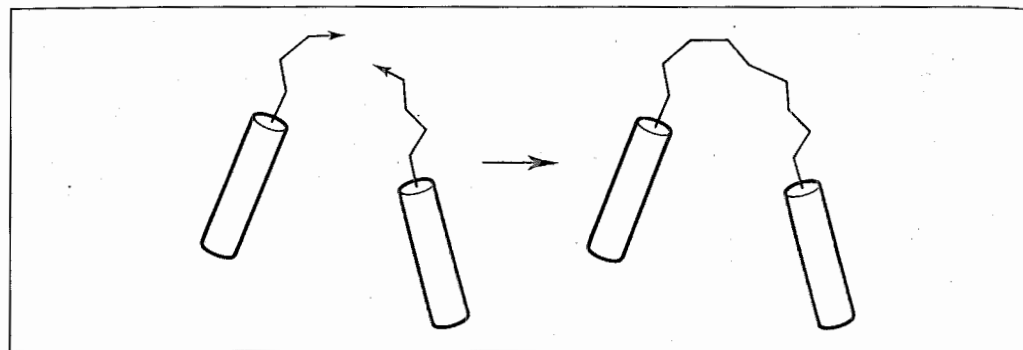


Fig. 10.24: An effective way to construct loops using a systematic search algorithm is to grow the two ends of the chain until they meet.

angles; later variants permit the bond angles to deviate slightly from their equilibrium values and so have a higher chance of finding an acceptable match [Brucoleri and Karplus 1985]. In the CONGEN program of Brucoleri and Karplus a systematic search is used to explore the space of $N - 6$ rotatable bonds (where N is the number of ϕ and ψ torsions in the loop) [Brucoleri and Karplus 1987]. For each conformation that is generated, the Go and Scheraga chain-closure algorithm is used to complete the structure. Purely systematic search methods can also be used to generate loop conformations. One interesting way to try to alleviate the combinatorial explosion is to construct the loop from both ends simultaneously; the half-complete loops are then joined in the middle (Figure 10.24).

Methods based on random algorithms have also been devised for modelling protein loops. One interesting method is the random tweak algorithm [Shenkin *et al.* 1987], which calculates the changes in the backbone ϕ and ψ torsion angles that will enable a randomly generated loop conformation to fit a set of distance constraints. An advantage of the random tweak procedure is that almost every chain can be 'tweaked' so that it satisfies these constraints; it is also extremely fast because it scales with the number of constraints rather than with the length of the chain. However, no information is included about the interactions with the rest of the protein in the calculation and this has to be checked once a loop conformation has been generated.

Having defined one or more backbone conformations for the protein, including the loop regions, it is then necessary to assign conformations to the side chains. In the core region there may be a high degree of sequence identity between the unknown protein and the template, and the side-chain conformations can often be transferred directly from the template. Changes in amino acids in the core are often very conservative (e.g. a change from a phenylalanine to a tyrosine) and it is also easy to model the side chain in such cases. Where there is less correspondence between the amino acid sequences (and especially for the loop regions) then the side chains must be added without reference to the template. A variety of systematic and random methods have been used to predict side-chain conformations; Monte Carlo, simulated annealing and genetic algorithm methods are particularly common [Vasquez 1996]. A popular tactic is to restrict the conformations of the side

chains to those that are observed in experimentally determined protein structures [Ponder and Richards 1987]; a further refinement of this approach recognises that side-chain conformations depend upon the conformation of the main chain [Dunbrack and Karplus 1993]. Side-chain prediction methods invariably keep the backbone fixed.

The initial structures obtained from a comparative modelling exercise can often be rather high in energy. Energy minimisation is thus often performed to refine the structure, though one should be careful to ensure that the minimisation does not cause any drastic changes and some practitioners deprecate its use.

Once a protein model has been constructed, it is important to examine it for flaws. Much of this analysis can be performed automatically using computer programs that examine the structure and report any significant deviations from the norm. A simple test is to generate a Ramachandran map, in order to determine whether the amino acid residues occupy the energetically favourable regions. The conformations of side chains can also be examined to identify any significant deviations from the structures commonly observed in X-ray structures. More sophisticated tests can also be performed. One popular approach is Eisenberg's '3D profiles' method [Bowie *et al.* 1991; Lüthy *et al.* 1992]. This calculates three properties for each amino acid in the proposed structure: the total surface area of the residue that is buried in the protein, the fraction of the side-chain area that is covered by polar atoms and the local secondary structure. These three parameters are then used to allocate the residue to one of eighteen environment classes. The buried surface area and fraction covered by polar atoms give six classes (Figure 10.25) for each of the three types of secondary structure (α -helix, β -sheet or coil). Each amino acid is given a score that reflects the compatibility of that amino acid for that environment, based upon a statistical analysis of known protein structures. Specifically, the score for a residue i in an environment j is calculated using:

$$\text{score} = \ln \left(\frac{P(i:j)}{P_i} \right) \quad (10.6)$$

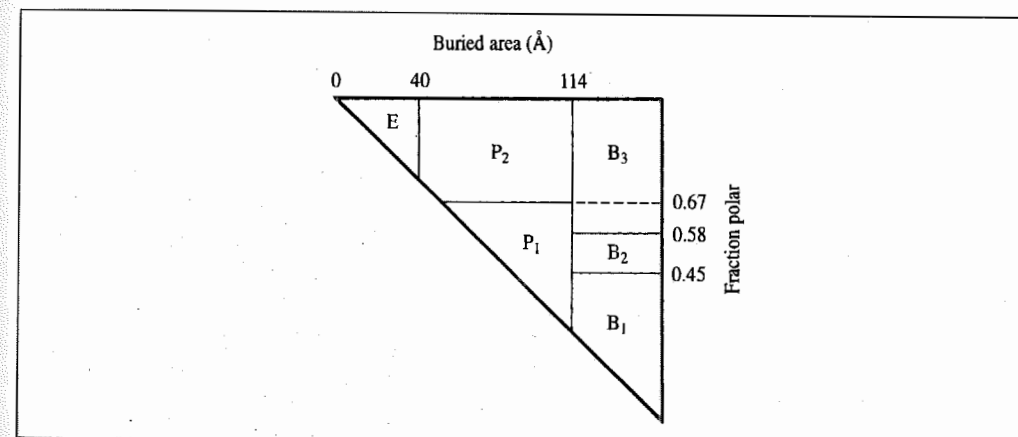


Fig. 10.25: The six environment categories used by the 3D profiles method. (Figure adapted from Bowie J U, R Lüthy and D Eisenberg 1991. A Method to Identify Protein Sequences That Fold into a Known Three-Dimensional Structure. *Science* 253:164-170.)

where $P(i:j)$ is the probability of finding residue i in environment j and P_i is the overall probability of finding residue i in any environment. For example, $P(i:j)$ is -0.45 for a valine residue in a partially buried environment with a high fraction ($>67\%$) of the surface covered by polar atoms in an α -helix. The negative number indicates that this environment is not favoured for valine. However, this environment is more favoured by arginine, for which the $P(i:j)$ value is 0.50 .

The 3D profiles method can be used to calculate an overall score for a protein model. It was found that deliberately misfolded protein models have low scores because they contain residues in environments with which they are not compatible. Such misfolded models often cannot be distinguished from the correct structures using molecular mechanics energies. The 3D profile can also be used to identify whether a generally correct model contains regions of incorrectly assigned residues. This is usually done by plotting the score as a function of the sequence, as shown in Figure 10.26. Any residues for which the score falls significantly below the average score should be investigated to check whether the model is faulty in that particular region.

Comparative modelling is a widely used method, with many models being published in the literature. Of particular importance are those papers which retrospectively compare a predicted model with the subsequent experimental structure. An early example was the comparison of a model of the aspartyl protease renin built using the Composer program [Frazao *et al.* 1994]. Renin is an important enzyme in the control of high blood pressure and so is of potential interest as a pharmaceutical target. Composer uses the rigid-body

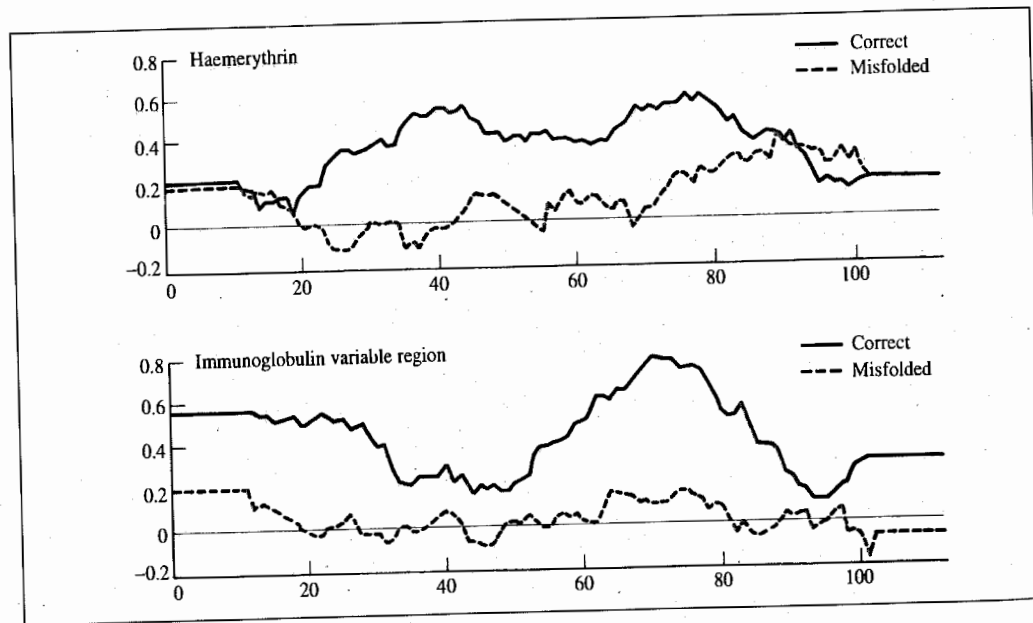


Fig. 10.26: The 3D profiles output for incorrect and partially incorrect protein models compared to the correct structures. The vertical axis gives the average profile score for a 21-residue window. (Figure redrawn from Lüthy R, J U Bowie and D Eisenberg 1992. *Assessment of Protein Models with Three-Dimensional Profiles*. Nature 356:83-85.)

approach; in this case, two homologous aspartyl proteases (pepsin and chymosin) were used as templates, followed by loop modelling using the 'spare parts' method. Side chains are assigned using a series of rules derived by examining topologically equivalent positions in homologous structures. This analysis leads to 1200 rules, one for each 20 by 20 amino acid replacement in each of the three types of secondary structure (α -helix, β -strand, neither). If no applicable rule exists then a conformation is chosen from a rotamer library. The model had an RMS fit for 280 alpha-carbon atoms of 0.84 \AA (the atom pairs were selected by applying a cutoff of 3.5 \AA). Of some interest was the fact that the model was closer to the X-ray structure than it was to either of the two structures used for its construction. However, the analysis did also highlight some areas for improvement, such as a proline-rich loop for which there were few representative examples in the database of known structures.

10.7 Predicting Protein Structures by 'Threading'

Threading (more formally known as 'fold recognition') is a method that may be used to suggest a general structure for a new protein [Jones *et al.* 1992; Jones and Thornton 1993]. The basic threading concept is very simple. Suppose we wish to predict the structure of an amino acid sequence, and that we have available a number of three-dimensional protein structures, typically chosen to represent common structural classes. We wish to know which structure is most compatible with the sequence of the unknown protein. This is done by 'threading' the sequence through each protein structure in turn (hence the name). Threading methods are closely related to *ab initio* approaches to protein structure prediction, but whereas the latter can effectively explore all of the conformational space (albeit a restricted lattice in many cases) threading methods inherently limit the search space to the conformations of known structures. As such, threading is doomed to fail for any protein which adopts a completely new fold. That threading can work is, of course, a consequence of the fact that there does appear to be a finite set of protein folds and that proteins with very weak sequence similarity can adopt very similar structures. Thus when two proteins have a sequence identity of more than 70% it should be straightforward to determine an alignment that leads to a reliable model. As the degree of similarity falls so the task becomes more difficult, and when one enters the so-called 'twilight zone' (corresponding to less than 20-30% sequence identity) then comparative modelling is often considered to be inappropriate (or at least, any model should be treated with caution). Threading should (in principle, at least) be particularly suited to such problems.

A naïve threading implementation involves advancing each amino acid to occupy the location occupied in the previous iteration by its predecessor. A score is calculated for each structure so generated and the process is repeated until the sequence has been entirely threaded through the structure. The output is the structure or structures that correspond to the lowest value of the scoring function. As can be imagined there are very many possibilities to consider and so threading programs use special searching methods such as double dynamic programming to efficiently find the best ways to match the sequence to the structure. Even so, finding the optimal alignment is a very complex problem (particularly as one needs to consider gaps) and has resulted in some useful approximations that can be used to

make the problem more manageable. For example, in the *frozen approximation* each residue from the target sequence is scored according to the residues present in the actual template [Godzik *et al.* 1992]. It is also possible to use very high (or infinite) gap penalties in certain regions of the structure such as elements of secondary structure.

A variety of different scoring functions have been used for threading [Bryant and Lawrence 1993; Maiorov and Crippen 1994; Jernigan and Bahar 1996; Jones and Thornton 1996], but most share some common features. A threading calculation can require a large number of possibilities to be considered, and so the scoring functions are usually quite simple. This also reflects the low-resolution nature of the problem, in that one is usually attempting just to predict the basic fold of the protein. Each amino acid is typically treated as a single interaction site. Many of the scoring functions used in threading algorithms are potentials of mean force that provide an estimate of the free energy of interaction between two residues as a function of their separation. These potentials of mean force are calculated from statistical analyses of known protein structures. For example, if one plots the distribution of distances observed in X-ray protein structures between amino acids that are separated by three other residues in the sequence (i.e. between residues i and $i + 4$) then a large peak is observed in the interval 5.9–6.5 Å and a broad shoulder at 11.4–13.3 Å. These reflect the presence of α -helix and β -strand conformations. It is from these distribution frequencies that one can determine the residue-residue potentials of mean force. Sippl has provided an example of the use of such potentials [Sippl 1990]. The pentapeptide sequence valine-asparagine-threonine-phenylalanine-valine (VNTFV in the one-letter amino acid code) adopts an α -helical conformation in the protein erythrocrucorin but a β -strand conformation in ribonuclease. The potential of a mean force suggests that the β -strand is the more stable conformation for the isolated pentapeptide, but that when it is flanked by aspartic acid at one end and alanine at the other (as in erythrocrucorin), the α -helix does indeed become the more stable conformation. For threading algorithms one is particularly interested in the interactions between amino acids that are close in three-dimensional space but far apart in the sequence, and the potentials used in such calculations are derived appropriately. In addition to the pairwise knowledge-based term a solvation contribution is often added. This measures the propensity of each amino acid for a certain degree of solvation, to help ensure that hydrophobic residues pack into the central core of the protein and hydrophilic residues are on the outside. As it tends to be just the cores that are conserved, the pairwise interaction term may be omitted for loop residues, which are therefore treated just by the solvation term.

Although knowledge-based potentials are most popular, it is also possible to use other types of potential function. Some of these are more firmly rooted in the fundamental physics of interatomic interactions whereas others do not necessarily have any physical interpretation at all but are able to discriminate the correct fold from 'decoy' structures. These decoy structures are generated so as to satisfy the basic principles of protein structure such as a close-packed, hydrophobic core [Park and Levitt 1996]. The fold library is also clearly important in threading. For practical purposes the library should obviously not be too large, but it should be as representative of the different protein folds as possible. To derive a fold database one would typically first use a relatively fast sequence comparison method in conjunction with cluster analysis to identify families of homologues, which are assumed to have the same fold. A sequence identity threshold of about 30% is commonly

applied, with one representative being chosen from each cluster. Each of these representative structures is then compared to all other structures to enable a smaller set of representative folds to be identified. The number of unique folds is usually found to be about $\frac{1}{20}$ th the total number of structures in the protein databank, as we saw for the SCOP database.

10.8 A Comparison of Protein Structure Prediction Methods: CASP

The utility of a protein model depends upon the use to which it is put. In some cases, one is only interested in the general fold that the protein adopts and so a relatively low-resolution structure is acceptable. For other applications, such as drug design, the model must be much more accurate, including the loops and side chains. In such cases, a poor model may often be far worse than no model at all, as it can be seriously misleading.

To evaluate the then available techniques for calculating protein models, a 'competition' was organised in 1994–1995. In this first CASP* challenge entrants were invited to predict the three-dimensional structures of seven proteins from their amino acid sequences [Mosimann *et al.* 1995]. The structures of these seven proteins were simultaneously solved using X-ray crystallography, but these structures were not made available to the modellers. A total of 43 separate structures were submitted by 13 research groups, each model then being compared with the X-ray structure. The quality of each model was also assessed using a variety of methods, including the calculation of Ramachandran maps and 3D profiles.

Each competitor first had to decide which of the known protein structures they wished to use as the template; they then had to construct a sequence alignment (possibly using other protein sequences from the same family), and finally had to construct the model. The degree of sequence identity for the seven proteins ranged from 22% to 77%. One reassuring conclusion was that in favourable cases very accurate models could be constructed; the 'best' structure had an RMS difference with the X-ray structure of just 0.6 Å (for the protein NM23, which not surprisingly had the highest sequence identity with its template and indeed the same number of amino acid residues). Overall, the accuracy of the models largely depended upon two factors: the percentage sequence identity and the presence of substantial insertions or deletions between the template and target structures. The need for an accurate sequence alignment was evident; an incorrect alignment almost always resulted in an incorrect structure. For those proteins where there were large insertion loops, these were invariably predicted incorrectly, demonstrating a clear need for new strategies to tackle this problem. In some cases, models of the same protein were generated using both 'hands-on' approaches (where the modeller directed the construction of the structure) and wholly automatic procedures. In all cases, the manual structure was superior to the automatic model.

One disturbing finding was that a significant proportion of the models contained errors, even including amino acids with the wrong stereochemistry. In addition, significant deviations

* CASP stands for Critical Assessment of techniques for protein Structure Prediction.

from planarity of the amide bond were noted in many structures, and the torsion angles of the side chains often deviated from the distributions observed in experimental protein structures. Some of the models contained high-energy steric interactions between non-bonded atoms and unlikely distributions of amino acids (i.e. hydrophilic residues on the inside and hydrophobic residues on the outside). Most, if not all, of these problems can be identified very simply using publicly available software, and the organisers of the competition suggested that any models submitted for publication should be accompanied by output from a structure-verification program to allow an objective assessment of the quality of the model. The use of energy minimisation to refine models was identified as the cause of many of these problems, thereby highlighting the need for alternative protocols for the refinement stage of a comparative modelling exercise.

The first CASP competition was generally regarded as a great success and subsequent rounds have been organised, with increasing numbers of research groups submitting entries. At the time of writing the final reports from CASP3 have now been published [Moult *et al.* 1999], and the arrangements for CASP4 are well in hand. Three general categories were identified for CASP3. These comprise comparative modelling (which relies upon a clear relationship between the target protein and a protein of known structure), fold recognition (of which threading is one example), and *ab initio* prediction methods (which do not rely directly upon knowledge of any complete structures). Of the targets in CASP3, 15 were considered to be in the first category, 22 in the second and 15 in the third. A high level of participation ensured the success of the enterprise, which culminated in a meeting at Asilomar, California, to discuss the results. Of especial value to the community is that participants are encouraged to assess not only what was successful but (more importantly) what lessons had been learned. Whilst it is difficult to identify any real trends from just three CASPs some of the key developments have been the continuing improvement of *ab initio* prediction methods, the introduction of advanced sequence-comparison methods such as PSI-BLAST and hidden Markov models, which often performed as well as more 'sophisticated' techniques for suggesting possible homologues, and a gradual improvement in comparative modelling methods. Perhaps the single most important message that has emerged from all of the CASP competitions (as well as other studies) is that the real key to comparative modelling is the quality of the alignment.

10.8.1 Automated Protein Modelling

As we mentioned in the introduction to this chapter, the Human Genome Project is generating thousands of protein sequences, at a far greater rate than the structures can be solved by experimental techniques. Given the close relationship between the three-dimensional structure and the function of a protein there is increasing interest in automating the process of predicting the structures of these proteins, as a prelude to assigning a tentative function. A number of the methods that we considered for comparative modelling and fold recognition can be run in an automated manner. Indeed, one subsection of CASP3 was an assessment of automated methods, many of which have been made available over the World Wide Web (one of the earliest being Swiss-Model [Peitsch 1996]). This assessment confirmed that the most accurate models are those which benefit from some human input (especially during

the alignment stage), but the automated methods clearly have huge potential in exploiting the data being generated by the genome project.

An automated modelling procedure obviously needs to generate its model without the need for intervention: identifying structural templates related to the target sequence, aligning the templates with the target, building the model and finally evaluating and assessing the model. Of course, not all of the unknown proteins in a given genome will necessarily have a known structure that can be used as a template. For example, Sánchez and Šali considered the baker's yeast genome (*Saccharomyces cerevisiae*) [Sánchez and Šali 1998]. Of 6218 open reading frames (ORFs; a region of DNA that is converted into RNA and thence into protein) there were related structures for 2256 (36.3%), with an average pairwise sequence identity of 27%. Model building was performed with the Modeller program [Šali and Blundell 1993]. The models were then assessed for their quality, with 1071 of the original ORFs being considered to have a reliable model (17.2%).

Such large-scale modelling experiments can require significant computational resources, but the main bottleneck is generally considered to be the absence of structurally defined members of many protein families and the difficulty in detecting weak similarities, which would enable the appropriate template structures to be identified for more detailed comparative modelling. Above all, it is important to remember that no one single theoretical or experimental technique can predict protein function from sequence; rather, it is the application of an appropriate combination of methods that is required. Moreover, although our emphasis has been on the importance of the three-dimensional structure, such information is only one part of the jigsaw. An illustration is provided by a study which compared all of the protein structures released in 1998 with all structures that were known by the end of 1997 [Koppensteiner *et al.* 2000]. Some 147 of the proteins (corresponding to 196 domains) solved in 1998 had no significant sequence similarity to any of the pre-1998 proteins. However, when the structures of these 196 domains were compared with the pre-1998 set it was found that 147 of the domains had significant structural similarity with a previously known protein fold. Moreover, in two-thirds of these cases the function was also the same. The implication from these and similar studies is that computational techniques can be very effective at processing and filtering the raw sequence information in order to identify proteins that may be of interest and thus to suggest what experiments should be performed in order to confirm the hypothesis.

10.9 Protein Folding and Unfolding

The mechanism by which a protein folds to its native state has long been a subject of considerable interest, from both experimental and theoretical perspectives. As we noted in the introduction to this chapter, proteins typically adopt a single structure, corresponding to the global minimum of free energy under physiological conditions. Moreover, protein sequences can generally fold into this unique state in just a few seconds (or less) from any arbitrary starting conformation. Whilst exceptions to both of these two facts can be found, they do often hold true for small, water-soluble proteins such as enzymes, which have been the focus of most of the studies to date. The mechanism by which a protein is able

to fold into its unique fold was considered by Levinthal, who showed that folding could not occur via a systematic enumeration of all possible conformations [Levinthal 1969]. If it is assumed that there are three conformations for each amino acid then a polypeptide chain with (say) 100 amino acids would have about 10^{48} conformations. If the interconversion between conformations required just 10^{-11} seconds then it would take about 10^{29} years to explore them all. Of course, this is for the most basic of grid search algorithms, but even the most advanced systematic conformational search would still require an inordinate amount of time to identify the global minimum energy conformation. This discrepancy between the time for the exhaustive search and the observed timescale of protein folding is popularly known as the *Levinthal paradox*.

Two general types of computational model have been used to investigate protein folding: simple lattice models and atomistic models. These two approaches are complementary; lattice models attempt to capture the essential physics of the problem but do not provide information about specific interactions at the atomic level. It is often possible to exhaustively enumerate all possible states on the entire energy surface of a lattice model, in contrast to an atomistic model. Another important feature of atomistic simulations is that most (to date) have considered protein *unfolding*, rather than folding. The two processes are obviously linked through the principle of microscopic reversibility, though it has been argued that an unfolding pathway obtained with the 'strong' unfolding conditions (such as high temperature) that are often used in the simulations may not necessarily correspond to the 'true' physiological folding pathway [Finkelstein 1997]. It is just becoming possible to directly simulate the folding of proteins which, although containing a very small number of amino acids, do have recognisable secondary structure in solution. In the remainder of this section, we will consider both types of model and what they can tell us about the nature of protein folding and how a combination of experiment and theory has led to a 'new view' of protein folding that appears to resolve the Levinthal paradox.

This 'new view' considers an ensemble of structures through which a protein can fold to the native conformation rather than a single pathway involving a number of distinct intermediates. As such, a statistical description of the energy surface can be applied [Bryngelson and Wolynes 1987; Bryngelson *et al.* 1995; Onuchic *et al.* 1997]. The resulting theory suggests that the energy landscape of protein folding (i.e. the variation of free energy with the protein conformation, and the form of that free energy) generally resembles a funnel, but one which is 'rough', with local minima where the protein can transiently reside. Most of the molecular organisation occurs early in the procedure and can be described using a few parameters. Later in the folding process the protein may become trapped in the local minima, which can give rise to the semblance of a discrete folding pathway that is sensitive to the amino acid sequence and three-dimensional structure. A schematic representation of a folding landscape that conforms to these ideas is shown in Figure 10.27.

The essential features of protein lattice models were described in Section 10.3.1. Protein folding is usually studied with a self-avoiding chain on a cubic lattice with one residue per vertex and a simple interaction model which only includes interactions between pairs of monomers that are in contact on the lattice but are not successive in the sequence. Polymers of length 27 that occupy all sites of a $3 \times 3 \times 3$ cube are particularly popular. There are nearly

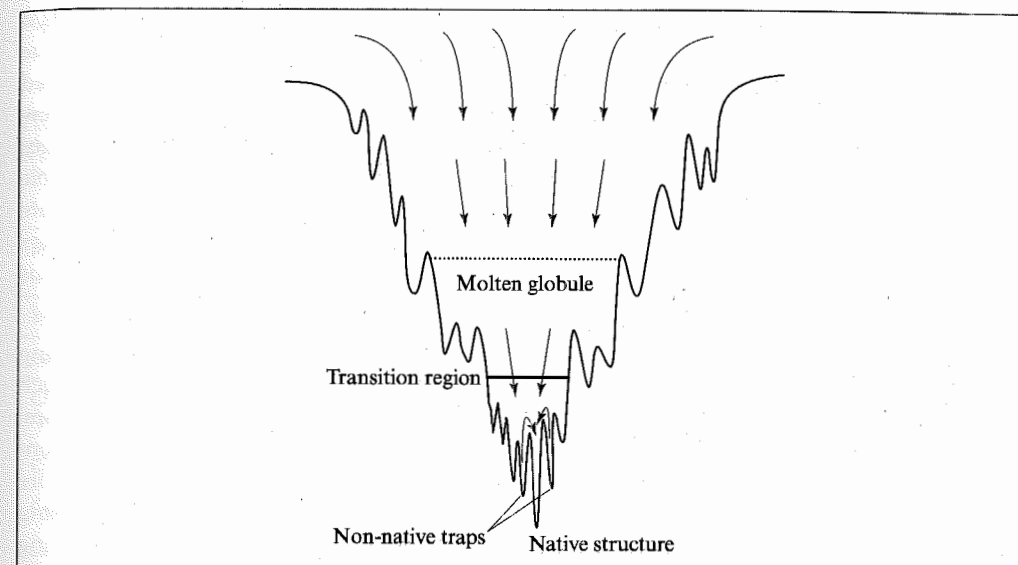


Fig. 10.27: Schematic representation of the energy landscape for protein folding. (Figure adapted from Onuchic J N, Z Luthey-Schulten and P Wolynes 1997. *Theory of Protein Folding: The Energy Landscape Perspective*. *Annual Reviews in Physical Chemistry* 48:545-600.)

5 million possible structures for this system, 51704 of which are unique (unrelated by rotation, reflection or reverse-labelling symmetries). In addition, it is estimated that there are about 10^{18} different non-compact arrangements.

In one study, Monte Carlo simulations were used to explore the conformational space of several hundred such sequences. The interactions between pairs of residues were selected at random from a Gaussian distribution [Sali *et al.* 1994a, b]. As such, this model corresponds to a heteropolymer with a random sequence of monomers of many different types. It was found that some sequences were able to find the global energy minimum (the native state) within a relatively short number of steps, whereas others did not. The key difference between the folding and non-folding sequences was the presence in the folding sequences of a pronounced global energy minimum, with a relatively large energy separation to the next lowest state. A three-stage folding pathway was suggested from these studies. The first stage involves a rapid collapse to a semi-compact random globule which contains about 30% of the contacts observed in the global minimum. In the second, rate-limiting, stage the protein searches for a transition state. There are about 1000 transition states which are structurally similar to the native state, having 80-95% of the native contacts. In the third stage, the chain rapidly progresses from one of the transition states to the native conformation. The transition region is key to the folding mechanism as it enables the search time to be reduced to a realistic value.

A related study was the exhaustive enumeration of the global minimum energy structures for all 2^{27} possible sequences of the 'HP' model described in Section 10.3.1 [Li *et al.* 1996]. This showed that 4.75% of these sequences have a unique ground state (i.e. just one

conformation of the polymer on the lattice gives rise to the minimum energy). From this data it was also possible to determine how many sequences had a given structure as their unique ground state. Some structures are adopted by many sequences (the best being represented by 3794 sequences) whereas other structures correspond to the ground state of just a few sequences. Of additional interest was the fact that these 'highly designable' structures tended to have a larger gap between the global energy minimum and the next most stable structure, as found by the earlier studies. Of course, one could argue that it is inappropriate to extrapolate from the results of a small $3 \times 3 \times 3$ cube to 'real' proteins, and some workers consider the HP interaction model to be too simplistic. Nevertheless, such investigations invariably have the benefit of stimulating debate, which often leads to further advances.

Atomistic simulations of protein unfolding have most frequently been performed using high-temperature molecular dynamics. As we have already noted, the results from such calculations might not always be directly relevant to the physiological folding mechanism. However, simulations under other denaturing conditions such as high or low pH or non-aqueous solvents or even high-concentration urea (a common laboratory denaturant) are also possible. Most simulations are for a time of at least 1 ns; some are performed using explicit solvent, whilst others use an implicit solvent model. Some of the most interesting and fruitful work in this field is on systems that have also been studied experimentally, with NMR spectroscopy being a particularly widely used technique. The NMR data from an unfolded protein is usually more difficult to interpret than for a folded protein, and simulations can be useful in helping its interpretation. Two examples where this was possible are the characterisation of partially unfolded states of ubiquitin in 60% methanol [Alonso and Daggett 1995] and in thermally denatured barnase [Bond *et al.* 1997].

The ultimate goal for some in this area is a full atomistic simulation of the folding process, starting from an arbitrary structure, with explicit representation of the solvent. Such simulations are currently at the very limit of what can be achieved due to the length of simulation required and the large number of particles involved. One example of the current state-of-the-art is the $1 \mu\text{s}$ simulation of a 36-residue peptide starting from a fully extended state [Duan and Kollman 1998]. This peptide is one of the smallest proteins that can fold autonomously, with folding estimated to take between $10 \mu\text{s}$ and $100 \mu\text{s}$. It contains three short α -helices. The simulation involved in addition to the protein about 3000 water molecules and was performed in a truncated octahedron simulation box with a time step of 2 fs. About four months of computing time on a 256 massively parallel supercomputer was required for the $1 \mu\text{s}$ simulation. Whilst the protein did not actually fold into the known experimental structure, a marginally stable state which showed significant resemblance to the native conformation was observed. This state had a lifetime of about 150 ns. A variety of metrics were used to monitor the simulation, including the RMS deviation to the experimental structure, the radius of gyration, the fraction of native contacts present and the solvation free energy. A high degree of fluctuation in all of these features was observed, characteristic of a relatively shallow free-energy landscape. Cluster analysis of the trajectory was used to identify the major conformations visited during the simulation and also to characterise the pathways between these conformations. Ready transitions were observed between the early states especially, giving a 'tangled' network of pathways. As computer power increases we are likely to see more studies of this type.

Appendix 10.1 Some Common Abbreviations and Acronyms Used in Bioinformatics

This is necessarily an incomplete list; more comprehensive glossaries can be found elsewhere (particularly on the World Wide Web).

A,G,C,T, (U)	Adenine, guanine, cytosine, thymine – the four bases present in DNA. Uracil replaces thymine in RNA
Bp	Base pair
cDNA	Complementary DNA, synthesised from messenger RNA
Chromosome	Discrete unit of the genome consisting of a single molecule of DNA that carries many genes.
Clone	Genetically identical copy (of a gene, cell or organism)
Codon	Sequence of three nucleotides that codes for a single amino acid (or a termination signal)
Contig	A group of pieces of DNA, derived from a cloning experiment (often a series of ESTs, see below), that represent overlapping regions of a chromosome
Deletion	One or more nucleotides that are not copied during DNA replication
DNA	Deoxyribose nucleic acid
Domain	Sequence of a polypeptide chain that can independently fold into a stable three-dimensional structure
Dynamic programming	Technique widely used in sequence alignment
EST	Expressed sequence tag. An EST is a partial sequence (typically less than 400 bases) selected from cDNA and used to identify genes expressed in a particular tissue
Eukaryote	Organism whose cells have a discrete nucleus and other subcellular compartments (<i>cf.</i> prokaryote)
Exon	Translated sequence of DNA
Gap	A break in a DNA or protein sequence which enables two or more sequences to be aligned
Gene	A sequence of DNA at a particular position on a specific chromosome that encodes a precise functional product (usually protein)
Genome	All of the genetic material in the chromosomes of an organism
Indel	Insertion or deletion required to optimise sequence alignment
Intron	Non-translated sequence of DNA
Kb	Kilobase – one thousand nucleotide bases
ktup	<i>k</i> -tuple. Parameter used in FASTA and FASTP sequence-alignment methods
Mb	Megabase – one million nucleotide bases

Appendix 10.1 Continued

mRNA	Messenger RNA
Mutation	A change in the DNA sequence
Nucleotide	Three components that make up the basic building block in DNA and RNA: a nitrogenous base (A, T, G, C, U), a phosphate and a sugar
Oligonucleotide	A molecule composed of a small number of nucleotides
Orthologue	Homologous proteins that perform the same function in different organisms
ORF	Open Reading Frame – region of DNA that is transcribed into RNA. Delineated by an initiator codon at one end and a stop codon at the other end
PAM	Point Accepted Mutation per 100 residues
Paralogue	Homologous proteins that perform different but related functions in one organism
PCR	Polymerase Chain Reaction. Widely used method for amplifying a DNA base sequence
Polymorphism	Differences in DNA sequence between individuals
Prokaryote	Organism lacking a nucleus and subcellular compartments (<i>cf.</i> eukaryote). Includes bacteria and viruses
RNA	Ribonucleic acid
SNP	Single Polynucleotide Polymorphism – single base-pair variations in DNA
STS	Sequence tagged site. A short DNA sequence that occurs just once in the human genome and whose location and base sequence are known
Transcription	First step in gene expression, corresponding to the generation of mRNA from the original DNA
Translation	Second step in gene expression, the synthesis of proteins from mRNA
tRNA	Transfer RNA

Appendix 10.2 Some of the Most Common Sequence and Structural Databases Used in Bioinformatics

GenBank (NCBI, USA) EMBL Nucleotide Sequence Database (Europe) DDBJ (Japan)	The three main nucleotide sequence databases, which are synchronised daily
PIR-International Protein Sequence Database Swiss-Prot, TrEMBL	Redundant protein sequence database Annotated non-redundant protein sequence database. TrEMBL is a computer-annotated supplement to Swiss-Prot. TrEMBL contains the translations of all coding sequences present in the EMBL Nucleotide Sequence Database which are not yet integrated into Swiss-Prot
GenPept	Compendium of amino acid translations derived from GenBank
PDB, NRL3D	Protein Data Bank – protein structures (mostly from X-ray crystallography). NRL3D is a derived sequence database in PIR format
SCOP	Structural Classification of Proteins. Hierarchical protein structure database
CATH, FSSP	Sequence-structure classification databases
Prosite	Motif database

Appendix 10.3 Mutation Probability Matrix for 1 PAM

	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val	
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
Ala	A	867	2	9	10	3	8	17	21	2	6	4	2	6	2	22	35	32	0	2	18
Arg	R	1	914	1	0	1	10	0	10	3	1	19	4	1	4	6	1	8	0	1	
Asn	N	4	1	822	36	0	4	6	6	21	3	1	13	0	1	2	20	9	1	4	1
Asp	D	6	0	42	859	0	6	53	6	4	1	0	3	0	0	1	4	3	0	0	1
Cys	C	1	1	0	0	973	0	0	0	1	1	0	0	0	0	1	5	1	0	3	2
Gln	Q	3	9	4	5	0	876	27	1	23	1	3	6	4	0	6	2	2	0	0	1
Glu	E	10	0	7	56	0	35	864	4	2	3	1	4	2	0	3	4	2	0	1	2
Gly	G	21	1	12	11	1	2	7	935	1	0	1	2	2	1	3	21	3	0	0	5
His	H	1	8	18	3	1	20	1	0	912	0	1	1	0	2	3	1	1	1	4	1
Ile	I	2	2	3	1	2	1	2	0	0	872	9	2	12	7	0	1	7	0	1	32
Leu	L	3	1	3	0	0	6	1	1	4	22	947	2	45	13	3	1	3	4	2	15
Lys	K	2	37	25	6	0	12	7	2	2	4	1	925	19	0	3	8	11	0	1	1
Met	M	1	1	0	0	0	2	0	0	5	8	4	875	1	0	1	2	0	0	4	
Phe	F	1	1	1	0	0	0	0	1	2	8	6	0	4	945	0	2	1	3	28	0
Pro	P	13	5	2	1	1	8	3	2	5	1	2	2	1	1	925	12	4	0	0	2
Ser	S	28	11	34	7	11	4	6	16	2	2	1	7	4	3	17	840	38	5	2	2
Thr	T	22	2	13	4	1	3	2	2	1	11	2	8	6	1	5	32	871	0	2	9
Trp	W	0	2	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	976	1	0
Tyr	Y	1	0	3	0	3	0	1	0	4	1	1	0	0	21	0	1	1	2	945	1
Val	V	13	2	1	1	3	2	2	3	3	57	11	1	17	1	3	2	10	0	2	902

Each element M_{ij} of this matrix corresponds to the probability that the amino acid in column i will mutate to the amino acid in row j after a period of 1 PAM. The values have been multiplied by 10 000. (Based on Dayhoff M O 1978. *Atlas of Protein Sequence and Structure* Volume 5 Supplement 3. Dayhoff M O (Editor) Georgetown University Medical Center, National Biomedical Research Foundation: Figure 82.)

Appendix 10.4 Mutation Probability Matrix for 250 PAM

	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
Ala	A	13	6	9	9	5	8	9	12	6	8	6	7	4	11	11	11	2	4	9
Arg	R	3	17	4	3	2	5	3	2	6	3	2	9	4	4	4	3	7	2	2
Asn	N	4	4	6	7	2	5	6	4	6	3	2	5	3	2	4	5	4	2	3
Asp	D	5	3	8	11	1	7	10	5	6	3	2	5	3	1	4	5	5	1	2
Cys	C	2	1	1	1	52	1	1	2	2	1	1	1	1	2	3	2	1	4	2
Gln	Q	3	5	5	6	1	10	7	3	8	2	3	5	3	1	4	3	3	1	2
Glu	E	5	4	7	11	1	9	12	5	6	3	2	5	3	1	4	5	5	1	2
Gly	G	12	5	10	10	4	7	9	27	5	5	4	6	5	3	8	11	9	2	3
His	H	2	5	5	4	2	7	4	2	15	2	2	3	2	3	3	2	2	3	2
Ile	I	3	2	2	2	2	2	2	2	10	6	2	6	5	2	3	4	1	3	9
Leu	L	6	4	4	3	2	6	4	3	5	15	34	4	20	13	5	4	6	6	7
Lys	K	6	18	10	8	2	10	8	5	8	5	4	24	9	2	6	8	8	4	3
Met	M	1	1	1	1	0	1	1	1	2	3	2	7	2	1	1	1	1	1	2
Phe	F	2	1	2	1	1	1	1	3	5	6	1	4	32	1	2	2	4	20	3
Pro	P	7	5	5	4	3	5	4	5	3	3	4	3	2	19	6	5	1	2	4
Ser	S	9	6	8	7	7	6	7	9	6	5	4	7	5	3	9	10	9	4	6
Thr	T	8	5	6	6	4	5	5	6	4	6	4	6	5	3	6	8	11	2	3
Trp	W	0	2	0	0	0	0	0	0	1	0	0	0	1	0	1	0	55	1	0
Tyr	Y	1	1	2	1	3	1	1	3	2	2	1	2	15	1	2	2	3	31	2
Val	V	7	4	4	4	4	4	4	5	4	15	10	4	10	5	5	7	2	4	17

Each element M_{ij} of this matrix corresponds to the probability that the amino acid in column i will mutate to the amino acid in row j after a period of 250 PAM. The values have been multiplied by 100. (Based on Dayhoff M O 1978. *Atlas of Protein Sequence and Structure* Volume 5 Supplement 3. Dayhoff M O (Editor) Georgetown University Medical Center, National Biomedical Research Foundation: Figure 83.)

Further Reading

- Altschul S F 1996. Sequence Comparison and Alignment. In Sternberg M E (Editor) *Protein Structure Prediction - A Practical Approach*. Oxford, IRL Press, pp. 137-167.
- Altschul S F, M S Boguski, W Gish and J C Wootton 1994. Issues in Searching Molecular Sequence Databases. *Nature Genetics* 6:119-129.
- Attwood T K and D J Parry-Smith 2000. *Introduction to Bioinformatics*. Harlow, Addison Wesley Longman.
- Barton G J 1996. Protein Sequence Alignment and Database Scanning. In Sternberg M E (Editor) *Protein Structure Prediction - A Practical Approach*. Oxford, IRL Press, pp. 31-63.
- Barton G J 1998. Protein Sequence Alignment Techniques. *Acta Crystallographica* D54:1139-1146.
- Blundell T L, B L Sibanda, M J E Sternberg and J M Thornton. Knowledge-based Prediction of Protein Structures and the Design of Novel Molecules. *Nature* 326:347-352.

- Branden C and J Tooze 1991. *Introduction to Protein Structure*.
 Chatfield C and A J Collins 1980. *Introduction to Multivariate Analysis*. London, Chapman & Hall.
 Dobson C M, A Šali and M Karplus 1998. Protein Folding: A Perspective from Theory and Experiment. *Angewandte Chemie International Edition* 37:868-893.
 Perutz M 1992. *Protein Structure. New Approaches to Disease And Therapy*. New York, W H Freeman.
 Schulz G E and R H Schirmer 1979. *Principles of Protein Structure*. New York, Springer-Verlag.

References

- Alonso D O V and V Daggett 1995. Molecular Dynamics Simulations of Protein Unfolding and Limited Refolding: Characterisation of Partially Unfolded States of Ubiquitin in 60% Methanol and in Water. *Journal of Molecular Biology* 247:501-520.
 Altschul S F, W Gish, W Miller, E W Myers and D J Lipman 1990. Basic Local Alignment Search Tool. *Journal of Molecular Biology* 215:403-410.
 Altschul S F, T L Madden, A A Schäffer, J Zhang, Z Zhang, W Miller and D J Lipman 1997. Gapped BLAST and PSI-BLAST: A New Generation of Protein Database Search Programs. *Nucleic Acids Research* 25:3389-3402.
 Bernstein F C, T F Koetzle, G J B Williams, E Meyer, M D Bryce, J R Rogers, O Kennard, T Shikhanouchi and M Tasumi 1977. The Protein Data Bank: A Computer-Based Archival File for Macromolecular Structures. *Journal of Molecular Biology* 112:535-542.
 Birktoft J J and D M Blow 1972. The structure of Crystalline Alpha-Chymotrypsin V. The Atomic Structure of Tosyl-Alpha-Chymotrypsin at 2 Ångstroms Resolution. *Journal of Molecular Biology* 68:187-240.
 Bond C J, K-B Wong, J Clarke, A R Ferscht and V Daggett 1997. Characterisation of Residual Structure in the Thermally Denatured State of Barnase by Simulation and Experiment: Description of the Folding Pathway. *Proceedings of the National Academy of Sciences USA* 94:13409-13413.
 Bowie J U, R Lüthy and D Eisenberg 1991. A Method to Identify Protein Sequences that Fold into a Known Three-Dimensional Structure: *Science* 253:164-170.
 Brenner S E, C Chothia and T J P Hubbard 1997. Population Statistics of Protein Structures: Lessons from Structural Classifications. *Current Opinion in Structural Biology* 7:369-376.
 Bruccoleri R E and M Karplus 1985. Chain Closure with Bond Angle Variations. *Macromolecules* 18:2767-1773.
 Bruccoleri R E and M Karplus 1987. Prediction of the Folding of Short Polypeptide Segments by Uniform Conformational Sampling. *Biopolymers* 26:137-168.
 Bryant S H and C E Lawrence 1993. An Empirical Energy Function for Threading Protein Sequences Through the Folding Motif. *Proteins: Structure, Function and Genetics* 16:92-112.
 Bryngelson J D, J N Onuchic, N D Socci and P G Wolynes 1995. Funnels, Pathways, and the Energy Landscape of Protein Folding: A Synthesis. *Proteins: Structure, Function and Genetics* 21:167-195.
 Bryngelson J D and P G Wolynes 1987. Spin Glasses and the Statistical Mechanics of Protein Folding. *Proceedings of the National Academy of Sciences USA* 84:7524-7528.
 Chan H S and K A Dill 1993. The Protein Folding Problem. *Physics Today* Feb:24-32.
 Chothia C and A M Lesk 1986. The Relation Between the Divergence of Sequence and Structure in Proteins. *EMBO Journal* 5:823-826.
 Chothia C, A M Lesk, A Tramontano, M Levitt, S J Smith-Gill, G Air, S Sheriff, E A Padlan and D Davies 1989. Conformations of Immunoglobulin Hypervariable Regions. *Nature* 342:877-883.
 Chou P Y and G D Fasman 1978. Prediction of the Secondary Structure of Proteins from Their Amino Acid Sequence. *Advances in Enzymology* 47:45-148.

- Cohen F E and S R Presnell 1996. The Combinatorial Approach. In Sternberg M J E (Editor) *Protein Structure and Prediction*. Oxford, IRL Press, pp. 207-227.
 Cohen F E, M J E Sternberg and W R Taylor 1982 Analysis and Prediction of the Packing of α -Helices against a β -Sheet in the Tertiary Structure of Globular Proteins. *Journal of Molecular Biology* 156:821-862.
 Cuff J A and G J Barton 1999. Evaluation and Improvement of Multiple Sequence Methods for Protein Secondary Structure Prediction. *Proteins: Structure, Function and Genetics* 34:508-519.
 Dayhoff M O 1978. A Model of Evolutionary Change. In Dayhoff M O (Editor) *Proteins in Atlas of Protein Sequence and Structure* Volume 5 Supplement 3. Georgetown University Medical Center, National Biomedical Research Foundation, pp. 345-358.
 Duan Y and P A Kollman 1998. Pathways to a Protein Folding Intermediate Observed in a 1-Microsecond Simulation in Aqueous Solution. *Science* 282:740-744.
 Dunbrack R L Jr and M Karplus 1993. Backbone-dependent Rotamer Library for Proteins. Application to Side-chain Prediction. *Journal of Molecular Biology* 230:543-574.
 Eddy S R 1996. Hidden Markov Models. *Current Opinion in Structural Biology* 6:361-365.
 Finkelstein A V 1997. Can Protein Unfolding Simulate Protein Folding? *Protein Engineering* 10:843-845.
 Frazao C, C Topham, V Dhanaraj and T L Blundell 1994. Comparative Modelling of Human Renin: A Retrospective Evaluation of the Model with Respect to the X-ray Crystal Structure. *Pure and Applied Chemistry* 66:43-50.
 Garnier J, D Osguthorpe and B Robson 1978. Analysis of the Accuracy and Implications of Simple Methods for Predicting the Secondary Structure of Globular Proteins. *Journal of Molecular Biology* 120:97-120.
 Gibson K D and H A Scheraga 1987. Revised Algorithms for the Build-up Procedure for Predicting Protein Conformations by Energy Minimization. *Journal of Computational Chemistry* 8:826-834.
 Go N and H A Scheraga 1970. Ring Closure and Local Conformational Deformations of Chain Molecules. *Macromolecules* 3:178-187.
 Godzik A, A Kolinski and J Skolnick 1993. De Novo and Inverse Folding Predictions of Protein Structure and Dynamics. *Journal of Computer-Aided Molecular Design* 7:397-438.
 Godzik A, J Skolnick and A Kolinski 1992. Simulations of the Folding Pathway of Triose Phosphate Isomerase-type α/β Barrel Proteins. *Proceedings of the National Academy of Sciences USA* 89:2629-2633.
 Gonnet G H, M A Cohen and S A Benner 1992. Exhaustive Matching of the Entire Protein Sequence Database. *Science* 256:1443-1445.
 Gribskov M, A D McLachlan and D Eisenberg 1987. Profile Analysis: Detection of Distantly Related Proteins. *Proceedings of the National Academy of Sciences USA* 84:4335-4358.
 Havelka W A, R Henderson and D Oesterhelt 1995. 3-Dimensional Structure of Halorhodopsin at 7-Ångstrom Resolution. *Journal of Molecular Biology* 247:726-738.
 Henderson R, J M Baldwin, T A Ceska, F Zemlin, E Beckmann and K H Downing 1990. Model for the Structure of Bacteriorhodopsin Based on High-resolution Electron Cryo-microscopy. *Journal of Molecular Biology* 213:899-929.
 Henikoff S and J G Henikoff 1992. Amino Acid Substitution Matrices from Protein Blocks. *Proceedings of the National Academy of Sciences USA* 89:10915-10919.
 Holm L and C Sander 1993. Protein Structure Comparison by Alignment of Distance Matrices. *Journal of Molecular Biology* 233:123-138.
 Holm L and C Sander 1994. The FSSP Database of Structurally Aligned Protein Fold Families. *Nucleic Acids Research* 22:3600-3609.
 Holm L and C Sander 1999. Protein Folds and Families: Sequence and Structure Alignments. *Nucleic Acids Research* 27:244-247.
 Jernigan R L and I Bahar 1996. Structure-derived Potentials and Protein Simulations. *Current Opinion in Structural Biology* 6:195-209.

- Jones D and J Thornton 1993. Protein Fold Recognition. *Journal of Computer-Aided Molecular Design* 7:439-456.
- Jones D T, W R Taylor and J M Thornton 1992. A New Approach to Protein Fold Recognition. *Nature* 358:86-89.
- Jones D T and J M Thornton 1996. Potential Energy Functions for Threading. *Current Opinion in Structural Biology* 6:210-216.
- Jones T A and S Thirup 1986. Using Known Substructures in Protein Model Building and Crystallography. *EMBO Journal* 5:819-822.
- King, R D, M Saqi, R Sayle and M J E Sternberg 1997. DSC: Public Domain Protein Secondary Structure Prediction. *Computer Applications in the Biosciences* 13:473-474.
- Koppensteiner W A, P Lackner, M Wiederstein and M J Sippl 2000. Characterization of Novel Proteins Based on Known Protein Structures. *Journal of Molecular Biology* 296:1139-1152.
- Kovacs H, A E Mark and W F van Gunsteren 1997. Solvent Structure at a Hydrophobic Protein Surface. *Proteins: Structure, Function and Genetics* 27:395-404.
- Krogh A, M Brown, S Mian, K Sjölander and D Haussler 1994. Hidden Markov Models in Computational Biology. Applications to Protein Modeling. *Journal of Molecular Biology* 235:1501-1531.
- Levinthal C 1969. In Debrunner P, J C M Tsibris and E Munck (Editors) *Mössbauer Spectroscopy in Biological Systems*, Proceedings of a Meeting held at Allerton House, Monticello, Illinois, University of Illinois Press, Urbana, p. 22.
- Levitt M 1976. A Simplified Representation of Protein Conformations for Rapid Simulation of Protein Folding. *Journal of Molecular Biology* 104:59-107.
- Levitt M 1992. Accurate Modeling of Protein Conformation by Automatic Segment Matching. *Journal of Molecular Biology* 226:507-533.
- Li H, R Helling, C Tang and N Wingreen 1996. Emergence of Preferred Structures in a Simple Model of Protein Folding. *Science* 273:666-669.
- Li Z Q and H A Scheraga 1987. Monte Carlo Minimization Approach to the Multiple Minima Problem in Protein Folding. *Proceedings of the National Academy of Sciences USA* 84:6611-6615.
- Lipman, D J and W R Pearson 1985. Rapid and Sensitive Protein Similarity Searches. *Science* 227:1435-1441.
- Lüthy R, J U Bowie and D Eisenberg 1992. Assessment of Protein Models with Three-Dimensional Profiles. *Nature* 356:83-85.
- Maiorov V N and G M Crippen 1994. Learning About Protein Folding via Potential Functions. *Proteins: Structure, Function and Genetics* 20:167-173.
- Mosimann S, S Meleshko and M N G Jones 1995. A Critical Assessment of Comparative Molecular Modeling of Tertiary Structures of Proteins. *Proteins: Structure, Function and Genetics* 23:301-317.
- Moult J, T Hubbard, K Fidelis and J T Pedersen 1999. Critical Assessment of Methods of Protein Structure Prediction (CASP): Round III. *Proteins: Structure, Function and Genetics* Suppl. 3:2-6.
- Murzin A G, S E Brenner, T Hubbard and C Chothia 1995. SCOP: A Structural Classification of Proteins Database for the Investigation of Sequences and Structures. *Journal of Molecular Biology* 247:536-540.
- Needleman S B and C D Wunsch 1970. A General Method Applicable to the Search for Similarities in the Amino Acid Sequences of Two Proteins. *Journal of Molecular Biology* 48:443-453.
- Ning Q and T J Sejnowski 1988. Predicting the Secondary Structure of Globular Proteins Using Neural Network Models. *Journal of Molecular Biology* 202:865-888.
- Noble M E M, R K Wierenga, A-M Lambeir, F R Opperdoes, W H Thunnissen, K H Kalk, H Groendijk and W G J Hol 1991. The Adaptability of the Active Site of Trypanosomal Triosephosphate Isomerase as Observed in the Crystal Structures of Three Different Complexes. *Proteins: Structure, Function and Genetics* 10:50-69.
- Novotny J, A A Rashin and R E Bruccoleri 1988. Criteria that Discriminate between Native Proteins and Incorrectly Folded Models. *Proteins: Structure, Function and Genetics* 4:19-30.
- Onuchic J N, Z Luthey-Schulten and P Wolynes 1997. Theory of Protein Folding: The Energy Landscape Perspective. *Annual Reviews in Physical Chemistry* 48:545-600.
- Orengo C A, N P Brown and W R Taylor 1992. Fast Structure Alignment for Protein Databank Searching. *Proteins: Structure, Function and Genetics* 14:139-167.
- Orengo C A and W R Taylor 1990. A Rapid Method of Protein Structure Alignment. *Journal of Theoretical Biology* 147:517-551.
- Orengo C A and W R Taylor 1993. A Local Alignment Method for Protein Structure Motifs. *Journal of Molecular Biology* 233:488-497.
- Orengo C A, T P Flores, W R Taylor and J M Thornton 1993. Identification and Classification of Protein Fold Families. *Protein Engineering* 6:485-500.
- Ortiz A R, A Kolinski and J Skolnick 1998. Fold Assembly of Small Proteins Using Monte Carlo Simulations Driven by Restraints Derived from Multiple Sequence Alignments. *Journal of Molecular Biology* 277:419-446.
- Park B and M Levitt 1996. Energy Functions that Discriminate X-ray and Near-native Folds from Well-constructed Decoys. *Journal of Molecular Biology* 258:367-392.
- Pauling L, R B Corey and H R Bronson 1951. The Structure of Proteins: Two Hydrogen-bonded Helical Configurations of the Polypeptide Chain. *Proceedings of the National Academy of Sciences USA* 37:205-211.
- Pearson W R 1990. Rapid and Sensitive Sequence Comparison with FASTP and FASTA. *Methods in Enzymology* 183:63-98.
- Pearson W R and D J Lipman 1988. Improved Tools for Biological Sequence Comparison. *Proceedings of the National Academy of Sciences USA* 85:2444-2448.
- Peitsch M C 1996. ProMod and Swiss-Model: Internet-based Tools for Automated Comparative Protein Modelling. *Biochemical Society Transactions* 24:274-279.
- Ponder J W and F M Richards 1987. Tertiary Templates for Proteins. Use of Packing Criteria in the Enumeration of Allowed Sequences for Different Structural Classes. *Journal of Molecular Biology* 193:775-791.
- Rabiner L R 1989. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE* 77:257-286.
- Ripoll D R and H A Scheraga 1988. On the Multiple-Minimum Problem in the Conformational Analysis of Polypeptides. II. An Electrostatically Driven Monte Carlo Method: Tests on Poly(L-Alanine). *Biopolymers* 27:1283-1303.
- Ripoll D R and H A Scheraga 1989. On the Multiple-Minimum Problem in the Conformational Analysis of Polypeptides. III. An Electrostatically Driven Monte Carlo Method: Tests on met-Enkephalin. *Journal of Protein Chemistry* 8:263-287.
- Rost B and C Sander 1993. Prediction of Protein Secondary Structure at Better than 70% Accuracy. *Journal of Molecular Biology* 232:584-599.
- Šali A 1995. Modelling Mutations and Homologous Proteins. *Current Opinion in Biotechnology* 6:437-451.
- Šali A and T L Blundell 1993. Comparative Protein Modelling by Satisfaction of Spatial Restraints. *Journal of Molecular Biology* 234:779-815.
- Šali A, E Shakhnovich and M Karplus 1994a. How Does a Protein Fold? *Nature* 369:248-251.
- Šali A, E Shakhnovich and M Karplus 1994b. Kinetics of Protein Folding. A Lattice Model Study of the Requirements for Folding to the Native State. *Journal of Molecular Biology* 235:1614-1636.
- Sánchez R and A Šali 1998. Large-scale Protein Structure Modelling of the *Saccharomyces cerevisiae* Genome. *Proceedings of the National Academy of Sciences USA* 95:13597-13602.
- Scheraga H A 1993. Searching Conformational Space. In van Gunsteren W F, P K Weiner and A J Wilkinson (Editors) *Computer Simulation of Biomolecular Systems* Volume 2. Leiden, ESCOM.

- Shenkin P S, D L Yarmusch, R M Fine, H Wang and C Levinthal 1987. Predicting Antibody Hypervariable Loop Conformation. I. Ensembles of Random Conformations for Ringlink Structures. *Biopolymers* **26**:2053–2085.
- Simons K T, R Bonneau, I Ruzinski and D Baker 1999b. *Ab Initio* Protein Structure Prediction of CASP III Targets Using ROSETTA. *Proteins: Structure, Function and Genetics Supplement* **3**:171–176.
- Simons K T, C Kooperberg, E Huang and D Baker 1997. Assembly of Protein Tertiary Structures from Fragments with Similar Local Sequences Using Simulated Annealing and Bayesian Scoring Functions. *Journal of Molecular Biology* **268**:209–225.
- Simons K T, I Ruzinski, C Kooperberg, B A Cox, C Bystroff and D Baker 1999a. Improved Recognition of Native-Like Protein Structures Using a Combination of Sequence-Dependent and Sequence-Independent Features of Proteins. *Proteins: Structure, Function and Genetics* **34**:82–95.
- Sippl M J 1990. Calculation of Conformational Ensembles from Potentials of Mean Force. An Approach to the Knowledge-Based Prediction of Local Structures in Globular Proteins. *Journal of Molecular Biology* **213**:859–883.
- Skolnick J, A Kolinski and A R Ortiz 1997. MONSSTER: A Method for Folding Globular Proteins with a Small Number of Distance Restraints. *Journal of Molecular Biology* **265**:217–241.
- Smith T F and M S Waterman 1981. Identification of Common Molecular Subsequences. *Journal of Molecular Biology* **147**:195–197.
- Srinivasan N, K Gurprasad and T L Blundell 1996. Comparative Modelling of Proteins. In Sternberg M E (Editor) *Protein Structure Prediction – A Practical Approach*. Oxford, IRL Press, pp. 111–140.
- Sternberg M J E, F E Cohen and W R Taylor 1982 A Combinatorial Approach to the Prediction of the Tertiary Fold of Globular Proteins. *Biochemical Society Transactions* **10**:299–301.
- Summers N L, W D Carlson and M Karplus 1987. Analysis of Side-Chain Orientations in Homologous Proteins. *Journal of Molecular Biology* **196**:175–198.
- Taylor W R and C A Orengo 1989. Protein Structure Alignment. *Journal of Molecular Biology* **208**:1–22.
- Thompson J D, D G Higgins and T J Gibson 1994. CLUSTAL W: Improving the Sensitivity of Progressive Multiple Sequence Alignment Through Sequence Weighting, Position-specific Gap Penalties and Weight Matrix Choice. *Nucleic Acids Research* **22**:4673–4680.
- Turk D, H W Hoeffken, D Grosse, J Stuerzebecher, P D Martin, B F P Edwards and W Bode 1992. Refined 2.3 Ångstroms X-Ray Crystal Structure of Bovine Thrombin Complexes Formed with the 3 Benzamidine and Arginine-Based Thrombin Inhibitors NAPAP, 4-TAPAP and MQPA: A Starting Point for Improving Antithrombotics. *Journal of Molecular Biology* **226**:1085–1099.
- Turk D, J Stuerzebecher and W Bode 1991. Geometry of Binding of the *N*-Alpha-Tosylated Piperidides of *meta*-Amidino-Phenylalanine, *Para* Amidino-Phenylalanine and *para*-Guanidino-Phenylalanine to Thrombin and Trypsin – X-ray Crystal Structures of Their Trypsin Complexes and Modeling of their Thrombin Complexes. *FEBS Letters* **287**:133–138.
- Vasquez M 1996. Modeling Side-chain Conformation. *Current Opinion in Structural Biology* **6**:217–221.
- Wilmot C M and J M Thornton 1988. Analysis and Prediction of the Different Types of β -turn in Proteins. *Journal of Molecular Biology* **203**:221–232.
- Zarembinski T I, L-W Hung, H-J Mueller-Dieckmann, K-K Kim, H Yokota, R Kim and S-H Kim 1998. Structure-based Assignment of the Biochemical Function of a Hypothetical Protein: A Test Case of Structural Genomics. *Proceedings of the National Academy of Sciences USA* **95**:15189–15193.

CHAPTER ELEVEN

Four Challenges in Molecular Modelling: Free Energies, Solvation, Reactions and Solid-state Defects

In this chapter we shall consider four important problems in molecular modelling. First, we discuss the problem of calculating free energies. We then consider continuum solvent models, which enable the effects of the solvent to be incorporated into a calculation without requiring the solvent molecules to be represented explicitly. Third, we shall consider the simulation of chemical reactions, including the important technique of *ab initio* molecular dynamics. Finally, we consider how to study the nature of defects in solid-state materials.

11.1 Free Energy Calculations

11.1.1 The Difficulty of Calculating Free Energies by Computer

The free energy is often considered to be the most important quantity in thermodynamics. The free energy is usually expressed as the Helmholtz function, A , or the Gibbs function, G . The Helmholtz free energy is appropriate to a system with constant number of particles, temperature and volume (constant NVT), whereas the Gibbs free energy is appropriate to constant number of particles, temperature and pressure (constant NPT). Most experiments are conducted under conditions of constant temperature and pressure, where the Gibbs function is the appropriate free energy quantity.

Unfortunately, the free energy is a difficult quantity to obtain for systems such as liquids or flexible macromolecules that have many minimum energy configurations separated by low-energy barriers. Associated quantities such as the entropy and the chemical potential are also difficult to calculate. As we showed in Section 6.3, the free energy cannot be accurately determined from a 'standard' molecular dynamics or Monte Carlo simulation, because such simulations do not adequately sample from those regions of phase space that make important contributions to the free energy. Specifically, we showed that the Helmholtz free energy is given by:

$$A = k_B T \ln \left(\iint d\mathbf{p}^N d\mathbf{r}^N \exp \left(-\frac{\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{k_B T} \right) \rho(\mathbf{p}^N, \mathbf{r}^N) \right) \quad (11.1)$$

The term $\exp[+\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)/k_B T]$ makes important contributions to the integral. However, a simulation using either Monte Carlo or molecular dynamics sampling seeks out the *lower-energy* regions of phase space. Such simulations will never adequately sample the important high-energy regions and so to calculate the free energy using a conventional simulation will lead to poorly converged and inaccurate values. The grand canonical and particle-insertion methods do provide a route to the free energy, but they are not applicable to many of the systems of interest, which contain complex molecules at high densities.

11.2 The Calculation of Free Energy Differences

Let us consider a closely related but slightly different problem: the calculation of the free energy difference of two states. As an example, we will consider the problem of calculating the free energy difference between ethanol ($\text{CH}_3\text{CH}_2\text{OH}$) and ethane thiol ($\text{CH}_3\text{CH}_2\text{SH}$) in water. As we shall see, this is a problem that can be tackled using methods that use Monte Carlo or molecular dynamics sampling. Three methods have been proposed for calculating free energy differences: thermodynamic perturbation, thermodynamic integration and slow growth. We shall consider each of these in turn.

11.2.1 Thermodynamic Perturbation

Consider two well-defined states X and Y. For example, X could be a system comprising a molecule of ethanol in a periodic box of water and Y could be ethane thiol in water. X contains N particles interacting according to the Hamiltonian \mathcal{H}_X . Y contains N particles interacting according to \mathcal{H}_Y . The free energy difference (ΔA) between the two states is as follows:

$$\Delta A = A_Y - A_X = -k_B T \ln \frac{Q_Y}{Q_X} \quad (11.2)$$

$$\Delta A = -k_B T \left\{ \frac{\iint d\mathbf{p}^N d\mathbf{r}^N \exp[-\mathcal{H}_Y(\mathbf{p}^N, \mathbf{r}^N)/k_B T]}{\iint d\mathbf{p}^N d\mathbf{r}^N \exp[-\mathcal{H}_X(\mathbf{p}^N, \mathbf{r}^N)/k_B T]} \right\} \quad (11.3)$$

Substituting 1 in the form $\exp[+\mathcal{H}_X(\mathbf{p}^N, \mathbf{r}^N)/k_B T] \exp[-\mathcal{H}_X(\mathbf{p}^N, \mathbf{r}^N)/k_B T]$ into the numerator gives:

$$\Delta A = -k_B T \left\{ \frac{\iint d\mathbf{r}^N d\mathbf{p}^N \exp\left(-\frac{\mathcal{H}_Y(\mathbf{r}^N, \mathbf{p}^N)}{k_B T}\right) \exp\left(+\frac{\mathcal{H}_X(\mathbf{r}^N, \mathbf{p}^N)}{k_B T}\right) \exp\left(-\frac{\mathcal{H}_X(\mathbf{r}^N, \mathbf{p}^N)}{k_B T}\right)}{\iint d\mathbf{r}^N d\mathbf{p}^N \exp\left(-\frac{\mathcal{H}_X(\mathbf{r}^N, \mathbf{p}^N)}{k_B T}\right)} \right\} \quad (11.4)$$

Equation (11.4) can be written in terms of an ensemble average, as follows:

$$\begin{aligned} \Delta A &= -k_B T \left\{ \frac{\iint d\mathbf{p}^N d\mathbf{r}^N \exp[-\mathcal{H}_Y(\mathbf{p}^N, \mathbf{r}^N)/k_B T] \exp[+\mathcal{H}_X(\mathbf{p}^N, \mathbf{r}^N)/k_B T] \exp[-\mathcal{H}_X(\mathbf{p}^N, \mathbf{r}^N)/k_B T]}{\iint d\mathbf{p}^N d\mathbf{r}^N \exp[-\mathcal{H}_X(\mathbf{p}^N, \mathbf{r}^N)/k_B T]} \right\} \\ &= -k_B T \langle \exp[-\mathcal{H}_Y(\mathbf{p}^N, \mathbf{r}^N) - \mathcal{H}_X(\mathbf{p}^N, \mathbf{r}^N)/k_B T] \rangle_0 \end{aligned} \quad (11.5)$$

The subscript 0 indicates averaging over the ensemble of configurations representative of the initial state X. If the averaging is over the ensemble corresponding to the final state Y (indicated by the subscript 1) then we are effectively simulating the reverse process, from which ΔA can be determined by:

$$\Delta A = k_B \ln \langle \exp[-(\mathcal{H}_X - \mathcal{H}_Y)/k_B T] \rangle_1 \quad (11.6)$$

This approach to the calculation of free energy differences, Equation (11.6), is generally attributed to Zwanzig [Zwanzig 1954]. To perform a thermodynamic perturbation calculation we must first define \mathcal{H}_Y and \mathcal{H}_X and then run a simulation at the state X, forming the ensemble average of $\exp[-(\mathcal{H}_Y - \mathcal{H}_X)/k_B T]$ as we proceed. Analogously, we could run a simulation at the state Y and obtain the ensemble average of $\exp[-(\mathcal{H}_X - \mathcal{H}_Y)/k_B T]$. Thus, if X corresponds to ethanol and Y to ethane thiol, the free energy difference could be obtained from a simulation of ethanol in a periodic box of water as follows. For each configuration we calculate the value of the energy for every instantaneous conformation of ethanol in which the oxygen atom is temporarily assigned the potential energy parameters of sulphur. Alternatively, we could simulate ethane thiol and for each configuration calculate the energy of the system in which the sulphur is 'mutated' into oxygen.

If X and Y do not overlap in phase space then the value of the free energy difference calculated using Equation (11.6) will not be very accurate, because we will not adequately sample the phase space of Y when simulating X. This problem arises when the energy difference between the two states is much larger than $k_B T$: $|\mathcal{H}_Y - \mathcal{H}_X| \gg k_B T$. How then can we obtain accurate estimates of the free energy difference under such circumstances? Consider what happens if we introduce a state that is intermediate between X and Y, with a Hamiltonian \mathcal{H}_1 and a free energy $A(1)$:

$$\begin{aligned} \Delta A &= A(Y) - A(X) \\ &= (A(Y) - A(1)) + (A(1) - A(X)) \\ &= -k_B T \ln \left[\frac{Q(Y)}{Q(1)} \cdot \frac{Q(1)}{Q(X)} \right] \\ &= -k_B T \ln \langle \exp[-(\mathcal{H}_Y - \mathcal{H}_1)/k_B T] \rangle - k_B T \ln \langle \exp[-(\mathcal{H}_1 - \mathcal{H}_X)/k_B T] \rangle \end{aligned} \quad (11.7)$$

If we define region 1 so that it overlaps with X and Y we may improve the sampling and obtain a more reliable value. This is shown in Figure 11.1.

An obvious extension is to use several different intermediate states in progressing from \mathcal{H}_X to \mathcal{H}_Y :

$$\begin{aligned} \Delta A &= A(Y) - A(X) \\ &= (A(Y) - A(N)) + (A(N) - A(N-1)) + \dots \\ &\quad + (A(2) - A(1)) + (A(1) - A(X)) \\ &= -k_B T \ln \left[\frac{Q(Y)}{Q(N)} \cdot \frac{Q(N)}{Q(N-1)} \cdot \frac{Q(N-1)}{Q(N-2)} \dots \frac{Q(2)}{Q(1)} \frac{Q(1)}{Q(X)} \right] \end{aligned} \quad (11.8)$$

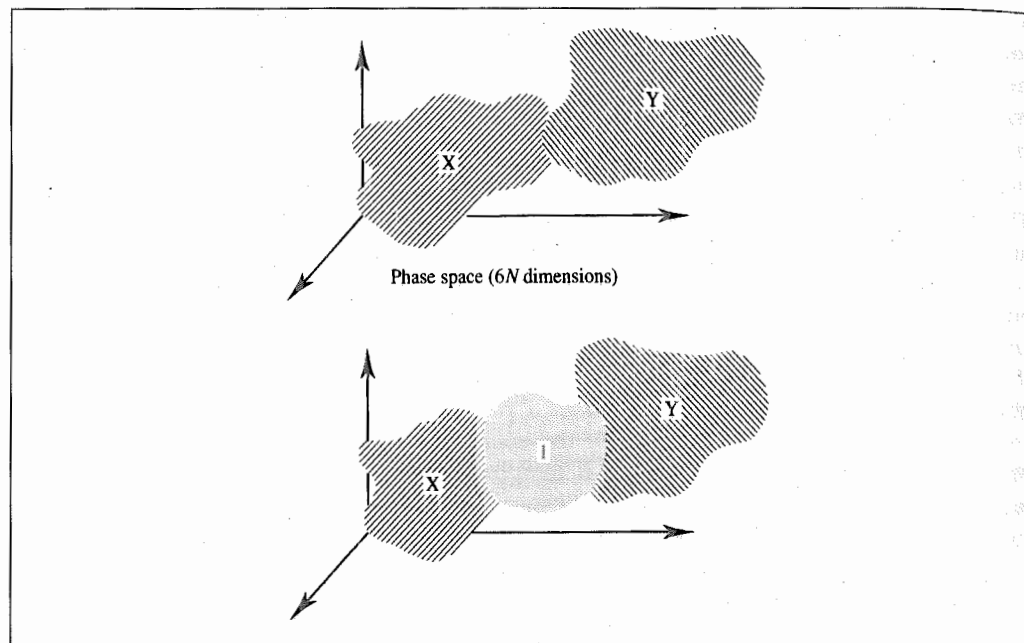


Fig. 11.1: An intermediate state (labelled I) can improve the degree of overlap in phase space and lead to improved sampling.

The important point to notice is that the intermediate terms cancel out and so we are free to choose as many intermediate states as are necessary to get good overlaps and thus reliable values of the free energy differences.

11.2.2 Implementation of Free Energy Perturbation

Suppose we are using an empirical energy function such as the following to describe the inter- and intramolecular interactions in our ethanol/ethane thiol system:

$$\begin{aligned} \mathcal{V}(\mathbf{r}^N) = & \sum_{\text{bonds}} \frac{k_i}{2} (l_i - l_{i,0})^2 + \sum_{\text{angles}} \frac{k_i}{2} (\theta_i - \theta_{i,0})^2 + \sum_{\text{torsions}} \frac{V_n}{2} (1 + \cos(n\omega - \gamma)) \\ & + \sum_{i=1}^N \sum_{j=i+1}^N \left(4\epsilon_{ij} \left[\left(\frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] + \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} \right) \end{aligned} \quad (11.9)$$

The force-field model for ethanol contains C–O and O–H bond-stretching contributions; in ethane thiol these are replaced by C–S and S–H parameters. Similarly, in ethanol there will be angle-bending terms due to C–O–H, C–C–O and H–C–O angles; in ethane thiol these will be C–S–H, C–C–S and H–C–S. The torsional contribution will be modified appropriately, as will the van der Waals and electrostatic interactions (both those within the

solute and between the solute and solvent). The partial atomic charges for all of the atoms in ethanol may all be different from those for ethane thiol.

The relationship between the initial, final and intermediate states is usefully described in terms of a *coupling parameter*, λ . As λ is changed from 0 to 1, the Hamiltonian varies from \mathcal{H}_X to \mathcal{H}_Y . Each of the terms in the force field for an intermediate state λ can be written as a linear combination of the values for X and Y:

$$1. \text{ Bonds:} \quad k_i(\lambda) = \lambda k_i(Y) + (1 - \lambda)k_i(X) \quad (11.10)$$

$$l_0(\lambda) = \lambda l_0(Y) + (1 - \lambda)l_0(X) \quad (11.11)$$

$$2. \text{ Angles:} \quad k_\theta(\lambda) = \lambda k_\theta(Y) + (1 - \lambda)k_\theta(X) \quad (11.12)$$

$$\theta_0(\lambda) = \lambda \theta_0(Y) + (1 - \lambda)\theta_0(X) \quad (11.13)$$

$$3. \text{ Dihedrals:} \quad v_\omega(\lambda) = \lambda v_\omega(Y) + (1 - \lambda)v_\omega(X) \quad (11.14)$$

$$4. \text{ Electrostatics:} \quad q_i(\lambda) = \lambda q_i(Y) + (1 - \lambda)q_i(X) \quad (11.15)$$

$$5. \text{ van der Waals:} \quad \epsilon(\lambda) = \lambda \epsilon(Y) + (1 - \lambda)\epsilon(X) \quad (11.16)$$

$$\sigma(\lambda) = \lambda \sigma(Y) + (1 - \lambda)\sigma(X) \quad (11.17)$$

For each value of λ (λ_i) a simulation is performed (using either Monte Carlo or molecular dynamics as appropriate) with the appropriate force field parameters. First, the system is equilibrated using the force field parameters appropriate to λ_i . A production phase is then performed during which the free energy difference $\Delta A(\lambda_i \rightarrow \lambda_{i+1})$ is accumulated as $-k_B T \ln \langle \exp(-\Delta \mathcal{H}_i / k_B T) \rangle$, where $\Delta \mathcal{H}_i = \mathcal{H}_{i+1} - \mathcal{H}_i$. The total free energy change for $\lambda = 0$ to $\lambda = 1$ is then the sum of the free energy changes for the various values of λ_i , as shown in Figure 11.2.

The approach that we have described so far is known as *forward sampling*, because the free energy is determined for $\lambda_i \rightarrow \lambda_{i+1}$. In *backward sampling*, the free energy difference between λ_i and λ_{i-1} is determined. The coupling parameter λ still increases from 0 to 1; it is just that the free energies are accumulated in a different manner. In *double-wide sampling*, the free energy differences for both $\lambda_i \rightarrow \lambda_{i+1}$ and $\lambda_i \rightarrow \lambda_{i-1}$ are obtained from a simulation as

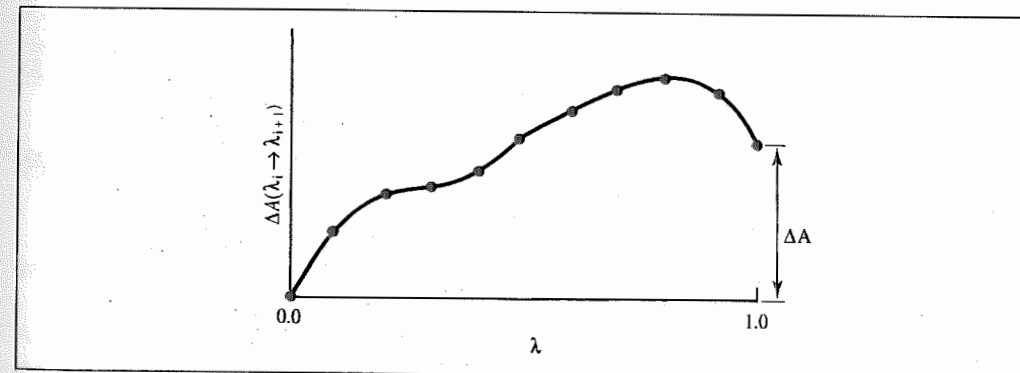


Fig. 11.2: Calculation of the free energy difference using perturbation.

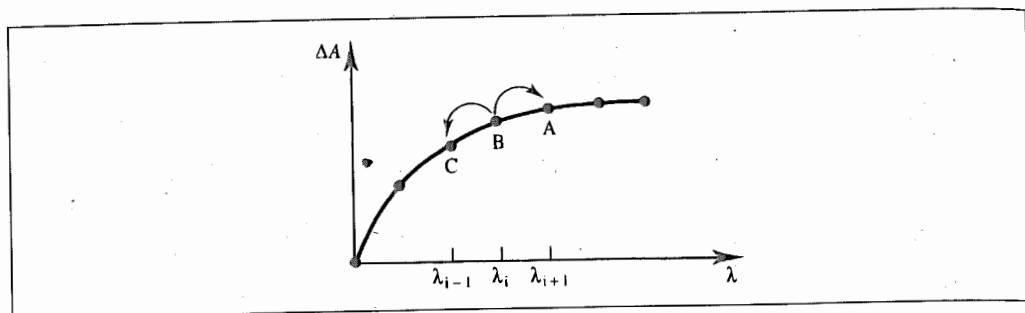


Fig. 11.3: Double wide sampling enables two free energies to be accumulated from a single simulation.

illustrated in Figure 11.3. Consider point B in Figure 11.3, which corresponds to a coupling parameter λ_i . A simulation performed using λ_i can be used to obtain values for both the free energy difference $\Delta A(\lambda_i \rightarrow \lambda_{i+1})$ and the free energy difference $\Delta A(\lambda_i \rightarrow \lambda_{i-1})$. This is clearly a more efficient way to obtain the desired free energy as twice as many free energy differences can be obtained from a single simulation.

11.2.3 Thermodynamic Integration

An alternative way to calculate the free energy difference uses thermodynamic integration. The formula for the free energy difference is derived in Appendix 11.1 and is:

$$\Delta A = \int_{\lambda=0}^{\lambda=1} \left\langle \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{\partial \lambda} \right\rangle_{\lambda} d\lambda \quad (11.18)$$

To calculate a free energy difference using thermodynamic integration we thus need to determine the integral in Equation (11.18). In practice, this is achieved by performing a series of simulations corresponding to discrete values of λ between 0 and 1. For each value of λ , the average of

$$\left\langle \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{\partial \lambda} \right\rangle_{\lambda} \quad (11.19)$$

is determined. These partial derivatives are calculated analytically in some programs but in others a finite difference approximation is used ($\partial \mathcal{H} / \partial \lambda \approx \Delta \mathcal{H} / \Delta \lambda$). The total free energy difference ΔA is then equal to the area under the graph of

$$\left\langle \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{\partial \lambda} \right\rangle_{\lambda} \quad (11.20)$$

versus λ , as illustrated in Figure 11.4.

11.2.4 The 'Slow Growth' Method

A third approach for the calculation of free energy differences from computer simulation is the slow growth method. Here, the Hamiltonian changes by a very small, constant amount at

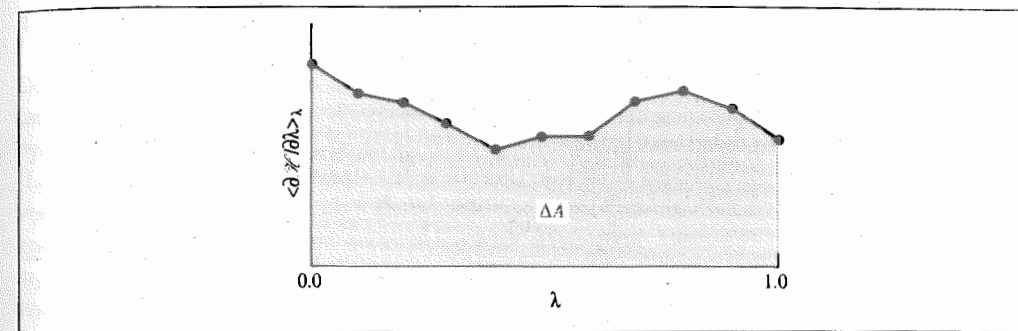


Fig. 11.4: Calculation of free energy differences by thermodynamic integration.

each step of the calculation. This means that at each stage the Hamiltonian $\mathcal{H}(\lambda_{i+1})$ is very nearly equal to $\mathcal{H}(\lambda_i)$. The free energy difference is given by:

$$\Delta A = \sum_{i=1; \lambda=0}^{i=N_{\text{step}}; \lambda=1} (\mathcal{H}_{i+1} - \mathcal{H}_i) \quad (11.21)$$

This expression is derived in Appendix 11.2.

In principle, all three methods for calculating the free energy difference should give the same result, as the free energy is a state function and so independent of path. However, there may be practical reasons for choosing one method or another, as we shall discuss in Section 11.6. The other point to note at this stage is that our formulation of the free energy has been in terms of the partition function Q and the Hamiltonian $\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)$, which have contributions from both kinetic and potential energies. When the kinetic energy contributions are integrated out they cancel and so the various equations can be written in terms of the difference between the potential function, $\mathcal{V}(\mathbf{r}^N)$, rather than the Hamiltonian, $\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)$. Q is then replaced by the configurational integral, Z , and the free energy values that are obtained are excess free energies, relative to an ideal gas.

Our discussion so far has considered the calculation of Helmholtz free energies, which are obtained by performing simulations at constant NVT . For proper comparison with experimental values we usually require the Gibbs free energy, G . Gibbs free energies are obtained from a simulation at constant NPT .

11.3 Applications of Methods for Calculating Free Energy Differences

11.3.1 Thermodynamic Cycles

An early application of the free energy perturbation method was the determination of the free energy required to create a cavity in a solvent. Postma, Berendsen and Haak determined the free energy to create a cavity ($\lambda = 1$) in pure water ($\lambda = 0$) using isothermal-isobaric

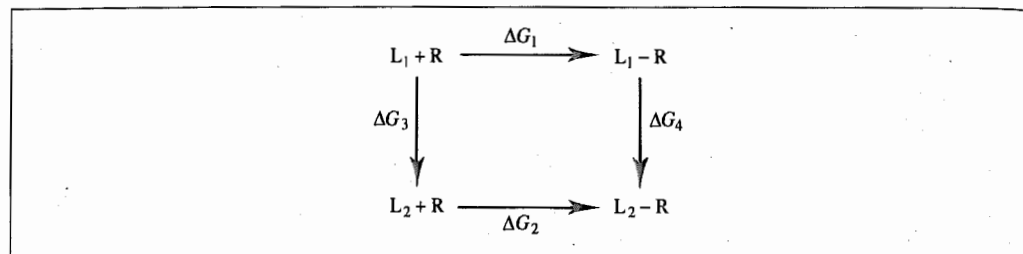


Fig. 11.5: Thermodynamic cycle for binding ligands L_1 and L_2 to receptor R .

molecular dynamics simulations [Postma *et al.* 1982]. Five different cavity sizes were considered and the results showed that, as expected, the free energy of cavity formation increased with the size of the cavity, the results being in good agreement with analytical theories. For small cavities ($< 1 \text{ \AA}$ radius) the results were inaccurate due to poor sampling. The calculations provided not only the free energy of cavity formation for the different cavity sizes but also structural and dynamic properties of the water molecules around the cavity. For example, the water structure varied with the cavity size. A cavity of radius 1.78 \AA had the most pronounced shell structure, with a high first-neighbour peak and a significant second-neighbour peak in the cavity-water pair distribution function.

Many processes of interest to the molecular modeller involve an equilibrium between molecules that interact via non-covalent forces, the free energy being related to the equilibrium constant by $\Delta G = -RT \ln K$. Let us consider the binding of two different ligands (L_1 and L_2) to a receptor molecule (R). L_1 and L_2 could be putative inhibitors of an enzyme R or two 'guests' for a host R . The thermodynamic cycle for the two binding processes is shown in Figure 11.5. The relative binding affinity of L_1 and L_2 equals $\Delta G_2 - \Delta G_1$ and is commonly written $\Delta\Delta G$. In principle, it would be possible to calculate values of ΔG_1 and ΔG_2 by simulating the actual association process. To do this we would bring the ligand and the receptor together from an initial large separation to gradually form the intermolecular complex. However, in most cases this would involve such a major reorganisation of the receptor, the ligand and the solvent that it would be difficult to ensure adequate sampling of phase space.

The free energy is a state function, and so its value round a thermodynamic cycle must be zero. Thus $\Delta G_2 - \Delta G_1 = \Delta G_4 - \Delta G_3$ (Figure 11.5). ΔG_3 corresponds to the free energy difference of the two ligands in solution; ΔG_4 is the free energy difference of the two intermolecular complexes. The changes ΔG_3 and ΔG_4 do not correspond to any transformation that can be performed in the laboratory, but they are quite feasible in the computer. The free energy difference only depends upon the endpoints, and so we are at liberty to change the Hamiltonians in any way we wish. The free energy differences obtained from such non-physical pathways are likely to be much more reliable than the 'physically plausible' processes as they should involve much less reorganisation of the system. This is particularly so if the two ligands L_1 and L_2 have similar structures. To calculate the relative free energy of binding of the two ligands we would therefore 'mutate' L_1 into L_2 in solution and L_1 to L_2 within the receptor. This is the *thermodynamic cycle perturbation approach* to calculating relative free energies.

11.3.2 Applications of the Thermodynamic Cycle Perturbation Method

One of the first applications of the thermodynamic cycle perturbation approach to the calculation of relative binding constants was the study by Lybrand, McCammon and Wipff of the synthetic macrocycle SC24, which, when protonated, can bind halide ions (Figure 11.6) [Lybrand *et al.* 1986]. SC24 binds Cl^- 4.30 kcal/mol more strongly than Br^- . Two simulations were performed to determine a theoretical value for this relative free energy using the free energy perturbation method with molecular dynamics. First, Cl^- was mutated to Br^- in aqueous solution, giving a free energy difference of 3.35 kcal/mol . The same mutation was then performed within the macrocycle, in a periodic box of water. The value obtained for this step was 7.50 kcal/mol , giving an overall relative free energy of binding of 4.15 kcal/mol . The experimental value was approximately 4.3 kcal/mol . Thus, although the free energy to desolvate Cl^- is unfavourable compared with Br^- , this is more than compensated for by favourable interactions between Cl^- and the host; Br^- is slightly too large to fit comfortably in the relatively inflexible SC24 molecule.

One of the most attractive applications of the free energy techniques is for predicting the relative free energies of binding of inhibitors of biological macromolecules such as proteins or DNA. If we know the binding constant of an inhibitor then we can, in principle at least, calculate the binding constant of a related inhibitor. The free energy cycle used to perform this calculation is analogous to that shown in Figure 11.5: we perform two separate free energy calculations: ligand L_1 is mutated to ligand L_2 in solution and within the binding site. An early calculation of this type was performed by Bash and colleagues, who studied two inhibitors of thermolysin (an enzyme which cleaves the amide bonds in peptides and proteins) [Bash *et al.* 1987]. The two inhibitors investigated had the general formula carboxybenzoxy-Gly^P(X)-L-Leu-L-Leu [Barlett and Marlowe 1987] (Figure 11.7). The experimentally determined binding constants (K_i) of the $X \equiv \text{NH}$ and $X \equiv \text{O}$ inhibitors were 9.1 nM and 9000 nM , i.e. the former binds 1000 times more strongly. This difference in binding constants is equivalent to 4.1 kcal/mol . X-ray crystallographic analysis showed that the two inhibitors bind in almost identical positions. The calculated free energy difference was determined to

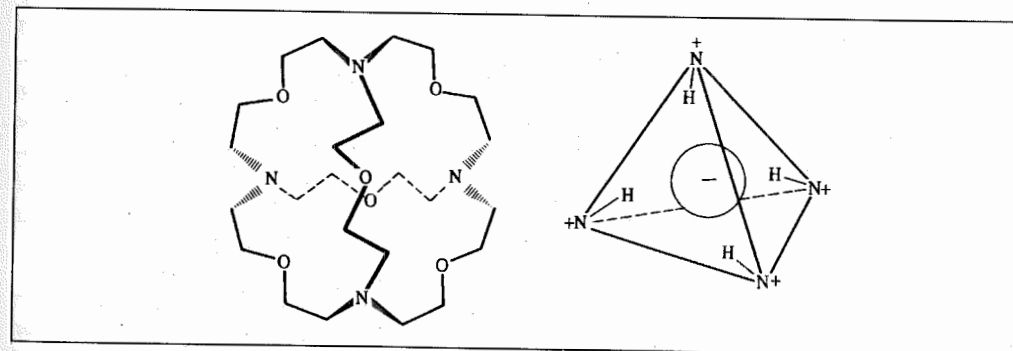


Fig. 11.6: The SC24/halide system. (Figure adapted from Lybrand T P, J A McCammon and G Wipff 1986. Theoretical Calculation of Relative Binding Affinity in Host-Guest Systems. Proceedings of the National Academy of Sciences USA 83:833-835.)

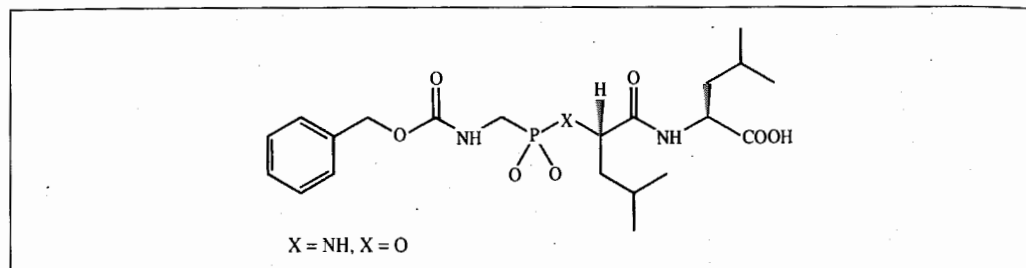


Fig. 11.7: Thermolysin inhibitors [Bartlett and Marlowe 1987].

be 4.2 ± 0.5 kcal/mol, in good agreement with the experimental result. In the active site of the enzyme the X group of the inhibitor interacts with the backbone carbonyl oxygen of one of the amino acids (Ala 113). The ester oxygen interacts unfavourably with this carbonyl, but the amide can form a hydrogen bond. The relative free energy of binding of the amide inhibitor to the protein was calculated to be 7.6 kcal/mol lower than the ester, but this was counteracted by the difference in the free energies of solvation, which was calculated to be 3.4 kcal/mol. The amide inhibitor thus incurs a greater desolvation penalty than the ester.

This study obviously gave very satisfactory agreement with the experimental data. However, a subsequent calculation by Merz and Kollman showed that the results were very sensitive to the charge model used for the inhibitor [Merz and Kollman 1989]. The charges for the inhibitor were obtained by electrostatic potential fitting in each case, though with different basis sets. This second calculation gave a free energy difference of 5.9 kcal/mol. Other studies have also shown that calculated free energies can be very sensitive to the charge model used; we will discuss some of the problems with performing free energy calculations in Section 11.6.

As one final example of the application of free energy calculations we will examine the determination of relative partition coefficients. The partition coefficient (P) is the equilibrium constant for the transfer of a solute between two solvents. The logarithm of the partition coefficient ($\log P$) for transfer between water and a variety of solvents (primarily 1-octanol) is widely used to derive structure-activity relationships (see Section 12.9), in which the biological activity of a molecule is correlated with its physicochemical properties. The thermodynamic cycle for the partition of two solutes, A and B, between two solvents is shown in Figure 11.8. If it were possible to calculate the free energy of transfer from one solvent to another (i.e. ΔG_1 or ΔG_2 in Figure 11.8) then this would give the partition coefficient directly. However, such a simulation would require an inordinate amount of time and probably be very inaccurate. A relative partition coefficient can be determined by mutating one solute into the other in the two separate solvents.

Calculations of relative partition coefficients have been reported using the free energy perturbation method with the molecular dynamics and Monte Carlo simulation methods. For example, Essex, Reynolds and Richards calculated the difference in partition coefficients of methanol and ethanol partitioned between water and carbon tetrachloride with molecular dynamics sampling [Essex *et al.* 1989]. The results agreed remarkably well with experiment

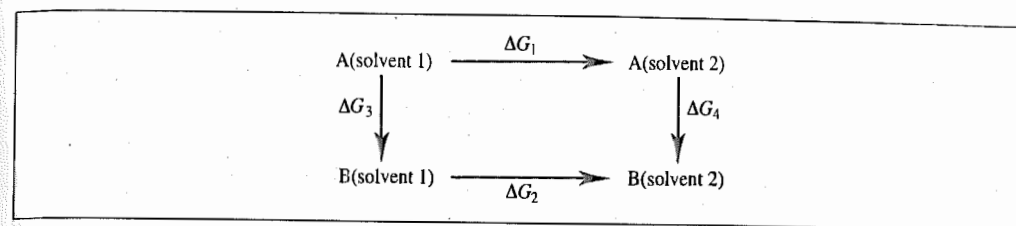


Fig. 11.8: Thermodynamic cycle for calculating relative partition coefficients.

(within 0.06 kcal/mol). Jorgensen, Briggs and Contreras used Monte Carlo methods to calculate the relative partition coefficients for eight pairs of solutes (including methanol/methylamine, acetic acid/acetamide and pyrazine/pyridine) between water and chloroform [Jorgensen *et al.* 1990]. For these eight systems good qualitative agreement with experimental data was obtained. However, the results involving acetic acid gave too broad a spread of values. This was traced to the relative free energies of hydration, which varied over too wide a range and indicated some areas for improvement in the force field model.

11.3.3 The Calculation of Absolute Free Energies

In some cases, it is possible to devise thermodynamic cycles which enable the absolute free energy of a change to be determined using free energy perturbation methods [Jorgensen *et al.* 1988]. Figure 11.9 shows a thermodynamic cycle for the association of L and R to give a complex LR in both the gas phase and in solution. ΔG_{ass} is the free energy of association in solution and is given by:

$$\Delta G_{\text{ass}} = \Delta G_{\text{gas}}(L + R \rightarrow LR) + \Delta G_{\text{sol}}(LR) - \Delta G_{\text{sol}}(L) - \Delta G_{\text{sol}}(R) \quad (11.22)$$

$\Delta G_{\text{sol}}(X)$ is the solvation free energy of the species X (the free energy of transfer from the gas phase to solvent). The solvation free energy can be written in terms of perturbations where the species disappear to nothing in the gas phase and in solution, $\Delta G_{\text{sol}}(X) = \Delta G_{\text{gas}}(X \rightarrow 0) - \Delta G_{\text{sol}}(X \rightarrow 0)$. The free energy of association, ΔG_{ass} , can then be written:

$$\begin{aligned} \Delta G_{\text{ass}} &= \Delta G_{\text{gas}}(L + R \rightarrow LR) - \Delta G_{\text{gas}}(L \rightarrow 0) + \Delta G_{\text{sol}}(L \rightarrow 0) \\ &\quad - \Delta G_{\text{gas}}(R \rightarrow 0) + \Delta G_{\text{sol}}(R \rightarrow 0) + \Delta G_{\text{gas}}(LR \rightarrow 0) \\ &\quad - \Delta G_{\text{sol}}(LR \rightarrow 0) \end{aligned} \quad (11.23)$$

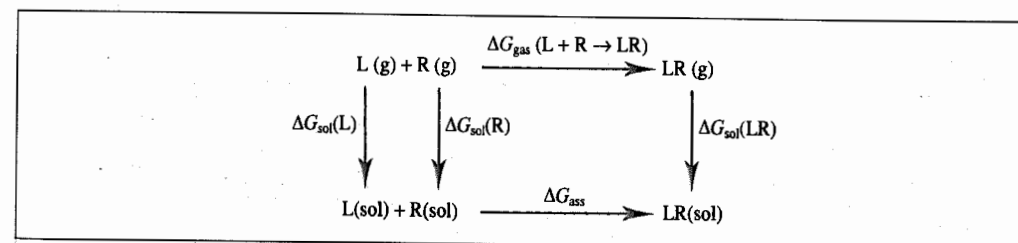


Fig. 11.9: Thermodynamic cycle used to calculate absolute free energies [Jorgensen *et al.* 1988].

The gas-phase terms cancel and $\Delta G_{\text{sol}}(\text{LR} \rightarrow 0)$ can be written as the sum of two separate calculations:

$$\Delta G_{\text{sol}}(\text{LR} \rightarrow 0) = \Delta G_{\text{sol}}(\text{LR} \rightarrow \text{R}) + \Delta G_{\text{sol}}(\text{R} \rightarrow 0) \quad (11.24)$$

Thus, the overall free energy change can be written:

$$\Delta G_{\text{ass}} = \Delta G_{\text{sol}}(\text{L} \rightarrow 0) - \Delta G_{\text{sol}}(\text{LR} \rightarrow \text{R}) \quad (11.25)$$

We thus need perform only two simulations, L to nothing in water and L to nothing in the LR complex. The first application of this approach was to the association of two methane molecules in water, where both species (L and R) are identical. In general, L should be chosen as the smaller component.

11.4 The Calculation of Enthalpy and Entropy Differences

Free energy changes can now be routinely calculated with errors of less than 1 kcal/mol in favourable cases. How does this compare with the error with which the enthalpy or entropy difference can be determined? One way to determine the enthalpy change would be to perform two separate simulations, one of the initial system and one of the final system. For example, the difference in the enthalpy of solvation of ethanol and ethane thiol in water could be determined by simulating the two species separately and then taking the difference in the total enthalpies of the two systems. These total energies are invariably large numbers, with relatively large errors. The error in the calculated enthalpy difference would be comparable in magnitude to the error in the energy of each system. By contrast, the free energy is determined solely in terms of interactions involving the solute. This means that the free energy can be calculated much more accurately. More efficient ways to calculate the enthalpy and entropy change have been proposed for use with both free energy perturbation and thermodynamic integration schemes [Fleischman and Brooks 1987; Yu and Karplus 1988]. The uncertainties in the enthalpies and entropies so calculated are better than would be obtained by subtracting the differences in total energies, but they are still about one order of magnitude larger than the corresponding free energies.

11.5 Partitioning the Free Energy

The overall free energy can be partitioned into individual contributions if the thermodynamic integration method is used [Boresch *et al.* 1994; Boresch and Karplus 1995]. The starting point is the thermodynamic integration formula for the free energy:

$$\Delta A = \int_{\lambda=0}^{\lambda=1} \left\langle \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{\partial \lambda} \right\rangle_{\lambda} d\lambda \quad (11.26)$$

The Hamiltonian can be written as a sum of contributions from bond stretching, angle bending, and so on:

$$\left\langle \frac{\partial \mathcal{H}(\lambda)}{\partial \lambda} \right\rangle_{\lambda} = \left\langle \frac{\partial \mathcal{H}_{\text{bonds}}(\lambda)}{\partial \lambda} + \frac{\partial \mathcal{H}_{\text{angles}}(\lambda)}{\partial \lambda} + \dots \right\rangle_{\lambda} \quad (11.27)$$

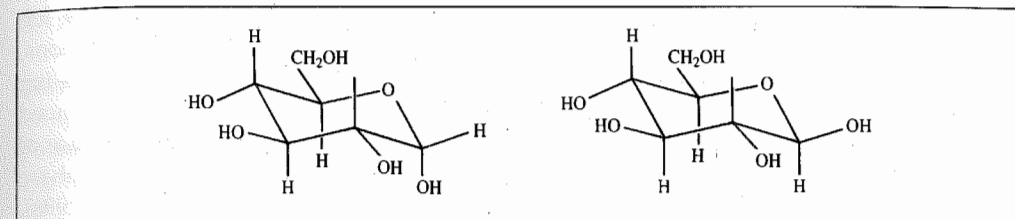


Fig. 11.10: The α and β anomers of D-glucose.

So the free energy is given by:

$$\begin{aligned} \Delta A &= \int_{\lambda=0}^{\lambda=1} \left\langle \frac{\partial \mathcal{H}_{\text{bonds}}(\lambda)}{\partial \lambda} \right\rangle_{\lambda} d\lambda + \int_{\lambda=0}^{\lambda=1} \left\langle \frac{\partial \mathcal{H}_{\text{angles}}(\lambda)}{\partial \lambda} \right\rangle_{\lambda} d\lambda + \dots \\ &= \Delta A_{\text{bonds}} + \Delta A_{\text{angles}} + \dots \end{aligned} \quad (11.28)$$

We should remember that only the sum of the contributions is truly meaningful, as the individual contributions are not state functions. This has led some to criticise any use of such partitioning schemes [Smith and van Gunsteren 1994b], though they may be useful to indicate which interactions contribute the most to the overall free energy, and may also suggest the source of most of the error in the calculation. It is not possible to perform such a partitioning using thermodynamic perturbation.

An example of this partitioning scheme is the study by Ha and colleagues of the anomeric equilibrium between the α and β anomers of D-glucose [Ha *et al.* 1991]. D-glucose can exist in two tautomeric forms: α -D-glucose, in which the C₁ hydroxyl group is axial; and β -D-glucose, in which it is equatorial (Figure 11.10). In the gas phase, the axial α isomer is more stable than the equatorial β isomer, due to the anomeric effect which is considered to arise from unfavourable dipole-dipole interactions and delocalisation of the lone pair on the ring oxygen into an anti-bonding σ^* orbital. However, the β -D (equatorial) anomer is more stable than the α -D (axial) anomer by 0.3 kcal/mol in aqueous solution. The free energy difference between the two isomers in water was calculated by Ha *et al.* using both free energy perturbation and thermodynamic integration to be -0.3 ± 0.4 kcal/mol for $\beta \rightarrow \alpha$. A partitioning of the free energy showed that this small difference arose from the cancellation of two large terms: the α isomer was predicted to be 3.6 kcal/mol more favourable than the β isomer in the gas phase, due mainly to electrostatic effects. However, the β isomer was favoured over the α isomer in aqueous solution, again due to electrostatic effects, such as the enhanced hydrogen-bonding capability of the β isomer with the solvent. The small free energy difference, the difficulties of obtaining a reliable force field model and the large number of accessible conformations makes this equilibrium particularly difficult to tackle. One way in which the sampling problem has been tackled is by the use of a method called locally enhanced sampling (LES), which uses multiple copies of those parts of the system that can exist in more than one conformation. In the case of glucose these are the hydroxyl hydrogens and the hydroxymethyl group. In LES, each of the copies does not interact with the other copies of the same group and each atom 'sees' the mean force from all the copies. The first application of the technique was the study of the diffusion of

carbon monoxide within the protein myoglobin [Elber and Karplus 1990], but there are many other potential applications. LES reduces the barriers to conformational transitions, which leads to more rapid transitions between the conformational minima [Simmerling and Elber 1995]. However, the free energies calculated for the LES energy surface do need to be corrected to give the corresponding result for the single-copy system. In the case of glucose, it was found that the α isomer was favoured by 0.5–1.0 kcal/mol in the gas phase but that the β isomer was favoured in solution by 0.2 kcal/mol [Simmerling *et al.* 1998]. The gas-phase result was suggested to be a consequence of the tendency of the O–C–O–C linkage to adopt a gauche conformation, whereas the solution result was due to solvation effects. Many quantum mechanical studies have also been performed on this system, some of which have also included solvation effects (using the methods to be discussed in Section 11.10.2). For example, Barrows and co-workers performed high-level *ab initio* calculations on 11 varied low-energy conformations on glucose using large basis sets and including the effects of electron correlation [Barrows *et al.* 1998]. From this data, they calculated a gas-phase equilibrium constant of 0.4 kcal/mol in favour of the α isomer (using a Boltzmann-weighted average), whereas the equivalent value in solution was 0.6 kcal/mol in favour of the β isomer.

Another common practice is to partition the free energy into contributions from van der Waals and electrostatic interactions. This can be achieved rather easily by first perturbing the electrostatic and then the van der Waals parameters. One system that has been studied in this way is the biotin/streptavidin complex. This protein–ligand complex is of particular interest due to the extremely strong association constant (–18.3 kcal/mol). The chemical structure of biotin is shown in Figure 11.11. The separate electrostatic and van der Waals free-energy calculations suggested that the largest contribution to the very negative free energy of binding was due to the non-polar van der Waals contribution rather than the electrostatic component [Miyamoto and Kollman 1993a, b]. Despite the presence of many hydrogen bonds between the ligand and the protein in the complex it was suggested that, whilst there was a large and favourable electrostatic interaction between biotin and streptavidin, this was almost cancelled by the free energy of interaction of biotin with water. By contrast, the van der Waals interaction gave a much greater contribution in the protein–ligand complex than for the ligand in water, so leading to its dominance. Indeed, the ligand is almost completely buried within the protein cavity, as can be seen in the structure of the intermolecular complex shown in Figure 11.12 (colour plate section).

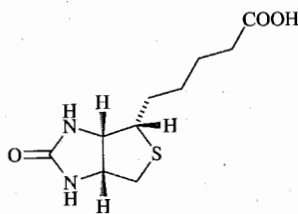


Fig. 11.11: Biotin.

11.6 Potential Pitfalls with Free Energy Calculations

There are two major sources of error associated with the calculation of free energies from computer simulations. Errors may arise from inaccuracies in the Hamiltonian, be it the potential model chosen or its implementation (the treatment of long-range forces, etc.). The second source of error arises from an insufficient sampling of phase space.

Unfortunately, there is no set recipe that guarantees adequate coverage of phase space and thus reliable free energy values [Mitchell and McCammon 1991]. The errors associated with inadequate sampling may be identified by running the simulation for longer periods of time (molecular dynamics) or for more iterations (Monte Carlo); the perturbation can be performed in both forward and reverse directions; a different scheme could be used to determine the free energy difference (e.g. thermodynamic perturbation and thermodynamic integration). At the very least, the simulation should be run in both directions; the difference in the calculated free energy values (often referred to as the *hysteresis*) gives a lower-bound estimate of the error in the calculation.

One possible pitfall to be aware of when estimating errors is that an excessively short simulation may give an almost zero difference between the forward and reverse directions. If the time of the simulation is much longer than the relaxation time of the system then the change can be performed reversibly. If the simulation time is of the same order of magnitude as the relaxation time then one would expect a significant degree of hysteresis. However, if the simulation is much shorter than the relaxation time then approximately zero hysteresis may result, due to the inability of the system to adjust to the changes. In such a situation, the free energies for both forward and reverse directions may be approximately the same, but quite likely incorrect.

11.6.1 Implementation Aspects

The allure of methods for calculating free energies and their associated thermodynamic values such as equilibrium constants has resulted in considerable interest in free energy calculations. A number of decisions must be made about the way that the calculation is performed. One obvious choice concerns the simulation method. In principle, either Monte Carlo or molecular dynamics can be used; in practice, molecular dynamics is almost always used for systems where there is a significant degree of conformational flexibility, whereas Monte Carlo can give very good results for small molecules which are either rigid or have limited conformational freedom.

One must choose from the thermodynamic perturbation, thermodynamic integration and slow growth methods. Each of these methods has been extensively used, but the slow growth method is not now recommended. This method suffers from a phenomenon known as 'Hamiltonian lag'; the system never has time to properly equilibrate for a given value of the coupling parameter, because the potential function changes at every step. An additional advantage of the integration and perturbation approaches is that, should one decide at the end of a simulation that more sampling needs to be done for particular values of λ , or that more λ values are required over a particular range, then this can

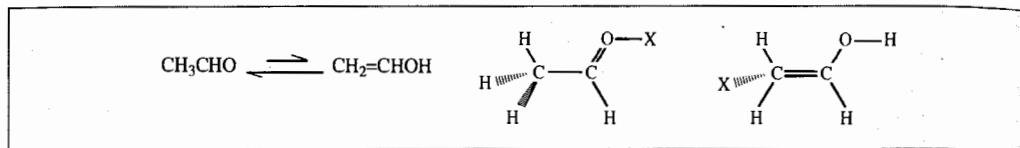


Fig. 11.13: To calculate the free energy difference between the aldehyde and enol forms of acetaldehyde, the single topology method uses dummy atoms (X).

easily be done without losing information from other parts of the calculation. With slow growth, one would have to redo the simulation from scratch.

Prior to the calculation, the increment $\delta\lambda$ in the coupling parameter must be specified. Traditionally, $\delta\lambda$ is set to a constant value before the simulation commences. It is important that there is enough overlap between successive states λ_i and λ_{i+1} so that reliable values can be obtained. An alternative approach is to use small changes in λ when the free energy is changing quickly and a larger change in λ when the free energy is changing more slowly. This is the basis of a method called *dynamically modified windows*, in which the slope of the free energy versus λ curve is used to determine the value of $\delta\lambda$ to use in the next iteration [Pearlman and Kollman 1989].

As the free energy is a thermodynamic state function the free energy difference between the initial and final states should be independent of the path along which the change is made, so long as it is reversible. It may be possible to proceed from the initial to the final state along more than one pathway. A change that involves high energy barriers will require much smaller increments to be made in the coupling parameter λ to ensure reversibility than a pathway that proceeds via a lower barrier.

Many free energy calculations involve changes in the molecular topologies of the species concerned; there are often different numbers of atoms in the initial and final states, and the atoms may be bonded in different ways. For example, suppose we wish to determine the free energy difference between acetaldehyde and its enol (Figure 11.13), in which a hydrogen atom migrates from the methyl carbon atom to the carbonyl oxygen. The system can be represented in the calculation using either a 'single' topology or a 'dual' topology. In the single-topology method, the molecular topology at all stages is the union of the initial and final states, using dummy atoms where necessary. A dummy atom does not interact with the other atoms in the system. Thus the hydrogen atom bonded to the oxygen in the enol form would be represented as a dummy atom when the simulation reached the endpoint corresponding to the aldehyde as shown in Figure 11.13.

The alternative to the single-topology representation is the dual-topology method. Here, both the molecular topologies are maintained during the entire simulation, such that both species 'exist' (in a topological sense) but do not interact with each other. The Hamiltonian that describes the interaction between these groups and the environment can be described in a number of ways, the simplest of which is the linear relationship:

$$\mathcal{H}(\lambda) = \lambda\mathcal{H}_Y + (1 - \lambda)\mathcal{H}_X \quad (11.29)$$

Many free energy calculations involve the creation or annihilation of atoms. A potential problem with such simulations is that a singularity may occur in the function for which

an ensemble average is to be formed. One way to try to deal with this is to scale the initial Hamiltonian by a factor λ^n (rather than just λ) and the final Hamiltonian by $(1 - \lambda)^n$ (rather than $(1 - \lambda)$). It can be shown that for Monte Carlo simulations the singularity problem can be dealt with, provided n is at least 4 [Buetler *et al.* 1994]. However, a molecular dynamics simulation requires not only the energies to be calculated but also the first and the second derivatives. If λ^n scaling is used then either a steadily decreasing time step must be used as λ approaches zero, or these regions of the simulation must be omitted altogether and their contributions estimated by extrapolation. An alternative approach is to replace the traditional Lennard-Jones interaction with a soft-core potential of the following form [Buetler *et al.* 1994; Liu *et al.* 1996]:

$$v_{ij}^{\text{LJ}} = 4\varepsilon_{ij} \left(\frac{\sigma_{ij}^{12}}{[\alpha_{\text{LJ}}\sigma_{ij}^6 + r_{ij}^6]^2} - \frac{\sigma_{ij}^6}{(\alpha_{\text{LJ}}\sigma_{ij}^6 + r_{ij}^6)} \right) \quad (11.30)$$

where ε_{ij} and σ_{ij} have their usual Lennard-Jones meanings. The parameter α_{LJ} determines the 'softness' of the interaction, which has the effect of making the interaction approach a finite value as the interatomic distance r_{ij} goes to zero. With a suitable choice of the parameter α_{LJ} it is possible to ensure that the position of the minimum in the soft-core potential coincides with that of the unscaled energy curve (Figure 11.14). When used to perturb the system from X at $\lambda = 0$ to Y at $\lambda = 1$ then the soft-core Lennard-Jones interaction between the particle i that is being perturbed and some other particle j at a distance r_{ij} varies as:

$$v_{ij}^{\text{LJ}}(\lambda) = 4(1 - \lambda)\varepsilon_X \left(\frac{\sigma_X^{12}}{[\alpha_{\text{LJ}}\lambda^2\sigma_X^6 + r_{ij}^6]^2} - \frac{\sigma_X^6}{[\alpha_{\text{LJ}}\lambda^2\sigma_X^6 + r_{ij}^6]} \right) + 4\lambda\varepsilon_Y \left(\frac{\sigma_Y^{12}}{[\alpha_{\text{LJ}}(1 - \lambda)^2\sigma_Y^6 + r_{ij}^6]^2} - \frac{\sigma_Y^6}{[\alpha_{\text{LJ}}(1 - \lambda)^2\sigma_Y^6 + r_{ij}^6]} \right) \quad (11.31)$$

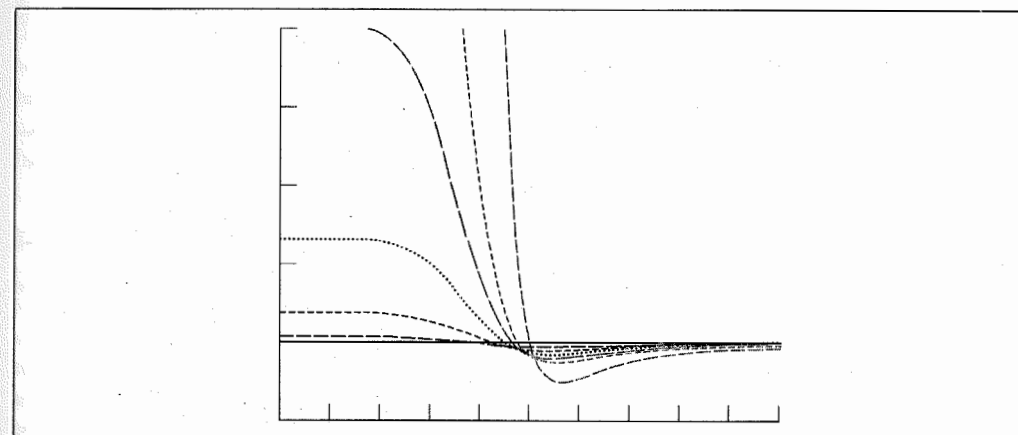


Fig. 11.14: Comparison of scaled and unscaled Lennard-Jones potentials (Equation (11.31)) for the case where a particle disappears at $\lambda = 0$. As λ decreases the curves get progressively closer to the x axis.

A similar soft-core expression can also be derived for the electrostatic interactions. In the situation where the particle disappears then only one of the terms remains (i.e. the second term if the particle disappears at $\lambda = 0$). The effect of using this type of soft-core potential can be seen in simulations of protein–ligand systems, where the ligand ‘disappears’ and is replaced by one or more solvent molecules. In a normal perturbation calculation, decreasing the effective radius of the ligand atoms will give rise to a collapse of the protein cavity and disruption to the surrounding protein structure. By contrast, placing the soft-core interaction sites at locations where the atoms are to be created or deleted maintains the protein cavity, because the solvent molecules are able to actually pass through the ligand as it is annihilated. This feature of soft-core potentials can also be useful for other kinds of free energy calculation where it is desired to try to simultaneously derive the relative free energies of binding of several ligands to a receptor and also in certain types of simulated annealing structure refinement.

11.7 Potentials of Mean Force

The free energy changes that we have considered so far correspond to chemical ‘mutations’. We may also be interested to know how the free energy changes as a function of some inter- or intramolecular coordinate, such as the distance between two atoms, or the torsion angle of a bond within a molecule. The free energy surface along the chosen coordinate is known as a *potential of mean force* (PMF). When the system is in a solvent, the potential of mean force incorporates solvent effects as well as the intrinsic interaction between the two particles. Potentials of mean force were introduced in our discussion of Langevin dynamics (Section 7.8), where we noted that the ratio of *trans* to *gauche* conformers of 1,2-dichloroethane was significantly different in the liquid than in an isolated molecule. Unlike the mutations so common in free energy perturbation calculations, which are often along non-physical pathways, the potential of mean force is calculated for a physically achievable process. Consequently, the point of highest energy on the free energy profile that is obtained from a PMF calculation corresponds to the transition state for the process, from which it is possible to derive kinetic quantities such as rate constants.

Various methods have been proposed for calculating potentials of mean force. The simplest type of PMF is the free energy change as the separation (r) between two particles is changed. We might anticipate that we could calculate the potential of mean force from the radial distribution function using the following expression for the Helmholtz free energy:

$$A(r) = -k_B T \ln g(r) + \text{constant} \quad (11.32)$$

The constant is often chosen so that the most probable distribution corresponds to a free energy of zero.

Unfortunately, the potential of mean force may vary by several multiples of $k_B T$ over the relevant range of the parameter r . The logarithmic relationship between the potential of mean force and the radial distribution function means that a relatively small change in the free energy (i.e. a small multiple of $k_B T$) may correspond to $g(r)$ changing by an order of magnitude from its most likely value. Unfortunately, standard Monte Carlo or molecular

dynamics simulation methods do not adequately sample regions where the radial distribution function differs drastically from the most likely value, leading to inaccurate values for the potential of mean force. The traditional way to avoid this problem uses a technique called *umbrella sampling*.

11.7.1 Umbrella Sampling

Umbrella sampling attempts to overcome the sampling problem by modifying the potential function so that the unfavourable states are sampled sufficiently. The method can be used with both Monte Carlo and molecular dynamics simulations. The modification of the potential function can be written as a perturbation:

$$\mathcal{V}'(\mathbf{r}^N) = \mathcal{V}(\mathbf{r}^N) + W(\mathbf{r}^N) \quad (11.33)$$

where $W(\mathbf{r}^N)$ is a weighting function, which often takes a quadratic form:

$$W(\mathbf{r}^N) = k_W (\mathbf{r}^N - \mathbf{r}_0^N)^2 \quad (11.34)$$

For configurations that are far from the equilibrium state \mathbf{r}_0^N the weighting function will be large and so a simulation using the modified energy function $\mathcal{V}'(\mathbf{r}^N)$ will be biased along some relevant ‘reaction coordinate’ away from the configuration \mathbf{r}_0^N . The resulting distribution will, of course, be non-Boltzmann. The corresponding Boltzmann averages can be extracted from the non-Boltzmann distribution using a method introduced by Torrie and Valleau [Torrie and Valleau 1977]. The result is:

$$\langle A \rangle = \frac{\langle A(\mathbf{r}^N) \exp[+W(\mathbf{r}^N)/k_B T] \rangle_W}{\langle \exp[+W(\mathbf{r}^N)/k_B T] \rangle_W} \quad (11.35)$$

The subscript W indicates that the average is based on the probability $P_W(\mathbf{r}^N)$, which in turn is determined by the modified energy function $\mathcal{V}'(\mathbf{r}^N)$. For example, to obtain the potential of mean force via the radial distribution function (Equation (11.32)) the distribution function with the forcing potential would be determined and then corrected to give the ‘true’ radial distribution function, from which the free energy can be calculated as a function of the separation. It is usual to perform an umbrella sampling calculation in a series of stages, each of which is characterised by a particular value of the coordinate and an appropriate value of the forcing potential $W(\mathbf{r}^N)$. However, if the forcing potential is too large, the denominator in Equation (11.35) is dominated by contributions from only a few configurations with especially large values of $\exp[W(\mathbf{r}^N)]$ and the averages take too long to converge.

To illustrate the use of umbrella sampling, let us consider how the technique has been used to determine the potential of mean force for rotation of the central C–C bond of butane in aqueous solution. The barrier between the *trans* and *gauche* conformations of butane is approximately 3.5 kcal/mol, which is sufficiently high to give sampling problems in simulations. For example, in the molecular dynamics simulation of Ryckaert and Bellemans the mean time between *gauche*–*trans* transitions was about 10 ps [Ryckaert and Bellemans 1978]. Jorgensen, Gao and Ravimohan used umbrella sampling with Monte Carlo simulations to calculate the potential of mean force as the central bond in butane is rotated in a periodic box of water molecules, to determine the effect of the solvent on the relative

populations of the different conformations [Jorgensen *et al.* 1985]. The results predicted a shift in the expected populations of *trans* and *gauche* isomers from 68% *trans* in the gas phase to 54% in aqueous solution, a change of 14%. In addition, the barrier height was reduced in solution. Jorgensen and colleagues performed many calculations on similar systems using umbrella sampling and Monte Carlo simulations; he recommended that to reduce the barriers to a value between 1 kcal/mol and 3 kcal/mol was appropriate. In some cases, it is possible to use a barrier height of zero, though the barriers cannot be reduced too severely as this makes the forcing potential too large.

It is also possible to calculate potentials of mean force using the free energy perturbation method with a molecular dynamics or Monte Carlo simulation. As usual, the calculation is broken into a series of steps that are characterised by a coupling parameter λ . With molecular dynamics, holonomic constraint methods are used to fix the desired coordinates without affecting the dynamic motion of the system. This is the essence of the extension of the SHAKE procedure by Tobias and Brooks to cope with general coordinate changes [Tobias and Brooks 1988] (see Section 7.5). In a Monte Carlo simulation the required coordinates are simply fixed at the desired value(s). This contrasts with umbrella sampling, in which the coordinate(s) of interest would be able to vary over their range of values throughout the simulation, subjected to a potential that has been modified using the forcing function. At each step of the perturbation calculation, the difference in the energy between the configuration and the configuration that corresponds to $\lambda + \delta\lambda$ is determined and the free energy accumulated in the appropriate way.

To compare the perturbation and umbrella sampling methods for calculating potentials of mean force, Jorgensen and Buckner repeated the PMF calculation for butane in water using the perturbation method [Jorgensen and Buckner 1987]. The *gauche* population was calculated to increase by 12.3% using this method, in accordance with the previous umbrella sampling calculations. Jorgensen put forward several arguments in favour of the perturbation approach. A major concern with umbrella sampling is that a proper sampling of the phase space may not be achieved. In some cases, the presence of bottlenecks in phase space may be identified if separate simulations starting from different configurations give different results, but even this approach is not fail-safe as all simulations may encounter the same problem. Indeed, Jorgensen suggested that just such a bottleneck may have occurred in a previous simulation of pentane in water using umbrella sampling (which involved 5 million Monte Carlo steps). The only real problem with the perturbation method is the need to choose an appropriate value of $\delta\lambda$ so that there is adequate overlap between the configuration corresponding to λ and that corresponding to $\lambda + \delta\lambda$. Jorgensen and Buckner varied the central torsion angle in their simulation of butane using 15° increments.

11.7.2 Calculating the Potential of Mean Force for Flexible Molecules

To calculate a potential of mean force using free energy perturbation (or indeed umbrella sampling) it is necessary to determine the pathway for the transition of interest. This is trivial for simple problems such as the separation of two particles or the rotation of butane but can be quite complicated for more detailed changes such as conformational interconversions.

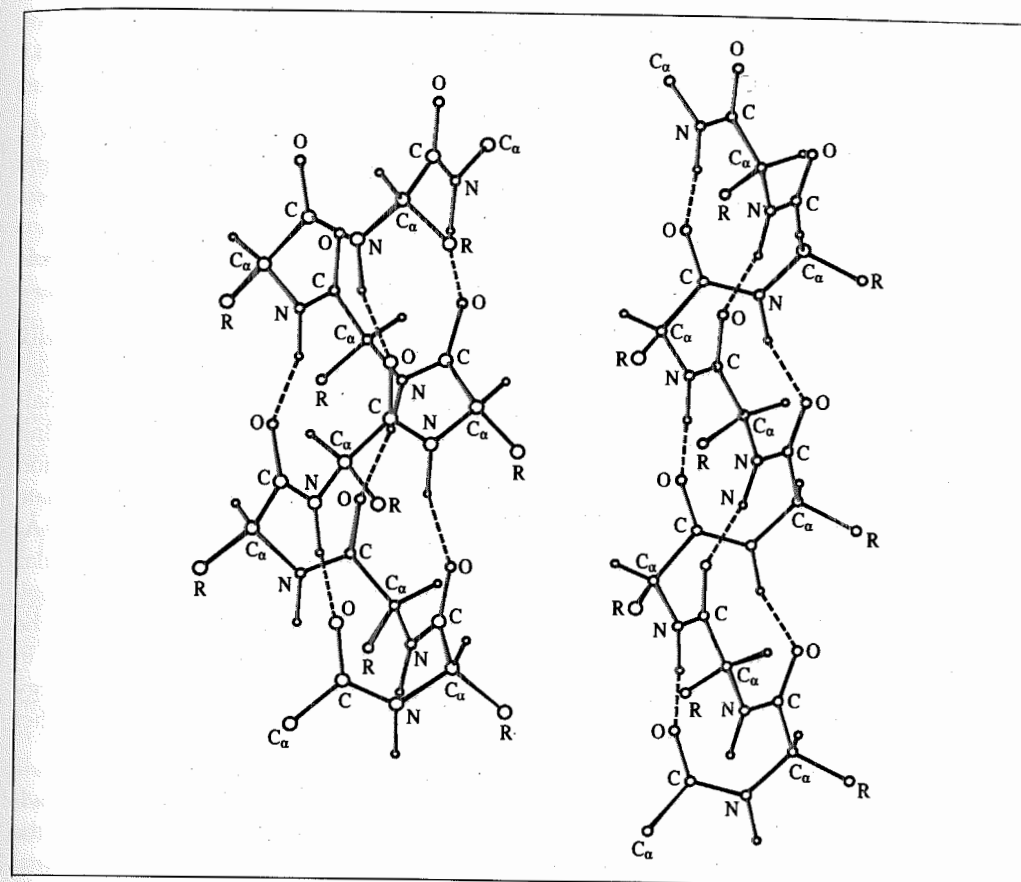


Fig. 11.15: The α -helix (left) and the 3_{10} -helix.

The reaction path methods discussed in Section 5.9.3 may be helpful in determining these pathways.

To illustrate the calculation of potentials of mean force for flexible systems we will consider helical conformations of polypeptide chains. We have already met the α -helix, which is commonly observed in protein structures (see Section 10.2). In this conformation, hydrogen bonds are formed between residues i and $i + 4$. Polypeptide chains can also form a different type of helix, called a 3_{10} -helix. Here, the hydrogen bonds are formed between residues i and $i + 3$. These two helices are compared in Figure 11.15. The backbone conformations of such helices do not differ significantly: the α -helix has backbone torsion angles ($\phi = -60^\circ$, $\psi = -50^\circ$) and the 3_{10} -helix has ($\phi = -50^\circ$, $\psi = -28^\circ$). The 3_{10} -helix is found to a small extent in protein structures, usually at the ends of α -helices. However, the 3_{10} -helix is much more common in peptides formed from α,α -dialkyl amino acids, which have two alkyl substituents at the α -carbon atom. The prototypical member of this class of amino acids is α -methylalanine (MeA; see Figure 11.16). Peptides containing this amino acid can

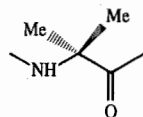


Fig. 11.16: Methylalanine.

form both α -helices and 3_{10} -helices with the actual conformation present being rather sensitive to the conditions; such peptides are 3_{10} -helical in CDCl_3 and α -helical in $(\text{CD}_3)_2\text{SO}$.

To calculate the potential of mean force for interconverting the α -helix to the 3_{10} -helix requires an appropriate reaction coordinate to be determined. Here we describe the calculations of three groups who all used different approaches. Smythe, Huston and Marshall studied a decamer of α -methylalanine, $\text{CH}_3\text{CO-MeA}_{10}\text{-NMe}$, using umbrella sampling [Smythe *et al.* 1993, 1995]. They used the self-penalty walk method described in Section 5.9.3 to determine the transition pathway and observed that the reaction coordinate correlated well with a smooth change in the end-to-end distance from the 3_{10} -helix (19 Å) to the α -helix (13 Å). Their umbrella sampling calculations were performed using molecular dynamics, with the end-to-end distance being subjected to a restraining potential. Simulations were performed in various solvents: in water, the free energy change for the α -helix \rightarrow 3_{10} -helix transition was calculated to be 7.6 kcal/mol, with the value in dichloromethane being 5.8 kcal/mol and *in vacuo* 3.2 kcal/mol. Although a distinct energy barrier was found for the vacuum calculations, no transition barrier was found for either solution calculation.

Zhang and Hermans studied a 10-residue alanine peptide as well as a 10-residue α -methylalanine peptide, *in vacuo* and in water [Zhang and Hermans 1994]. In their calculations, the transition from one conformation to the other was performed using a restraining potential that forced the structure to exchange one set of hydrogen bonds to the set of hydrogen bonds appropriate to the other structure. This additional potential function could be used to drive the molecule back and forth between the two conformations by varying a coupling parameter, λ , between 0 and 1. Free energy profiles were determined for the α -helix to 3_{10} -helix transition using molecular dynamics and the slow growth method. The results showed that the alanine peptide had a clear preference for the α -helix both *in vacuo* and in water but that the free energy change for the MeA peptide was approximately zero in water and that the 3_{10} -helix was preferred *in vacuo*. It was proposed by Zhang and Hermans that the discrepancy between these results for the α -methylalanine peptide and those obtained by Smythe, Huston and Marshall was probably due to the different force field models employed; Smythe *et al.* used a united atom model, whereas Zhang and Hermans used an all-atom model.

Tirado-Reeves, Maxwell and Jorgensen used yet another approach for calculating the potential of mean force, this time for an undecaalanine peptide in water [Tirado-Reeves *et al.* 1993]. The free energy profile was calculated using the perturbation method and Monte Carlo simulations by gradually varying the ψ backbone torsion angles, keeping the ϕ torsion angles fixed at -60° . The free energy difference between the two conformations

was calculated to be 10.6 kcal/mol in favour of the α -helix, with a small activation barrier of 2.8 kcal/mol for the 3_{10} - to α -helical transition. *In vacuo*, a larger free energy difference was predicted (13.6 kcal/mol).

These three studies have been described at some length, in part to illustrate the different approaches available for calculating thermodynamic properties of complex systems but also to emphasise the fact that different methods can give quite different (and sometimes contradictory) results. Such comparative studies serve to highlight the fact that it is necessary to examine critically the methods and models used in a calculation. All three studies were in part prompted by experimental electron spin resonance results that suggested that a 16-residue alanine-based peptide adopted a 3_{10} -helical conformation in water [Miick *et al.* 1992]. These results were contradicted by all the simulations, and indeed prompted Smythe and Marshall to undertake similar experiments on their conformationally constrained peptides, experiments which showed that these peptides were α -helical, in agreement with the calculations.

11.8 Approximate/'Rapid' Free Energy Methods

Free energy calculations are notoriously time-consuming to perform. Whilst one might have anticipated that ever faster computers would have made significant inroads on this problem, in some respects the opposite has happened, as researchers are now able to more fully quantify the need for sufficient sampling of phase space and to attain better convergence. In addition, of course, there is a natural desire to investigate ever larger systems. A practical illustration of the dilemma facing the proponents of free energy methods as a predictive tool, at least in an industrial environment, is that, if the calculation takes longer to perform than a candidate molecule can be synthesised and tested, then there is little practical benefit from attempting the calculation. There has thus been continued interest in the development of alternative methods which, whilst still being based upon 'exact' statistical mechanics, are intended to provide answers with less computational effort than a full-blown free energy calculation. These methods tend to approach the problem from one of two perspectives. Some, such as the λ -dynamics method, enable a single simulation to provide information on a number of molecules. Others, such as the linear response method, aim to limit the amount of simulation that needs to be performed.

In a traditional free energy calculation a coupling parameter, λ , provides the link between the initial and final systems. In most free energy calculations λ varies uniformly from 0 to 1 (or from 1 to 0), one exception being the dynamically modified windows technique discussed in Section 11.6.1. By contrast, the λ -dynamics technique considers λ to be another 'particle' in the simulation, with its own fictitious mass. As such, λ -dynamics is similar in some respects to those charge calculation schemes where the charges can vary as a dynamic variable (see Section 4.9.6). Specifically, λ corresponds to the reaction coordinate along which the potential would be modified in an umbrella sampling calculation. The advantage of making this association is that the biasing potentials that are used in the umbrella sampling method can be used in the λ -dynamics technique to provide enhanced sampling in relevant regions of configurational space. Indeed, it is possible to

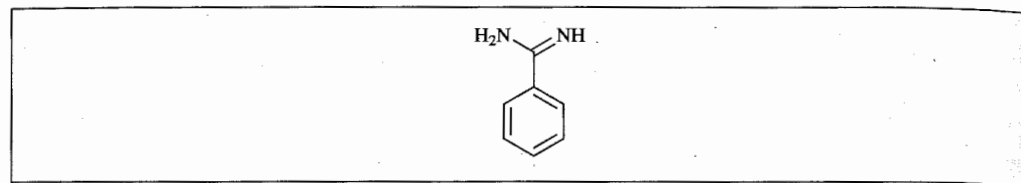


Fig. 11.17: Benzamidine.

use a set of coupling variables, λ_i , $i = 1, \dots, n$. In the case where we have just one molecule being changed into another, then these different λ_i represent changes in different components of the interaction potential (for example, the Lennard-Jones and Coulombic interactions). If λ_1 and λ_2 refer to the Coulombic and the van der Waals interactions, respectively, then the potential function can be written:

$$\mathcal{V}(\mathbf{r}^N, \lambda_1, \lambda_2) = (1 - \lambda_1)\mathcal{V}_A^{\text{coul}} + \lambda_1\mathcal{V}_B^{\text{coul}} + (1 - \lambda_2)\mathcal{V}_A^{\text{L-J}} + \lambda_2\mathcal{V}_B^{\text{L-J}} + \mathcal{V}_{\text{env}}(\mathbf{r}^N) \quad (11.36)$$

Here, the term $\mathcal{V}_{\text{env}}(\mathbf{r}^N)$ corresponds to all interactions involving those parts of the system that are not changing (i.e. the solvent and the unchanging part of the solute). The lambda variables move under the influence of a specific term which serves to limit their absolute extent (i.e. between 0 and 1) and which can be used to restrict their value to particular ranges during the simulation in order to provide enhanced sampling at particular points.

The basic λ -dynamics scheme can be used to perform a 'regular' type of free energy calculation in which one solute is perturbed into another such as the perturbation of methanol to ethane or to methane thiol [Kong and Brooks 1996]. However, it can also be used to investigate a number of perturbations simultaneously. As such, it provides a route to assess several free energies from a single simulation. One published example concerns the binding of benzamidine derivatives to the enzyme trypsin [Guo and Brooks 1998; Guo *et al.* 1998]. Benzamidine is shown in Figure 11.17; this molecule binds relatively strongly to the enzyme because the positively charged amidine group interacts with a negatively charged aspartate residue in the protein. However, substitution at the para position can affect the strength of binding, with *p*-amino benzamidine binding slightly more strongly, *p*-methyl slightly more weakly and *p*-chloro more weakly still than the parent molecule. When λ -dynamics is applied to this problem, each of the L ligands ($L = 4$ in this case) is represented by a different value of λ_i . Initially, all values of λ_i are set to $1/L$ and their velocities set to zero. This means that each molecule is set on an equal footing at the beginning of the calculation. The system then evolves under the influence of the following hybrid potential:

$$\mathcal{V}(\mathbf{r}^N, \lambda_i) = \sum_{i=1}^L \lambda_i^2 (\mathcal{V}_i(\mathbf{r}^{\text{int}}) - F_i) + \mathcal{V}_{\text{env}}(\mathbf{r}^N) \quad (11.37)$$

As in Equation (11.36), $\mathcal{V}_{\text{env}}(\mathbf{r}^N)$ corresponds to those interactions concerning all atoms not directly involved in the perturbations, whereas $\mathcal{V}_i(\mathbf{r}^{\text{int}})$ concerns those atoms associated with the group being perturbed in ligand i (for which the associated lambda parameter is λ_i). F_i is a reference free energy and can serve two purposes. If F_i equals the solvation/desolvation free energy of the relevant ligand then the free energy value obtained from

the calculation corresponds to the free energy change for the full cycle. F_i can also be used as a biasing potential to control the sampling in particular regions of phase space. Finally, there is a constraint on the values of λ_i :

$$\sum_{i=1}^L \lambda_i^2 = 1 \quad (11.38)$$

As the simulation proceeds, the values of λ_i fluctuate, subject to the constraint in Equation (11.38). The free energy difference between two molecules i and j can be determined by identifying the probability that each molecule occupies the state $\lambda_i = 1$ or $\lambda_j = 1$, respectively. Thus:

$$\Delta\Delta A_{ij} = -\frac{1}{k_B T} \ln \left[\frac{P(\lambda_i = 1, \lambda_{m \neq i} = 0)}{P(\lambda_j = 1, \lambda_{n \neq j} = 0)} \right] \quad (11.39)$$

These relative probabilities can be easily determined by simply counting the number of times during the simulation that the relevant value of lambda reaches unity. In the case of the para-substituted benzamidines it was possible after only a relatively short simulation (110 ps) to observe that the *p*-chloro and *p*-methyl derivatives were significantly weaker than the *p*-amino and the parent compound (Figure 11.18). In this particular case, all four

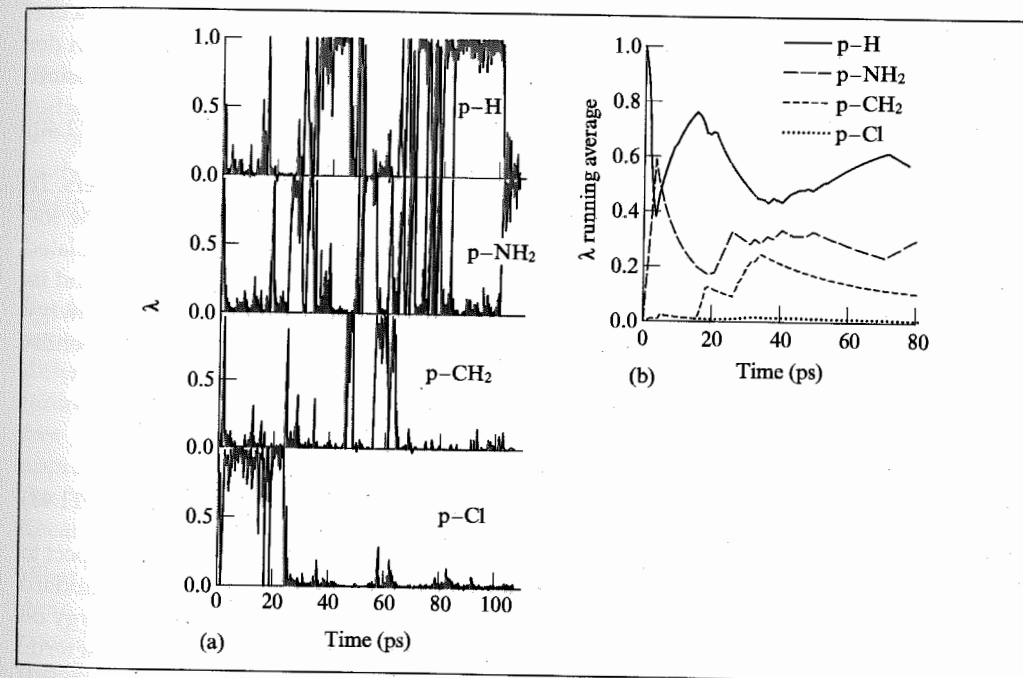


Fig. 11.18: λ -dynamics simulation of benzamidine derivatives binding to trypsin. (a) The larger the value of λ the stronger the interaction with the protein at that instant. (b) Running average of each value of λ over the course of the simulation. (Figure redrawn from Guo Z and C L Brooks III 1998. Rapid Screening of Binding Affinities: Application of the λ -Dynamics Method to a Trypsin-Inhibitor System. Journal of the American Chemical Society 120:1920-1921.)

inhibitors have rather similar binding affinities (within 1 kcal/mol), thus requiring a relatively long simulation to separate them. In general, a compound with a binding affinity more than 3 kcal/mol worse than the most favourable molecule should be screened out within a few tens of picoseconds, though longer simulation times would be required to provide a correct rank ordering.

Conceptually similar to the lambda-dynamics approach is the so-called 'chemical-Monte Carlo/molecular dynamics' method [Pitera and Kollman 1998; Eriksson *et al.* 1999], which also considers many molecules simultaneously. In this approach, molecular dynamics is used to sample the coordinate space with Monte Carlo moves sampling the various chemical states. To avoid possible problems associated with hybrid states the chemical sampling is restricted to jumps between the relevant end-states. At the end of the simulation the relative free energies of the various chemical states is given by the ratio of the populations. Both host-guest and protein-ligand systems have been successfully investigated with the method, which, like the other methods discussed in this section, is designed to rapidly identify which candidates look most promising for further investigation.

The linear response (LR) method was originally devised by Åqvist and co-workers [Åqvist *et al.* 1994] for estimating the binding affinities of ligands binding to proteins. Also known as the linear interaction energy (LIE) approach, it is a semi-empirical method for estimating absolute binding free energies and requires just two simulations, one of the solvated ligand-protein system and one of the ligand alone in solution. In both cases, the interaction between the ligand and its environment is broken down into the electrostatic and van der Waals contributions. The free energy of binding is then given by the following expression:

$$\Delta G = \beta(\langle \mathcal{V}_{\text{ligand-protein}}^{\text{el}} \rangle - \langle \mathcal{V}_{\text{ligand-solvent}}^{\text{el}} \rangle) + \alpha(\langle \mathcal{V}_{\text{ligand-protein}}^{\text{vdw}} \rangle - \langle \mathcal{V}_{\text{ligand-solvent}}^{\text{vdw}} \rangle) \quad (11.40)$$

As usual, the angle brackets $\langle \rangle$ indicate ensemble averages. α and β are two parameters. To determine ΔG one thus needs to perform just two simulations, one of the ligand in the solvent and the other of the ligand bound to the protein. The interactions that are accumulated consist solely of the electrostatic and van der Waals interactions between the ligand and its environment. First we consider an expansion of the Zwanzig expression for the free energy difference between two states X and Y (Equation (11.6)). The result obtained (see Appendix 11.3) is:

$$\Delta A = \frac{1}{2}[(\Delta \mathcal{H})_0 + \langle \Delta \mathcal{H} \rangle_1] - \frac{1}{4k_B T} [(\langle (\Delta \mathcal{H} - \langle \Delta \mathcal{H} \rangle_0)^2 \rangle_0 - \langle (\Delta \mathcal{H} - \langle \Delta \mathcal{H} \rangle_1)^2 \rangle_1) + \dots] \quad (11.41)$$

where $\Delta \mathcal{H} = \mathcal{H}_Y - \mathcal{H}_X$

For the electrostatic component, the free energy varies in a harmonic fashion with respect to deviations from equilibrium with a constant force constant (Figure 11.19). This is a standard result from dielectric theory and means that the mean square fluctuations of the energy on the two surfaces (the second terms in Equation (11.41)) will cancel, leaving just the first term. This leads to a value of $\frac{1}{2}$ for the electrostatic component (i.e. $\beta = 0.5$). A simple test of this theory is to calculate the electrostatic contribution to solvation free energies. Here, state X corresponds to the situation where all of the solvent-solvent and intramolecular solute interactions are present but the interaction between the solute and solvent is only described

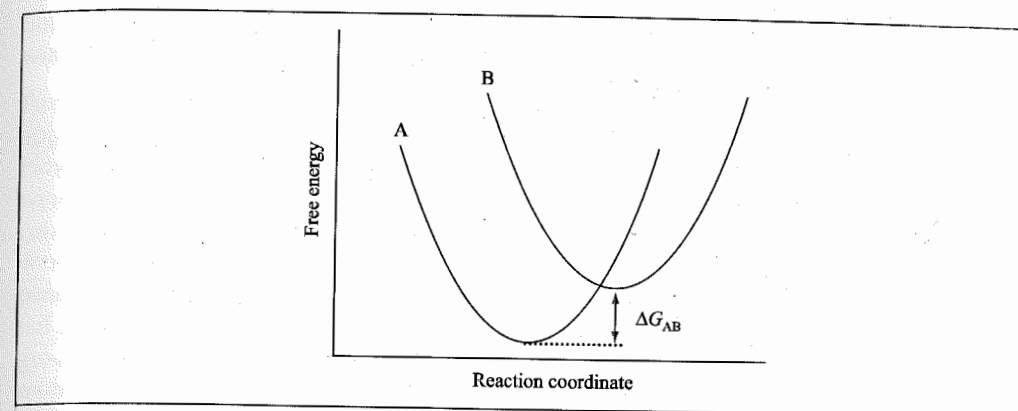


Fig. 11.19: Representation of the harmonic variation of the electrostatic component of the free energy according to the linear response approximation.

by a Lennard-Jones potential; the solute-solvent electrostatic interactions are missing. In state Y all interactions are included. The only difference between X and Y is thus the presence of the solute-solvent electrostatic terms, and so $\Delta \mathcal{H} (= \mathcal{H}_Y - \mathcal{H}_X)$ in Equation (11.41) is equal to $\mathcal{H}^{\text{el}}(\text{ligand-solvent})$. Thus:

$$\Delta A_{\text{sol}}^{\text{el}} = \frac{1}{2} \langle \mathcal{H}_{\text{ligand-solvent}}^{\text{el}} \rangle \quad (11.42)$$

The validity of this result has been confirmed by for example comparing the free energy perturbation result for charging Na^+ and Ca^{2+} ions in water with the ensemble average value of $\mathcal{V}^{\text{el}}(\text{ion-solvent})$, giving factors of 0.49 and 0.52. Moreover, one can apply the same arguments to the case of a ligand in a protein environment, leading to the following expression for the electrostatic contribution to the free energy of binding (i.e. the first term in Equation (11.40)):

$$\Delta A_{\text{binding}}^{\text{el}} = \frac{1}{2} (\langle \mathcal{H}_{\text{ligand-protein}}^{\text{el}} \rangle - \langle \mathcal{H}_{\text{ligand-solvent}}^{\text{el}} \rangle) \quad (11.43)$$

For the van der Waals component no such analytical theory exists. Åqvist and co-workers assumed that a similar linear treatment would work for these interactions but with a different empirical factor, to be determined from calibration experiments. There was some indirect evidence that this approach would be reasonable. For example, the experimental free energies of solvation for various hydrocarbons (e.g. *n*-alkanes) depend in an approximately linear fashion on the length of the carbon chain. In addition, the mean van der Waals solute-solvent energies from molecular dynamics simulations did show a linear variation with chain length (the slope of the line varying according to the solvent).

What remains is to determine a value of the parameter α . In the original publication this was done using a series of ligands which bind to endothiapepsin, an enzyme for which various crystal structures are known. Some of these ligands are shown in Figure 11.20; as can be seen, they are quite substantial. Molecular dynamics simulations of four ligands were performed within the enzyme binding site and in water and accumulating the required average interaction energies. Assuming the factor $\frac{1}{2}$ for the electrostatic contribution and comparing with the experimental binding affinities gave $\alpha = 0.161$. When a fifth ligand was evaluated,

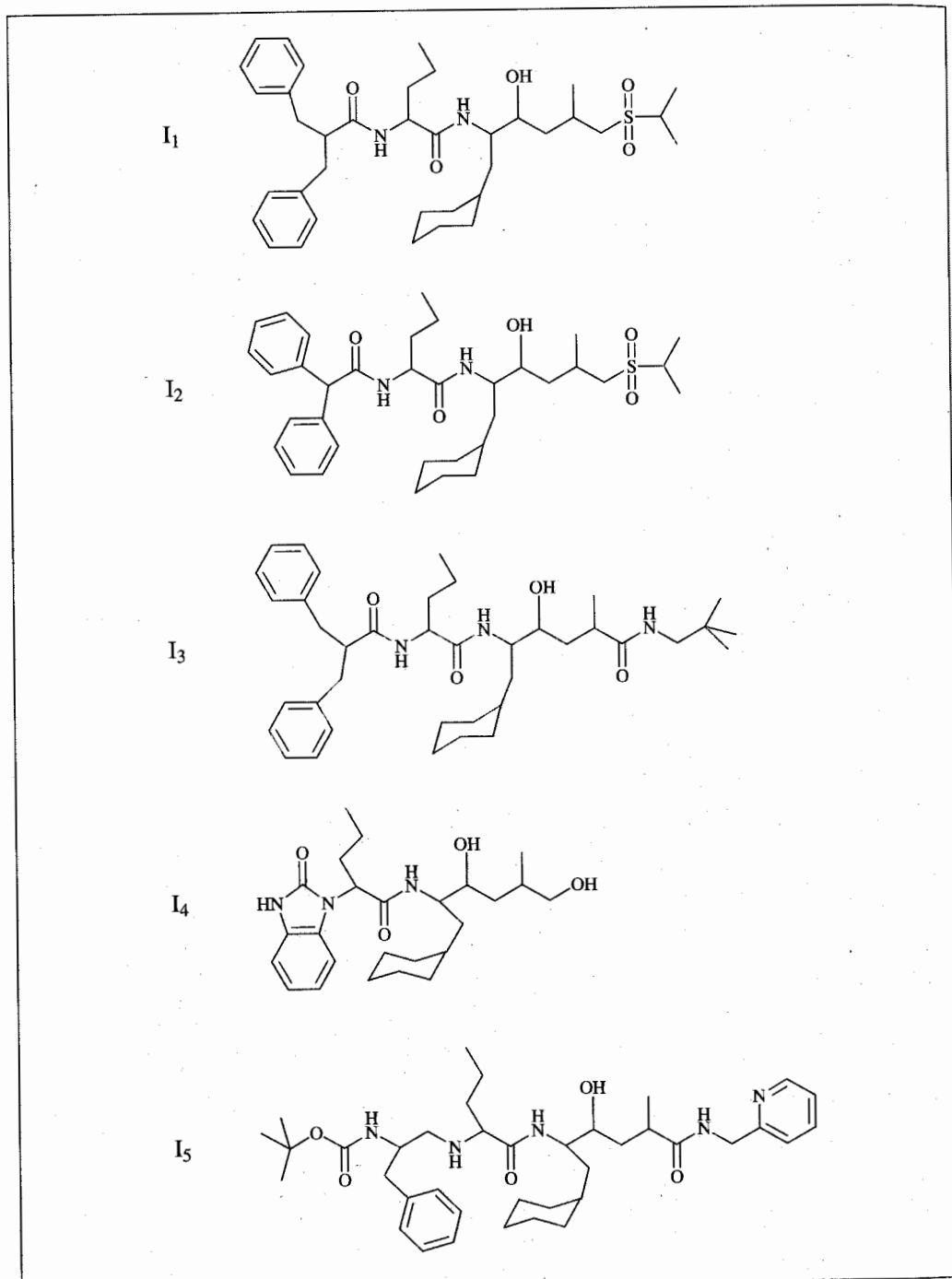


Fig. 11.20: Structures of the endothiapepsin ligands used to calibrate the LIE approach.

not present in the calibration set, rather remarkably the predicted free energy of binding was within 0.2 kcal/mol of the experimental result.

Further studies for different ligands and different enzymes appeared to support the approach, and also the constants α and β [Hansson and Åqvist 1995, Hansson *et al.* 1998]. However, other groups found that different values of the van der Waals parameter, α , were required. One possible reason for this discrepancy could be due to the different protocols (for example, different force fields), but some groups found that different parameters were required for different systems, even when the same protocols were employed. One possible explanation for this is that α depends on the nature of the binding site. This is not an unreasonable conclusion, given the different distributions of polar and non-polar groups in different binding sites. Wang and co-workers investigated this variation in more detail and showed that there appeared to be a correlation between the value of α and the 'weighted non-polar desolvation ratio', which is a measure of the hydrophobicity of the binding site [Wang *et al.* 1999]. It was found that, whilst it is generally more accurate to calibrate α for each system if experimental binding data for similar ligands is available, choosing a value based on the weighted non-polar desolvation ratio could give better results for dissimilar compounds.

Other groups have applied the linear response method to problems other than protein-ligand binding. A good problem for any new free energy approach is to predict the free energies of hydration of small organic molecules. Accurate hydration data are available for a wide variety of systems, and the calculations can usually be run relatively quickly. One immediate problem with the two-parameter linear response method is that, as α and β are both positive, it is not possible for any solute to have a positive hydration free energy (both the electrostatic and van der Waals interactions between solutes and water give negative solute-solvent energies). To deal with this problem, Carlsen and Jorgensen introduced an additional term which was related to the penalty for forming a solute cavity [Carlsen and Jorgensen 1995]. This third term was proportional to the solvent-accessible surface area:

$$\Delta G_{\text{hyd}} = \beta \langle \psi^{\text{el}} \rangle + \alpha \langle \psi^{\text{vdw}} \rangle + \gamma \text{SASA} \quad (11.44)$$

In their work on hydration, Carlsen and Jorgensen attempted to fit the three coefficients α , β and γ , obtaining the best fit for $\alpha = 0.4$, $\beta = 0.45$ and $\gamma = 0.03$ kcal/(mol Å²). In a subsequent study of the binding of a series of sulphonamide inhibitors to the enzyme thrombin, however, these parameter values were found to be ineffective and that new values were required to give an acceptable fit to the experimental data, with β now being much reduced in value (0.146) [Jones-Hertzog and Jorgensen 1997]. As more variables are considered for inclusion in an LIE-like relationship, it is important that a statistically correct strategy is employed to ensure that the 'optimal' equation is derived (the one with the most predictive power). The techniques to derive such equations are discussed in Section 12.12 on quantitative structure-activity relationships; one study where they were successfully employed considered the binding of a series of inhibitors to the enzyme neuraminidase [Wall *et al.* 1999].

Other attempts to predict free energies from a single simulation have explored the relationship between the coupling parameter, λ , and the free energy. Specifically, the free energy is

expressed as a Taylor series expansion in terms of λ around the point $\lambda = 0$. This expansion is [Smith and van Gunsteren 1994a; Liu *et al.* 1996] (showing just the first two terms explicitly):

$$\begin{aligned} A(\lambda) &= A(\lambda) - A(0) = A'_{\lambda=0}\lambda + \frac{1}{2!}A''_{\lambda=0}\lambda^2 + \frac{1}{3!}A'''_{\lambda=0}\lambda^3 + \dots \\ &= \left\langle \frac{\partial \mathcal{H}}{\partial \lambda} \right\rangle_{\lambda} + \frac{1}{k_B T} \left\langle \left(\frac{\partial \mathcal{H}}{\partial \lambda} - \left\langle \frac{\partial \mathcal{H}}{\partial \lambda} \right\rangle_0 \right)^2 \right\rangle_0 + \dots \end{aligned} \quad (11.45)$$

Truncating this series after the first derivative and integrating provides the basis for the thermodynamic integration approach. Moreover, if the Taylor series expansion is continued until it converges then Equation (11.45) is equivalent to the thermodynamic perturbation formula, so providing a link between the two approaches. In practice, it is always necessary to truncate the series; the problem then is whether it is appropriate to assume that the discarded higher-order terms are zero. A good way to test this approach is to consider model systems where the free energy change is known to be zero. One such system involves a simple diatomic molecule in a box of water. Each atom of the diatomic molecule is assigned a charge, equal in magnitude (0.25) but of opposite signs. The state to which this system is 'perturbed' corresponds to simply switching the charges (i.e. the start and final states are equivalent). For this system, a standard free energy perturbation calculation can give an answer very close to zero. The series expansion was not able to reproduce this result well, even from a 1 ns simulation. However, if one considered an alternative problem involving a change from $\lambda = 0$ to $\lambda = 0.5$ (which corresponds to decharging the system) the series expansion did give a very good result. Nevertheless, all these methods were found to fail for calculations involving the creation or deletion of atoms (a problem we discussed above). In the same paper, Liu and colleagues suggested an interesting method that could not only overcome this problem but also enable many free energy values for a series of related ligands to be obtained from a single simulation. At those positions where atoms are created or deleted, soft-core interaction sites are used of the form in Equation (11.30). A single long simulation of this (non-physical) reference state is performed. The soft-core potential has a functional form such that solvent molecules can sometimes penetrate 'within' the usual van der Waals radius. This extends the configurational space accessible to the system. Estimates of free energy differences can be obtained by running through the trajectory, substituting the soft-core sites for the appropriate 'real' atoms and calculating the energy for incorporation into the free-energy perturbation formula. In a 'proof-of-concept' illustration, the free energies of hydration for a series of small molecules were calculated from a single simulation consisting of a soft-core cavity in water [Schäfer *et al.* 1999]. These calculations suggested that the efficiency gains over conventional free energy calculations could reach 2–3 orders of magnitude but that the method did require some further development for certain types of system.

11.9 Continuum Representations of the Solvent

Most chemical processes take place in a solvent and so it is clearly important to consider how the solvent affects the behaviour of a system. In some cases, solvent molecules are directly

involved, as in ester hydrolysis reactions or in systems where solvent molecules are so tightly bound that they are effectively an integral part of the solute. Such solvent molecules should be modelled explicitly. In other systems, the solvent does not directly interact with the solute but it provides an environment that strongly affects the behaviour of the solute. For example, the highly anisotropic environment in a liquid crystal or lipid bilayer strongly influences the conformations of dissolved solutes. Here, it may not be necessary to explicitly model the solvent molecules, though special treatments such as mean field theories (see Section 7.10) may be required. In the third case, the solvent merely acts as a 'bulk medium' but can still significantly affect solute behaviour, with the dielectric properties of the solvent often being particularly important. In this case, it would clearly be useful not to have to explicitly include every single solvent molecule in the system, to enable us to concentrate on the behaviour of the solute(s). The solvent acts as a perturbation on the gas-phase behaviour of the system. This is the purpose of the 'continuum' solvent models [Smith and Pettitt 1994]. A considerable variety of such models have been proposed, for use with both quantum mechanics and empirical models [Cramer and Truhlar 1992]. Our discussion will be restricted to a few of the more widely used methods.

11.9.1 Thermodynamic Background

The solvation free energy (ΔG_{sol}) is the free energy change to transfer a molecule from vacuum to solvent. The solvation free energy can be considered to have three components:

$$\Delta G_{\text{sol}} = \Delta G_{\text{elec}} + \Delta G_{\text{vdw}} + \Delta G_{\text{cav}} \quad (11.46)$$

where ΔG_{elec} is the electrostatic component. This contribution is particularly important for polar and charged solutes due to the polarisation of the solvent, which we will model as a uniform medium of constant dielectric ϵ . ΔG_{vdw} is the van der Waals interaction between the solute and solvent; this may in turn be divided into a repulsive term, ΔG_{rep} , and an attractive dispersion term, ΔG_{disp} . ΔG_{cav} is the free energy required to form the solute cavity within the solvent. This component is positive and comprises the entropic penalty associated with the reorganisation of the solvent molecules around the solute together with the work done against the solvent pressure in creating the cavity. In addition to the above three components, an explicit hydrogen-bonding term, ΔG_{hb} , may be added for those systems where there is localised hydrogen bonding between the solute and solvent. Initially, we will discuss the electrostatic contribution to the free energy of solvation. We will then consider the van der Waals and cavity contributions.

11.10 The Electrostatic Contribution to the Free Energy of Solvation: The Born and Onsager Models

Two important contributions to the study of solvation effects were made by Born (in 1920) and Onsager (in 1936). Born derived the electrostatic component of the free energy of solvation for placing a charge within a spherical solvent cavity [Born 1920], and Onsager extended this to a dipole in a spherical cavity (Figure 11.21) [Onsager 1936]. In the Born

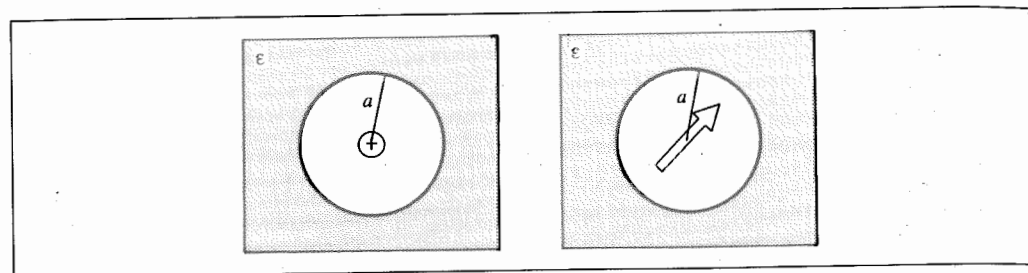


Fig. 11.21: The Born and Onsager models.

model, ΔG_{elec} of an ion is equal to the work done to transfer the ion from vacuum to the medium. This in turn is equal to the difference in the electrostatic work to charge the ion in the two environments. The work to charge an ion in a medium of dielectric constant ϵ equals $q^2/2\epsilon a$, where q is the charge on the ion and a is the radius of the cavity. The electrostatic contribution to the solvation free energy is thus the difference in the work done in charging the ion in the dielectric and *in vacuo*:

$$\Delta G_{\text{elec}} = -\frac{q^2}{2a} \left(1 - \frac{1}{\epsilon}\right) \quad (11.47)$$

Note that in this equation, as throughout our discussion, we have used reduced electrostatic units, in which the factor $4\pi\epsilon_0$ is ignored. This is common practice in the literature. The Born model is very simple yet can be quite successful. It is necessary to choose a set of cavity radii. Traditionally, ionic radii from crystal structures are used. However, for the alkali halides it is found that adding 0.1 Å to the radii of anions and 0.85 Å to the radii of cations gives much better agreement with experimental data. Justification for this adjustment was provided by Rashin and Honig, who examined electron density distributions in crystals and concluded that the ionic radii are reasonably good indicators of cavity size for anions but that for cations it is more appropriate to use covalent radii [Rashin and Honig 1985]. They subsequently suggested that the optimal agreement with experiment could be obtained by increasing these radii by an empirical factor of 7%.

11.10.1 Calculating the Electrostatic Contribution via Quantum Mechanics

The Born model is obviously only appropriate to species with a formal charge. Onsager's dipole model is relevant to many more molecules (in fact, the Onsager model is a special case of the result derived by Kirkwood [Kirkwood 1934], who considered an arbitrary distribution of charges within a spherical cavity). The solute dipole within the cavity induces a dipole in the surrounding medium, which in turn induces an electric field within the cavity (the *reaction field*). The reaction field then interacts with the solute dipole, so providing additional stabilisation of the system. The magnitude of the reaction field was determined by Onsager to be:

$$\phi_{\text{RF}} = \frac{2(\epsilon - 1)}{(2\epsilon + 1)a^3} \mu \quad (11.48)$$

where μ is the dipole moment of the solute; a and ϵ are the radius of the cavity and the dielectric constant of the medium, as before. The energy of a dipole in an electric field ϕ_{RF} is $-\phi_{\text{RF}}\mu$, but for a polarisable dipole it is necessary to add an additional term which represents the work done assembling the charge distribution within the cavity. This additional term has magnitude $\phi_{\text{RF}}\mu/2$ and so the electrostatic contribution to the free energy of solvation in this model is:

$$\Delta G_{\text{elec}} = -\frac{\phi_{\text{RF}}\mu}{2} \quad (11.49)$$

If the species is charged then an appropriate Born term must also be added. The reaction field model can be incorporated into quantum mechanics, where it is commonly referred to as the *self-consistent reaction field* (SCRF) method, by considering the reaction field to be a perturbation of the Hamiltonian for an isolated molecule. The modified Hamiltonian of the system is then given by:

$$\mathcal{H}_{\text{tot}} = \mathcal{H}_0 + \mathcal{H}_{\text{RF}} \quad (11.50)$$

where \mathcal{H}_0 is the Hamiltonian of the isolated molecule and \mathcal{H}_{RF} is the perturbation, given by [Tapia and Goscinski, 1975]:

$$\mathcal{H}_{\text{RF}} = -\hat{\mu}^T \frac{2(\epsilon - 1)}{(2\epsilon + 1)a^3} \langle \Psi | \hat{\mu} | \Psi \rangle \quad (11.51)$$

where $\hat{\mu}$ is the dipole moment operator written in matrix form and $\hat{\mu}^T$ is its transpose. The wavefunction Ψ for the modified Hamiltonian is determined and the electrostatic contribution to the solvation free energy is then given by:

$$\Delta G_{\text{elec}} = \langle \Psi | \mathcal{H}_{\text{tot}} | \Psi \rangle - \langle \Psi_0 | \mathcal{H}_0 | \Psi_0 \rangle + \frac{1}{2} \frac{2(\epsilon - 1)}{(2\epsilon + 1)a^3} \mu^2 \quad (11.52)$$

The third term in Equation (11.52) is the correction factor corresponding to the work done in creating the charge distribution of the solute within the cavity in the dielectric medium. Ψ_0 is the gas-phase wavefunction.

A drawback of the SCRF method is its use of a spherical cavity; molecules are rarely exactly spherical in shape. However, a spherical representation can be a reasonable first approximation to the shape of many molecules. It is also possible to use an ellipsoidal cavity; this may be a more appropriate shape for some molecules. For both the spherical and ellipsoidal cavities analytical expressions for the first and second derivatives of the energy can be derived, so enabling geometry optimisations to be performed efficiently. For these cavities it is necessary to define their size. In the case of a spherical cavity a value for the radius can be calculated from the molecular volume:

$$a^3 = 3V_m/4\pi N_A \quad (11.53)$$

The molecular volume V_m can in turn be obtained by dividing the molecular weight by the density or from refractivity measurements; N_A is Avogadro's number. The cavity radius can also be estimated from the largest interatomic distance within the molecule. A third approach is to calculate the 'volume' of the molecule from a suitable electron density contour. The radii obtained by these procedures are often adjusted by adding an empirical constant to give the 'true' cavity radius. This extra value accounts for the fact that solvent

molecules cannot approach right up to the molecule. An additional extension to the simple SCRF procedure is the use of a multipolar expansion to represent the solute [Rinaldi *et al.* 1983]. This overcomes a drawback of the basic model in which a molecule with a zero dipole would have zero solvation energy.

A yet more realistic cavity shape is that obtained from the van der Waals radii of the atoms of the solute. This is the approach taken in the *polarisable continuum* method (PCM) [Miertus *et al.* 1981], which has been implemented in a variety of *ab initio* and semi-empirical quantum mechanical programs. Due to the non-analytical nature of the cavity shapes in the PCM approach, it is necessary to calculate ΔG_{elec} numerically. The cavity surface is divided into a large number of small surface elements, and there is a point charge associated with each surface element. This system of point charges represents the polarisation of the solvent, and the magnitude of each surface charge is proportional to the electric field gradient at that point. The total electrostatic potential at each surface element equals the sum of the potential due to the solute and the potential due to the other surface charges:

$$\phi(\mathbf{r}) = \phi_{\rho}(\mathbf{r}) + \phi_{\sigma}(\mathbf{r}) \quad (11.54)$$

where $\phi_{\rho}(\mathbf{r})$ is the potential due to the solute and $\phi_{\sigma}(\mathbf{r})$ is the potential due to the surface charges. The PCM algorithm is as follows. First, the cavity surface is determined from the van der Waals radii of the atoms. That fraction of each atom's van der Waals sphere which contributes to the cavity is then divided into a number of small surface elements of calculable surface area. The simplest way to do this is to define a local polar coordinate frame at the centre of each atom's van der Waals sphere and to use fixed increments of $\Delta\theta$ and $\Delta\phi$ to give rectangular surface elements (Figure 11.22). The surface can also be divided using tessellation methods [Paschual-Ahuir *et al.* 1987]. An initial value of the point charge for each surface element is then calculated from the electric field gradient due to the solute alone:

$$q_i = - \left[\frac{\epsilon - 1}{4\pi\epsilon} \right] E_i \Delta S \quad (11.55)$$

where ϵ is the dielectric constant of the medium, E_i is the electric field gradient and ΔS is the area of the surface element. The contribution $\phi_{\sigma}(\mathbf{r})$ due to the other point charges can then be calculated using Coulomb's law. These charges are modified iteratively until they are

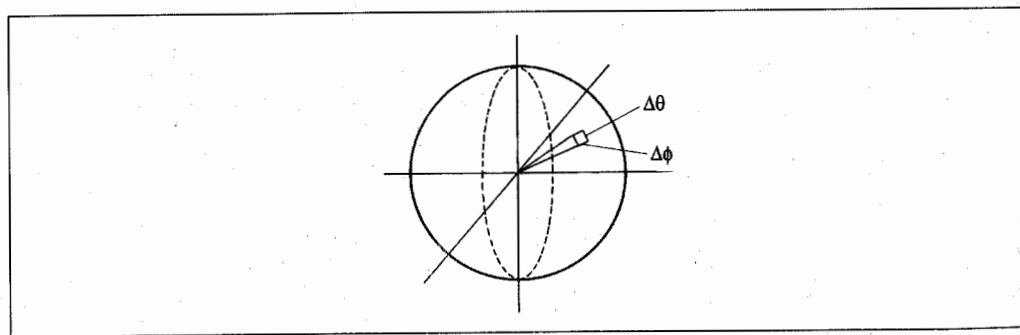


Fig. 11.22: Small surface elements can be created on the van der Waals surface of an atom using constant increments of the polar angles, θ and ϕ .

self-consistent. The potential $\phi_{\sigma}(\mathbf{r})$ from the final part of the charge is then added to the solute Hamiltonian ($\mathcal{H} = \mathcal{H}_0 + \phi_{\sigma}(\mathbf{r})$) and the SCF calculation initiated. After each SCF calculation new values of the surface charges are calculated from the current wavefunction to give a new value of $\phi_{\sigma}(\mathbf{r})$ which is used in the next iteration until the solute wavefunction and the surface charges are self-consistent.

To calculate ΔG_{elec} , we must take account of the work done in creating the charge distribution within the cavity in the dielectric medium. This is equal to one-half of the electrostatic interaction energy between the solute charge distribution and the polarised dielectric, and so:

$$\Delta G_{\text{elec}} = \int \Psi \mathcal{H} \Psi d\tau - \int \Psi_0 \mathcal{H}_0 \Psi_0 d\tau - \frac{1}{2} \int \phi(\mathbf{r}) \rho(\mathbf{r}) d\mathbf{r} \quad (11.56)$$

where $\rho(\mathbf{r})$ is the charge distribution of the surface elements.

There are two slight complications with the PCM approach. The first of these arises as a consequence of representing a continuous charge distribution over the cavity surface as a set of single point charges. When calculating the electrostatic potential due to the charges on the surface elements one must exclude the charge for the current surface element. To include it would cause the charges to diverge rather than converge. The contribution of the charge on that surface element is therefore determined separately using the Gauss theorem. The second complication arises because the wavefunction of the solute extends beyond the cavity. Thus the sum of the charges on the surface is not equal and opposite to the charge of the solute. This problem can be easily overcome by scaling the charge distribution on the surface so that it is equal and opposite to the charge of the solute.

COSMO is an interesting variant on the PCM method (COSMO stands for 'conductor-like screening model') [Klamt and Schüürmann 1993; Klamt 1995; Klamt *et al.* 1998]. The cavity is considered to be embedded in a conductor with an infinite dielectric constant. The advantage of this is that screening effects in an infinitely strong dielectric (i.e. a conductor) are much easier to handle. A small correction to the results for this conductor can provide the appropriate value for water, which of course has a high dielectric constant. On the surface of a conductor the potential due to the solute and due to the surface charges is set to zero, which gives rise to a convenient boundary condition when determining the surface charges. For an alternative dielectric these charges are scaled by a factor:

$$q' = q \frac{\epsilon_r - 1}{\epsilon_r + 0.5} \quad (11.57)$$

The SCRF and PCM models have been used to investigate the effect of solvent upon energetics and equilibria. For example, Wong, Wiberg and Frisch used the SCRF method to investigate the effect of different solvents upon the tautomeric equilibria of 2-pyridone (Figure 11.23) [Wong *et al.* 1992]. Geometry optimisations were performed for the various tautomeric species at high levels of theory, and vibrational frequencies were calculated. Results were reported for the gas phase, for a non-polar solvent (cyclohexane, $\epsilon = 2.0$) and for an aprotic polar solvent (acetonitrile, $\epsilon = 35.9$). The calculated free energy changes in the gas phase, cyclohexane and acetonitrile were -0.64 kcal/mol, 0.36 kcal/mol and 2.32 kcal/mol, respectively, which compared favourably with the experimental values of -0.81 kcal/mol, 0.33 kcal/mol and 2.96 kcal/mol. The dielectric medium was found to have a much more pronounced effect

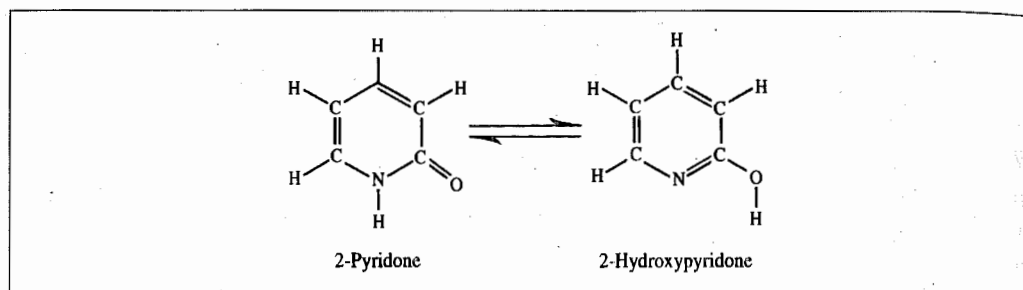


Fig. 11.23: Tautomers of 2-pyridone.

on the structure, charge distribution and vibrational frequencies of the keto form than of the enol form. This was ascribed to the more polar nature of the keto tautomer.

11.10.2 Continuum Models for Molecular Mechanics

In many cases, solvent effects can be incorporated into a force field model using one of the theories that we have just examined. It is possible to study larger systems with the empirical models, in which case it is necessary to take account of the dielectric properties of the solute as well as those of the solvent. Before reading this section it may be useful to revise the definitions of molecular surface and accessible surface given in Section 1.5, as these are widely referenced.

The *boundary element method* of Rashin is similar in spirit to the polarisable continuum model, but the surface of the cavity is taken to be the molecular surface of the solute [Rashin and Namboodiri 1987; Rashin 1990]. This cavity surface is divided into small boundary elements. The solute is modelled as a set of atoms with point polarisabilities. The electric field induces a dipole proportional to its polarisability. The electric field at an atom has contributions from dipoles on other atoms in the molecule, from polarisation charges on the boundary, and (where appropriate) from the charges of electrolytes in the solution. The charge density is assumed to be constant within each boundary element but is not reduced to a single point as in the PCM model. A set of linear equations can be set up to describe the electrostatic interactions within the system. The solutions to these equations give the boundary element charge distribution and the induced dipoles, from which thermodynamic quantities can be determined.

The *generalised Born equation* has been widely used to represent the electrostatic contribution to the free energy of solvation [Constanciel and Contreras 1984]. The model comprises a system of particles with radii a_i and charges q_i . The total electrostatic free energy of such a system is given by the sum of the Coulomb energy and the Born free energy of solvation in a medium of relative permittivity ϵ :

$$G_{\text{elec}} = \sum_{i=1}^N \sum_{j=i+1}^N \frac{q_i q_j}{\epsilon r_{ij}} - \frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^N \frac{q_i^2}{a_i} \quad (11.58)$$

The first term in Equation (11.58) can be written as the sum of a Coulomb interaction *in vacuo* and a second term in $(1 - 1/\epsilon)$:

$$\sum_{i=1}^N \sum_{j=i+1}^N \frac{q_i q_j}{\epsilon r_{ij}} = \sum_{i=1}^N \sum_{j=i+1}^N \frac{q_i q_j}{r_{ij}} - \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^N \sum_{j=i+1}^N \frac{q_i q_j}{r_{ij}} \quad (11.59)$$

In the generalised Born approach the total electrostatic energy is written as a sum of three terms, the first of which is the Coulomb interaction between the charges *in vacuo*:

$$G_{\text{elec}} = \sum_{i=1}^N \sum_{j=i+1}^N \frac{q_i q_j}{r_{ij}} - \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^N \sum_{j=i+1}^N \frac{q_i q_j}{r_{ij}} - \frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^N \frac{q_i^2}{a_i} \quad (11.60)$$

where ΔG_{elec} equals the difference between G_{elec} and the Coulomb energy *in vacuo*. This is the generalised Born (GB) equation:

$$\Delta G_{\text{elec}} = - \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^N \sum_{j=i+1}^N \frac{q_i q_j}{r_{ij}} - \frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^N \frac{q_i^2}{a_i} \quad (11.61)$$

The generalised Born equation has been incorporated into both molecular mechanics calculations (by Still and co-workers [Still *et al.* 1990; Qiu *et al.* 1997]) and semi-empirical quantum mechanics calculations (by Cramer and Truhlar, in an ongoing series of models called SM1, SM2, SM3, etc. [Cramer and Truhlar 1992; Chambers *et al.* 1996]). In these treatments, the two terms in Equation (11.61) are combined into a single expression of the following form:

$$\Delta G_{\text{elec}} = - \frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^N \sum_{j=1}^N \frac{q_i q_j}{f(r_{ij}, a_{ij})} \quad (11.62)$$

where $f(r_{ij}, a_{ij})$ depends upon the interparticle distances r_{ij} and the Born radii a_i . A variety of forms are possible for the function f ; that proposed by Still and colleagues was:

$$f(r_{ij}, a_{ij}) = \sqrt{(r_{ij}^2 + a_{ij}^2) e^{-D}} \quad \text{where } a_{ij} = \sqrt{a_i a_j} \quad \text{and } D = r_{ij}^2 / (2a_{ij})^2 \quad (11.63)$$

This form of the function f can be justified for the following reasons. When $i = j$, the equation returns the Born expression; for two charges close together (i.e. a dipole, in which r_{ij} is small compared to a_i and a_j) the expression is close to the Onsager result; and for two charges separated by a significant distance ($r_{ij} \gg a_i, a_j$) the result is very close to the sum of the Coulomb and Born expressions. A further advantage of this functional form is that the expression can be differentiated analytically, thereby enabling the solvation term to be included in gradient-based optimisation methods and molecular dynamics simulations.

A rather complex procedure is used to determine the Born radii a_i , values of which are calculated for each atom in the molecule that carries a charge or a partial charge. The Born radius of an atom (more correctly considered to be an 'effective' Born radius) corresponds to the radius that would return the electrostatic energy of the system according to the Born equation if all other atoms in the molecule were uncharged (i.e. if the other atoms only acted to define the dielectric boundary between the solute and the solvent). In Still's force field implementation, atomic radii from the OPLS force field are assigned to each

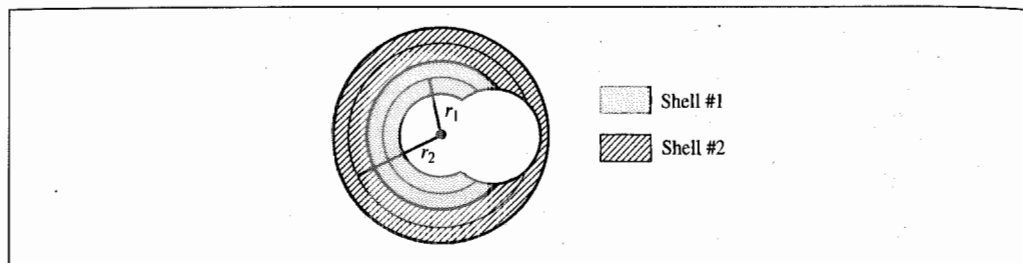


Fig. 11.24: Calculation of the effective Born radius in the generalised Born model. Shells are constructed until they contain the entire molecule. For each shell the amount of exposed surface area is determined for the middle of the shell. (Figure adapted from Still W C, A Tempczyk, R C Hawley and T Hendrickson 1990. *Semianalytical Treatment of Solvation for Molecular Mechanics and Dynamics*. Journal of the American Chemical Society 112:6127-6129.)

atom and amended by an empirically determined offset of -0.09 \AA to define the dielectric boundary. In the quantum mechanical approach of Cramer and Truhlar, the radius of each atom is a function of the charge on the atom. The dielectric boundary is then taken to be the union of the relevant radii.

The electrostatic energy of an atom i is calculated numerically by constructing a series of spherical shells until the outer shell (shell M) entirely contains the entire van der Waals surface of the molecule, as shown in Figure 11.24. The Born electrostatic energies of the dielectric in these shells are determined using the following equation:

$$\Delta G_{\text{elec}} = -\frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) q_i^2 \left\{ \sum_{k=1}^M \frac{A_k}{4\pi r_k^2} \left[\left(\frac{1}{r_k - 0.5T_k}\right) - \left(\frac{1}{r_k + 0.5T_k}\right) \right] + \frac{1}{r_{M+1} - 0.5T_{M+1}} \right\} \quad (11.64)$$

The radius of the k th shell (r_k) is measured at the middle of the shell. A_k is the amount of surface area of a sphere of radius r_k that is not contained within the van der Waals surface of the molecule. The thickness of each shell, T_k , increases with distance from the atom as follows:

$$T_{k+1} = (1 + F)T_k \quad (11.65)$$

where F (the expansion factor) and T_1 (the radius of the first shell) are parameters, chosen by Still to be 0.5 and 0.1 \AA , respectively. The shells are constructed from where the dielectric boundary commences (i.e. in Still's case the shells start 0.9 \AA inside the van der Waals radii). The final term in Equation (11.65) is the contribution due to the dielectric that lies beyond the van der Waals surface of the molecule. The effective Born radius is then given by equating Equation (11.65) with the Born equation for the atom, and so:

$$\frac{1}{a_i} = \sum_{k=1}^M \frac{A_k}{4\pi r_k^2} \left[\left(\frac{1}{r_k - 0.5T_k}\right) - \left(\frac{1}{r_k + 0.5T_k}\right) \right] + \frac{1}{r_{M+1} - 0.5T_{M+1}} \quad (11.66)$$

The effective Born radii do not change very much and so are recalculated whenever the non-bonded list is updated. The Still formulation of the generalised Born equation requires the surface areas A_k of the spherical shells that are exposed to solvent to be calculated. For

this, a fast numerical method devised by Wodak and Janin [Wodak and Janin 1980] is used in which the accessible surface area is given by:

$$A_i = S_i \prod_j (1.0 - b_{ij}/S_i) \quad (11.67)$$

where S_i is the total accessible surface area of an atom i with radius r_i as defined with a solvent probe of radius r_s . b_{ij} is the amount of surface area removed due to overlap with an atom j which is a distance d_{ij} from atom i :

$$S_i = 4\pi(r_i + r_s)^2 \quad (11.68)$$

$$b_{ij} = \pi(r_i + r_s)(r_j + r_i - 2r_s - d_{ij})[1.0 + (r_i - r_j)/d_{ij}] \quad (11.69)$$

The Wodak-Janin method is only approximate for more than two spheres. Exact values of A_i can be calculated, but only with a significant computational effort. A comparison of results obtained with the approximate and exact methods showed that, for molecules significantly larger than the probe, the approximation was valid (Wodak and Janin's expression was intended to be used to study solvent effects in proteins). Still showed that it was also possible to reduce the b_{ij} term by an empirical constant and obtain accurate results for smaller systems.

The generalised Born model has been incorporated into a number of quantum mechanical and molecular mechanical programs. Still and his group have made extensive use of the model in conformational searching and for calculating relative free energies of binding with free energy perturbation methods. For example, the relative free energies of binding of D- and L-enantiomeric α -amino acid-derived substrates to a podand ionophore (**1**; Figure 11.25) were calculated to be in good agreement with experiment using a mixed Monte Carlo/dynamics method and the generalised Born model for chloroform [Burger *et al.* 1994]. Similar calculations were then used to predict which of a variety of substituted ionophores (varying the group X in Figure 11.25) would be expected to show the greatest selectivity for the guest ($Y = \text{NHMe}$). The derivative **2** was predicted to show the greatest enantioselectivity. Unfortunately, this particular compound was too insoluble to measure binding affinities, but a related compound, **3**, did show the desired selectivity. Moreover, when the calculations were repeated on **3** the predicted difference in binding affinity was within 0.3 kcal/mol of the experimental result. Such studies clearly illustrate the potential applicability of such calculations, but it should be noted that this system was carefully chosen to minimise any errors associated with the force field parameters (through the use of enantiomeric guests) and sampling (as the host is locked into a single binding conformation). Even so, to achieve accurate results it was usually necessary to perform simulations of the order of 10 ns ; at the time such simulations could only be realistically achieved using a continuum model of the solvent.

11.10.3 The Langevin Dipole Model

The Langevin dipole method of Warshel and Levitt [Warshel and Levitt 1976] is intermediate between a continuum and an explicit solvation model. A three-dimensional

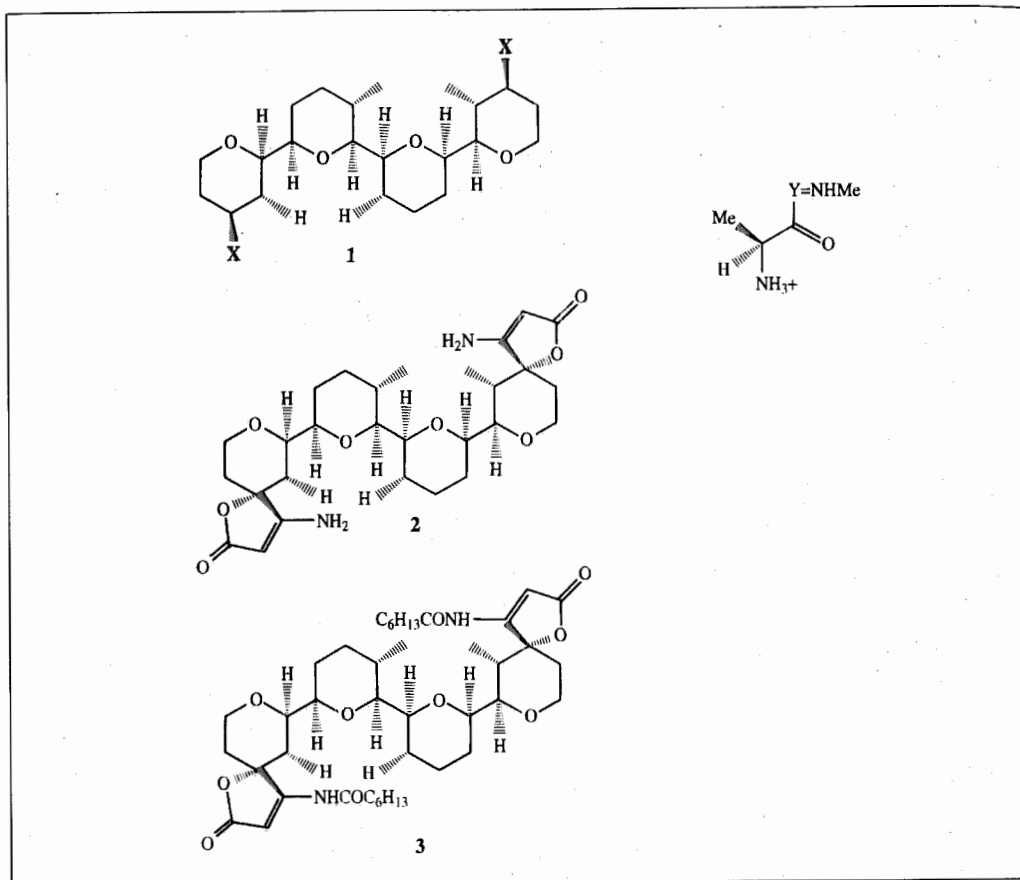


Fig. 11.25: Ionophores that selectively bind amino acids [Burger et al. 1994].

grid of rotatable point dipoles is established in the region beyond the boundary (which can be of arbitrary shape; Figure 11.26). For macromolecules the boundary corresponds to the solvent accessible surface. These dipoles represent the molecular dipoles of the solvent molecules in the outer region, and the separation between them is chosen accordingly. The electric field E_i at each dipole has a contribution from the solute and from other solvent dipoles. The size and direction of each dipole is determined using the Langevin equation:

$$\mu_i = \mu_0 \frac{E_i}{|E_i|} \left[\frac{\exp[C\mu_0|E_i|/k_B T] + \exp[-C\mu_0|E_i|/k_B T]}{\exp[C\mu_0|E_i|/k_B T] - \exp[-C\mu_0|E_i|/k_B T]} - \frac{1}{C\mu_0|E_i|/k_B T} \right] \quad (11.70)$$

where μ_0 is the size of the dipole moment of a solvent molecule and C is a parameter that represents the degree to which the dipoles resist reorientation; its value may be obtained from a separate simulation using explicit solvent. Converged values of the dipoles are usually obtained within a few iterations. The free energy of the Langevin dipoles is then

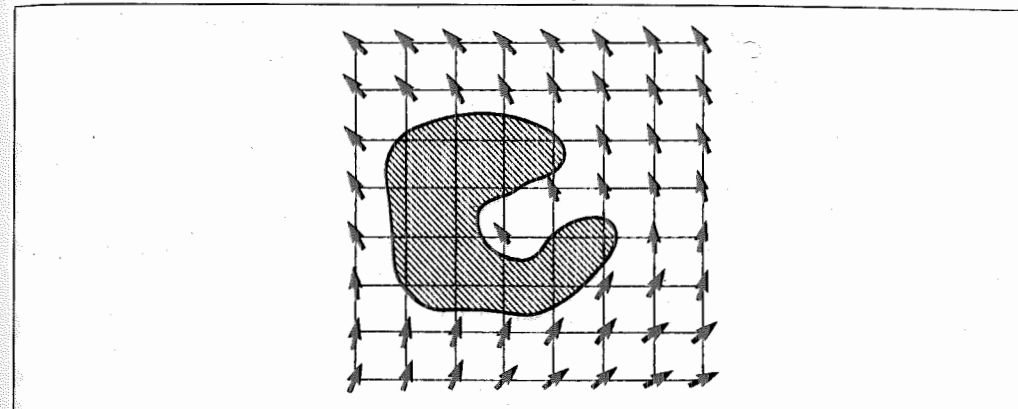


Fig. 11.26: The Langevin dipole model.

given by:

$$\Delta G_{\text{sol}} = -\frac{1}{2} \sum_i \mu_i \cdot E_i^0 \quad (11.71)$$

E_i^0 is the field due to the solute charges alone. The Langevin dipole method has been widely used by Warshel in his studies of enzyme reactions (see Section 11.13.3).

11.10.4 Methods Based upon the Poisson–Boltzmann Equation

The final class of methods that we shall consider for calculating the electrostatic component of the solvation free energy are based upon the Poisson or the Poisson–Boltzmann equations. These methods have been particularly useful for investigating the electrostatic properties of biological macromolecules such as proteins and DNA. The solute is treated as a body of constant low dielectric (usually between 2 and 4), and the solvent is modelled as a continuum of high dielectric. The Poisson equation relates the variation in the potential ϕ within a medium of uniform dielectric constant ϵ to the charge density ρ :

$$\nabla^2 \phi(\mathbf{r}) = -\frac{\rho(\mathbf{r})}{\epsilon_0 \epsilon} \quad (11.72)$$

In reduced electrostatic units, the factor $4\pi\epsilon_0$ is eliminated and the Poisson equation becomes:

$$\nabla^2 \phi(\mathbf{r}) = -\frac{4\pi\rho(\mathbf{r})}{\epsilon} \quad (11.73)$$

The charge density is simply the distribution of charge throughout the system and has SI units of C m^{-3} . The Poisson equation is thus a second-order differential equation (∇^2 is the usual abbreviation for $(\partial^2/\partial x^2) + (\partial^2/\partial y^2) + (\partial^2/\partial z^2)$). For a set of point charges in a constant dielectric the Poisson equation reduces to Coulomb's law. However, if the dielectric

is not constant but varies with position, then Coulomb's law is not applicable and the Poisson equation is:

$$\nabla \cdot \epsilon(\mathbf{r}) \nabla \phi(\mathbf{r}) = -4\pi\rho(\mathbf{r}) \quad (11.74)$$

The Poisson equation must be modified when mobile ions are present, to account for their redistribution in the solution in response to the electric potential. The ions are prevented from congregating at the locations of extreme electrostatic potential due to repulsive interactions with other ions and their natural thermal motion. The ion distribution is described by a Boltzmann distribution of the following form:

$$n(\mathbf{r}) = \mathcal{N} \exp(-\mathcal{V}(\mathbf{r})/k_B T) \quad (11.75)$$

where $n(\mathbf{r})$ is the number density of ions at a particular location \mathbf{r} , \mathcal{N} is the bulk number density and $\mathcal{V}(\mathbf{r})$ is the energy change to bring the ion from infinity to the position \mathbf{r} . When these effects are incorporated into the Poisson equation the result is the *Poisson-Boltzmann equation*:

$$\nabla \cdot \epsilon(\mathbf{r}) \nabla \phi(\mathbf{r}) - \kappa' \sinh[\phi(\mathbf{r})] = -4\pi\rho(\mathbf{r}) \quad (11.76)$$

κ' is related to the Debye-Hückel inverse length, κ , by:

$$\kappa^2 = \frac{\kappa'^2}{\epsilon} = \frac{8\pi N_A e^2 I}{1000 \epsilon k_B T} \quad (11.77)$$

where e is the electronic charge, I is the ionic strength of the solution and N_A is Avogadro's number. This is a non-linear differential equation that can be written as an alternative form by expanding the hyperbolic sine function as a Taylor series:

$$\nabla \cdot \epsilon(\mathbf{r}) \nabla \phi(\mathbf{r}) - \kappa' \phi(\mathbf{r}) \left[1 + \frac{\phi(\mathbf{r})^2}{6} + \frac{\phi(\mathbf{r})^4}{120} + \dots \right] = -4\pi\rho(\mathbf{r}) \quad (11.78)$$

The linearised Poisson-Boltzmann equation is obtained by taking only the first term in the expansion, giving:

$$\nabla \cdot \epsilon(\mathbf{r}) \nabla \phi(\mathbf{r}) - \kappa' \phi(\mathbf{r}) = -4\pi\rho(\mathbf{r}) \quad (11.79)$$

How can Equation (11.79) be solved? Before computers were available only simple shapes could be considered. For example, proteins were modelled as spheres or ellipses (Tanford-Kirkwood theory); DNA as a uniformly charged cylinder; and membranes as planes (Gouy-Chapman theory). With computers, numerical approaches can be used to solve the Poisson-Boltzmann equation. A variety of numerical methods can be employed, including finite element and boundary element methods, but we will restrict our discussion to the finite difference method first introduced for proteins by Warwicker and Watson [Warwicker and Watson 1982]. Several groups have implemented this method; here we concentrate on the work of Honig's group, whose DelPhi program has been widely used.

A cubic lattice is superimposed onto the solute(s) and the surrounding solvent. Values of the electrostatic potential, charge density, dielectric constant and ionic strength are assigned to each grid point. The atomic charges do not usually coincide with a grid point and so the

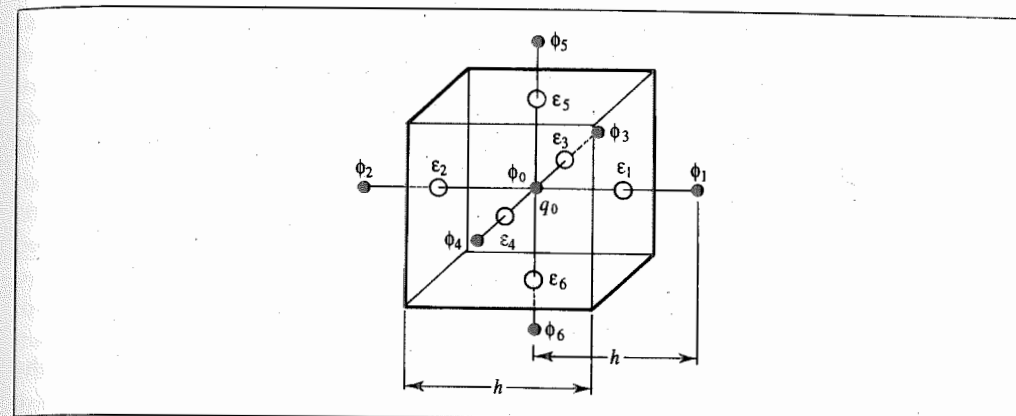


Fig. 11.27: The cube used in the finite difference method for solving the Poisson-Boltzmann equation. (Figure adapted from Klapper I, R Hagstrom, R Fine, K Sharp and B Honig 1986. *Focusing of Electric Fields in the Active Site of Cu-Zn Superoxide Dismutase: Effects of Ionic Strength and Amino-Acid Substitution*. *Proteins: Structure, Function and Genetics* 1:47-59.)

charge is allocated to the eight surrounding grid points in such a way that the closer the charge to the grid point, the greater the proportion of its total charge that is allocated. The derivatives in the Poisson-Boltzmann equation are then determined by a finite difference formula. Consider the cube of side h surrounding the grid point shown in Figure 11.27. A charge q_0 is associated with the grid point; this is equivalent to a uniform charge density of q_0/h^3 within the cube (i.e. $\rho_0 = q_0/h^3$). The potential at the grid point is given by:

$$\phi_0 = \frac{\sum \epsilon_i \phi_i + 4\pi \frac{q_0}{h}}{\sum \epsilon_i + \kappa'^2 f(\phi_0)} \quad (11.80)$$

The summations are over the potentials ϕ_i at the six adjoining grid points and the dielectric constants ϵ_i which are associated with the midpoints of the lines between the grid points. The function $f(\phi_0)$ in the denominator has the value 1 for the linear Poisson-Boltzmann equation, and is equivalent to the series expansion $(1 + \phi_0^2/6 + \phi_0^4/120 + \dots)$ for the non-linear case. κ'^2 is obtained from the ionic strength at the grid point. The crucial feature is that the potential at each grid point influences the potential at the neighbouring grid points, and so by iteratively repeating the calculation converged values will eventually be obtained.

To perform a Poisson-Boltzmann calculation it is necessary to allocate a value for the dielectric constant to each grid point, which requires us to decide which grid points lie within the solute(s) and which are in the solvent. The boundary between the solute and solvent is defined as either the molecular surface or the accessible surface. All grid points outside this surface are assigned a high dielectric constant (80 for water) and an ionic strength value. Grid points within the surface are assigned the dielectric constant of the

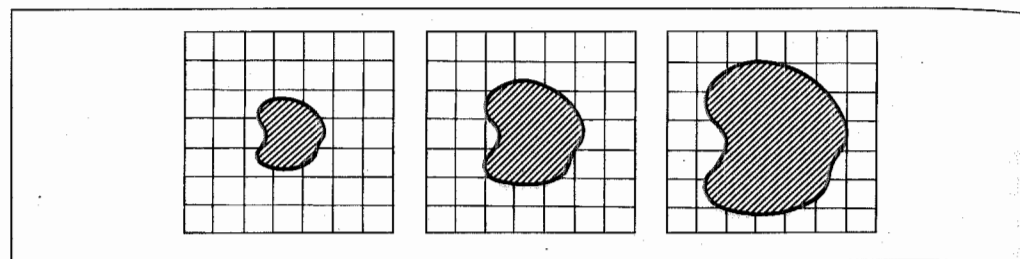


Fig. 11.28: Focusing can improve the accuracy of finite difference Poisson-Boltzmann calculations.

macromolecule, which is usually considered to lie between 2 and 4. This value of the dielectric constant is justified by the following arguments. The dielectric constant of a material is due to several factors, including its inherent polarisability and its ability to reorient internal dipoles within a changing electric field. A molecule that is fixed in conformation will not be able to change the orientations of its dipolar groups and so the only contribution to the dielectric constant will be due to polarisation effects. Polarisation effects alone lead to a dielectric constant of about 2 for organic liquids. If the conformation of the molecule can change then the dipolar effects should be taken into account, leading to an increased dielectric of 4. The atomic charges and van der Waals radii are often taken directly from an existing force field, though parameter sets designed specifically for use with the Poisson-Boltzmann method have been developed [Sitkoff *et al.* 1994].

The correct choice of grid size can be crucial to the success of a finite difference Poisson-Boltzmann calculation. The finer the lattice, the more accurate the results, though more computer time will be required. A grid size of 65^3 has been widely used. A technique called *focusing* can help alleviate some of these problems. In this method, a series of calculations are performed with the system occupying an ever greater fraction of the total grid box at each step. The boundary points in each new grid are internal points from its predecessor, as shown in Figure 11.28. Focusing enables better estimates of the potential values at the boundary to be obtained. The results can also depend upon the orientation of the solute(s) within the grid. The error associated with this can be reduced by performing a series of calculations on randomly translated and rotated copies of the system and then averaging the results.

11.10.5 Applications of Finite Difference Poisson-Boltzmann Calculations

A wide variety of problems have been studied using the finite difference Poisson-Boltzmann (FDPB) method. In addition to the numerical values that the method can provide, significant insights can often be gleaned by graphical examination of the electrostatic potential around the molecule [Honig and Nicholls 1995]. It is often found that the electrostatic potential around a protein calculated using the FDPB method differs significantly from that obtained with a uniform dielectric model. The location of the charged and polar groups in the protein and the shape of the molecule (which determines the shape of the boundary between the regions of high and low dielectric) significantly influence the shape of the potential. This

can be seen in Figure 11.29 (colour plate section), which shows the electrostatic potential around the enzyme trypsin. The activity of this enzyme is regulated *in vivo* by trypsin inhibitor, which is a smaller protein that binds strongly to trypsin. However, both trypsin and trypsin inhibitor have net positive charges. How then do the two molecules associate? If the electrostatic potential around trypsin is calculated assuming a uniform dielectric constant of 80 then, as expected, the potential is positive everywhere. However, when the effects of the dielectric boundary are included then a region of negative electrostatic potential appears in the region where the inhibitor binds. A second example is provided by the enzyme Cu-Zn superoxide dismutase, which catalyses the conversion of O_2^- radicals to O_2 and H_2O_2 . The rate constant for the reaction is high, being only about one order of magnitude smaller than the expected collision rate of the substrate with the entire enzyme. However, the active site constitutes a very small proportion of the surface, and uniform collisions of the substrate over the protein surface would not explain the observed kinetics. It has therefore been suggested that the substrate is 'steered' into the active site by the electric field of the protein. Figure 11.30 (colour plate section) shows the electrostatic potential around the enzyme (which is a dimer) with the active sites at the top left and bottom right. As can be seen, a concentrated region of positive electrostatic potential extends from the active site into solution [Klapper *et al.* 1986]. The cleft-like nature of the protein around the active site enhances the positive electrostatic potential by focusing electric field lines out into the solvent.

The finite difference Poisson-Boltzmann method can be used to calculate the electrostatic contribution to various processes such as solvation and the formation of intermolecular complexes. The electrostatic component of the solvation free energy equals the change in electrostatic energy for transfer from vacuum to the solvent where the electrostatic energy of a charge q_i in a potential ϕ_i equals $q_i\phi_i$. The solvation free energy is determined by performing two separate calculations using the same grids and the same solute dielectric but exterior dielectrics of 80 (when the solvent is water) and 1 (for the vacuum). Then ΔG_{elec} is given by:

$$\Delta G_{elec} = \frac{1}{2} \sum_i q_i (\phi_i^{80} - \phi_i^1) \quad (11.81)$$

The summation in Equation (11.81) is over all charges in the solute.

The change in free energy for the association of two molecules (assumed to have the same internal dielectric constant, ϵ_m) can be calculated using the finite difference Poisson-Boltzmann method. This problem is usefully discussed by dividing the free energy of association into a series of steps, as shown in Figure 11.31 [Gilson and Honig 1988]. First, the free energy associated with the transfer of the two isolated species from the solvent (dielectric constant ϵ_s) to a medium of dielectric ϵ_m is calculated in the same manner as for the solvation free energy (but here the transfer is from the solvent to a medium of dielectric ϵ_m , not to a vacuum). The free energy to bring the two molecules together is calculated using Coulomb's law in a medium of dielectric ϵ_m . Finally, the energy to transfer the complex from the medium of dielectric ϵ_m to the solvent is determined. The same procedure can be applied to other processes, such as the calculation of the free energy difference between two conformations in solution.

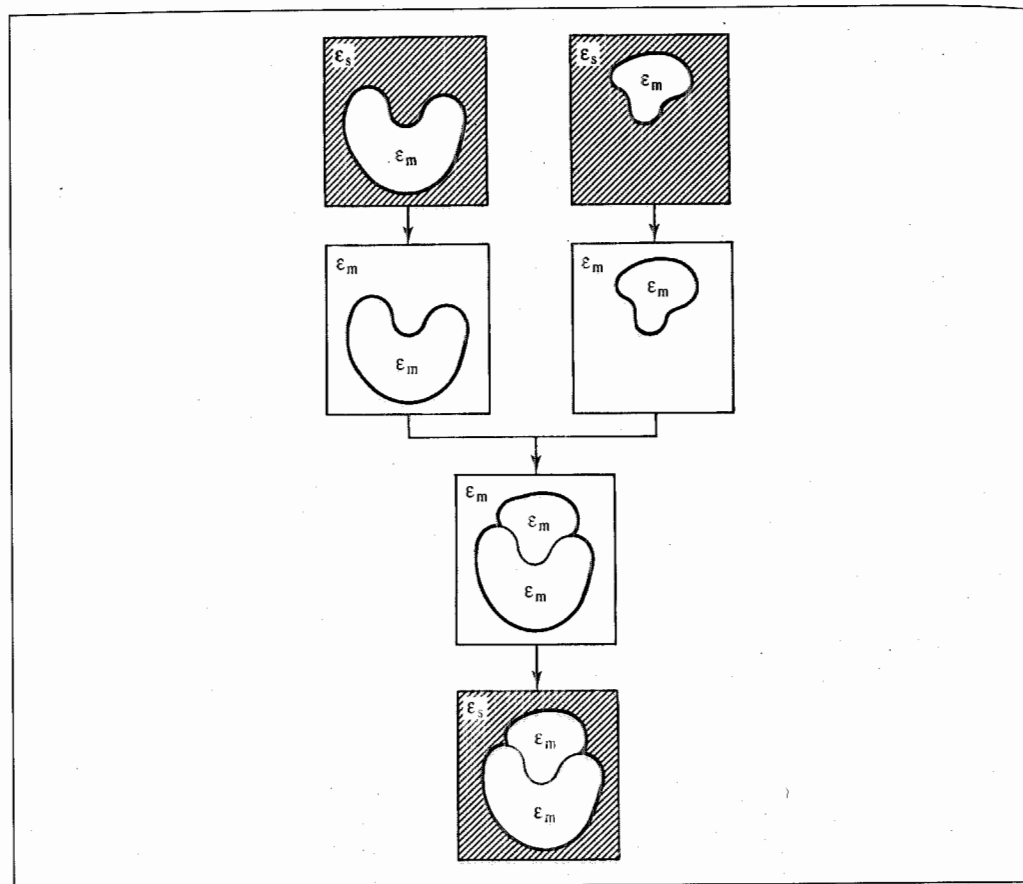


Fig. 11.31: Calculation of the electrostatic free energy of association of two molecules. (Figure adapted from Gilson M K and B Honig 1988. Calculation of the Total Electrostatic Energy of a Macromolecular System: Solvation Energies, Binding Energies and Conformational Analysis. Proteins: Structure, Function and Genetics 4:7-18.)

11.11 Non-electrostatic Contributions to the Solvation Free Energy

So far, we have only considered the electrostatic contribution to the free energy of solvation. Important though this is, there are some additional factors that contribute to the overall free energy of solvation, as shown in Equation (11.46). These extra contributions can be especially significant for solutes that are neither charged nor highly polar. The cavity and van der Waals terms are often combined and represented using an equation of the following form:

$$\Delta G_{\text{cav}} + \Delta G_{\text{vdw}} = \gamma A + b \quad (11.82)$$

where A is the total solvent accessible area and γ and b are constants. This linear dependence upon the area A can be explained as follows. The cavity term equals the work to create the cavity against the solvent pressure and the entropy penalty associated with the reorganisation of solvent molecules around the solute. The solvent molecules most affected by this reorganisation are those in the first solvation shell. The number of solvent molecules in the first solvation shell is approximately proportional to the accessible surface area of the solute. The solute-solvent van der Waals interaction energy would also be expected to be dependent primarily upon the number of solvent molecules in the first solvent shell, as van der Waals interactions fall off rapidly with distance. Hence both the cavity and van der Waals terms should be approximately proportional to the solvent accessible surface area. The parameters γ and b in Equation (11.82) are usually taken from experimentally determined free energies for the transfer of alkanes from vacuum to water. The parameter b is commonly set to zero, making the cavity plus van der Waals terms directly proportional to the solvent accessible surface area. Still's generalised Born/surface area model (GB/SA) uses the generalised Born approach for the electrostatic contribution together with a cavity and van der Waals surface area term in which the surface area is calculated using a variant of the Wodak and Janin algorithm, the constant γ having the value $7.2 \text{ cal}/(\text{mol } \text{\AA}^2)$ [Hasel *et al.* 1988]. As we have already stated, analytical first and second derivatives of the surface area with respect to the atomic coordinates can be rapidly determined using the Wodak-Janin method, so enabling the GB/SA method to be incorporated into energy-minimisation and molecular dynamics calculations.

The cavity and van der Waals contributions may also be modelled as separate terms. In some implementations an estimate of the cavity term may be obtained using scaled particle theory [Pierotti 1965; Claverie *et al.* 1978], which uses an equation of the form:

$$\Delta G_{\text{cav}} = K_0 + K_1 a_{12} + K_2 a_{12}^2 \quad (11.83)$$

The constants K depend upon the volume of the solvent molecule (assumed to be spherical in shape) and the number density of the solvent. a_{12} is the average of the diameters of a solvent molecule and a spherical solute molecule. This equation may be applied to solutes of a more general shape by calculating the contribution of each atom and then scaling this by the fraction of that atom's surface that is actually exposed to the solvent. The dispersion contribution to the solvation free energy can be modelled as a continuous distribution function that is integrated over the cavity surface [Floris and Tomasi 1989].

11.12 Very Simple Solvation Models

Some particularly simple solvation models include all contributions to the solvation free energy (including the electrostatic contribution) in an equation of the following form:

$$\Delta G_{\text{sol}} = \sum_i a_i S_i \quad (11.84)$$

where S_i is the exposed solvent accessible surface area of atom i , and the summation is over all atoms in the solute. a_i is a parameter that depends upon the nature of atom i . Despite the

obvious assumptions inherent in such an approach, it does have the advantage of providing an extremely rapid way to calculate a solvation contribution. Eisenberg and McLachlan developed such a model to study proteins, with the parameters a_i being derived by considering just five classes of atom (carbon, neutral oxygen and nitrogen, charged oxygen, charged nitrogen, and sulphur) [Eisenberg and McLachlan 1986]. The values themselves were obtained by fitting to experimentally determined free energies of transfer. Eisenberg and McLachlan applied their solvation model to a variety of problems, such as the recognition of misfolded protein structures and ligand binding.

11.13 Modelling Chemical Reactions

It is obviously important to be able to model chemical reactions, as these lie at the heart of chemistry and biochemistry. Most reactions of interest do not take place in the gas phase but in some medium, be it in a solvent, in an enzyme or on the surface of a catalyst. The environment can have a significant impact upon the reaction by speeding it up or slowing it down or even changing the reaction pathway. Good agreement can sometimes be obtained for calculations performed on isolated systems (i.e. in the gas phase), but to model the system properly the environment must be taken into account.

The preferred technique for modelling chemical reactions is usually considered to be quantum mechanics. Unfortunately, if one wishes to represent the whole system explicitly, the large number of atoms that must be considered means that *ab initio* quantum mechanics is rarely practical. Here we will consider three methods that have been used to study chemical reactions involving large systems. One strategy is to use a purely empirical approach. An alternative is to divide the system into two and treat the 'reaction region' using quantum mechanics, with the rest of the system being modelled using molecular mechanics. Third, we shall consider techniques such as the Car-Parrinello method and density functional theory, which, when allied to extremely powerful computers, can enable the entire reacting system to be simulated using quantum mechanics.

11.13.1 Empirical Approaches to Simulating Reactions

Despite the often-held belief that reactions can only be studied using quantum mechanics, this is by no means the case. Many research groups have developed force field models for studying reactions, which can provide very satisfactory results. Such force fields are used to estimate the activation energies of possible transition states to explain and to predict the stereo- and regioselectivity of the reaction. The force field model is usually derived by extending an existing force field to enable the structures and relative energies of transition structures to be determined.

Here we will illustrate the method using a single example. The aldol reaction between an enol boronate and an aldehyde can lead to four possible stereoisomers (Figure 11.32). Many of these reactions proceed with a high degree of diastereoselectivity (i.e. *syn* : *anti*) and/or enantioselectivity (*syn*-I : *syn*-II and *anti*-I : *anti*-II). Bernardi, Capelli, Gennari,

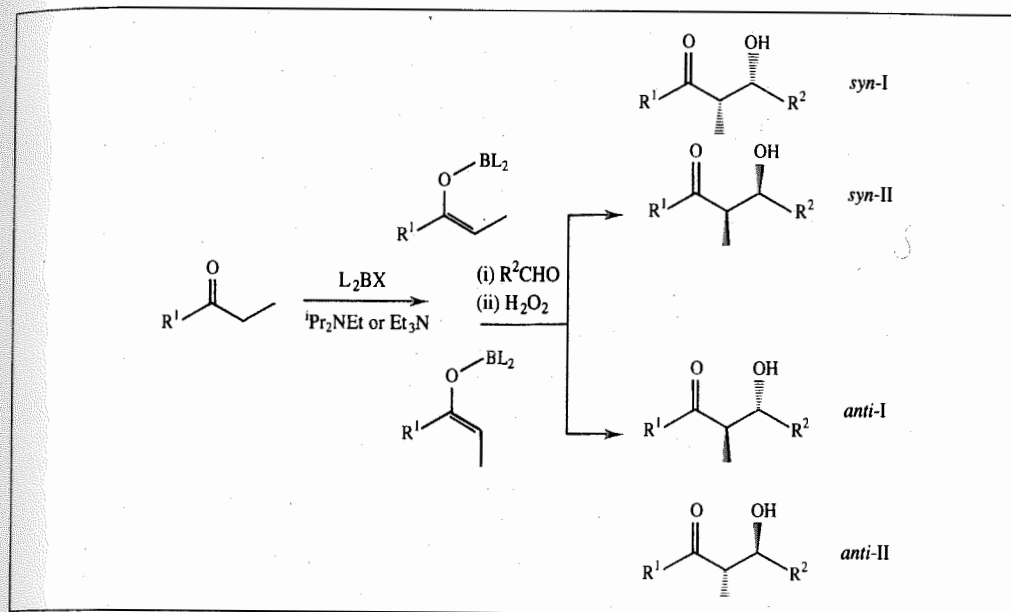


Fig. 11.32: The aldol reaction between an enol boronate and an aldehyde leads to four possible stereoisomers.

Goodman and Paterson studied this reaction using a force field based on MM2 [Bernardi *et al.* 1990]. The force field was parametrised to reproduce the geometries and relative energies of the chair and twist-boat transition structures with unsubstituted reactants, previously determined using *ab initio* methods (see Figure 11.33). It was assumed that the stereoselectivity was determined by the relative energies of the various possible transition structures (i.e. the reaction is assumed to be kinetically controlled).

The force field was then used to predict the results for the addition of the *E* and *Z* isomers of the enol boronate of butanone ($R^1 = \text{Me}$) to ethanol ($R^2 = \text{Me}$). The relevant transition structures are shown in Figure 11.34. A Boltzmann distribution, calculated at the temperature of the reaction (-78°C), predicted that the *Z* isomer would show almost complete *syn* selectivity (*syn* : *anti* = 99 : 1) and that the *E* isomer would be selective for the *anti* product (*anti* : *syn* = 86 : 14). These results were in good agreement with the experimental

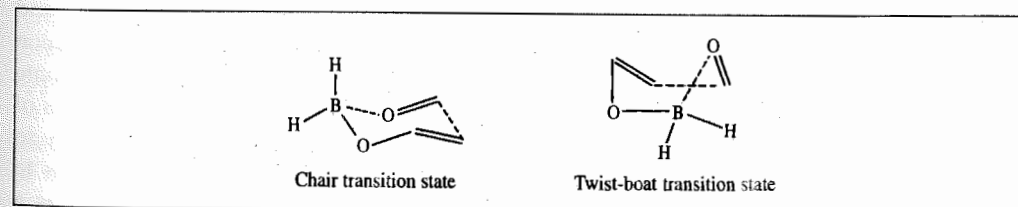


Fig. 11.33: Transition structures for the enol boronate/aldehyde reaction.

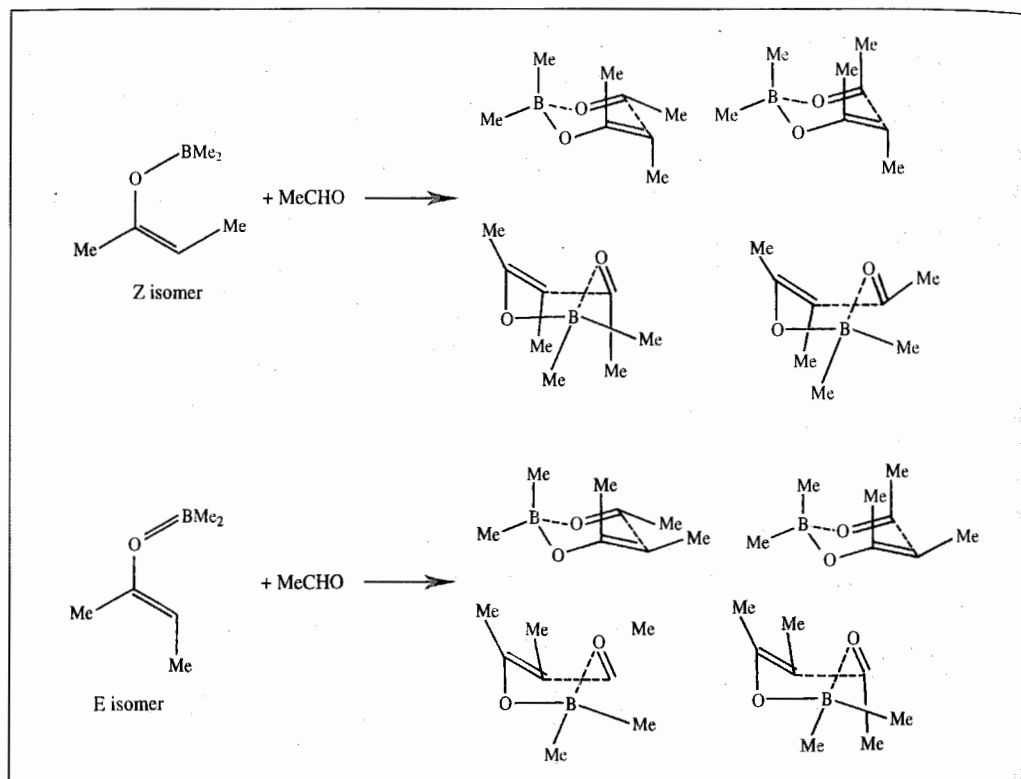


Fig. 11.34: Transition states for aldol reaction between butanone and ethanol.

observations. The major product in each case was obtained from a chair-like transition structure, but the reduced fraction of the *anti* product for the *E*-isomer was due to a significant contribution from the boat pathway, which leads to the *syn* product.

11.13.2 The Potential of Mean Force of a Reaction

A complete description of a chemical reaction needs to take account of solvent effects. The most realistic way to achieve this is by including explicit solvent molecules. A classic example of how to tackle this problem is Jorgensen's study of the nucleophilic attack of the chloride anion on methyl chloride [Chandrasekhar *et al.* 1985; Chandrasekhar and Jorgensen 1985]. This reaction proceeds via the S_N2 reaction, in which the chloride anion approaches along the carbon-chlorine bond of methyl chloride to give a five-coordinate transition state, which then collapses to give the products. We considered some aspects of the energy surface for this system in Section 5.9, though there we were interested only in the energy change for the gas-phase reaction. The aim of Jorgensen's calculation was to obtain a potential of mean force for the reaction (i.e. the change in the free energy as a function of the reaction coordinate) in a variety of solvents.

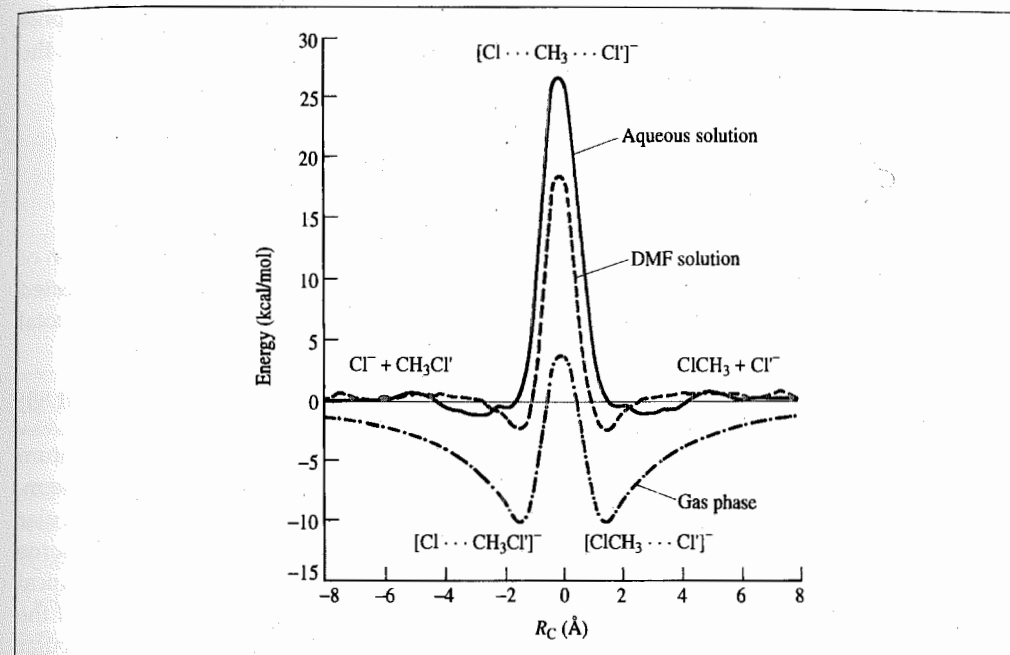


Fig. 11.35: Potential of mean force for the Cl⁻ + MeCl reaction in various solvents. (Figure redrawn from Chandrasekhar J and W L Jorgensen 1985. Energy Profile for a Nonconcerted S_N2 Reaction in Solution. *Journal of the American Chemical Society* 107:2974-2975.)

The first step was to determine the quantum mechanical reaction pathway; a series of geometries along the path were determined using the path-following method of Gonzalez and Schlegel [Gonzalez and Schlegel 1988]. The solute-solvent interactions were modelled using Lennard-Jones and electrostatic terms in which the parameters smoothly varied with the reaction coordinate. To perform the Monte Carlo simulations, umbrella sampling was employed to constrain the geometry of the solute to a series of windows along the pathway, and thus calculate the potential of mean force. Preferential sampling methods were used, so that solvent molecules near the solute were sampled more often than solvent molecules further away.

The results are summarised in Figure 11.35, which shows how the potential of mean force varies for the reaction in the gas phase, in water and in dimethyl formamide (DMF). The results exhibit a number of interesting features. In the gas phase an ion-dipole complex forms, giving a minimum in the free energy profile. There is then an activation barrier of approximately 13.9 kcal/mol to reach the pentagonal transition state. In aqueous solution, no ion-dipole minimum is observed. This is because any favourable contribution due to the formation of the ion-dipole is compensated for by the energy lost in the desolvation of the chloride ion. There is then a large activation free energy barrier of approximately 26.3 kcal/mol from the ion-dipole pair to the transition state. This barrier is much larger than in the gas phase because of the poorer solvation of the transition state relative to the

ion-dipole complex. In DMF (a solvent with smaller anion-solvating ability), the ion-dipole complex is at a minimum in the free energy as less energy is required to desolvate the chloride anion in this solvent.

11.13.3 Combined Quantum Mechanical/Molecular Mechanical Approaches

One approach to the simulation of chemical reactions in solution is to use a combination of quantum mechanics and molecular mechanics. The 'reacting' parts of the system are treated quantum mechanically, with the remainder being modelled using the force field. The total energy E_{TOT} for the system can be written:

$$E_{TOT} = E_{QM} + E_{MM} + E_{QM/MM} \quad (11.85)$$

where E_{QM} is the energy of those parts of the system treated exclusively with quantum mechanics, and E_{MM} is the energy of the purely molecular mechanical parts of the system. $E_{QM/MM}$ is the energy of interaction between the quantum mechanical and molecular mechanical parts of the system. This is described by a Hamiltonian $\mathcal{H}_{QM/MM}$. In some cases, $E_{QM/MM}$ is due entirely to non-bonded interactions between the quantum mechanical and molecular mechanical atoms. An example where this could arise would be if all of the atoms in the reacting species were treated quantum mechanically, with molecular mechanics being used exclusively for the solvent. For example, Cl^- and MeCl could be treated using quantum mechanics and solvent with molecular mechanics. In this case, the Hamiltonian $\mathcal{H}_{QM/MM}$ can be written:

$$\mathcal{H}_{QM/MM} = - \sum_i \sum_M \frac{q_M}{r_{i,M}} + \sum_\alpha \sum_M \frac{Z_\alpha q_M}{R_{\alpha,M}} + \sum_\alpha \sum_M \left(\frac{A_{\alpha,M}}{R_{\alpha,M}^{12}} - \frac{C_{\alpha,M}}{R_{\alpha,M}^6} \right) \quad (11.86)$$

The subscript i in Equation (11.86) refers to a quantum mechanical electron and the subscript α to a quantum mechanical nucleus. The subscript M indicates a molecular mechanical nucleus and q_M is its partial atomic charge. There are thus electrostatic interactions between the electrons of the quantum mechanical region and the molecular mechanical nuclei, electrostatic interactions between quantum mechanical and molecular mechanical nuclei, and van der Waals interactions between the quantum mechanical and molecular mechanical atoms. The second and third terms in Equation (11.86) do not involve electronic coordinates and so can be calculated in a straightforward way (i.e. they are constant for a given nuclear configuration). The first term must be incorporated into the quantum mechanical calculation via one-electron integrals added to the one-electron matrix, H^{core} . These one-electron integrals have the form:

$$\int \phi_\mu(1) \frac{1}{r_{1,M}} \phi_\nu(1) d\nu(1) \quad (11.87)$$

In some cases, the quantum mechanical and molecular mechanical regions are in the same molecule and so there are bonds between atoms from each region. The energy $E_{QM/MM}$ must now contain terms that describe this interaction. This can be done by adding a molecular mechanical-like energy which contains bond-stretching, angle-bending

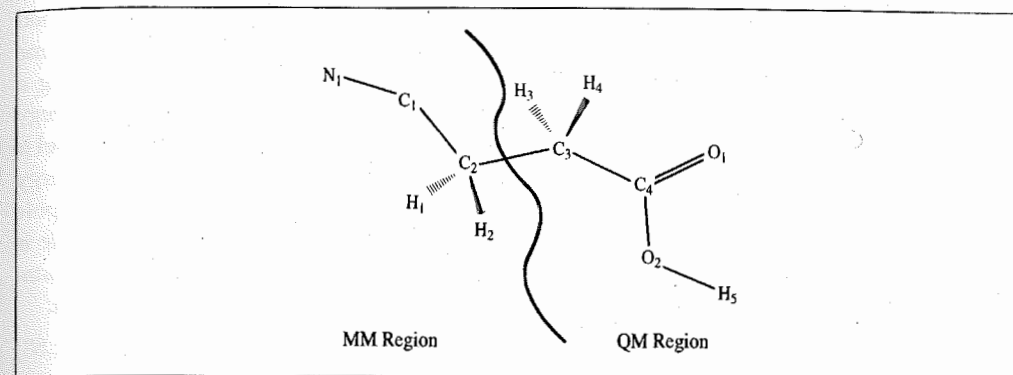


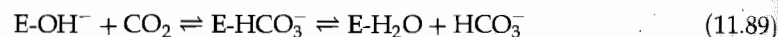
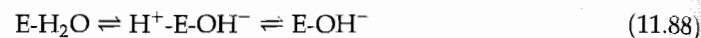
Fig. 11.36: The division of a molecule into quantum mechanical and molecular mechanical regions, with the molecular mechanical contributions as indicated.

and torsional terms for atoms from both the quantum mechanical and molecular mechanical sets. This is illustrated in Figure 11.36, which shows which terms would be included in $E_{QM/MM}$.

Various combined quantum mechanical/molecular mechanical implementations have been described [Warshel and Levitt 1976; Singh and Kollman 1986; Field *et al.* 1990; Maseras and Morokuma 1995]. These implementations differ in the quantum mechanical theory that is used (semi-empirical, *ab initio*, valence bond or density functional theory), the molecular mechanical model, and the way in which the solvent is represented (either explicitly or using a simplified model). Another important difference is the way in which the junction between the QM/MM regions is handled. In particular, one must avoid half-filled orbitals for the quantum mechanical region, which would arise if the connection bonds were simply truncated. Two general approaches to this problem have been developed. In one approach, a hybrid sp^2 orbital containing one electron is established along the QM-MM [Warshel and Levitt 1976]. The alternative method simply includes 'link' atoms (typically hydrogen atoms), which ensure that valency is maintained. Interactions between these link atoms and the molecular mechanical region is reduced in magnitude or completely neglected. Comparisons of the two approaches on simple model systems suggest that neither is systematically better than the other provided care is taken in the formulation [Reuter *et al.* 2000].

Combined quantum mechanical/molecular mechanical methods are not, of course, restricted to studies of reactions but can also be used to study association processes and conformational transitions. Most implementations use a two-zone model as described above, but Morokuma and colleagues have described a multilayered approach called ONIOM [Svensson *et al.* 1996]. ONIOM is a particularly apt name given that a typical calculation is constructed from a series of layers. For example, a three-layer ONIOM calculation on the Diels-Alder reaction involved an inner core treated with the B3LYP density functional approach, the intermediate layer with a Hartree-Fock level of theory and the outer layer with MM3. A particular feature of ONIOM and its related methods is that they provide rigorous gradients and second derivatives, so enabling properties such as vibrational frequencies to be calculated [Dapprich *et al.* 1999]

The objective of many of the research groups involved in the development of combined quantum mechanical/molecular mechanical models has been the simulation of enzyme reactions. Warshel has reported studies in which the reaction centre is treated using a valence bond model [Warshel 1991; Åqvist and Warshel 1993]. The first part of his strategy is a calibration of the valence bond model for the reference reaction in solution. This model is then used to simulate the enzyme reaction using molecular dynamics and free energy perturbation methods, with solvent effects being treated using the Langevin dipole model. Warshel has extensively studied a wide range of enzyme systems. One example is his study of the enzyme carbonic anhydrase, which is a zinc-containing enzyme that catalyses the reversible hydration of carbon dioxide according to the following mechanism:



where E represents the enzyme. In the first step, a bound water molecule is proteolysed and the protein is transferred to solution. This is the rate-determining step of the reaction. In the second step, CO_2 is converted to HCO_3^- . Åqvist, Fothergill and Warshel have examined both steps; here we concentrate on their results for the nucleophilic attack on the carbon dioxide (Equation (11.89)) [Åqvist *et al.* 1993]. Simulations of the hydration reaction of CO_2 in water were performed to find valence bond parameters which reproduced the experimentally observed value. This valence bond model was then used to simulate the same reaction in the enzyme. The resulting free energy profiles for the reference reaction and the enzyme reaction are shown in Figure 11.37. These results suggest that the enzyme markedly lowers the activation barrier for the reaction and that the reaction is less exothermic in the enzyme than in water ($\Delta G^\ddagger = 6.3 \text{ kcal/mol}$ versus 11.9 kcal/mol ; $\Delta G^0 = -4.8 \text{ kcal/mol}$ versus -10.5 kcal/mol). The experimental values were estimated to be $\Delta G^\ddagger = 7.1 \text{ kcal/mol}$ and $\Delta G^0 = -4.1 \text{ kcal/mol}$. The enzyme thus speeds up the reaction by a factor of about 10^3 compared with aqueous solution. The transition-state geometry obtained from the simulation was found to be similar to geometries obtained using gas-phase *ab initio* calculations.

11.13.4 *Ab Initio* Molecular Dynamics and the Car-Parrinello Method

The 'ideal' way to simulate reactions (and indeed many other processes where we might wish to derive properties dependent upon the electronic distribution) would of course be to use a fully quantum mechanical approach.

In principle, it would be relatively straightforward to use a quantum mechanical model to determine the forces required by molecular dynamics or the energies for a Monte Carlo simulation algorithm. Hartree-Fock calculations are normally solved using iterative matrix diagonalisation techniques, as discussed in Chapter 2. Density functional calculations can also be tackled using such methods. However, for systems with many atoms and/or basis functions such calculations can be very time-consuming, and it may also be difficult to achieve convergence. Even with pseudopotentials, the number of plane-wave basis functions that can be required for density functional calculations may be very large, and as the number of occupied orbitals is often considerable, it can be a major task to solve the Kohn-Sham equations and determine the energy of a given configuration of atoms. In

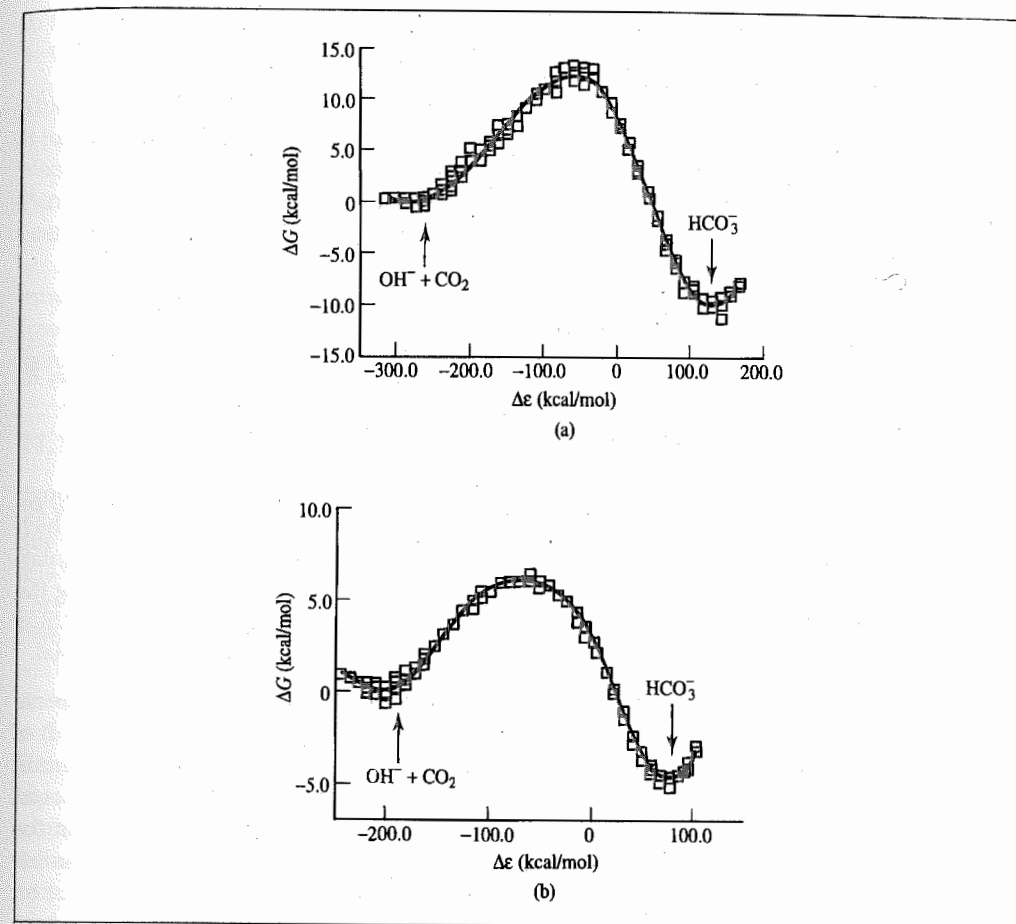


Fig. 11.37: Free energy profile for the nucleophilic attack of water on CO_2 (a) in aqueous solution and (b) in the enzyme carbonic anhydrase. (Graphs redrawn from Åqvist J, M Fothergill and A Warshel 1993. *Computer Simulation of the $\text{CO}_2/\text{HCO}_3^-$ Interconversion Step in Human Carbonic Anhydrase I*. Journal of the American Chemical Society 115:631-635.)

1985, Car and Parrinello described a method that brought together a number of key concepts that we have considered in earlier chapters [Car and Parrinello 1985; Remler and Madden 1990]. They were primarily concerned with the problem of performing *ab initio* simulations involving both the electronic and the nuclear motions ('total energy' simulations or '*ab initio* molecular dynamics'). However, their scheme can be used to perform energy minimisation or simply to determine the basis set coefficients for a fixed atomic configuration.

A key feature of the Car-Parrinello proposal was the use of molecular dynamics and simulated annealing to search for the values of the basis set coefficients that minimise the electronic energy. In this sense, their approach provides an alternative to the traditional matrix diagonalisation methods. In the Car-Parrinello scheme, 'equations of motion' for

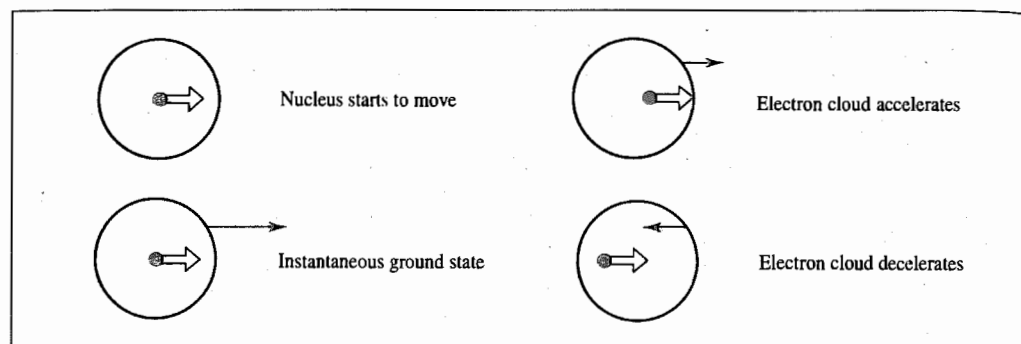


Fig. 11.38: Lag effects in *ab initio* molecular dynamics. (Figure redrawn from Payne M C, M P Teter, D C Allan, R A Arias and D J Joannopoulos 1992. *Iterative Minimisation Techniques for Ab Initio Total-Energy Calculations: Molecular Dynamics and Conjugate Gradients*. *Reviews of Modern Physics* 64:1045–1097.)

the coefficients are set up, and then molecular dynamics is used to move the system through the space of the basis set coefficients. Starting from a random set of coefficients (which correspond to a high energy) the system moves downhill on the energy surface, accumulating 'kinetic energy'. Simulated annealing was proposed as a mechanism for preventing the system becoming trapped in a local minimum. The SHAKE algorithm (see Section 7.5) is used to impose constraints on the system to ensure that the orbitals remain orthonormal.

To perform *ab initio* molecular dynamics, Car and Parrinello suggested that the electronic and nuclear dynamics could be performed simultaneously. Somewhat surprisingly, it is found that in the Car-Parrinello scheme it is not necessary for the electronic configuration to be at a minimum in coefficient space for each molecular dynamics time step, even though this gives errors in the forces on the nuclei. It can be shown that the errors in the nuclear forces are cancelled by the associated errors in the electronic motion. One possible explanation for this rather strange (and fortuitous!) result is to consider the motion of an atom with a single occupied molecular orbital. If the nucleus starts to move with a constant velocity, then the orbital will initially lag behind the nucleus. The orbital starts to accelerate until it eventually overtakes the nucleus. Having overtaken the nucleus, the orbital starts to slow down, until the nucleus overtakes the orbital, and so on, as illustrated in Figure 11.38.* An important practical feature of this molecular dynamics approach is that the fictitious masses assigned to the coefficients must be chosen so that the frequencies of electron motion are higher than those of the nuclei to avoid energy exchange. This can have the practical consequence of requiring a smaller time step, which adds to the computational cost.

An alternative to the Car-Parrinello method is the following scheme, which separates the electronic and nuclear motions:

1. Calculate the forces on the nuclei.
2. Move the nuclei according to the molecular dynamics integration scheme.

* It should be pointed out that not all workers in the field are convinced of this particular explanation, appealing though it is.

3. Optimise the electronic configuration for the new nuclear configuration.
4. Go to step 1.

This algorithm alternates between the electronic structure problem and the nuclear motion. It turns out that to generate an accurate nuclear trajectory using this decoupled algorithm the electrons must be fully relaxed to the ground state at each iteration, in contrast to the Car-Parrinello approach, where some error is tolerated. This need for very accurate basis set coefficients means that the minimum in the space of the coefficients must be located very accurately, which can be computationally very expensive. However, conjugate gradients minimisation is found to be an effective way to find this minimum, especially if information from previous steps is incorporated [Payne *et al.* 1992]. This reduces the number of minimisation steps required to locate accurately the best set of basis set coefficients.

11.13.5 Examples of *Ab Initio* Molecular Dynamics Simulations

As we have already observed, liquid water is one of the most challenging systems to model due to its almost unique properties such as hydrogen bonding and high dielectric constant. It is therefore not surprising that it was one of the first systems to be considered for a true Car-Parrinello *ab initio* molecular dynamics simulation [Laasonen *et al.* 1993; Sprik *et al.* 1996]. As the calculations were performed using density functional methods the two studies were thus not only designed to actually perform the simulation but also to investigate a variety of DFT models. Specifically, a variety of the gradient-corrected functionals discussed in Section 3.7.3 were considered, together with two different pseudopotential schemes. These latter differed in the number of plane waves needed; in the first study, a so-called supersoft pseudopotential was employed, which required fewer plane waves than the more conventional pseudopotential used in the second study. This was largely driven by the available computational resources; there were some shortcomings with the supersoft pseudopotential, but it did enable the simulation to be run on what would now be considered a very modest computer. It is worth recalling from Section 3.8.6 that the combination of plane waves and density functional theory provides a very natural and appealing way to tackle periodic systems. However, there we were concerned with naturally periodic systems, whereas for *ab initio* molecular dynamics it is the use of periodic boundary conditions which gives rise to the periodicity.

The simulations were restricted to a relatively small number of molecules (32) under periodic boundary conditions. Only rather short simulations were possible (of the order of 5 ps), but it was still possible to determine many of the standard structural properties such as the radial distribution function together with properties such as the vibrational spectra, which provide information on hydrogen bonding within the system. More recent simulations using a larger number of molecules and for a longer time focused on the molecular charge distribution and polarisation effects [Silvestrelli and Parrinello 1999]. A broad distribution was found for the dipole moment around an average value of 3.0 D. This may have important implications for empirical potentials, which are often parametrised to reproduce a lower value around 2.6 D. In addition, the anisotropy of the electronic charge distribution in the water molecule was found to be reduced in the liquid.

Building upon the earlier simulations, a subsequent study investigated systems containing hydronium and hydroxyl ions in water [Tuckerman *et al.* 1995a, b]. Protons show exceptionally high mobilities that are far in excess of the values expected from a straightforward diffusion process. A model that accounts for this (the *Grotthuss mechanism*) involves the proton jumping from the oxygen atom of one water molecule to another. This simple picture works very well for proton conduction in ice, but the situation is more complicated for the liquid species. The simulation of a single hydronium ion (H_3O^+) in water showed that for about 60% of the time the proton is associated with a single water molecule, with the three protons making hydrogen bonds to three neighbouring molecules, giving an H_9O_4^+ complex. For the remaining 40% of the time the proton could not be assigned to a unique oxygen atom but was shared between the oxygen atoms of two water molecules, to give an H_5O_2^+ structure. Close examination of these two structures indicated that they were, in fact, part of the same fluctuating complex. Much less experimental information is available about the OH^- ion, which the simulation suggests is coordinated to four water molecules, each pointing one OH bond towards it. This H_9O_5^- species remains intact for about 2–3 ps before one of the hydrogen bonds breaks, giving a transient tetrahedral H_7O_4^- complex.

Liquid hydrogen fluoride is another fluid of interest due to its strong hydrogen-bonding potential. Experimental data suggest the existence of chain-like structures, each containing between six and eight HF molecules held together by hydrogen bonds. In the liquid these chains adopt a zig-zag conformation and are significantly entangled. In addition, there is the possibility of branched structures forming, but the relative importance of these is a matter of debate. The structure of the liquid is very sensitive to the nature of the potential model. The *ab initio* molecular dynamics simulations used a density functional approach, and it was necessary to use a gradient-corrected functional in order to describe the system correctly. The simulation contained 54 molecules, with the production phase lasting 0.8 ps [Rothlisberger and Parrinello 1997]. Although the data from the simulation were rather noisy due to the short simulation time, a number of features were apparent. For example, a small degree of branching was observed, with a difference between the likelihood of branching at the hydrogen (1%) and fluorine atoms (6%).

Ab initio molecular dynamics has been applied to many 'materials science' problems. One interesting early application was the *ab initio* molecular dynamics simulation of the reaction between a chlorine molecule and a silicon surface [Stich *et al.* 1994]. This reaction is particularly important in silicon chip manufacture, where the dissociative chemisorption of chlorine (and other halogens) is widely used for processes such as dry etching and surface cleaning. A series of simulations was performed, in each of which a chlorine molecule was 'fired' towards the silicon surface. The subsequent motion and reaction was then determined using the *ab initio* molecular dynamics approach based upon conjugate gradients minimisation. The motions of the nuclei were determined using the Verlet algorithm with a time step of approximately 0.5 fs, and each simulation was performed for a total time of between 200 and 400 fs.

The silicon surface contains chains of atoms that are formally bonded to just three other atoms. These atoms compensate for the lack of a full valence complement of bonds by π

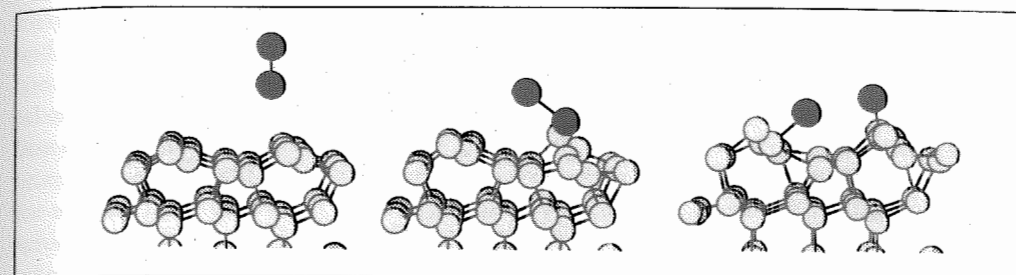


Fig. 11.39: Structural changes observed during the reaction of a chlorine molecule with a silicon surface. (Figure redrawn from Stich I, A De Vita, M C Payne, M J Gillan and L J Clarke 1994. *Surface Dissociation from First Principles: Dynamics and Chemistry*. *Physical Review B* 49:8076–8085.)

bonding along the chains (Figure 11.39). These π -bonded chains represent regions of high electron density, with the valleys between the chains being of relatively low density. Despite this difference in electron density, the chlorine molecule dissociated when it was directed towards either of these two regions. Bonds form between the chlorine atoms and π -bonded silicon atoms, which in turn causes a change in the local hybridisation from sp^2 to sp^3 . This then leads to a large local deformation, which lifts the silicon atoms involved above the π -bonded chains. Many other processes of significant commercial interest are now within the scope of the *ab initio* simulation technique, a more recent example being the Ziegler–Natta catalysed polymerisation of ethylene [Boero *et al.* 1999].

A somewhat more unusual illustration of the use of *ab initio* molecular dynamics was of the effect of a knot on the breaking strength of a polymer strand [Saitta *et al.* 1999]. A simple linear alkane formed the basis for the work; the initial calculations involved stretching *n*-decane until one of the bonds broke, to give two radicals. An analogous calculation involving a polyethylene chain with a trefoil knot was then performed. In this case the chain broke at the entrance to the knot (Figure 11.40). Of some interest is the fact that the presence of the knot significantly weakens the strand, as measured by the strain energy per C–C bond in the chain (12.7 kcal/mol for the knotted strand and 16.2 kcal/mol for the linear unknotted case).

Our small selection of examples has tended to concentrate on those which involve the making or breaking of bonds, but this is not of course a requirement for using *ab initio* molecular dynamics. Systems of particular interest are those for which it is difficult to generate empirical force-field models. One such example is the study by de Wijs and colleagues of the viscosity of liquid iron under the conditions believed to exist at the Earth's core [de Wijs *et al.* 1998]. The Earth's magnetic field is believed to arise from the convection of this liquid, an understanding of which is clearly dependent upon knowledge of the viscosity of the medium. Estimates of this viscosity vary over many orders of magnitude; for obvious reasons, it is unlikely that this uncertainty will be resolved by experimental measurements. Two regions were of particular interest: the boundary between the solid inner core and the molten outer core, and the boundary between the core and the mantle. The temperatures of these two regions are somewhat uncertain; for the inner core boundary a temperature of 6000 K was assumed and for the core–mantle boundary two temperatures (4300 K and

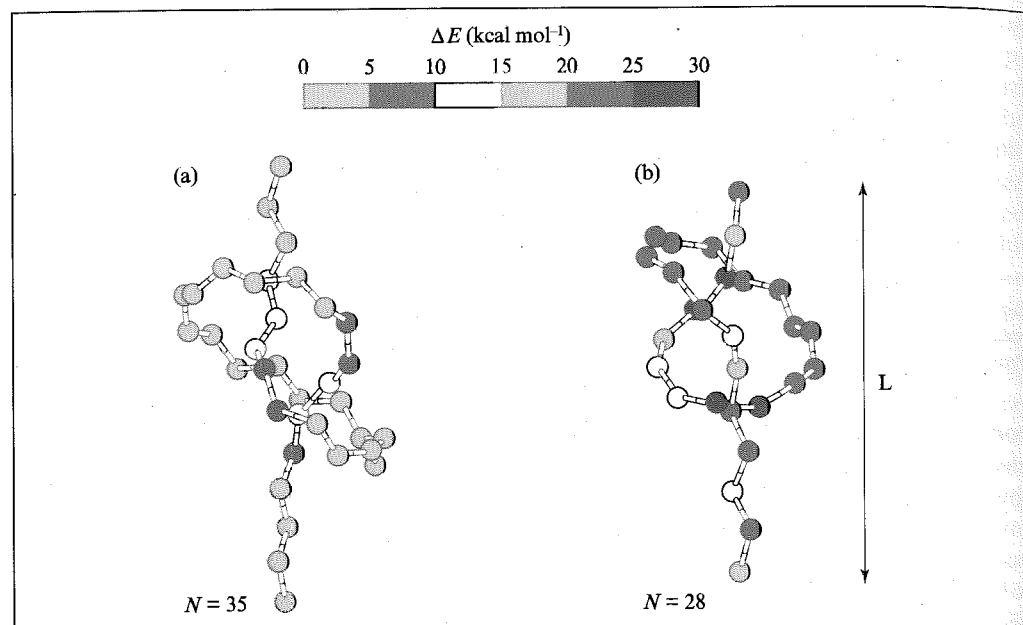


Fig. 11.40: Distribution of strain energy in two knotted polymer chains containing 35 (left) and 28 (right) carbon atoms. The strain energy is localised and most of the bonds immediately outside the entrance point to the knot. (Figure redrawn from Saitta A M, P D Sooper, E Wasserman and M L Klein 1999. Influence of a knot on the strength of a polymer strand. *Nature* 399:46–48.)

3500 K) were investigated. From an equation of state for iron the densities at these temperatures could be predicted to enable the simulations to be performed. A periodic system containing 64 atoms was used and the simulation run for 2 ps after equilibration. The calculated pressure agreed within 10% with the 'experimental' values (330 GPa at the inner core boundary and 135 GPa at the core–mantle boundary). Additional parameters could also be calculated, including the viscosity, the values for which were at the low end of previous suggestions.

11.14 Modelling Solid-state Defects

Materials that contain defects and impurities can exhibit some of the most scientifically interesting and economically important phenomena known. The nature of disorder in solids is a vast subject and so our discussion will necessarily be limited. The smallest degree of disorder that can be introduced into a perfect crystal is a *point defect*. Three common types of point defect are vacancies, interstitials and substitutionals. *Vacancies* form when an atom is missing from its expected lattice site. A common example is the Schottky defect, which is typically formed when one cation and one anion are removed from the bulk and placed on the surface. Schottky defects are common in the alkali halides. *Interstitials* are due to the presence of an atom in a location that is usually unoccupied. A

Frenkel defect arises when an ion (usually the smaller cation) is removed from its regular site and is placed in an interstitial position. Frenkel defects are common when there is a significant difference in size between the cation and the anion (such as AgBr). *Substitutionals* occur when a foreign ion occupies a regular lattice site. The presence of additional atoms may be due either to accidental impurities or to deliberate doping. These impurities can occupy interstitial sites or they may substitute for an existing atom. The substitution of one atom for another is often difficult due to the tightly packed nature of most solids, but it can be achieved if the atomic sizes are approximately equal. An *aliovalent* substituent is one which has a different valence state from the host. An example of this is the introduction of magnesium (as Mg^{2+} ions) into NaCl (the term *heterovalent* is also used). An additional consequence of this is the formation of cation vacancies, which neutralise the extra charge of the impurity. Defects must always be present in any crystalline solid above absolute zero, purely on entropic grounds (though the concentration can still be very small). These are known as intrinsic defects. Extrinsic defects, by contrast, arise from the accidental or deliberate incorporation of impurities.

Two point defects may aggregate to give a defect pair (such as when the two vacancies that constitute a Schottky defect come from neighbouring sites). Clusters of defects can also form. These defect clusters may ultimately give rise to a new periodic structure or to an extended defect such as a dislocation. Increasing disorder may alternatively give rise to a random, amorphous solid. As the properties of a material may be dramatically altered by the presence of defects it is obviously of great interest to be able to understand these relationships and ultimately predict them. However, we will restrict our discussion to small concentrations of defects.

The most direct effect of defects on the properties of a material usually derive from the altered ionic conductivity and diffusion properties. So-called superionic conductors are materials which have an ionic conductivity comparable to that of molten salts. This high conductivity is due to the presence of defects, which can be introduced thermally or via the presence of impurities. Diffusion affects important processes such as corrosion and catalysis. The specific heat capacity is also affected; near the melting temperature the heat capacity of a defective material is higher than for the equivalent ideal crystal. This reflects the fact that the creation of defects is enthalpically unfavourable but is more than compensated for by the increase in entropy, so leading to an overall decrease in the free energy.

Energy minimisation, molecular dynamics and Monte Carlo simulations have all been used to study the nature of defects and their influence on the material's properties. Special treatments are required for defects, because they can lead to very long-range perturbations. This is particularly the case when the defect has a net positive or negative charge. The calculation of defect energies using energy minimisation is commonly performed using a two-region strategy, based upon a paper published by Mott and Littleton [Mott and Littleton 1938]. The ions in the inner region are fully and explicitly affected by the presence of the defect, in contrast to the ions in the second region (which extends to infinity). Labelling the inner region as 1 and the outer region as region 2 leads to the following expression for the total energy of the system:

$$E = E_1(\mathbf{x}) + E_{12}(\mathbf{x}, \mathbf{y}) + E_2(\mathbf{y}) \quad (11.90)$$

where E_1 is the energy of region 1 (dependent on the coordinates \mathbf{x} of the ions within region 1), E_2 is the energy of region 2 (dependent on the displacements \mathbf{y} of the ions in region 2) and E_{12} is the energy of interaction between the two regions. It is assumed that E_2 is a quadratic function of the displacements, which means that it can be written as follows:

$$E_2(\mathbf{y}) = \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} \quad (11.91)$$

where \mathbf{A} is a force constant matrix. The harmonic-well assumption for the ions in region 2 is appropriate provided the perturbations are small. In practice, it also requires the bulk lattice to be optimised before the defect calculation is performed. At equilibrium, the derivative of the energy with respect to these coordinates \mathbf{y} is zero, from which we can derive:

$$(\partial E / \partial \mathbf{y})_{\mathbf{x}} = (\partial E_{12}(\mathbf{x}, \mathbf{y}) / \partial \mathbf{y})_{\mathbf{x}} + \mathbf{A} \cdot \mathbf{y} = 0 \quad (11.92)$$

This can be used to eliminate the energy E_2 from Equation (11.90), giving the following expression for the total energy:

$$E = E_1(\mathbf{x}) + E_{12}(\mathbf{x}, \mathbf{y}) - \frac{1}{2} (\partial E_{12}(\mathbf{x}, \mathbf{y}) / \partial \mathbf{y})_{\mathbf{x}} \cdot \mathbf{y} \quad (11.93)$$

In order to determine the energy it would thus seem that it is necessary merely to minimise E with respect to the positions \mathbf{x} and the displacements \mathbf{y} . However, a complication arises due to the fact that the displacements in the outer region are themselves a function of the inner-region coordinates. The solution to this problem is to require that the forces on the ions in region 1 are zero, rather than that the energy should be at a minimum (for simple problems the two are synonymous, but in practice there may still be some non-zero forces present when the energy minimum is considered to have been located). An additional requirement is that the ions in region 2 need to be at equilibrium.

Implementation of the two-region method requires calculation of the interaction between the ions in region 1 and region 2. For short-range potentials (e.g. the van der Waals contribution) it is only the inner part of region 2 that contributes significantly to the energy, E , and the forces on ions in region 1. Thus in current practical implementations of the method the outer region is subdivided into two regions, 2a and 2b (Figure 11.41). In region 1, an atomistic representation is used with full relaxation of the ions. Region 2a also contains explicit ions, whereas in region 2b it is assumed that the only effect of the defect is to change the polarisation of the ions. An iterative approach is used to identify the configuration in which the forces on the ions in region 1 are zero and the ions in region 2a are at equilibrium. The displacements of the ions in region 2a are commonly determined using just the electrostatic force from the defect species alone and equals the force due to any interstitial species less the force due to any vacancies (based on

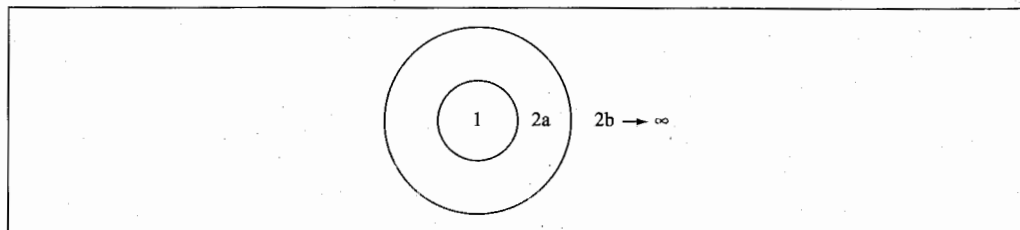


Fig. 11.41: Two region scheme used in Mott-Littleton calculations.

the position of the original vacancy site). A Newton-Raphson approach is employed, wherein the displacement \mathbf{y} of an ion from its current position is given by:

$$\mathbf{y} = -\mathcal{Y}' \cdot \mathcal{Y}''^{-1} \quad (11.94)$$

From these calculated displacements, the contributions to the energy from the ions in region 2a can be determined. Finally, it is necessary to determine the contribution from the ions in region 2b. As we have mentioned, these are not included explicitly but are considered to polarise due to the electrostatic field from the total charge on the defect. This contribution is given by the following summation:

$$E_{2b} = -Q \sum_j \frac{q_j (\mathbf{y}_j \cdot \mathbf{R}_j)}{|\mathbf{R}_j|^3} \quad (11.95)$$

where Q is the total effective charge on the defect and q_j is the charge on ion j in region 2b, with \mathbf{y}_j and \mathbf{R}_j being its displacement and equilibrium position. The Mott-Littleton approximation provides a formula for the displacement, leading to the following expression for the energy (for an isotropic medium):

$$E_{2b} = -\frac{Q^2 V_m}{8\pi\epsilon_0} \sum_j \frac{M_i q_j}{|\mathbf{R}_j|^4} \quad (11.96)$$

where V_m is the unit cell volume and M_i is the Mott-Littleton factor, which is related to the polarisabilities of the ions and the dielectric constant by:

$$M_i = \left(\frac{\alpha_i}{\sum_j \alpha_j} \right) \left(1 - \frac{1}{\epsilon} \right) \quad (11.97)$$

The summation is over the different types of ion in the unit cell. The summation can be written as an analytical expression, depending upon the lattice structure (the original Mott-Littleton paper considered the alkali halides, which form simple cubic lattices) and evaluated in a manner similar to the Ewald summation; this typically involves a summation over the complete lattice from which the explicit sum for the inner region is subtracted.

The defect energy equals the difference in the total energies for the defective and the perfect lattice, corrected for the intrinsic energy of the defective species (interstitial and/or vacancy) at infinite separation. In a modern Mott-Littleton calculation the inner region may contain up to a few hundred atoms; ideally, a series of calculations is performed with increasing numbers of explicit ions until the defect energy converges. Incorporation of polarisability is usually important and is often handled using a shell model (Section 4.22.2). In addition to the energy of defect formation it is also possible to calculate the associated entropy change. This requires a consideration of the effect of the defect on the lattice phonon spectrum (Section 5.10). One technical point that is important to note is that these calculations are performed at constant volume and should be corrected prior to comparison with the constant pressure values typically obtained by experiment. This is particularly important for phenomena which occur at high temperatures, where the difference between the constant volume and constant pressure results can be significant.

Table 11.1 gives the results of Mott-Littleton calculations on some simple ionic systems [Mackrodt 1982]. By comparing the relative energies of the various types of defect, it is

	Theory (eV)	Experiment (eV)
LiF (Schottky)	2.37	2.34–2.68
NaCl (Schottky)	2.22	2.20–2.75
KBr (Schottky)	2.27	2.37–2.53
RbI (Schottky)	2.16	2.1
MgF ₂ (anion Frenkel)	3.12	—
CaF ₂ (anion Frenkel)	2.75	2.7
BaF ₂ (anion Frenkel)	1.98	1.91
CaCl ₂ (anion Frenkel)	4.7	—
MgO (Schottky)	7.5	5–7

Table 11.1: Defect energies for various materials. Data from [Mackrodt 1982].

possible to predict which types of defect might be expected in a particular material. For example, the energies to form Schottky defects in the alkali halides are 1–2 eV lower than to form Frenkel defects. By contrast, the dominant type of defect in the alkaline earth fluorides is the anion Frenkel defect.

An alternative to the Mott–Littleton method is to use a so-called supercell calculation, wherein the defect is located within a lattice that is subjected to periodic boundary conditions. The main difficulty with this approach is that, when a charged defect is present, the Ewald summation that is used to determine the Coulombic contribution diverges. This is dealt with by compensating for the net charge of the cell by a uniform background charge density. In addition, the energies of defect formation must be corrected for the interactions between defects in different cells. Supercells are often considered to be simpler for the calculation of defect entropies (and hence free energies) through the use of lattice statics and lattice dynamics. Full free energy minimisation can be performed on cells containing up to 1000 atoms under conditions of either constant volume or constant pressure. Calculations on such large systems are facilitated by the calculation of analytical derivatives of the vibrational frequencies with respect to all external and internal variables, an example being the study of defects in MgO [Taylor *et al.* 1997].

Defect calculations are traditionally performed using an empirical potential function, but there are some types of problem for which a quantum mechanical model is required, such as when the defect formation is accompanied by a transition to an excited electronic state. The obvious drawback to this is that the quantum mechanical method is computationally more expensive than the empirical potentials that are typically used in Mott–Littleton or supercell calculations. As a consequence, the number of atoms that can be treated quantum mechanically is often limited to the defect and its immediate neighbours. It is then necessary in some way to incorporate the effects of the surrounding region. In the embedded cluster approach this outer region provides a representation of the electrostatic potential due to the surrounding lattice, most easily simulated using point charges placed at the appropriate lattice sites. In more sophisticated approaches, the influence of the defect on the surrounding region can be taken into account in a manner similar to the Mott–Littleton approach [Grimes *et al.* 1989; Pisani 1999].

The most basic data that the Mott–Littleton and supercell methods provide are the energies and entropies of defect formation. Nevertheless, despite the fact that these techniques are essentially static approaches it can also be possible to deduce information on the ‘dynamic’ processes of diffusion and conductivity. These two processes are related by the Nernst–Einstein relationship:

$$\frac{\sigma}{D} = \frac{Nq^2}{fk_B T} \quad (11.98)$$

where σ is the electrical conductivity, D is the diffusion coefficient, N is the number of particles per unit volume and q is the charge on the mobile species. f is a correlation factor whose value depends upon the underlying migration mechanism. f may deviate from unity if the atomic movements affect the migration of charge and mass in different ways. For example, if charge transport is caused by a vacancy mechanism in which atoms jump into vacant sites then this is effectively a random process. However, after the atom has jumped into the vacancy it is possible for it to jump back to the original site. Mass transport is thus a correlated process. Three different defect migration mechanisms are shown in Figure 11.42. Of these, the vacancy mechanism dominates in most close-packed crystal structures. In the interstitial mechanism, the interstitial atom jumps from one site to another, an example being the diffusion of carbon in iron. The interstitialcy mechanism involves an interstitial atom displacing a lattice atom onto a new interstitial site. An example of this is the motion of silver ions in silver halides.

If the transport is due to discrete jumps of atoms then the diffusion coefficient D is related to the concentration of the jumping species (x), the jumping frequency (ν) and the distance over which the jump occurs (d):

$$D = \frac{1}{6} x \nu d^2 \quad (11.99)$$

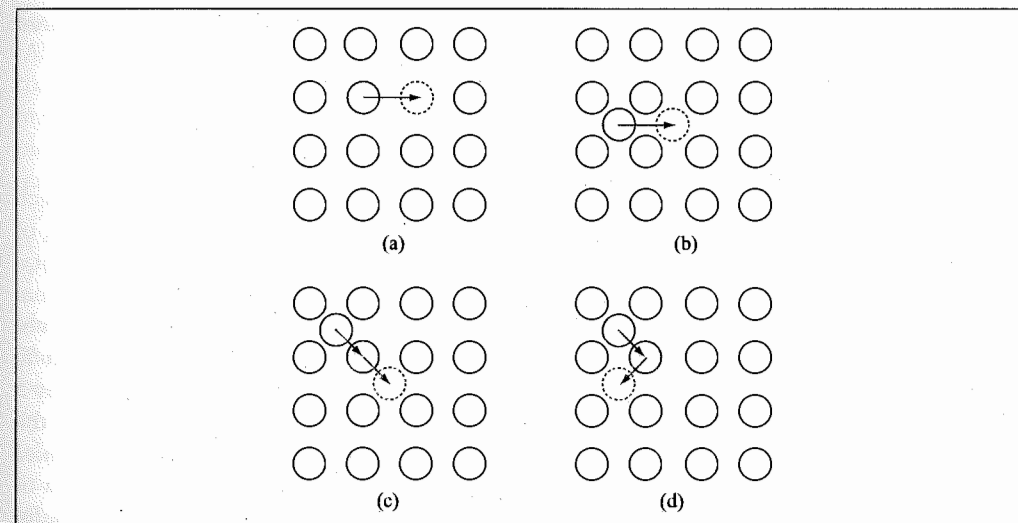


Fig. 11.42. Three different defect migration mechanisms: (a) vacancy, (b) interstitial, (c) collinear interstitialcy and (d) non-collinear interstitialcy. (Figure redrawn from Chadwick A V and J Corish 1997. *Defects and Matter Transport in Solid Materials*. In NATO ASI Series C 498 (New Trends in Materials Chemistry), pp. 285–318.)

The jump frequency is related in an exponential fashion to the free energy of activation between the ground state and the saddle point:

$$\nu = \nu_0 \exp(-\Delta G_{\text{act}}/k_B T) \quad (11.100)$$

ΔG_{act} in turn is composed of the enthalpy and entropy of activation, quantities which can in principle be calculated using other methods, such as those discussed above. The concentration of the jumping species is also predicted to vary in an exponential manner. It is thus expected that transport coefficients will follow an Arrhenius-like behaviour, and plotting the logarithm of the diffusion coefficient or the conductivity against $1/T$ will give a straight line (in fact, in the case of conductivity it is usual to plot $\log(\sigma T)$ against $1/T$ in accordance with the Nernst-Einstein relationship). It is indeed quite common to observe a series of linear regions, each corresponding to different types of defect population. Some typical activation energies are 0.66 eV (cation vacancy migration in NaCl), 0.35 eV (anion vacancy migration in CaF₂) and 2.0 eV (cation vacancy migration in MgO).

Molecular dynamics and Monte Carlo simulations can also be used to investigate systems with defects. In many respects these simulation methods are complementary to static techniques for the study of diffusion and conductivity. As we have discussed, calculation of transport coefficients using static methods is based upon the random jump model. This model is most appropriate when there are relatively high energy barriers involved. These high energy barriers make such systems less appropriate for simulation methods due to sampling difficulties. Simulation methods are most applicable to systems with facile diffusion (i.e. low activation energy barriers to transport), where the random jump model is less valid. Of course, one advantage of molecular dynamics is that diffusion coefficients can be calculated directly. The early molecular dynamics simulations concentrated on superionic materials such as SrCl₂, CaF₂ and Li₃N; the latter has a layered structure with much higher conductivity parallel to the layers, leading to very different mean squared displacements (Figure 6.10).

11.14.1 Defect Studies of the High- T_c Superconductor YBa₂Cu₃O_{7-x}

The discovery of materials which exhibit 'high'-temperature superconductivity (for which Bednorz and Müller were awarded the Nobel Prize for physics in 1987) led to a frenzy of activity to discover similar materials. This activity was not restricted to purely experimental considerations, as various theories were proposed to explain the reasons for this abnormal behaviour, with particular emphasis on variants on the so-called BCS theory, which involves the formation of pairs of electrons (Cooper pairs). One of the most studied of these high- T_c superconductors is the Y-Ba-Cu-O system, which was the first material found to display a transition temperature at liquid nitrogen temperatures (~ 90 K). The formula of the pertinent material is best written as YBa₂Cu₃O_{6+x}, with the superconducting properties being very sensitive to the value of x (rather than arising from aliovalent substitution of the Y³⁺ ions, as occurs in some other materials). High- T_c behaviour does not generally occur when x is less than approximately 0.3. The two related 'parent' molecules are YBa₂Cu₃O₆ and YBa₂Cu₃O₇, whose structures are shown in Figure 11.43. These materials contain two different types of copper atom and a variety of oxygen sites. In YBa₂Cu₃O₆,

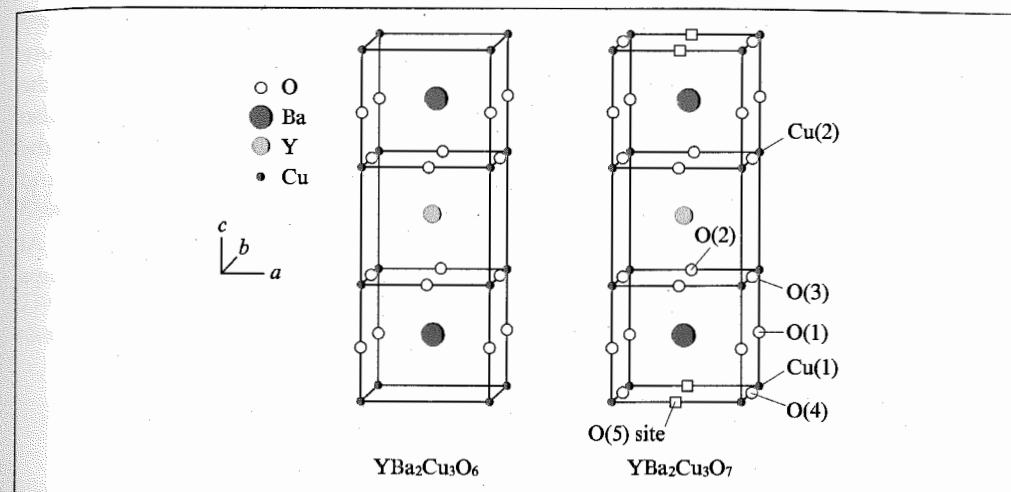


Fig. 11.43: The structures of YBa₂Cu₃O₆ and YBa₂Cu₃O₇.

the copper in the Cu(2) site is five-fold coordinate to oxygen, whereas in the Cu(1) site the copper is two-fold coordinated. The oxidation states of the copper in these two sites are Cu(II) and Cu(I), respectively. YBa₂Cu₃O₇ is derived from YBa₂Cu₃O₆ by introducing oxygen into the vacant O(4) sites; the copper in the Cu(1) sites becomes four-fold coordinate with the formation of the so-called CuO₂ basal plane.

These materials have been the subject of considerable computational investigations, involving both static and molecular dynamics methods. For example, static methods have been used to calculate the energies of formation of various vacancy and interstitial defects [Allan and Mackrodt 1994]. From these calculations, it was suggested that the lack of superconductivity for $x < 0.3$ was because the holes resulting from the excess oxygen are trapped in the basal plane adjacent to the O(4) position. In this case, we use the word 'hole' in the physicists' sense to denote a species (in this case Cu) from which an electron has been removed. Thus the addition of oxygen to YBa₂Cu₃O₆ causes Cu(I) to be oxidised to Cu(II), with this extra positive charge (i.e. the hole) being trapped by the negative charge of the now negatively charged oxygen ion. Superconductivity is associated with oxidation of the Cu(II) in the CuO₂ planes to Cu(III); this process does not occur until $x > 0.3$.

Many of the high- T_c superconductors also appear to be good oxygen ion conductors at high temperatures, and these effects have also been studied using computational methods. For example, oxygen diffusion in the YBa₂Cu₃O_{6.9} system has been studied using molecular dynamics [Zhang and Catlow 1992]. This particular composition has been shown experimentally to have the fastest oxygen diffusion. The system was set up by removing three O(4) atoms from 32 YBa₂Cu₃O₇ units, giving three oxygen vacancies, and then assigning three of the remaining O(4) species a charge of -2 to maintain charge neutrality, for a total of 413 atoms. The simulations were performed at quite a high temperature (higher than most relevant experimental studies) in order to obtain reasonable statistics. It was found that the oxygen diffuses in three directions, but that the diffusion coefficients were larger in the a

and b directions than in the c direction (see Figure 11.43). The overall values were in very nice agreement with those extrapolated from experiment data, although some low-temperature experiments had suggested that the b direction diffusion was significantly larger than in the a direction, in contrast to the simulations. The mechanism of oxygen transport was also investigated by analysing the trajectories using molecular graphics. This showed that the oxygen vacancies migrated between the O(4), O(1) and O(5) sites but not to the O(2) and O(3) sites. It was suggested that at lower temperatures the oxygen vacancies mostly followed a O(4)-O(1)-O(4) jump mechanism with a small proportion of O(4)-O(5)-O(4) jumps. At higher temperatures, the vacant O(5) sites were more easily occupied and O(4)-O(5) jumps occurred with a frequency comparable to that of the O(4)-O(1) mechanism.

Appendix 11.1 Calculating Free Energy Differences Using Thermodynamic Integration

If the free energy, A , is a continuous function of λ then we can write:

$$\Delta A = \int_0^1 \frac{\partial A(\lambda)}{\partial \lambda} d\lambda \quad (11.101)$$

Now

$$A(\lambda) = -k_B T \ln Q(\lambda) \quad (11.102)$$

Thus

$$\Delta A = -k_B T \int_0^1 \left[\frac{\partial \ln Q(\lambda)}{\partial \lambda} \right] d\lambda = \int_0^1 \frac{-k_B T}{Q(\lambda)} \frac{\partial Q(\lambda)}{\partial \lambda} d\lambda \quad (11.103)$$

From the definition of Q (Section 6.1.1):

$$Q_{NVT} = \frac{1}{N!} \frac{1}{h^{3N}} \iint d\mathbf{p}^N d\mathbf{r}^N \exp \left[-\frac{\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{k_B T} \right] \quad (11.104)$$

we can write the following for $\partial Q(\lambda)/\partial \lambda$:

$$\frac{\partial Q(\lambda)}{\partial \lambda} = \frac{1}{N!} \frac{1}{h^{3N}} \iint d\mathbf{p}^N d\mathbf{r}^N \frac{\partial}{\partial \lambda} \exp \left[-\frac{\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{k_B T} \right] \quad (11.105)$$

Applying the chain rule:

$$\frac{\partial Q(\lambda)}{\partial \lambda} = -\frac{1}{N!} \frac{1}{h^{3N}} \frac{1}{k_B T} \iint d\mathbf{p}^N d\mathbf{r}^N \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{\partial \lambda} \exp \left[-\frac{\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{k_B T} \right] \quad (11.106)$$

Substituting back into the expression for $\partial A/\partial \lambda$ gives:

$$\begin{aligned} \frac{\partial A(\lambda)}{\partial \lambda} &= \frac{1}{N!} \frac{1}{h^{3N}} \frac{1}{Q(\lambda)} \iint d\mathbf{p}^N d\mathbf{r}^N \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{\partial \lambda} \exp \left[-\frac{\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{k_B T} \right] \\ &= \iint d\mathbf{p}^N d\mathbf{r}^N \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)}{\partial \lambda} \left\{ \frac{\exp[-\mathcal{H}(\mathbf{p}^N, \mathbf{r}^N)/k_B T]}{Q(\lambda)} \right\} = \left\langle \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N, \lambda)}{\partial \lambda} \right\rangle_\lambda \end{aligned} \quad (11.107)$$

Thus

$$\Delta A = \int_{\lambda=0}^{\lambda=1} \left\langle \frac{\partial \mathcal{H}(\mathbf{p}^N, \mathbf{r}^N, \lambda)}{\partial \lambda} \right\rangle_\lambda d\lambda \quad (11.108)$$

Appendix 11.2 Using the Slow Growth Method for Calculating Free Energy Differences

The slow growth expression can be derived from the thermodynamic perturbation expression (Equation (11.7)) if it is written as a Taylor series:

$$\Delta A = -k_B T \sum_{i=0}^{N_{\text{step}}-1} \ln \langle \exp(-[\mathcal{H}(\lambda_{i+1}) - \mathcal{H}(\lambda_i)]/k_B T) \rangle_{NVT} \quad (11.109)$$

$$\Delta A \approx -k_B T \sum_{i=0}^{N_{\text{step}}-1} \ln \langle 1 - [\mathcal{H}(\lambda_{i+1}) - \mathcal{H}(\lambda_i)]/k_B T + \dots \rangle_{NVT} \quad (11.110)$$

$$\Delta A \approx -k_B T \sum_{i=0}^{N_{\text{step}}-1} \ln \left\{ 1 - \frac{1}{k_B T} \langle [\mathcal{H}(\lambda_{i+1}) - \mathcal{H}(\lambda_i)] \rangle_{NVT} + \dots \right\} \quad (11.111)$$

$$\Delta A \approx \sum_{i=0}^{N_{\text{step}}-1} \langle [\mathcal{H}(\lambda_{i+1}) - \mathcal{H}(\lambda_i)] \rangle_{NVT} \quad (11.112)$$

Appendix 11.3 Expansion of Zwanzig Expression for the Free Energy Difference for the Linear Response Method

The starting point is the standard expression for the free energy difference, Equation (11.6):

$$\Delta A = -k_B T \ln \langle \exp[-(\mathcal{H}_Y - \mathcal{H}_X)/k_B T] \rangle_0 \quad (11.113)$$

We expand the exponential:

$$\begin{aligned} \Delta A &= -k_B T \ln \left\langle 1 - \frac{(\mathcal{H}_Y - \mathcal{H}_X)}{k_B T} + \frac{(\mathcal{H}_Y - \mathcal{H}_X)^2}{2(k_B T)^2} - \dots \right\rangle_0 \\ &= -k_B T \ln \left[1 - \frac{\langle \mathcal{H}_Y - \mathcal{H}_X \rangle_0}{k_B T} + \frac{\langle (\mathcal{H}_Y - \mathcal{H}_X)^2 \rangle_0}{2(k_B T)^2} - \dots \right] \end{aligned} \quad (11.114)$$

Using the series expansion of $\ln(1+x)$ gives:

$$\begin{aligned} \Delta A &= -k_B T \left\{ -\frac{\langle \mathcal{H}_Y - \mathcal{H}_X \rangle_0}{k_B T} + \frac{\langle (\mathcal{H}_Y - \mathcal{H}_X)^2 \rangle_0}{2(k_B T)^2} \right. \\ &\quad \left. - \frac{1}{2} \left[\left(\frac{\langle \mathcal{H}_Y - \mathcal{H}_X \rangle_0}{k_B T} \right)^2 - \frac{\langle \mathcal{H}_Y - \mathcal{H}_X \rangle_0 \langle (\mathcal{H}_Y - \mathcal{H}_X)^2 \rangle_0}{2(k_B T)^3} + \left(\frac{\langle (\mathcal{H}_Y - \mathcal{H}_X)^2 \rangle_0}{2(k_B T)^2} \right)^2 \right] \right\} \end{aligned} \quad (11.115)$$

This can be rearranged to:

$$\Delta A = \langle \mathcal{H}_Y - \mathcal{H}_X \rangle_0 - \frac{1}{2k_B T} \langle [(\mathcal{H}_Y - \mathcal{H}_X) - \langle \mathcal{H}_Y - \mathcal{H}_X \rangle_0]^2 \rangle_0 + \dots \quad (11.116)$$

A similar procedure applied to the result from averaging at Y gives:

$$\Delta A = \langle \mathcal{H}_Y - \mathcal{H}_X \rangle_1 + \frac{1}{2k_B T} \langle [(\mathcal{H}_Y - \mathcal{H}_X) - \langle \mathcal{H}_Y - \mathcal{H}_X \rangle_1]^2 \rangle_1 + \dots \quad (11.117)$$

When Equations (11.116) and (11.117) are added together and we substitute $\Delta \mathcal{H}$ for $\mathcal{H}_Y - \mathcal{H}_X$ then we obtain:

$$\Delta A = \frac{1}{2} [\langle \Delta \mathcal{H} \rangle_0 + \langle \Delta \mathcal{H} \rangle_1] - \frac{1}{4k_B T} [\langle (\Delta \mathcal{H} - \langle \Delta \mathcal{H} \rangle_0)^2 \rangle_0 + \langle (\Delta \mathcal{H} - \langle \Delta \mathcal{H} \rangle_1)^2 \rangle_1] + \dots \quad (11.118)$$

Further Reading

- Allan N L and W C Mackrodt 1997. High- T_c Superconductors in Computer Modelling. In Catlow C R A (Editor) *Inorganic Crystallography*, pp. 241–268.
- Amara P and M J Field 1998. Combined Quantum Mechanical and Molecular Mechanical Potentials. In Schleyer, P v R, N L Allinger, T Clark, J Gasteiger, P A Kollman H F Schaefer III and P R Schreiner (Editors). *The Encyclopedia of Computational Chemistry*. Chichester, John Wiley & Sons.
- Beveridge D L and F M DiCapua 1989. Free Energy via Molecular Simulation: A Primer. In van Gunsteren W F and P K Weiner (Editors) *Computer Simulation of Biomolecular Systems*. Leiden, ESCOM, pp. 1–26.
- Catlow C R A 1994. An Introduction to Disorder in Solids. In NATO ASI Series C 418 (*Defects and Disorder in Crystalline and Amorphous Solids*), pp. 1–23.
- Catlow C R A 1994. Molecular Dynamics Studies of Defects in Solids. In NATO ASI Series C 418 (*Defects and Disorder in Crystalline and Amorphous Solids*), pp. 357–373.
- Catlow C R A, R G Bell and J D Gale 1994. Computer Modelling as a Technique in Materials Chemistry. *Journal of Materials Chemistry* 4:781–792.
- Catlow C R A and W C Mackrodt 1982. Theory of Simulation Methods for Lattice and Defect Energy Calculations in Crystals. In *Lecture Notes in Physics* 166 (Comput. Simul. Solids), pp. 3–20.
- Chadwick A V and J Corish 1997. Defects and Matter Transport in Solid Materials. In NATO ASI Series C 498 (*New Trends in Materials Chemistry*), pp. 285–318.
- Cramer C J and Truhlar D G 1995. Continuum Solvation Models: Classical and Quantum Mechanical Implementations. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 6. New York, VCH Publishers, pp. 1–72.
- Gale J 1999. *General Utility Lattice Program Manual*, Imperial College, London.
- Gao J 1995. Methods and Applications of Combined Quantum Mechanical and Molecular Mechanical Potentials. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 7. New York, VCH Publishers, pp. 119–185.
- Gillan M J 1989. *Ab Initio* Calculation of the Energy and Structure of Solids. *Journal of the Chemical Society Faraday Transactions* 2 85:521–536.
- Gillan M J 1997. The Virtual Matter Laboratory. *Contemporary Physics* 38:115–130.
- Harding J H 1997. Defects, Surfaces and Interfaces. In Catlow C R A (Editor) *Inorganic Crystallography*, pp. 185–199.

- Jorgensen W L 1983. Theoretical Studies of Medium Effects on Conformational Equilibria. *Journal of Physical Chemistry* 87:5304–5314.
- King P M 1993. Free Energy via Molecular Simulation: A Primer. In van Gunsteren W F, P K Weiner and A J Wilkinson (Editors) *Computer Simulation of Biomolecular Systems* Volume 2. Leiden, ESCOM, pp. 267–314.
- Kollman P A 1993. Free Energy Calculations: Applications to Chemical and Biochemical Phenomena. *Chemical Reviews* 93:2395–2417.
- Lybrand T P 1990. Computer Simulation of Biomolecular Systems Using Molecular Dynamics and Free Energy Perturbation Methods. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 1. New York, VCH Publishers, pp. 295–320.
- Mark A E and van Gunsteren W F 1995. Free Energy Calculations in Drug Design: A Practical Guide. In Dean P M, G Jolles and C G Newton (Editors) *New Perspectives in Drug Design*. London, Academic Press, pp. 185–200.
- Mezei M and D L Beveridge 1986. Free Energy Simulations. In Beveridge D L and W L Jorgensen (Editors) *Computer Simulation of Chemical and Biomolecular Systems*. *Annals of the New York Academy of Sciences* 482:1–23.
- Sandre E and A Pasturel 1997. An Introduction to *Ab-Initio* Molecular Dynamics Schemes. *Molecular Simulation* 20:63–77.
- Straatsma T P 1996. Free Energy by Molecular Simulation. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 9. New York, VCH Publishers, pp. 81–127.
- van Gunsteren W F 1989. Methods for Calculation of Free Energies and Binding Constants: Successes and Problems. In van Gunsteren and P K Weiner (Editors) *Computer Simulation of Biomolecular Systems*. Leiden, ESCOM, pp. 27–59.

References

- Allan N L and W C Mackrodt 1994. Oxygen Interstitial Defects in High- T_c Oxides. *Molecular Simulation* 12:89–100.
- Åqvist J, C Medina and J-E Samuelsson 1994. A New Method for Predicting Binding Affinity in Computer-aided Drug Design. *Protein Engineering* 7:385–391.
- Åqvist J and A Warshel 1993. Simulation of Enzyme Reactions Using Valence Bond Force Fields and Other Hybrid Quantum/Classical Approaches. *Chemical Reviews* 93:2523–2544.
- Åqvist J, M Fothergill and A Warshel 1993. Computer Simulation of the $\text{CO}_2/\text{HCO}_3^-$ Interconversion Step in Human Carbonic Anhydrase I. *Journal of the American Chemical Society* 115:631–635.
- Barrows S E, J W Storer, C J Cramer, A D French and D G Truhlar 1998. Factors Controlling Relative Stability of Anomers and Hydroxymethyl Conformers of Glucopyranose. *Journal of Computational Chemistry* 19:1111–1129.
- Bartlett P A and C K Marlowe 1987. Evaluation of Intrinsic Binding Energy from a Hydrogen-bonding Group in an Enzyme Inhibitor. *Science* 235:569–571.
- Bash P A, U C Singh, F K Brown, R Langridge and P A Kollman 1987. Calculation of the Relative Change in Binding Free-Energy of a Protein-Inhibitor Complex. *Science* 235:574–576.
- Bernardi A, A M Capelli, A Comotti, C Gannari, J M Goodman and I Paterson 1990. Transition-State Modeling of the Aldol Reaction of Boron Enolates: A Force Field Approach. *Journal of Organic Chemistry* 55:3576–3581.
- Boero M, M Parrinello and K Terakura 1999. Ziegler–Natta Heterogeneous Catalysis by First Principles Computer Experiments. *Surface Science* 438:1–8.
- Boresch S, G Archontis and M Karplus 1994. Free Energy Simulations: The Meaning of the Individual Contributions from a Component Analysis. *Proteins: Structure, Function and Genetics* 20:25–33.

- Boresch S and M Karplus 1995. The Meaning of Component Analysis: Decomposition of the Free Energy in Terms of Specific Interactions. *Journal of Molecular Biology* **254**:801-807.
- Born M 1920. Volumen and Hydratationswärme der Ionen. *Zeitschrift für Physik* **1**:45-48.
- Buetler T C, A E Mark, R C van Schaik, P R Gerber and W F van Gunsteren 1994. Avoiding Singularities and Numerical Instabilities in Free Energy Calculations Based on Molecular Simulations. *Chemical Physics Letters* **222**:529-539.
- Bürger M T, A Armstrong, F Guarnieri, D Q McDonald and W C Still 1994. Free Energy Calculations in Molecular Design: Predictions by Theory and Reality by Experiment with Enantioselective Podand Ionophores. *Journal of the American Chemical Society* **116**:3593-3594.
- Car R and M Parrinello 1985. Unified Approach for Molecular Dynamics and Density Functional Theory. *Physical Review Letters* **55**:2471-2474.
- Carlson H A and W L Jorgensen 1995. An Extended Linear Response Method for Determining Free Energies of Hydration. *Journal of Physical Chemistry* **99**:10667-10673.
- Chambers C C, G D Hawkins, C J Cramer and D G Truhlar 1996. Model for Aqueous Solvation Based on Class IC Atomic Charges and First Solvation Shell Effects. *Journal of Physical Chemistry* **100**:16385-16398.
- Chandrasekhar J and W L Jorgensen 1985. Energy Profile for a Nonconcerted S_N2 Reaction in Solution. *Journal of the American Chemical Society* **107**:2974-2975.
- Chandrasekhar J, S F Smith and W L Jorgensen 1985. Theoretical Examination of the S_N2 Reaction Involving Chloride Ion and Methyl Chloride in the Gas Phase and Aqueous Solution. *Journal of the American Chemical Society* **107**:154-163.
- Claverie P, J P Daudey, J Langlet, B Pullman, D Piazzola and M J Huron 1978. Studies of Solvent Effects. I. Discrete, Continuum and Discrete-Continuum Models and Their Comparison for Some Simple Cases: NH_4^+ , CH_3OH and substituted NH_4^+ . *Journal of Physical Chemistry* **82**:405-418.
- Constanciel R and R Contreras 1984. Self-Consistent Field Theory of Solvent Effects Representation by Continuum Models - Introduction of Desolvation Contribution. *Theoretica Chimica Acta* **65**:1-11.
- Cramer C J and D G Truhlar 1992. AM1-SM2 and PM3-SM3 Parametrized SCF Solvation Models for Free Energies in Aqueous Solution. *Journal of Computer-Aided Molecular Design* **6**:629-666.
- Dapprich S, I Komirovi, K S Byun, K Morokuma and M J Frisch 1999. A New ONIOM Implementation in Gaussian '98. Part I. The Calculation of Energies, Gradients, Vibrational Frequencies and Electric Field Derivatives. *THEOCHEM* **461-462**:1-21.
- de Wijs G A, G Kresse, L Vočadlo, D Dobson, D Alfè, M J Gillan and G D Price 1998. The Viscosity of Liquid Iron at the Physical Conditions of the Earth's Core. *Nature* **392**:805-807.
- Eisenberg D and A D McLachlan 1986. Solvation Energy in Protein Folding and Binding. *Nature* **319**:199-203.
- Elber R and M Karplus 1990. Enhanced Sampling in Molecular Dynamics: Use of the Time-Dependent Hartree Approximation for a Simulation of Carbon Monoxide Diffusion through Myoglobin. *Journal of the American Chemical Society* **112**:9161-9175.
- Eriksson M A L, J Pitera and P A Kollman 1999. Prediction of the Binding Free Energies of New TIBO-like HIV-1 Reverse Transcriptase Inhibitors Using a Combination of PROFEC, PB/SA, CMC/MD, and Free Energy Calculations. *Journal of Medicinal Chemistry* **42**:868-881.
- Essex J W, C A Reynolds and W G Richards 1989. Relative Partition Coefficients from Partition Functions: A Theoretical Approach to Drug Transport. *Journal of the Chemical Society Chemical Communications* 1152-1154.
- Field M J, P A Bash and M Karplus 1990. A Combined Quantum Mechanical and Molecular Mechanical Potential for Molecular Dynamics Simulations. *Journal of Computational Chemistry* **11**:700-733.
- Fleischman S H and C L Brooks III 1987. Thermodynamics of Aqueous Solvation - Solution Properties of Alcohols and Alkanes. *Journal of Chemical Physics* **87**:3029-3037.

- Floris F and J Tomasi 1989. Evaluation of the Dispersion Contribution to the Solvation Energy - A Simple Computational Model in the Continuum Approximation. *Journal of Computational Chemistry* **10**:616-627.
- Freitag S, I Le Trong, P S Stayton and R E Stenkamp 1997. Structural Studies of the Streptavidin Binding Loop. *Protein Science* **6**:1157.
- Gilson M K and B Honig 1988. Calculation of the Total Electrostatic Energy of a Macromolecular System: Solvation Energies, Binding Energies and Conformational Analysis. *Proteins: Structure, Function and Genetics* **4**:7-18.
- Gonzalez C and H B Schlegel 1988. An Improved Algorithm for Reaction Path Following. *Journal of Chemical Physics* **90**:2154-2161.
- Grimes R W, C R A Catlow and A M Stoneham 1989. Quantum-mechanical Cluster Calculations and the Mott-Littleton Methodology. *Journal of the Chemical Society, Faraday Transactions* **85**:485-495.
- Guo Z and C L Brooks III 1998. Rapid Screening of Binding Affinities: Application of the λ -Dynamics Method to a Trypsin-Inhibitor System. *Journal of the American Chemical Society* **120**:1920-1921.
- Guo Z, C L Brooks III and X Kong 1998. Efficient and Flexible Algorithm for Free Energy Calculations using the λ -Dynamics Approach. *Journal of Physical Chemistry* **B102**:2032-2036.
- Ha S, J Gao, B Tidor, J W Brady and M Karplus 1991. Solvent Effect on the Anomeric Equilibrium in D-Glucose: A Free Energy Simulation Analysis. *Journal of the American Chemical Society* **113**:1553-1557.
- Hansson T and J Åqvist 1995. Estimation of Binding Free Energies for HIV Proteinase Inhibitors by Molecular Dynamics Simulations. *Protein Engineering* **8**:1137-1144.
- Hansson T, J Marelis and J Åqvist 1998. Ligand Binding Affinity Prediction by Linear Interaction Energy Methods. *Journal of Computer-Aided Molecular Design* **12**:27-35.
- Hasel W, T F Hendrickson and W C Still 1988. A Rapid Approximation to the Solvent Accessible Surface Areas of Atoms. *Tetrahedron Computer Methodology* **1**:103-116.
- Honig B and A Nicholls 1995. Classical Electrostatics in Biology and Chemistry. *Science* **268**:1144-1149.
- Jones-Hertzog D K and W L Jorgensen 1997. Binding Affinities for Sulphonamide Inhibitors with Human Thrombin Using Monte Carlo Simulations with a Linear Response Method. *Journal of Medicinal Chemistry* **40**:1539-1549.
- Jorgensen W L, J M Briggs and M L Contreras 1990. Relative Partition Coefficients for Organic Solutes from Fluid Simulations. *Journal of Physical Chemistry* **94**:1683-1986.
- Jorgensen W L and J K Buckner 1987. Use of Statistical Perturbation Theory for Computing Solvent Effects on Molecular Conformation. Butane in Water. *Journal of Physical Chemistry* **91**:6083-6085.
- Jorgensen W L, J K Buckner, S Boudon and J Tirado-Reeves 1988. Efficient Computation of Absolute Free Energies of Binding by Computer Simulations - Applications to the Methane Dimer in Water. *Journal of Chemical Physics* **89**:3742-3746.
- Jorgensen W L, J Gao and C Ravimohan 1985. Monte Carlo Simulations of Alkanes in Water: Hydration Numbers and the Hydrophobic Effect. *Journal of Physical Chemistry* **89**:3470-3473.
- Kirkwood J G 1934. Theory of Solutions of Molecules Containing Widely Separated Charges with Special Application to Zwitterions. *Journal of Chemical Physics* **2**:351-361.
- Klamt A 1995. Conductor-like Screening Model for Real Solvent: A New Approach to the Quantitative Calculation of Solvation Phenomena. *Journal of Physical Chemistry* **99**:2224-2235.
- Klamt A, V Jonas, T Bürger and J C W Lohrenz 1998. Refinements and Parametrisation of COSMO-RS. *Journal of Physical Chemistry* **102**:5074-5085.
- Klamt A and G Schüürmann 1993. COSMO: A New Approach to Dielectric Screening in Solvents with Explicit Expressions for the Screening Energy and its Gradient. *Journal of the Chemical Society, Perkin Transactions* **2**:799-805.
- Klapper J, R Hagstrom, R Fine, K Sharp and B Honig 1986. Focusing of Electric Fields in the Active Site of CuZn Superoxide Dismutase: Effects of Ionic Strength and Amino-Acid Substitution. *Proteins: Structure, Function and Genetics* **1**:47-59.

- Kong X and C L Brooks III 1996. λ -Dynamics: A New Approach to Free Energy Calculations. *Journal of Chemical Physics* **105**:2414–2423.
- Laasonen, M Sprik and M Parrinello 1993. 'Ab Initio' Liquid Water. *Journal of Chemical Physics* **99**:9080–9089.
- Liu H, A E Mark and W F van Gunsteren 1996. Estimating the Relative Free Energy of Different Molecular States with Respect to a Single Reference State. *Journal of Physical Chemistry* **100**:9485–9494.
- Lybrand T P, J A McCammon and G Wipff 1986. Theoretical Calculation of Relative Binding Affinity in Host–Guest Systems. *Proceedings of the National Academy of Sciences USA* **83**:833–835.
- Mackrodt W C 1982. Defect Calculations for Ionic Materials. *Lecture Notes in Physics* **166** (Computer Simulation of Solids):175–194.
- Marquart M, J Walter, J Deisenhofer, W Bode and R Huber 1983. The Geometry of the Reactive Site and of the Peptide Groups in Trypsin, Trypsinogen and its Complexes with Inhibitors. *Acta Crystallographica B* **39**:480–490.
- Maseras F and K Morokuma 1995. IMOMM: A New Integrated *Ab Initio* + Molecular Mechanics Geometry Optimisation Scheme of Equilibrium Structures and Transition States. *Journal of Computational Chemistry* **16**:1170–1179.
- McRee D E, S M Redford, E D Getzoff, J R Lepock, R A Hallewell and J A Tainer 1990. Changes in Crystallographic Structure and Thermostability of a Cu, Zn Superoxide Dismutase Mutant Resulting from the Removal of Buried Cysteine. *Journal of Biological Chemistry* **265**:14234–14241.
- Merz K M Jr and P A Kollman 1989. Free Energy Perturbation Simulations of the Inhibition of Thermolysin: Prediction of the Free Energy of Binding of a New Inhibitor. *Journal of the American Chemical Society* **111**:5649–5658.
- Miertus S, E Scrocco and J Tomasi 1981. Electrostatic Interaction of a Solute with a Continuum – A Direct Utilization of *Ab Initio* Molecular Potentials for the Provision of Solvent Effects. *Chemical Physics* **55**:117–129.
- Miick S M, G V Martinez, W R Fiori, A P Todd and G L Millhauser 1992. Short Alanine-based Peptides May Form 3(10)-Helices and not Alpha-helices in Aqueous Solution. *Nature* **359**:653–655.
- Mitchell M J and J A McCammon 1991. Free Energy Difference Calculations by Thermodynamic Integration: Difficulties in Obtaining a Precise Value. *Journal of Computational Chemistry* **12**:271–275.
- Miyamoto S and P A Kollman 1993a. Absolute and Relative Binding Free Energy Calculations of the Interaction of Biotin and its Analogues with Streptavidin Using Molecular Dynamics/Free Energy Perturbation Approaches. *Proteins: Structure, Function and Genetics* **16**:226–245.
- Miyamoto S and P A Kollman 1993b. What Determines the Strength of Noncovalent Association of Ligands to Proteins in Aqueous Solution? *Proceedings of the National Academy of Sciences USA* **90**:8402–8406.
- Mott N F and M J Littleton 1938. Conduction in Polar Crystals. I. Electrolytic Conduction in Solid Salts. *Transactions of the Faraday Society* **34**: 485–499.
- Onsager L 1936. Electric Moments of Molecules in Liquids. *Journal of the American Chemical Society* **58**:1486–1493.
- Paschual-Ahuir J L, E Silla, J Tomasi and R Bonaccorsi 1987. Electrostatic Interaction of a Solute with a Continuum. Improved Description of the Cavity and of the Surface Cavity Bound Charge Distribution. *Journal of Computational Chemistry* **8**:778–787.
- Payne M C, M P Teter, D C Allan, R A Arias and D J Joannopoulos 1992. Iterative Minimisation Techniques for *Ab Initio* Total-Energy Calculations: Molecular Dynamics and Conjugate Gradients. *Reviews of Modern Physics* **64**:1045–1097.
- Pearlman D A and P A Kollman 1989. A New Method for Carrying Out Free-Energy Perturbation Calculations – Dynamically Modified Windows. *Journal of Chemical Physics* **90**:2460–2470.
- Pierotti R 1965. Aqueous Solutions of Nonpolar Gases. *Journal of Physical Chemistry* **69**:281–288.

- Pisani C 1999. Software for the Quantum-mechanical Simulation of the Properties of Crystalline Materials: State of the Art and Prospects. *THEOCHEM* **463**:125–137.
- Pitera J and P Kollman 1998. Designing an Optimum Guest for a Host Using Multimolecule Free Energy Calculations: Predicting the Best Ligand for Rebek's 'Tennis Ball'. *Journal of the American Chemical Society* **120**:7557–7567.
- Postma J P M, H J C Berendsen and J R Haak 1982. Thermodynamics of Cavity Formation in Water. *Faraday Symposium of the Chemical Society* **17**:55–67.
- Qiu D, P S Shenkin, F P Hollinger and W C Still 1997. The GB/SA Continuum Model for Solvation. A Fast Analytical Method for the Calculation of Approximate Born Radii. *Journal of Physical Chemistry* **101**:3005–3014.
- Rashin A A 1990. Hydration Phenomena, Classical Electrostatics, and the Boundary Element Method. *Journal of Physical Chemistry* **94**:1725–1733.
- Rashin A A and B Honig 1985. Reevaluation of the Born Model of Ion Hydration. *Journal of Physical Chemistry* **89**:5588–5593.
- Rashin A A and K Nambodiri 1987. A Simple Method for the Calculation of Hydration Enthalpies of Polar Molecules with Arbitrary Shapes. *Journal of Physical Chemistry* **91**:6003–6012.
- Remler D K and P A Madden 1990. Molecular Dynamics without Effective Potentials via the Car-Parrinello Approach. *Molecular Physics* **70**:921–966.
- Reuter N, A Dejaegere, B Maigret and M Karplus 2000. Frontier Bonds in QM/MM Methods: A Comparison of Different Approaches. *Journal of Physical Chemistry A* **104**:1720–1733.
- Rinaldi D, M F Ruiz-Lopez and J L Rivail 1983. *Ab Initio* SCF Calculations on Electrostatically Solvated Molecules Using a Deformable Three Axes Ellipsoidal Cavity. *Journal of Chemical Physics* **78**:834–838.
- Röthlisberger and M Parrinello 1997. *Ab Initio* Molecular Dynamics Simulation of Liquid Hydrogen Fluoride. *Journal of Chemical Physics* **106**:4658–4664.
- Ryckaert J-P and A Bellemans 1978. *Molecular Dynamics of Liquid Alkanes*, *Faraday Discussions* **20**:95–106.
- Saitta A M, P D Sooper, E Wasserman and M L Klein 1999. Influence of a Knot on the Strength of a Polymer Strand. *Nature* **399**:46–48.
- Schäfer H, W F van Gunsteren and A E Mark 1999. Estimating Relative Free Energies from a Single Ensemble: Hydration Free Energies. *Journal of Computational Chemistry* **20**:1604–1617.
- Silvestrelli P L and M Parrinello 1999. Structural, Electronic and Bonding Properties of Liquid Water from First Principles. *Journal of Chemical Physics* **111**:3572–3580.
- Simmerling C and R Elber 1995. Computer Determination of Peptide Conformations in Water: Different Roads to Structure. *Proceedings of the National Academy of Sciences USA* **92**:3190–3193.
- Simmerling C, T Fox and P A Kollman 1998. Use of Locally Enhanced Sampling in Free Energy Calculations: Testing and Application to the $\alpha \rightarrow \beta$ Anomerisation of Glucose. *Journal of the American Chemical Society* **120**:5771–5782.
- Singh U C and P A Kollman 1986. A Combined *Ab Initio* Quantum Mechanical and Molecular Mechanical Method for Carrying out Simulations on Complex Molecular Systems: Applications to the $\text{CH}_3\text{Cl} + \text{Cl}^-$ Exchange Reaction and Gas Phase Protonation of Polyethers. *Journal of Computational Chemistry* **7**:718–730.
- Sitkoff D, K A Sharp and B Honig 1994. Accurate Calculation of Hydration Free Energies Using Macroscopic Solvent Models. *Journal of Physical Chemistry* **98**:1978–1988.
- Smith P E and B M Pettitt 1994. Modeling Solvent in Biomolecular Systems. *Journal of Physical Chemistry* **98**:9700–9711.
- Smith P E and W F van Gunsteren 1994a. Predictions of Free Energy Differences from a Single Simulation of the Initial State. *Journal of Chemical Physics* **100**:577–585.
- Smith P E and W F van Gunsteren 1994b. When Are Free Energy Components Meaningful? *Journal of Physical Chemistry* **98**:13735–13740.

- Smythe M L, S E Huston and G R Marshall 1993. Free Energy Profile of a 3_{10} to α -Helical Transition of an Oligopeptide in Various Solvents. *Journal of the American Chemical Society* **115**:11594–11595.
- Smythe M L, S E Huston and G R Marshall 1995. The Molten Helix: Effects of Solvation on the α - to 3_{10} -Helical Transition. *Journal of the American Chemical Society* **117**:5445–5452.
- Sprink M, J Hutter and M Parrinello 1996. *Ab Initio* Molecular Dynamics Simulation of Liquid Water: Comparison of Three Gradient-corrected Density Functionals. *Journal of Chemical Physics* **105**:1142–1152.
- Stich I, A De Vita, M C Payne, M J Gilland and L J Clarke 1994. Surface Dissociation from First Principles: Dynamics and Chemistry. *Physical Review* **B49**:8076–8085.
- Still W C, A Tempczyk, R C Hawley and T Hendrickson 1990. Semianalytical Treatment of Solvation for Molecular Mechanics and Dynamics. *Journal of the American Chemical Society* **112**:6127–6129.
- Svensson M, S Humbel, R D J Froese, T Matsubara, S Sieber and K Morokuma 1996. ONIOM: A Multilayered Integrated MO+MM Method for Geometry Optimisations and Single Point Energy Predictions. A Test for Diels–Alder Reactions and Pt(P(t-Bu)₃)₂ + H₂ Oxidative Addition. *Journal of Physical Chemistry* **100**:19357–19363.
- Tapia O and O Goscinski 1975 Self-Consistent Reaction Field Theory of Solvent Effects. *Molecular Physics* **29**:1653–1661.
- Taylor M B, G D Barrera, N L Allan, T H K Barron and W C Mackrodt 1997. Free Energy of Formation of Defects in Polar Solids. *Faraday Discussions* **106**:377–387.
- Tirado-Reeves J, D S Maxwell and W L Jorgensen 1993. Molecular Dynamics and Monte Carlo Simulations Favor the α -Helical Form for Alanine-Based Peptides in Water. *Journal of the American Chemical Society* **115**:11590–11593.
- Tobias D J and C L Brooks III 1988. Molecular Dynamics with Internal Coordinate Constraints. *Journal of Chemical Physics* **89**:5115–5126.
- Torrie G M and J P Valleau 1977. Nonphysical Sampling Distributions in Monte Carlo Free-Energy Estimation: Umbrella Sampling. *Journal of Computational Physics* **23**:187–199.
- Tuckerman M, K Laasonen, M Sprink and M Parrinello 1995a. *Ab Initio* Molecular Dynamics Simulation of the Solvation and Transport of Hydronium and Hydroxyl Ions in Water. *Journal of Chemical Physics* **103**:150–161.
- Tuckerman M, K Laasonen, M Sprink and M Parrinello 1995b. *Ab Initio* Molecular Dynamics Simulation of the Solvation and Transport of H₃O⁺ and OH⁻ Ions in Water. *Journal of Physical Chemistry* **99**:5749–5752.
- Wall I D, A R Leach, D W Salt, M G Ford and J W Essex 1999. Binding Constants of Neuraminidase Inhibitors: An Investigation of the Linear Interaction Energy Method. *Journal of Medicinal Chemistry* **42**:5142–5152.
- Wang W, J Wang and P A Kollman 1999. What Determines the van der Waals Coefficient β in the LIE (Linear Interaction Energy) Method to Estimate Binding Free Energies Using Molecular Dynamics Simulations? *Proteins: Structure, Function and Genetics* **34**:395–402.
- Warshel A 1991. *Computer Modelling of Chemical Reactions in Enzymes and Solutions*. New York, John Wiley & Sons.
- Warshel A and M Levitt 1976. Theoretical Studies of Enzymic Reactions: Dielectric, Electrostatic and Steric Stabilization of the Carbonium Ion in the Reaction of Lysozyme. *Journal of Molecular Biology* **103**:227–249.
- Warwicker J and H C Watson 1982. Calculation of the Electric Potential in the Active-Site Cleft Due to Alpha-Helix Dipoles. *Journal of Molecular Biology* **157**:671–679.
- Wodak S J and J Janin 1980. Analytical Approximation to the Solvent Accessible Surface Area of Proteins. *Proceedings of the National Academy of Sciences USA* **77**:1736–1740.
- Wong M W, K B Wiberg and M J Frisch 1992. Solvent Effects. 3. Tautomeric Equilibria of Formamide and 2-Pyridone in the Gas Phase and Solution. An *Ab Initio* SCRF Study. *Journal of the American Chemical Society* **114**:1645–1652.

- Yu H-A and M Karplus 1988. A Thermodynamic Analysis of Solvation. *Journal of Chemical Physics* **89**:2366–2379.
- Zhang L and J Hermans 1994. 3_{10} -Helix versus α -Helix: A Molecular Dynamics Study of Conformational Preferences of Aib and Alanine. *Journal of the American Chemical Society* **116**:11915–11921.
- Zhang X and C R A Catlow 1992. Molecular Dynamics Study of Oxygen Diffusion in YBa₂Cu₃O_{6.19}. *Physical Review* **B46**:457–462.
- Zwanzig R W 1954. High-temperature Equation of State by a Perturbation Method. I. Nonpolar Gases. *Journal of Chemical Physics* **22**:1420–1426.

CHAPTER TWELVE

The Use of Molecular Modelling and Chemoinformatics to Discover and Design New Molecules

Molecular modelling techniques are widely used in the chemical, pharmaceutical and agrochemical industries. Much of this modelling activity employs the tools that we have discussed in earlier chapters, such as energy minimisation, molecular dynamics and Monte Carlo simulations and conformational analysis. In this chapter, we will discuss a number of methods that do not fit naturally into any of these categories or which bring together several of these tools to create a new approach. Other techniques have been developed as a consequence of, or in conjunction with, some new technological advance such as combinatorial chemistry or high-throughput screening. Our discussion will often but not exclusively use examples drawn from the pharmaceutical industry, though many of the techniques are also applicable to molecular design in other areas.

12.1 Molecular Modelling in Drug Discovery

Most drugs produce their effect by interacting with a biological macromolecule such as an enzyme, DNA, glycoprotein or receptor. The interaction between a ligand and its target* may be due entirely to non-bonded forces, but in some cases a covalent interaction may be involved. Drugs which interact with receptor proteins can be classified as *agonists*, *antagonists* or *inverse agonists*. Agonists produce the same or elevated effect as the natural substrate or effector molecule, whereas antagonists inhibit the effect of the natural ligand. Inverse agonists create an effect which appears opposite to that of the agonist. Tight-binding ligands often have a high degree of complementarity with the target. This complementarity can be assessed and measured in various ways. Many ligands show significant shape complementarity with the region of the macromolecule where they bind (the binding

* We shall use the generic term 'ligand' to indicate the inhibitor or substrate and the term 'receptor' to indicate the macromolecule to which it binds, be it an enzyme, a gene or a receptor protein.

site). This can be observed by constructing the molecular surfaces as illustrated in Figure 11.12 (colour plate section), which shows the molecular surfaces of biotin bound to streptavidin. The ligand often forms hydrogen bonds with the receptor. Some receptors have hydrophobic 'pockets', formed by groups of non-polar amino acids, into which the ligand can place a hydrophobic group of an appropriate size. It is also crucial to remember that a good drug does more than simply bind tightly to its target. After administration, a drug must get to the site of action. This transport process often requires the drug to pass through cell membranes. A cell membrane is a hydrophobic environment and so the drug must be sufficiently lipophilic (lipid loving) to partition into the membrane but not so lipophilic that it stays there. Once inside the cell, the drug must access its target. During this process, the molecule may also be removed from the body by metabolism, excretion and other pathways.

Discovering and developing any new medicine is a long and expensive process. A new compound must not only produce the desired response with minimal side-effects but must also be demonstrably better than existing therapies. Two key steps in many drug discovery programmes are the identification of hit molecules ('hits') and lead series ('leads'). A hit is a molecule that has some reproducible activity in a biological assay. A lead series comprises a set of related molecules that usually share some common structural feature, and which show some variation in the activity as the structure is modified. This gives one confidence that further synthetic modification to the lead series (termed *lead optimisation*) has a good chance of resulting in a drug candidate with the desired potency and selectivity, lack of toxicity and the appropriate characteristics to enable it to reach its target *in vivo*. Such a drug candidate will then enter the early stages of development, where further large-scale investigations are undertaken.

Finding novel lead series can be a difficult problem. Serendipity often played an important role in the past, a classic example being the discovery of penicillin by Alexander Fleming. For many years pharmaceutical companies have screened soil and other biological samples to find new leads, but it can be difficult to extract and purify the bioactive ingredient. The 1990s saw the widespread adoption of high-throughput screening (HTS), which enables large numbers of compounds to be screened using highly automated, robotic techniques. The molecules used for HTS come from compounds synthesised in earlier medicinal chemistry programmes, from compounds that can be purchased from chemical suppliers and from combinatorial chemistry (discussed in Section 12.14). However, although HTS makes it possible in principle to test every available compound against every biological assay, there are a number of practical reasons why this is not necessarily feasible, let alone desirable. The first reason is financial: although robotics and miniaturisation have significantly reduced the unit cost, the sheer number of samples now available in many companies means that the overall expense can be significant. A second reason is that some assays cannot be converted to a high-throughput format and so have to be conducted using more traditional techniques. Third, a significant proportion of the available samples might not be considered appropriate structures, suitable for taking forward to the next stage. For example, some molecules may contain functional groups which are known to react in a non-specific manner with biological targets. Other molecules might interfere with the proper interpretation of the assay, such as a strong fluorophore. Yet more molecules may just be considered 'inappropriate', or not sufficiently 'drug-like'.

For these and other reasons, it is often necessary to identify subsets of compounds. Computational techniques have a significant role to play in the ways in which such subsets can be constructed, with various techniques being available depending upon the type of molecule that one wishes to screen, what kind of information is available to assist the selection and what properties one wishes to take into account. Sections 12.2–12.11 describe a wide variety of methods that can be used either individually or in combination to select compounds. Some of these methods only use information about the underlying chemical structure of the molecule. These are often referred to as '2D' properties, as distinct from '3D' methods, which take into account the three-dimensional nature of a molecule (i.e. its conformation and properties dependent upon the conformation). Some of the methods can take into account information about the target protein or about other molecules that are known to be active at the target, whereas other methods are designed to produce 'diverse' collections of compounds for more general screening.

Having tested a number of compounds, it is then usually desired to construct a model which relates the observed activity to the molecular structure. The model can then be used in the next iteration of the process. Many different kinds of model are possible. A popular approach is to use statistical techniques to derive the model. Such statistical techniques are discussed in Sections 12.12–12.13.

12.2 Computer Representations of Molecules, Chemical Databases and 2D Substructure Searching

Substructure searching is probably the most basic approach to identifying compounds of interest. It is widely used for all kinds of problems. Most chemists take substructure searching for granted, testament to the decades of effort that has gone into the development of extremely powerful algorithms and database systems.

Many organisations maintain databases of chemical compounds. Some of these databases are publicly accessible; others are proprietary. A database may contain an extremely large number of compounds; several hundred thousand is common, and the database maintained by the American Chemical Society contains more than 18 million compounds. A more recent development involves the creation of databases containing *virtual* molecules. These are compounds that do not yet exist but which could be synthesised readily, typically using combinatorial chemistry techniques. We do not have space to consider in any detail the nature of chemical database systems, save for a few key points. The first issue concerns the representation of molecular structures in a computer. We are all familiar with the chemical diagrams in journals and lab-books, but simply storing the chemical diagram itself (as an image) is of little value. Rather, most systems represent molecules as *molecular graphs*. A graph contains *nodes*, which are connected by *edges*. Two examples are shown in Figure 12.1. In a molecular graph, the nodes correspond to the atoms and the edges to the bonds, as shown for acetic acid in Figure 12.2. The locations of the nodes and edges of a graph on the page are irrelevant; only the way in which the nodes are connected together matters. The conformational search trees that we met in Section 9.2 are a special kind of graph. A *subgraph* is a subset of the nodes and edges of a graph; thus the graph for CH₃ is

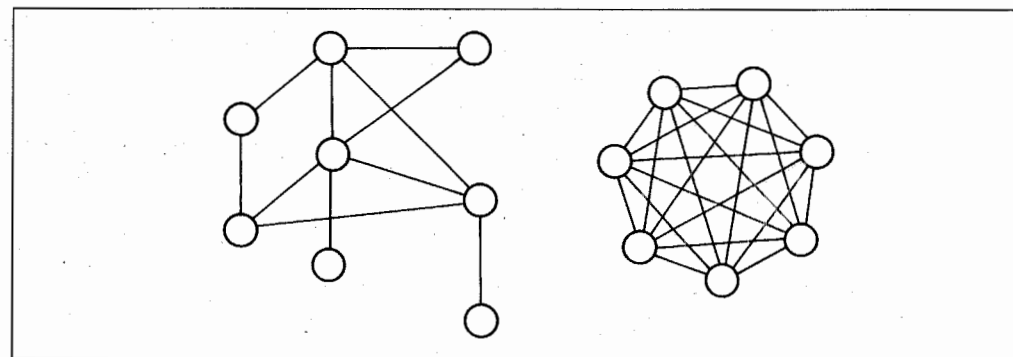


Fig. 12.1: Graphs contain nodes connected by edges. A completely connected graph (right) has an edge between all pairs of nodes.

a subgraph of the graph of acetic acid. A graph is said to be *completely connected* if there is an edge between all pairs of nodes. Only in rare cases is the molecular graph a completely connected graph, one example being the P₄ form of elemental phosphorus.

There are a number of different ways that the molecular graph can be communicated between the computer and the end-user. One common representation is the *connection table*, of which there are various flavours, but most provide information about the atoms present in the molecule and their connectivity. The most basic connection tables simply indicate the atomic number of each atom and which atoms form each bond; others may include information about the atom hybridisation state and the bond order. Hydrogens may be included or they may be implied. In addition, information about the atomic coordinates (for the standard two-dimensional chemical drawing or for the three-dimensional conformation) can be included. The connection table for acetic acid in one of the most popular formats, the Molecular Design mol format [Dalby *et al.* 1992], is shown in Figure 12.3.

An alternative way to represent molecules is to use a linear notation. A linear notation uses alphanumeric characters to code the molecular structure. These have the advantage of being much more compact than the connection table and so can be particularly useful for transmitting information about large numbers of molecules. The most famous of the early line notations is the Wiswesser line notation [Wiswesser 1954]; the SMILES notation is a more recent example that is increasingly popular [Weininger 1988]. To construct the Wiswesser

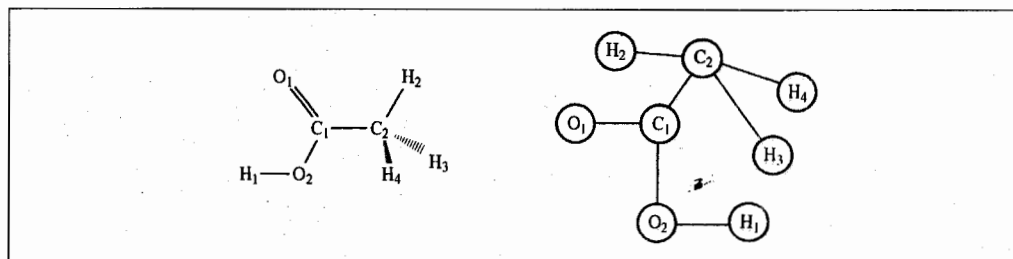


Fig. 12.2: The molecular graph of acetic acid.

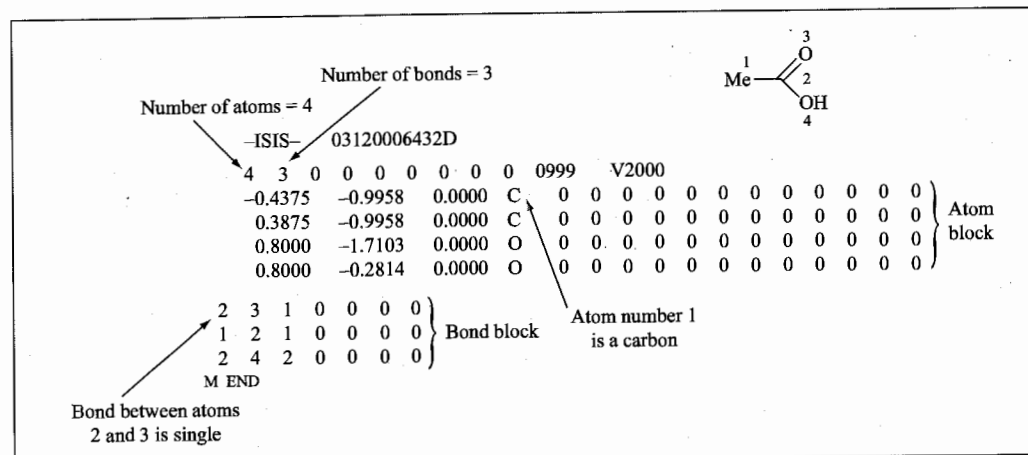


Fig. 12.3: MDL mol file for acetic acid, in the hydrogen-suppressed form.

notation for a molecule requires the application of a complex series of rules. SMILES is rather simpler, and with just a few rules one can write and interpret most SMILES strings. Thus atoms are represented by their atomic symbol. Hydrogens are not explicitly represented, except in special cases, as it is a *hydrogen-suppressed* notation. Upper case is used for aliphatic atoms, lower case for aromatic atoms. Single or aromatic bonds are not explicitly written but are assumed. Double bonds are represented by '=', triple bonds by '#'. The SMILES is constructed by 'walking' through the chemical diagram from one end to the other such that all atoms are visited just once. Rings are dealt with by 'breaking' one of the ring bonds, which are then indicated by appending an integer to the relevant atoms. Branching is indicated using brackets; any level of nesting if possible. Thus the simplest SMILES is probably C (methane). Ethane is CC, propane is CCC, 2-methyl propane is CC(C)C. Cyclohexane is C1CCCCC1 (note the use of the integer to indicate the ring bond). Benzene is c1ccccc1. Acetic acid is CC(=O)O. The SMILES for ranitidine (see Figure 9.26) is CNC(=CN(=O)=O)NCCSCc1ccc(CN(C)C)o1. Information about stereochemistry and geometrical isomerism can also be included in the SMILES notation.

One feature of both the connection table and the SMILES string formats is that there may be many different ways to represent the same molecule. Thus one may choose to number the atoms in the connection table in a different order or to write the SMILES differently (for example, acetic acid can be represented as OC(=O)C, O=C(C)O, O=C(O)C, etc.). A key requirement for any chemical database system is that it can determine whether or not a new molecule is already present in the system. This is typically done by generating some form of *canonical representation* of the structure. The canonical representation is unique irrespective of the numbering of the atoms in a mol file or the order of the atoms in a SMILES string. A popular method for doing this is the Morgan algorithm [Morgan 1965], which considers the properties of each atom together with those of its neighbours; this would enable the methyl carbon in acetic acid to be differentiated from the carboxyl carbon. It is also possible to generate a unique SMILES string for each molecule [Weininger

et al. 1989] (the canonical SMILES for acetic acid is CC(=O)O). Some algorithms can also incorporate information about stereochemistry and chirality into the canonicalisation. One immediate use of the canonical representation is that it can provide a very quick method to retrieve information about the compound. This is often done by generating a *hash key* from the structure. The hash key is typically an integer that is often used to indicate the location within a file where the requisite data is stored, so enabling the information to be retrieved very rapidly. The generation of the hash key can be done using well-established computer algorithms which can take a string of characters and generate the requisite integer.

A substructure search retrieves all the molecules from the database that contain the substructure. For example, we might wish to identify all compounds containing a carboxylic acid group. More complex queries are also possible in most systems; these would, for example, permit a query atom to match groups of atoms (e.g. 'any halogen') or features such as ring bonds or to specify stereochemistry. In the language of graph theory, substructure searching is known as *subgraph isomerism* - determining whether one graph is entirely contained within another. Even with the most efficient algorithms this is a relatively time-consuming process and so chemical database systems commonly use some form of screening method to rapidly eliminate molecules that cannot match the query. Such screens are often implemented using binary representations (a *bitstring*) and so operate very rapidly, especially if held in memory. There are two types of binary screen in common use. In a *structural key*, each position in the bitstring corresponds to a particular substructure. If that substructure is present in the molecule then the relevant bit in the molecule's key is set to 1. A predefined fragment dictionary is used to specify the substructures. As each molecule is added to the database a substructure search is performed for each fragment and the relevant bit assigned. Many different types of substructure can be incorporated, such as the presence or absence of particular elements, rings and common functional groups. It is also possible to assign bits which encode how many occurrences of a particular feature there are, such as 'at least two methyl groups'. The features in the fragment dictionary are defined to give optimal performance in 'typical' searches, depending upon the type of molecules in the database. This has the advantage that, when an effective choice of screens is used, the performance of the search should be very efficient, but if the dictionary is not appropriate then many molecules will pass through the screen and be subjected to the slow atom-based substructure search. Thus, a dictionary designed for typical 'organic' or 'drug-like' molecules might be inappropriate for a database containing just hydrocarbon molecules. The structural keys used by the MACCS and Isis systems from Molecular Design are probably the best known of this type of bitstring.

The alternative is to use a *hashed fingerprint*, which does not require a predefined fragment dictionary. Rather, an algorithmic approach is used to derive the bitstring, which initially contains all zeros. This method generates all possible linear paths of connected atoms through the molecule containing between 1 and a pre-defined number of atoms (e.g. 8). For example, in acetic acid the paths of length zero are just the atoms C and O, the paths of length 1 are CC, C=O and CO, and of length 2 are CCO and CC=O. Each path defines a pattern of atoms and bonds which serves as the input to a pseudo-random number generator, which produces a set of bits which are then set to the value 1. The hashing process typically sets 4 or 5 bits per pattern. A bitstring might contain 1024 bits, and after all paths

have been examined a typical organic, drug-like molecule might have a total of 200–300 bits set to 1. Obviously, the greater the number of different paths in the molecule the more bits (on average) that are set. Note that the use of the hashing algorithm means that it is possible (indeed, quite likely for typical molecules) that any one bit could be set by more than one pattern. However, it is much less likely (though nevertheless still possible) that the same set of bits would be set by different patterns. Hashed fingerprints are used in a number of database systems and are particularly associated with the systems from Daylight Chemical Information Systems.

When using a bitstring screen, the first operation is to calculate the corresponding bitstring for the substructure query. This query bitstring is then compared with the bitstrings for all the molecules in the database. A molecule can only possibly match the query if it contains a '1' for every position in the bitstring where the query also has a '1'. This comparison can be performed very quickly and so the database can be screened very rapidly. Well-designed screens can eliminate up to 99% of the molecules during this phase. The presence of clashes in hashed fingerprints does not affect the final results of a substructure search, though they

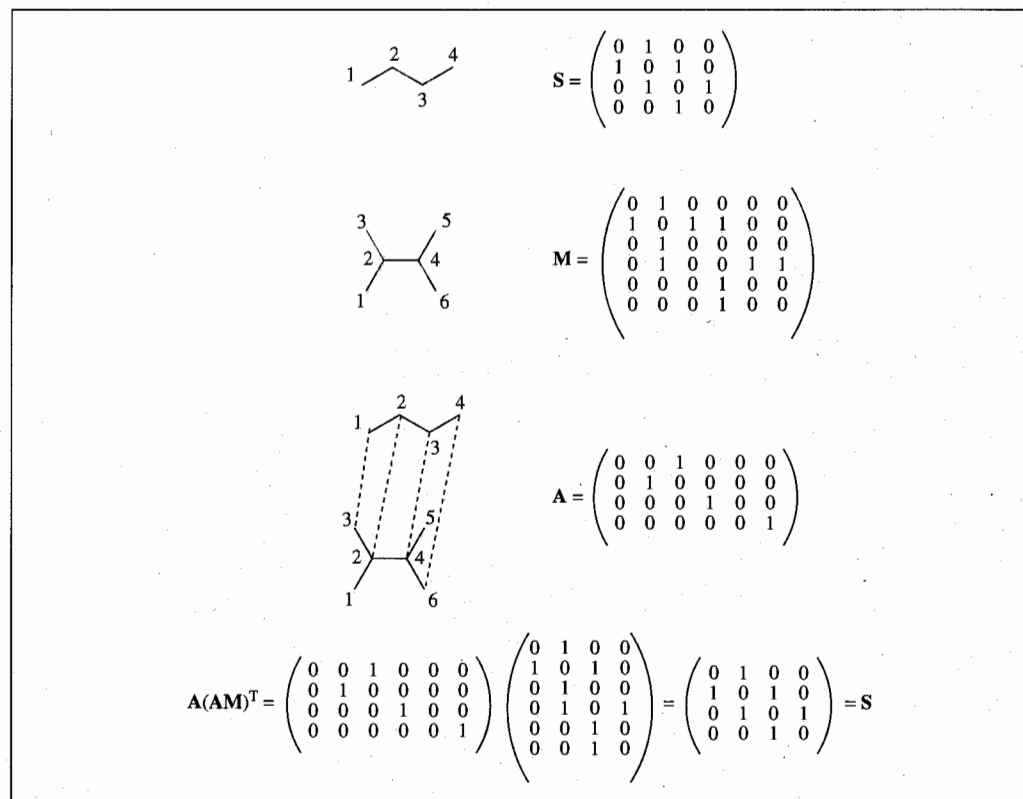


Fig. 12.4: Illustration of the operation of the Ullmann algorithm using a 4-atom substructure and a 6-atom 'molecule'. The proposed match is shown in the bottom figure, together with the relevant matrices used in the calculation.

might have an impact upon the screening efficiency, because more molecules need to be considered for the full substructure search.

Having eliminated molecules that could not match the query using the bitstring screen it is then necessary to undertake the more time-consuming atom-by-atom search for the molecules that remain. One commonly used method for the subgraph isomorphism problem was described by Ullmann [Ullmann 1976]. This algorithm represents the molecular graphs of both the query substructure and the potential molecular match by an *adjacency matrix*, which is a square, symmetric matrix such that the element (ij) has the value 1 if atoms i and j are bonded, and zero otherwise. Sample adjacency matrices are shown in Figure 12.4. If there are N_m atoms in the database molecule and N_s atoms in the substructure then the Ullmann algorithm tries to find matrices A such that $A(AM)^T$ is identical to S , where M is the adjacency matrix of the molecule and S is the adjacency matrix of the substructure. The matrix A has N_m columns and N_s rows such that each row contains just one 1 and each column contains no more than one 1. This matrix represents a possible match between the substructure and the molecule such that if an element A_{ij} is set to 1 then the atom numbered j in the substructure matches atom i in the molecule. Figure 12.4 shows an example of a matrix A which does indeed correspond to a successful match. The simplest implementation of the Ullmann algorithm is to generate all possible matrices A systematically, testing each to see whether they meet the requirements and represent a match. However, refinements of this simplistic (and time-consuming) algorithm are possible. Indeed, in his original paper Ullmann showed that a consideration of the neighbours of each potential match could dramatically improve its performance. A query atom cannot match a database atom unless each of the neighbour atoms of the query atom also matches a neighbour of the database atom.

12.3 3D Database Searching

2D substructure searching is a very powerful and widely used technique for identifying molecules with some particular feature (or combination of features, as the substructure can contain disconnected fragments). However, it does have some serious limitations if we wish to discover novel molecules with the desired biological activity. Key to understanding these limitations is the fact that receptors do not recognise substructures – rather, it is the three-dimensional stereoelectronic features of a molecule that are important for molecular recognition. In a 3D database search one tries to find molecules that satisfy the chemical and geometric requirements of the receptor. As such, a 3D database contains information about the conformational properties and functionality features of the molecules contained within it. Moreover, in contrast to a 2D search, 3D searching can enable lead series to be identified that are structurally quite different from those already known. There are two general types of 3D database search, the choice of which to use being dependent on the information available about the target receptor. In the first case, detailed structural information about the target receptor is not available, but it may be possible to derive an abstract model called a *pharmacophore* that indicates the key features of a series of active molecules. In the second case, a three-dimensional structure of the target macromolecule is available from X-ray crystallography or NMR or comparative modelling.

12.4 Deriving and Using Three-dimensional Pharmacophores

In drug design, the term 'pharmacophore' refers to a set of features that is common to a series of active molecules. Hydrogen-bond donors and acceptors, positively and negatively charged groups, and hydrophobic regions are typical features. We will refer to such features as 'pharmacophoric groups'. These groupings can be considered an illustration of the important concept of *bioisosteres*, which are atoms, functional groups or molecules with similar physical and chemical properties such that they produce generally similar biological properties [Thornber 1979; Patani and LaVoie 1996]. Some common bioisosteric groups are shown in Figure 12.5. A *three-dimensional (3D) pharmacophore* specifies the spatial relationships between the groups. These relationships are often expressed as distances or distance ranges but may also include other geometric measures such as angles and planes. For example, a commonly used 3D pharmacophore for antihistamines contains two aromatic rings and a tertiary nitrogen distributed as shown in Figure 12.6. The development of methods for studying the conformations of ligands has stimulated an interest in the influence of the three-dimensional structures of molecules on their chemical and biological activity. The objective of a procedure known as *pharmacophore mapping* is to determine possible 3D pharmacophores for a series of active compounds and is usually used when an experimental structure of the target macromolecule is not available. Once a pharmacophore has been developed, it can then be used to find or suggest other active molecules.

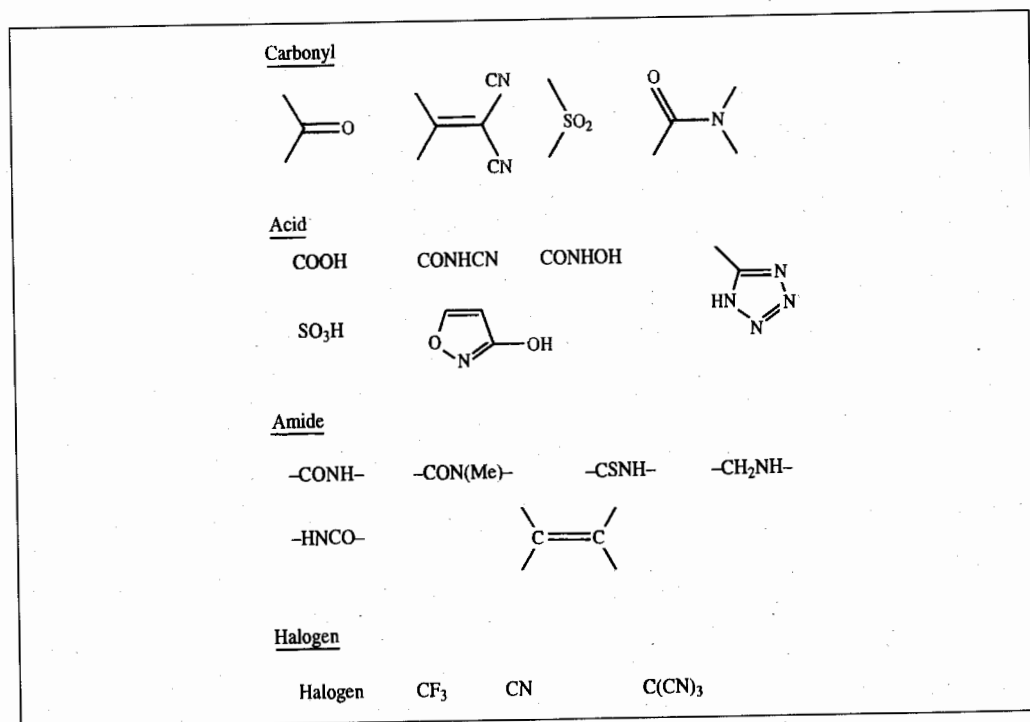


Fig. 12.5: Some common bioisosteric groups.

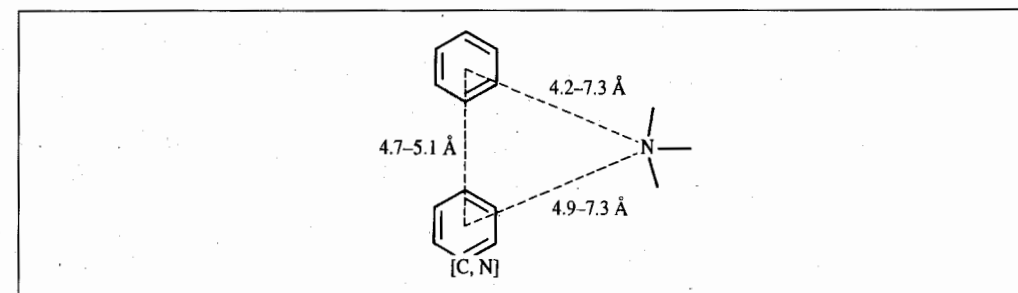


Fig. 12.6: Antihistamine 3D pharmacophore.

There are two problems to consider when calculating 3D pharmacophores. First, unless the molecules are all completely rigid, one must take account of their conformational properties. The second problem is to determine which combinations of pharmacophoric groups are common to the molecules and can be positioned in a similar orientation in space. More than one pharmacophore may be possible; indeed, some algorithms can generate hundreds of possible pharmacophores, which must then be evaluated to determine which best fits the data. It is important to realise that all of these approaches to finding 3D pharmacophores assume that all of the molecules bind in a common manner to the macromolecule.

12.4.1 Constrained Systematic Search

In some cases, it is relatively straightforward to deduce which features are required for activity. A well-known example is the pharmacophore for the angiotension-converting enzyme (ACE), which is involved in regulating blood pressure. Four typical ACE inhibitors are shown in Figure 12.7, including captopril, which is widely used to treat hypertension. Angiotension-converting enzyme is a zinc metalloprotease whose X-ray structure has not yet been solved. Three features within the class of inhibitors such as captopril are required for activity: a terminal carboxyl group (believed to interact with an arginine residue in the enzyme), an amido carbonyl group (which hydrogen bonds to a hydrogen-bond donor in the enzyme), and a zinc-binding group. The problem is to determine conformations in which the inhibitors can position these three pharmacophoric groups in the same relative position in space.

One of the most widely used methods for tackling this problem is the constrained systematic search method of Dammkoehler, Motoc and Marshall [Dammkoehler *et al.* 1989]. At first sight, it would appear that a systematic search over 20–30 molecules would greatly magnify the combinatorial explosion associated with a systematic conformational analysis. In fact, one can significantly reduce the scale of the problem by making use of information about molecules whose conformational space has already been considered. Thus, we are only interested in those conformations that would enable the current molecule's pharmacophoric groups to be positioned in the same locations that have already been found for previous molecules. Dammkoehler and colleagues showed that it is possible to determine what torsion angles of the rotatable bonds will enable conformations consistent with the previous

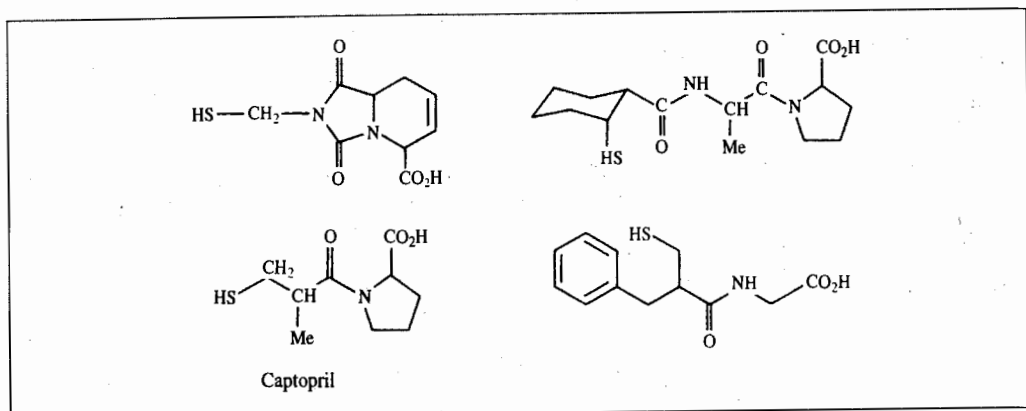


Fig. 12.7: Four typical ACE inhibitors.

results to be obtained. It is best to choose the most conformationally restricted molecules first, as these will have a reduced conformational space.

To derive an ACE pharmacophore, four points were defined for each molecule. The derivation of these four points for captopril is shown in Figure 12.8. Five distances (also shown in Figure 12.8) were defined between these four points. Note that one of the points corresponds to the presumed location of the enzyme's zinc atom. The number of rotatable bonds in each inhibitor varied between 3 and 9 and the molecules were considered in order of increasing number of rotatable bonds. The entire conformational space was explored for the first (most

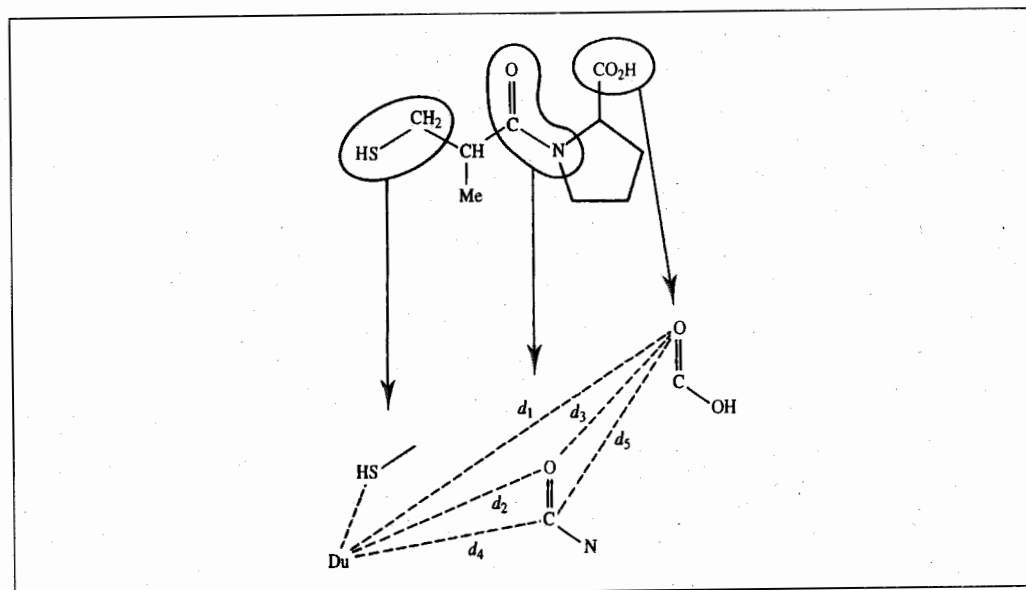


Fig. 12.8: Four points and five distances define the ACE pharmacophore.

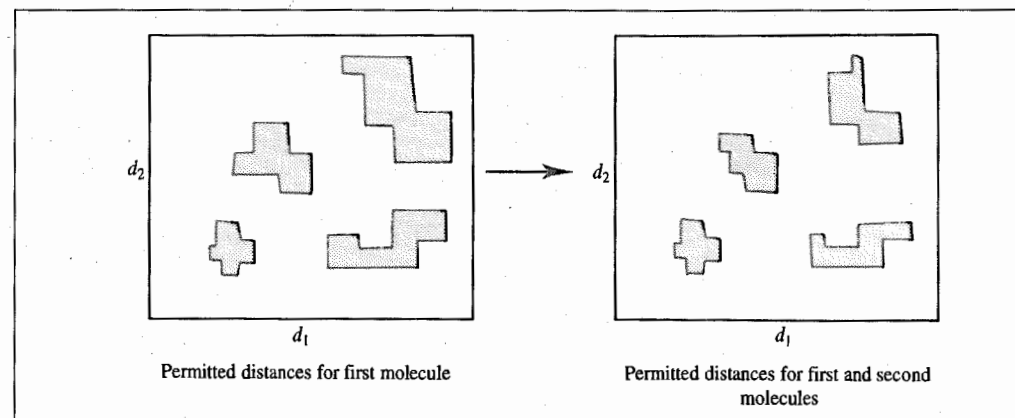


Fig. 12.9: A distance map indicates the distances available to specified groups. As more molecules are considered, the permitted regions get smaller.

inflexible) molecule. For each conformation, a point was registered in a five-dimensional hyperspace that corresponded to that particular combination of the five distances. When the second molecule was considered, only those torsion angles that would enable these distances to be achieved were permitted to the rotatable bonds. As more molecules were examined, so the common regions in the five-dimensional hypersurface were reduced, as illustrated schematically for a two-dimensional example in Figure 12.9. Two distinct 3D pharmacophores were obtained from the search, shown in Figure 12.10. The constrained search can be performed three orders of magnitude faster than the approach involving a separate systematic search on all the molecules.

12.4.2 Ensemble Distance Geometry, Ensemble Molecular Dynamics and Genetic Algorithms

A variant of distance geometry called *ensemble distance geometry* [Sheridan *et al.* 1986] can be used to simultaneously derive a set of conformations with a previously defined set of pharmacophoric groups overlaid. Ensemble distance geometry uses the same steps as

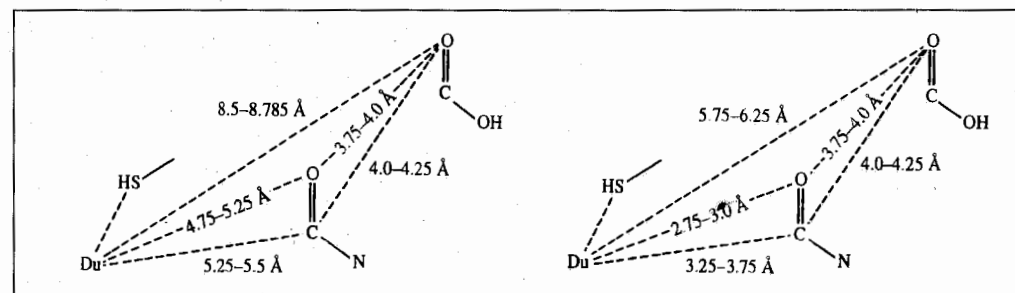


Fig. 12.10: Two ACE pharmacophores identified by the constrained systematic search [Dammkoehler *et al.* 1989].

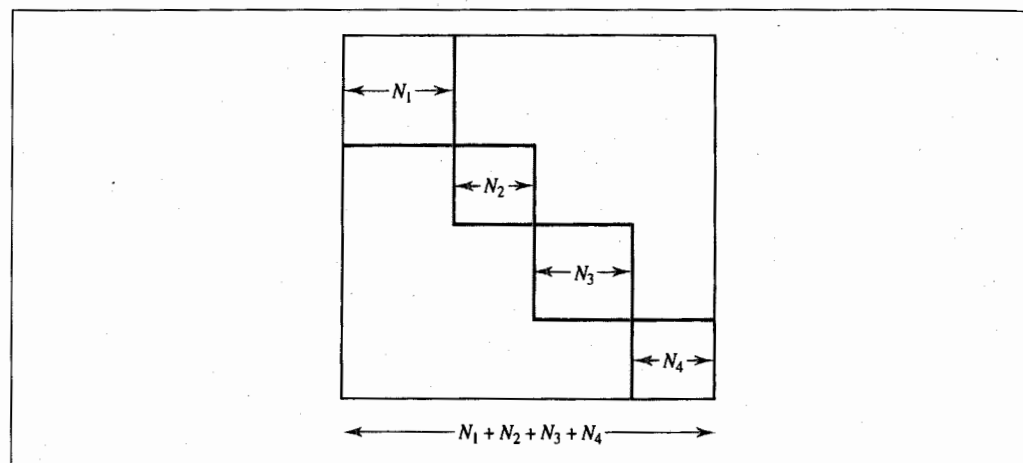


Fig. 12.11: Distance matrix used in ensemble distance geometry. There are N_1 atoms in the first molecule, N_2 in the second, and so on.

standard distance geometry, with the special feature that the conformational spaces of all the molecules are considered simultaneously. This is done using much larger bounds and distance matrices, with dimensions equal to the sum of the atoms in all the molecules. In these matrices, elements 1 to N_1 correspond to the N_1 atoms of molecule 1, elements $N_1 + 1$ to $N_1 + N_2$ to the N_2 atoms of molecule 2, and so on (Figure 12.11). Elements (i, j) and (j, i) of the bounds matrix thus represent the upper and lower bounds between atoms i and j (which may or may not be in the same molecule). The upper and lower bound distances between two atoms that are in the same molecule are set in the usual way. The lower bounds for atoms that are in different molecules are set to zero. This enables the molecules to be overlaid in three-dimensional space. The upper bounds for pairs of atoms that are in different molecules are set to a large value, except for those atoms that need to be superimposed in the pharmacophore, which are set to a small tolerance parameter. Having defined the bounds matrix, the usual distance geometry steps are followed: smoothing, assignment of random distances, and optimisation against the initial bounds.

The first application of ensemble distance geometry was to derive a model of the nicotinic pharmacophore using the four nicotinic agonists shown in Figure 12.12. Three sets of atoms were selected as the pharmacophoric groups labelled A, B and C in Figure 12.12. The ensemble distance geometry algorithm generated several different solutions, but after eliminating those that contained distorted bond lengths or angles or unfavourable van der Waals contacts the remaining solutions corresponded to a single pharmacophore. This pharmacophore can be represented as a triangle (Figure 12.8). Note that the B-C distance is fixed at the length of the C=O bond. Also note that in (-)-nicotine the centroid of the pyridine ring is defined as one of the pharmacophoric points. The pharmacophore was then validated by confirming that low-energy conformations could be generated for other known nicotinic agonists that were consistent with the distance constraints of the pharmacophore.

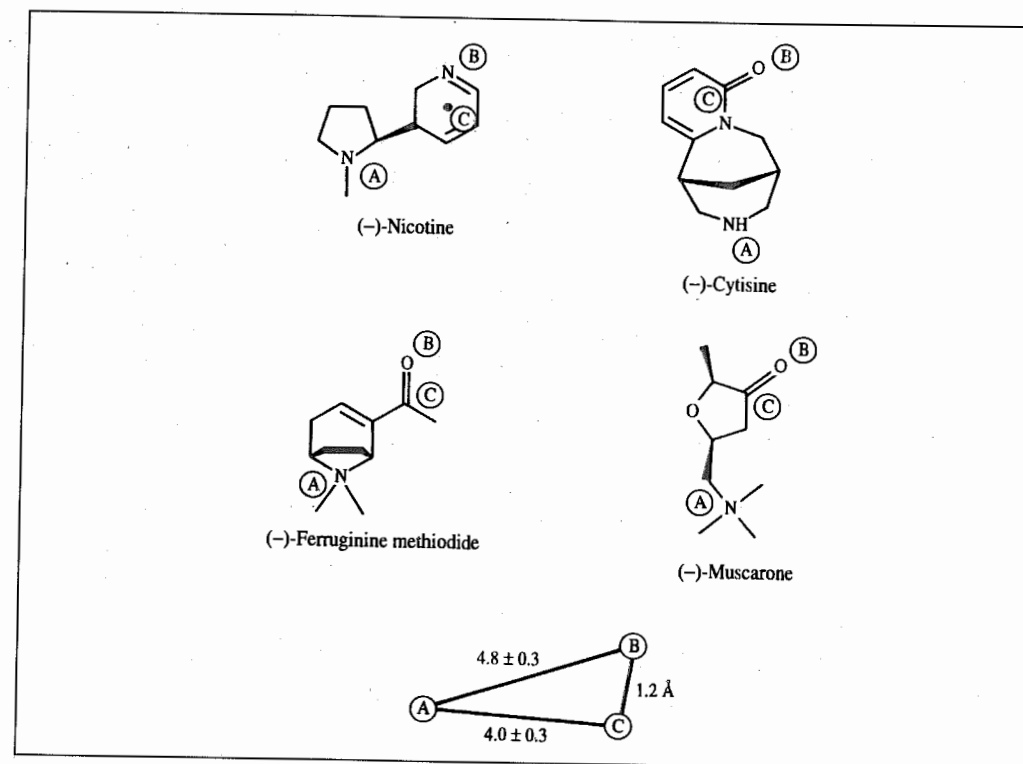


Fig. 12.12: Four molecules used to derive the nicotinic pharmacophore by distance geometry and the pharmacophore obtained.

In a related way, *ensemble molecular dynamics* derives a pharmacophore using restrained molecular dynamics for a collection of molecules. A force field model is set up so that none of the atoms in each molecule 'sees' the atoms in any other molecule. This enables the molecules to be overlaid in space. A restraint term is included in the potential, which forces the appropriate atoms or functional groups to be overlaid in space.

In the ensemble distance geometry and ensemble molecular dynamics methods together with the constrained systematic search it is necessary to provide the sets of matching atoms. These are then used to constrain the conformational search space. By contrast, the genetic algorithm method [Jones *et al.* 1995a] explores not only the conformational degrees of freedom of the various molecules but also the possible feature matches. These are thus all encoded within the chromosome. A standard genetic algorithm (see Section 9.9.1) is then applied to generate possible pharmacophores.

12.4.3 Clique Detection Methods for Finding Pharmacophores

When many pharmacophoric groups are present in the molecule it may be very difficult to identify all possible combinations of the functional groups (there may be thousands of

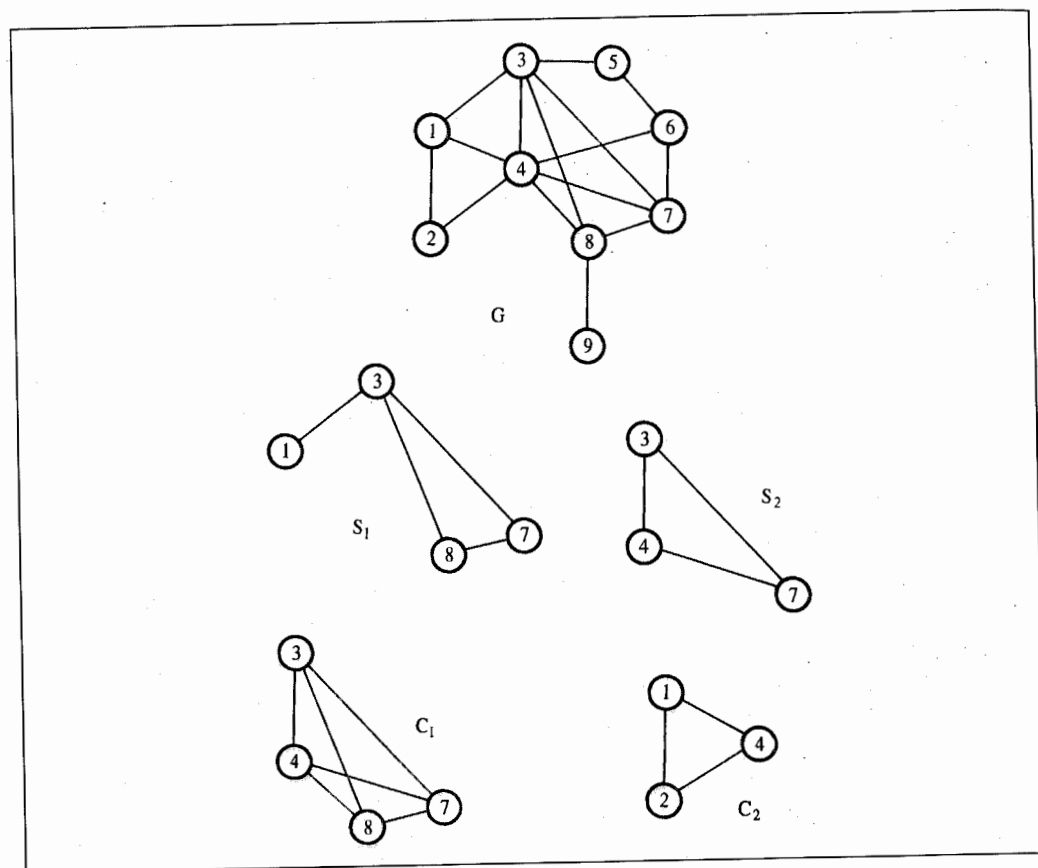


Fig. 12.13: Identifying cliques in a graph.

possible pharmacophores). To tackle this problem, *clique detection* algorithms can be applied to a set of precalculated conformations of the molecules. Cliques are based upon the graph-theoretical approach to molecular structure that we discussed above.

A clique is defined as a 'maximal completely connected subgraph'. This definition is best understood by considering a simple example. Consider the graph G in Figure 12.13, together with various subgraphs. G is not a completely connected graph, because there is not an edge between all the nodes. The subgraph S₁ is not a completely connected subgraph, because there is no edge between nodes 1 and 8. The subgraph S₂ is a completely connected subgraph, because there are edges between all the nodes. However, S₂ is not a clique, because it is not a maximal completely connected subgraph; it is possible to add node 8 in order to obtain the clique C₁. A graph may contain many cliques; thus Figure 12.13 also shows a second clique, C₂. Finding the cliques in a graph belongs to a class of problems that are known as NP-complete. This means that the computational time required to find an exact solution increases in an exponential fashion with the size of the problem. Many algorithms

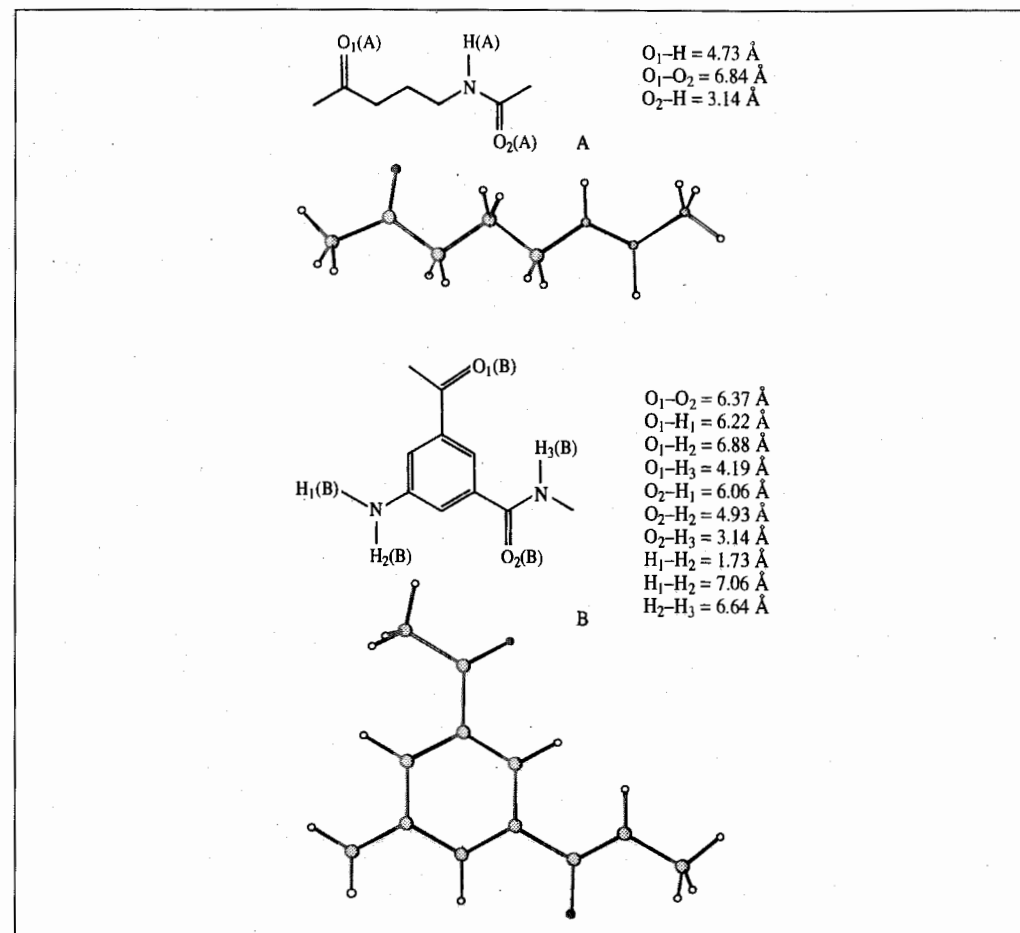


Fig. 12.14: Two molecules used to illustrate clique detection.

have been devised for finding cliques; the method of Bron and Kerbosch has been found to be suitably efficient for pharmacophore identification [Bron and Kerbosch 1973].

How is clique detection related to the identification of pharmacophores [Martin *et al.* 1993]? Let us suppose that we are comparing two conformations of two molecules, A and B (Figure 12.14). We construct a graph in which there is a node for every pair of matching pharmacophoric groups in the two structures. The two hydrogen-bond acceptors in molecule A (O₁(A) and O₂(A)) and the two in molecule B (O₁(B) and O₂(B)) give rise to four nodes in the joint graph. There is one hydrogen-bond donor in molecule A (H(A)) but three in molecule B (H₁(B), H₂(B) and H₃(B)), giving rise to three nodes in the graph. The intramolecular distances between these groups are indicated in Figure 12.14. An edge is drawn between each pair of nodes when the distance between the corresponding groups in the two

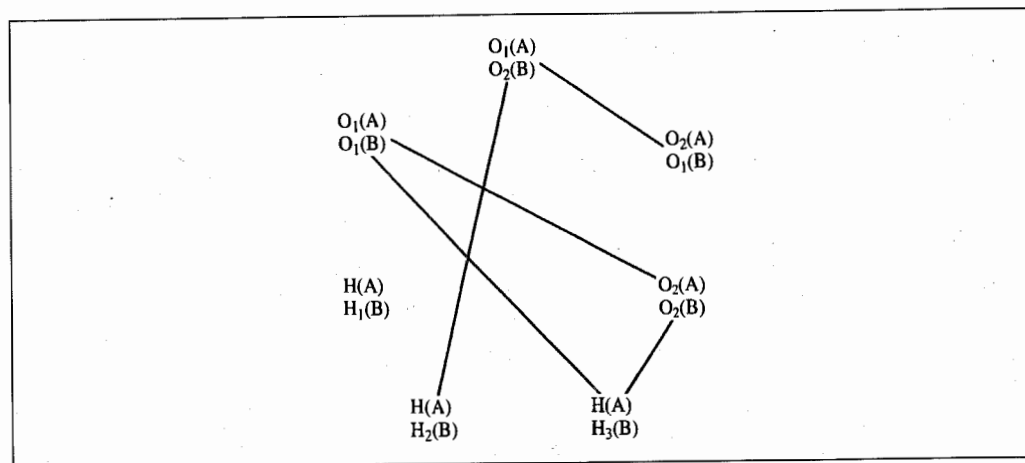


Fig. 12.15: The matching graph for the two molecules in Figure 12.14.

molecules is the same, within some tolerance. For example, the distance between $O_1(A)$ and $O_2(A)$ is 4.73 Å and the distance between $O_1(B)$ and $O_2(B)$ is 4.19 Å. If the tolerance is at least 0.54 Å then the two distances would be considered equal and so an edge is drawn between the corresponding pairs of nodes in the graph, as shown in Figure 12.15. The full graph is shown in Figure 12.15, where we have assumed a tolerance of 0.6 Å. Clique detection is used to find maximal sets of matching groups for the two molecules; in this simple example there are three cliques, two containing just two overlapping atoms and one containing three matching atoms (Table 12.1).

In the clique detection approach, the first step is to generate a family of low-energy conformations for the molecules. The molecule with the smallest number of conformations is used as the starting point, with each of its conformations being used in turn as the reference structure. Each conformation of every other molecule is then compared with the reference conformations and the cliques identified. The cliques for each molecule are obtained by combining the results for each of its conformations. Those cliques that are common to at least one conformation from each molecule can then be combined to give a possible 3D pharmacophore for the entire set.

Clique number	Atom from A	Atom from B
1	O_1	O_2
	H	H_2
2	O_1	O_2
	O_2	O_1
3	O_1	O_1
	O_2	O_2
	H	H_3

Table 12.1: Cliques found when matching the molecules in Figure 12.14.

12.4.4 Maximum Likelihood Method

One limitation of clique detection is that it needs to be run repeatedly with different reference conformations and the run-time scales with the number of conformations per molecule. The maximum likelihood method [Barnum *et al.* 1996] eliminates the need for a reference conformation, effectively enabling every conformation of every molecule to act as the reference. Despite this, the algorithm scales linearly with the number of conformations per molecule, so enabling a larger number of conformations (up to a few hundred) to be handled. In addition, the method scores each of the possible pharmacophores based upon the extent to which it fits the set of input molecules and an estimate of its 'rarity'. It is not required that every molecule has to be able to match every feature for the pharmacophore to be considered.

Prior to the pharmacophore identification phase, a set of conformations is generated for each molecule. Typically, the 'poling' method (see Section 9.14) is used to produce a reasonably small but representative set of low-energy conformations. First, all possible combinations of pharmacophore features (e.g. donor-donor-acceptor, aromatic ring-donor-acceptor-hydrophobic region) are exhaustively considered. Possible geometric arrangements of the features in 3D space are identified by taking each molecule to be the reference structure and examining its conformations. These configurations are scored and ranked according to how well they describe the set of active molecules. Each configuration is considered an 'hypothesis' which can be used to assign a probability that a molecule is active depending on whether or not it matches the pharmacophore. If there are K features in the pharmacophore then the algorithm defines $K + 2$ possible ways in which a molecule can match the pharmacophore on a scale from 0 to $K + 1$. If the molecule matches all K features it is placed in the class $x = K + 1$. If it does not match all K features or any of the subsets obtained by removing one of the features (there are K such subsets, each containing $K - 1$ features) then the molecule is placed in the class $x = 0$. A molecule which can match one of the subsets with $K - 1$ features is assigned to an intermediate class between $x = 1$ and $x = K$, depending upon which of the subsets (ordered by selectivity) it matches. Both the full match and each of the K partial matches are assigned a 'rarity value' depending upon the type of the features present and on their relative disposition. Pharmacophores that contain 'rare' features (such as positively ionisable groups) are scored more highly than those that contain more common features (such as hydrophobic regions). In addition, the greater the distribution of the features (as measured by the squared distance between the feature and the common centroid) the higher the score.

This rarity value is equated with the fraction of hits that would be returned by searching a large database of diverse molecules with the full pharmacophore (all K features) or the subset (with $K - 1$ features) as appropriate. Labelling this fraction of hits as $p(x)$ we now define $q(x)$ as the fraction of the M active molecules (i.e. the molecules originally supplied as input to the procedure) which match each of the $K + 1$ possible classes. The overall configuration is scored using:

$$\text{score} = M \sum_x q(x) \log_2 \left(\frac{q(x)}{p(x)} \right) \quad (12.1)$$

Higher scores are achieved by pharmacophores that have large values of $q(x)$ (i.e. are matched by more of the active molecules in the initial set) and lower values of $p(x)$ (i.e. it is less likely that a 'random' molecule would match). The actual numerical values $p(x)$ are obtained from a predefined regression equation.

The pharmacophores generated by this approach are typically expressed in terms of 'location constraints' rather than inter-feature distance ranges [Greene *et al.* 1994]. These location constraints are typically expressed as a point in 3D space surrounded by a spherical region. A molecule must be able to place the relevant features within the appropriate spheres (see Figure 12.16, colour plate section). The spheres can be of different sizes to reflect the fact that the various interactions have differing sensitivities to changes in distance. For example, the tolerance on a charge interaction might be smaller than on a hydrogen-bonding interaction to reflect the fact that the energy of an ionic interaction is more sensitive to changes in relative position than for the hydrogen bond. Another useful aspect of this approach is its requirement that the features must be 'surface accessible' in the molecule for it to be considered a match.

12.4.5 Incorporating Additional Geometric Features Into a 3D Pharmacophore

The features used to define a 3D pharmacophore are most easily derived from the positions of specific atoms within each molecule. It may be more appropriate to consider locations around the molecule where the receptor might position its functional groups. This is especially relevant for hydrogen-bond donors and acceptors; two ligands may be able to hydrogen bond to the same protein atom with the ligand atoms being in a completely different location in the binding site, as illustrated in Figure 12.17. 3D pharmacophores may also be defined in terms of specific geometrical relationships between the pharmacophoric groups, such as the angle between the planes of two aromatic rings. The 3D pharmacophore may also contain features that are designed to mimic the presence of the receptor. These are commonly represented as *exclusion spheres*, which indicate locations within the 3D pharmacophore where no part of a ligand is permitted to be positioned. Some of these additional features are illustrated in Figure 12.18.

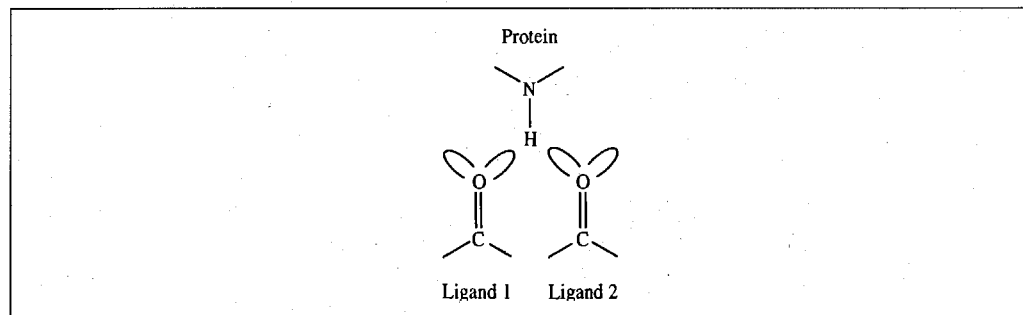


Fig. 12.17: Two ligands may be able to position a hydrogen-bond acceptor in different locations in space yet still interact with the same hydrogen-bond donor.

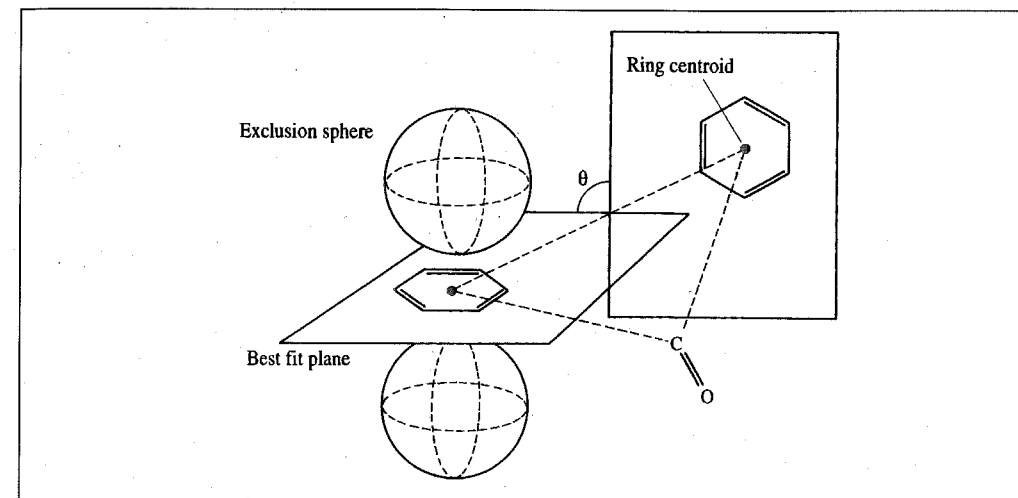


Fig. 12.18: Features that can be incorporated into 3D pharmacophores.

12.5 Sources of Data for 3D Databases

In a 3D database search one needs to consider the three-dimensional structures of the molecules. Where does such structural information come from? An obvious source is the Cambridge Structural Database, which contains experimental X-ray structures of more than 150 000 compounds. However, for most of the compounds in a typical compound database no crystal structure is available. Structure generation programs are designed to produce one or more low-energy conformations solely from the molecular graph. As the number of compounds may be very large, such programs must be able to operate automatically, rapidly and with little or no user intervention (i.e. without crashing!). The two most widely used structure generators to date are the CONCORD [Rusinko *et al.* 1988] and CORINA programs [Gasteiger *et al.* 1990], which both use a knowledge-based approach combined with energy minimisation.

Most structure generation algorithms produce only a single conformation for each molecule. With the possible exception of wholly rigid molecules, there is no guarantee that this structure corresponds to the conformation adopted when the molecule binds. We therefore need some way to take conformational flexibility into account during the 3D database search. The simplest way to do this is to store information about many conformations. To store conformations explicitly would usually require a large amount of disk space and so the information is usually compressed into a more compact form. To make the 3D search more efficient screening methods are commonly used, similar to those employed for '2D' substructure searches. One straightforward way to do this is to derive a set of distance 'keys' for the pairs of pharmacophoric groups in the molecule. Each key is a binary number in which each bit corresponds to a distance range between the appropriate pair of groups (donor-donor, donor-acceptor, etc.). For example, the first bit could correspond

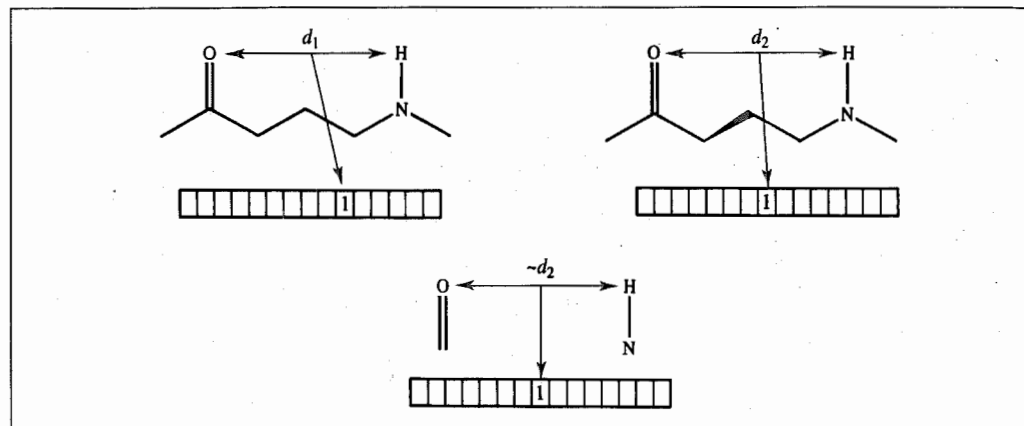


Fig. 12.19: 3D database searching. As each conformation is generated an appropriate bit is set in the binary key. At search time, the binary key appropriate to the pharmacophore is set up and compared with the keys in the database.

to a distance in the range 2.0–2.5 Å, the second bit to the range 2.5–3.0 Å, and so on. In fact, it is more efficient to use smaller bin sizes for the more common distances and larger bins for less common distances, because the distribution of distances between such groups in molecules is not uniform. The key initially contains all zeros. As each conformation is generated, the distances between the pharmacophoric groups are calculated and the appropriate bits in the relevant keys are changed to 1s (Figure 12.19). To search the database, the keys corresponding to the pharmacophore are calculated. The pharmacophore's keys are then compared with each molecular key, so identifying all molecules which could match the pharmacophore. Separate substructure-like screens are also used; these contain information about the number of each type of feature (e.g. number of donors). If the molecule does not contain the minimal number of groups in the pharmacophore, then it obviously cannot match and so can be discarded before its conformational properties need to be considered.

An alternative strategy is to explore the conformational space for each molecule during the database search. Systems which employ such an approach rely heavily upon screens which identify and reject molecules that could not satisfy the requirements of the pharmacophore before their conformational space is explored. These screens can be determined solely from the molecular graph and are typically represented as distance ranges. Triangle smoothing (Section 9.5) is one way in which such distance screens can be calculated; it provides the upper and lower bounds on interatomic distances. However, the distance ranges provided by triangle smoothing can be much wider than the actual distances that are observed in real structures. A simple example suffices to illustrate this point: the distance obtained when triangle smoothing is used to calculate the lower bound distance between the amide nitrogen and the carbonyl oxygen of the carboxylic acid group in 4-acetamido benzoic acid (Figure 12.20) is equal to the sum of the van der Waals radii (approximately 3.3 Å, depending upon the van der Waals radii used), compared to a distance of about 6.4 Å in all the accessible conformations of this molecule.

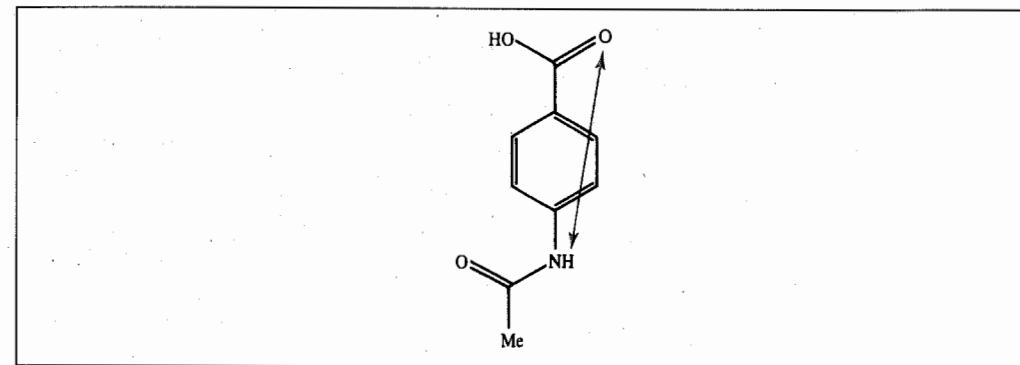


Fig. 12.20: 4-Acetamido benzoic acid. Triangle smoothing predicts that the lower bound distance between the amide nitrogen and the carbonyl oxygen is equal to the sum of the van der Waals radii. The actual distance is about 6.4 Å.

Having eliminated those molecules that could not possibly satisfy the geometric and chemical requirements of the pharmacophore, the program must explore the conformational degrees of freedom of the molecules that remain. This is done using methods to rapidly identify one or more conformations which satisfy the constraints of the pharmacophore. A natural method to use would be distance geometry, in which the pharmacophoric constraints would be incorporated into the bounds matrix, thereby leading to the generation of conformations that satisfy the constraints. However, distance geometry is rather too slow for this purpose. An alternative strategy is to 'adjust' or 'tweak' the conformation by rotating about single bonds, to force it to fit the pharmacophore. Adjustment is usually performed in torsional space (i.e. only the torsion angles are varied) by minimising an appropriate potential function expressed in terms of distances.

12.6 Molecular Docking

In molecular docking, we attempt to predict the structure (or structures) of the intermolecular complex formed between two or more molecules. Docking is widely used to suggest the binding modes of protein inhibitors. Most docking algorithms are able to generate a large number of possible structures, and so they also require a means to score each structure to identify those of most interest. The 'docking problem' is thus concerned with the generation and evaluation of plausible structures of intermolecular complexes [Blaney and Dixon 1993].

The docking problem involves many degrees of freedom. There are six degrees of translational and rotational freedom of one molecule relative to the other as well as the conformational degrees of freedom of each molecule. The docking problem can be tackled manually, using interactive computer graphics. This 'hands-on' approach can be very effective if we have a good idea of the expected binding mode, for example because we already know the binding mode of a closely related ligand. However, even in such cases one must be wary; X-ray crystallographic experiments have revealed that even very similar inhibitors may adopt

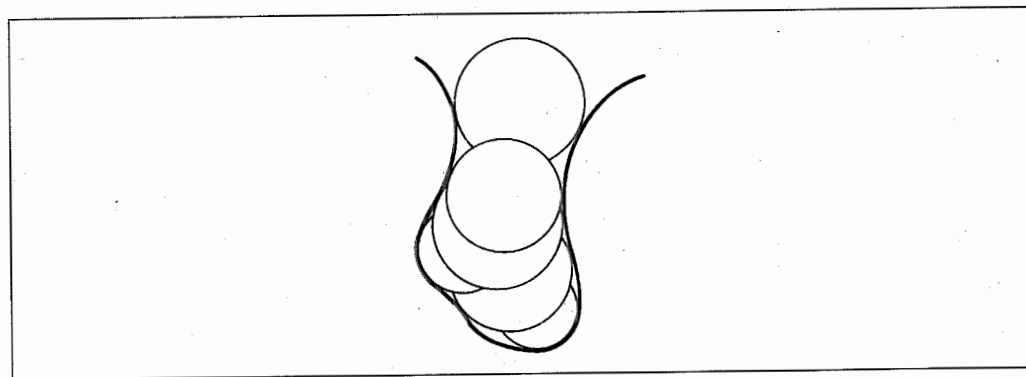


Fig. 12.21: A binding site represented as a collection of overlapping spheres.

quite different binding modes. Automatic docking algorithms can be less biased than human modellers and usually consider many more possibilities.

Various algorithms have been developed to tackle the docking problem. These can be characterised according to the number of degrees of freedom that they ignore. Thus, the simplest algorithms treat the two molecules as rigid bodies and explore only the six degrees of translational and rotational freedom. The earliest algorithms for docking small molecule ligands into the binding sites of proteins and DNA used this approximation. A well-known example of such an algorithm is the DOCK program of Kuntz and co-workers [Kuntz *et al.* 1982]. DOCK is designed to find molecules with a high degree of shape complementarity to the binding site. The program first derives a 'negative image' of the binding site from the molecular surface of the macromolecule. This negative image consists of a collection of overlapping spheres of varying radii, each of which touches the molecular surface at just two points, as shown schematically in Figure 12.21. Ligand atoms are then matched to the sphere centres to find matching sets (cliques) in which all the distances between the ligand atoms in the set are equal to the corresponding sphere centre–sphere centre distances (within some user-specified tolerance). The ligand can then be oriented within the site by performing a least-squares fit of the atoms to the sphere centres, as shown in Figure 12.22. The orientation is checked to ensure there are no unacceptable steric interactions between the ligand and the receptor. If the orientation is acceptable then an interaction energy is computed to give the 'score' for that binding mode. New orientations are generated by matching different sets of atoms and sphere centres. The top-scoring orientations are retained for subsequent analysis.

To perform conformationally flexible docking the conformational degrees of freedom need to be taken into account. Most of the methods that attempt to include the conformational degrees of freedom only consider the conformational space of the ligand; the receptor is invariably assumed to be rigid. All of the common methods for searching conformational space have been incorporated at some stage into a docking algorithm. For example, Monte Carlo methods have been used to perform molecular docking, often in conjunction with simulated annealing [Goodsell and Olson 1990]. At each iteration of the Monte Carlo procedure the internal conformation of the ligand is changed (by rotating about a bond) or the entire molecule is randomly translated or rotated. The energy of the ligand within

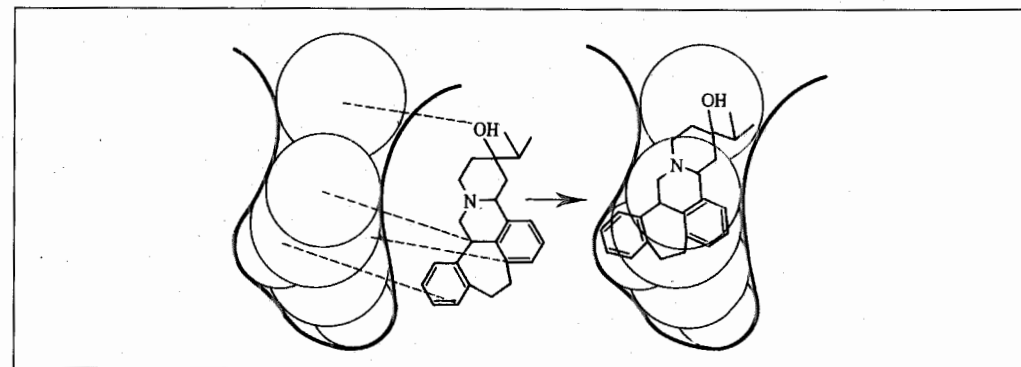


Fig. 12.22: The DOCK algorithm [Kuntz *et al.* 1982]. Atoms are matched to sphere centres and then the molecule is positioned within the binding site.

the binding site is calculated using molecular mechanics and the move is then accepted or rejected using the standard Metropolis criterion. An interesting variant on the basic Monte Carlo approach is the tabu search [Baxter *et al.* 1998]. This maintains a record of those regions of the search space that have already been visited, so ensuring that the method is encouraged to explore more of the binding site.

Genetic algorithms can also be used to perform molecular docking [Judson *et al.* 1994; Jones *et al.* 1995b; Oshiro *et al.* 1995]. Each chromosome codes not only for the internal conformation of the ligand as described in Section 9.9.1 but also for the orientation of the ligand within the receptor site. Both the orientation and the internal conformation will thus vary as the populations evolve. The score of each docked structure within the site acts as the fitness function used to select the individuals for the next iteration.

Distance geometry can be used to perform molecular docking. The major problem to be addressed with this method is to find a way to generate conformations of the ligand within the binding site. One way to achieve this is by using a modified penalty function that forces the ligand conformation to remain within the binding site. For example, an additional penalty term can be added which has the effect of forcing the ligand to lie in the DOCK-derived cluster of spheres that represents the binding site.

An approach that is used by a number of programs involves the incremental construction of the ligand [Leach and Kuntz 1990; Welch *et al.* 1996; Rarey *et al.* 1996]. This is similar in spirit to the depth-first systematic conformational search described in Section 9.2. The main difference, of course, is that in docking the conformational search is performed within the binding site. A typical incremental construction algorithm first identifies one or more 'base fragments' within the ligand. These base fragments are often chosen to be a reasonably significant, fairly rigid part of the molecule such as a ring system. The base fragment(s) are docked into the binding site and may then be clustered to remove similar orientations. Each docked orientation of the base fragment(s) then represents the starting point for the conformational analysis of the rest of the ligand. One might anticipate that such an approach would be very time-consuming, as it is in effect necessary to perform the conformational

analysis for each orientation of the base fragments. However, it is often found that the protein provides a particularly useful constraint, enabling the search tree to be pruned very effectively.

The ideal docking method would allow both ligand and receptor to explore their conformational degrees of freedom. Perhaps the most 'natural' way to incorporate the flexibility of the binding site is via a molecular dynamics simulation of the ligand-receptor complex. However, such calculations are computationally very demanding and are in practice only useful for refining structures produced using other docking methods; molecular dynamics does not explore the range of binding modes very well except for very small, mobile ligands. For many systems, the energy barriers that separate one binding mode from another are often too large to be overcome. Some other attempts have been made to incorporate protein flexibility (at least at the level of the side chains [Leach 1994]) but these methods are generally in their infancy and take much longer than rigid-protein docking.

When the first docking methods were developed the speed of the typical computer was such that only rigid-body docking of single molecules was feasible. As computational performance increased it was recognised that rigid-body docking could be used to examine large numbers of molecules from a database. At approximately the same time, algorithms that addressed the conformational flexibility of the ligand were devised. With the passage of time, it is now possible to search databases using a flexible-ligand algorithm. However, there is still a clear distinction between the use of docking to predict the binding mode of a single active molecule, where one can afford to use a particularly thorough search, and the use of docking for searching databases for possible lead compounds.

12.6.1 Scoring Functions for Molecular Docking

Most docking algorithms are capable of generating a large number of potential solutions. Some of these can be rejected immediately because they have a high-energy clash with the protein. The remainder must be assessed using some scoring function. When we are only interested in how a single ligand binds to the protein then the scoring function need only be able to identify the docked orientation that most closely corresponds to the 'true' structure of the intermolecular complex. However, when docking a database of molecules then not only should the scoring function be able to identify the 'true' docking mode of a given ligand but it also needs to be able to rank one ligand relative to another. Moreover, the large number of orientations that may be generated during a docking run means that it must be possible to calculate the scoring function rapidly.

Many of the scoring functions in common use attempt to approximate the binding free energy for the ligand binding to the receptor. We have previously encountered a number of ways in which simulation techniques can be used to predict (relative) free energies of binding (see Chapter 11), but these are far too slow to be of value in docking calculations. Faster, more approximate methods tend to be used. In contrast to the free energy perturbation approach these alternatives tend to consider that the free energy of binding can be written as an additive equation of various components to reflect the various contributions to binding [Bohm and Klebe 1996]. A complete equation of this kind would have the

following contributions [Ajay and Murcko 1995]:

$$\Delta G_{\text{bind}} = \Delta G_{\text{solvent}} + \Delta G_{\text{conf}} + \Delta G_{\text{int}} + \Delta G_{\text{rot}} + \Delta G_{\text{t/r}} + \Delta G_{\text{vib}} \quad (12.2)$$

where $\Delta G_{\text{solvent}}$ is the contribution due to solvent effects, arising from the balance of interactions between the solvent and the ligand, protein and intermolecular complex. Various methods can be used to determine these contributions. ΔG_{conf} arises from conformational changes in the protein and in the ligand. In many cases, the protein does not change much on binding (which is fortunate, given that most docking methods assume a rigid receptor). By contrast, the ligand changes from an ensemble of conformations in solution to what is often assumed to be a single dominant conformation in the bound state. Various analyses have been performed to try to determine the size of this energetic penalty for the ligand. When measured relative to the most significant conformation in solution, an average penalty of 3 kcal/mol was found [Bostrom *et al.* 1998]. ΔG_{int} is the free energy due to specific protein-ligand interactions. ΔG_{rot} is the free energy loss associated with freezing internal rotations of the protein and the ligand. This is mostly due to the entropic contribution. The simplest way to calculate this penalty is to assume that there are three states per rotatable bond (trans and \pm gauche) of equal energy, thus leading to a free energy loss of $RT \ln 3$ (~ 0.7 kcal/mol) per rotatable bond. $\Delta G_{\text{t/r}}$ is the loss in translational and rotational free energy caused by the association of two bodies (the ligand and the receptor) to give a single body (the intermolecular complex). This is often assumed to be constant for all ligands and so is ignored when one is interested in the relative binding strengths of different ligands. ΔG_{vib} is the free energy due to changes in vibrational modes. This contribution is difficult to calculate and is usually ignored.

Each of the terms in Equation (12.2) has been the subject of considerable discussion in the literature, and for some of them there may be a number of different approaches to their estimation. However, many of these methods are unsuitable for docking, due to the calculation time required. Some very simple functions have been employed for docking, such as that originally used in the DOCK program (illustrated in Figure 12.23 together with another similar form, the piecewise linear potential [Gelhaar *et al.* 1995]). Despite their apparent simplicity, such functions continue to rate well in comparisons of different functional

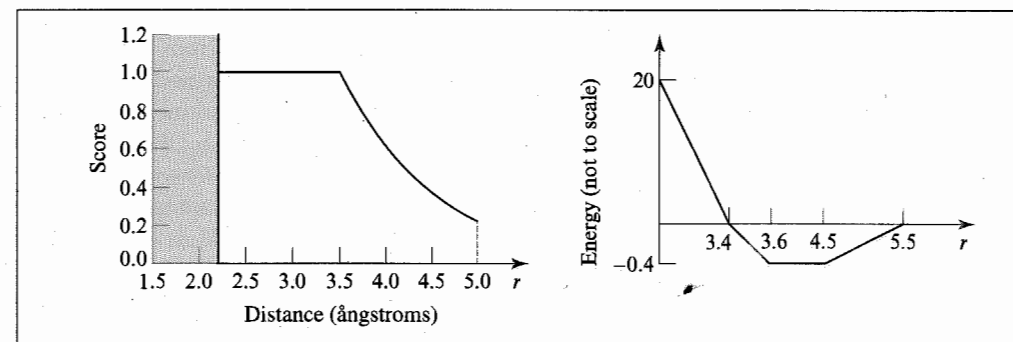


Fig. 12.23: Two simple scoring functions used in docking. On the left is the basic scoring scheme used by the DOCK program [Desjarlais *et al.* 1988]. On the right is the 'piecewise linear potential' [Gelhaar *et al.* 1995].

forms. Molecular mechanics is also widely used to calculate the energy of interaction; one way in which such a calculation can be speeded up is to pre-calculate electrostatic and van der Waals 'potentials' on a regular grid which covers the binding site [Meng *et al.* 1992]. The computational effort required to calculate the energy of interaction between ligand and protein is then linear in the number of atoms in the ligand, rather than being proportional to the product of the number of ligand atoms multiplied by the number of protein atoms.

Simple molecular mechanics scoring functions are popular, but we can see from Equation (12.2) that they provide only part of the overall free energy of binding. Thus whilst they have proved successful in some cases (such as a study of HIV-protease inhibitors [Holloway *et al.* 1995]), one should not be surprised if they do not always work. An interesting approach to this problem was suggested by Böhm. He tried to find a simple linear relationship between the free energy of binding and a variety of parameters which it was anticipated would be relevant to the overall free energy of binding and which could also be calculated rapidly [Böhm 1994]. The terms in this original formulation related to hydrogen bonding, ionic interactions, lipophilic interactions and the loss of internal degrees of freedom of the ligand:

$$\Delta G_{\text{bind}} = \Delta G_0 + G_{\text{hb}} \sum_{\text{h-bonds}} f(\Delta R, \Delta \alpha) + \Delta G_{\text{ionic}} \sum_{\text{ionic interactions}} f(\Delta R, \Delta \alpha) + \Delta G_{\text{lipo}} |A_{\text{lipo}}| + \Delta G_{\text{rot}} NROT \quad (12.3)$$

where ΔG_0 is a constant term, independent of the system, which was interpreted to correspond to the overall change in translational/rotational free energy ($\Delta G_{t/r}$ in Equation (12.2)). ΔG_{hb} corresponds to the contribution from an ideal hydrogen bond. This contribution is multiplied by a penalty function $f(\Delta R, \Delta \alpha)$ which accounts for large deviations of the hydrogen bond from the ideal geometry; ΔR is the deviation of the hydrogen-bond distance from its ideal value of 1.9 Å, and $\Delta \alpha$ is the deviation from the ideal angle of 180°. The same geometric dependency is applied to the ionic interactions. ΔG_{lipo} is a contribution from lipophilic interactions, which are assumed to be proportional to the lipophilic contact surface (i.e. involving non-polar atoms) between the protein and the ligand, A_{lipo} . ΔG_{rot} is the loss of free energy due to freezing a rotatable bond in the ligand upon binding. It is thus multiplied by the number of rotatable bonds in the ligand, $NROT$.

Experimental binding data on 45 protein–ligand complexes was extracted from the literature and then a multiple linear regression analysis (see Section 12.12.2) was performed to derive the parameters in this equation (i.e. the various ΔG values). The values of the parameters obtained from this analysis ($\Delta G_{\text{hb}} = -1.2$ kcal/mol, $\Delta G_{\text{ionic}} = -2.0$ kcal/mol, $\Delta G_{\text{lipo}} = -0.04$ kcal/mol Å², $\Delta G_{\text{rot}} = +0.3$ kcal/mol, $\Delta G_0 = +1.3$ kcal/mol) mostly correspond reasonably closely to values estimated from other approaches, with the exception of the constant term, ΔG_0 , for which a value between 7 and 11 kcal/mol is generally agreed [Ajay and Murcko 1995]. The model reproduced the experimental binding data with a standard deviation of 1.7 kcal/mol. The exponential relationship between the binding free energy and the equilibrium constant means that a change of just 1.4 kcal/mol in the free energy corresponds to a ten-fold change in affinity. This work has spawned a number of related studies, which differ in the terms included. For example, the surface area is commonly divided into polar and non-polar regions,

with different parameters for polar/polar, polar/non-polar and non-polar/non-polar interactions. Various statistical techniques have been used to derive the equation and various sources of data used to derive the function [Head *et al.* 1996; Böhm 1998; Eldridge *et al.* 1997]. One possible problem with such functions is that they are typically derived from ligands that bind very tightly to their receptor, whereas docking is increasingly used to identify ligands of only modest potency from a large database. For this particular problem, combining the results from more than one scoring function has been shown to give better results than just using individual scoring functions on their own, an approach referred to as 'consensus scoring' [Charifson *et al.* 1999].

12.7 Applications of 3D Database Searching and Docking

There are now a number of published studies that demonstrate the utility of 3D database searching in drug design, using both docking and pharmacophore searching. Kuntz's group has used the DOCK program against a number of targets, including HIV protease, DNA, thymidylate synthase and haemagglutinin [Kuntz 1992; Kuntz *et al.* 1994]. In each case, one or more inhibitors of modest potency were discovered. Information about hits from the first generation was then used to perform more exhaustive database searches to identify yet more potent compounds. The structures of some of the 'hits' were determined by X-ray crystallography, revealing that not all of the ligands bound in the same way as predicted by the docking algorithm. A degree of serendipity is still important even with automated docking methods. For this reason, it is important to assess the performance of any new docking methods against as many experimentally determined protein–ligand complexes as possible. The much larger number of X-ray structures now available means that it is possible to choose at least one hundred ligands that vary in size, shape, flexibility and functionality (charged, polar, hydrophobic) and which dock into many different proteins. Two good examples of the kind of analysis that is now possible are those evaluating the GOLD program [Jones *et al.* 1997] and the FlexX program [Kramer *et al.* 1999]. GOLD uses a genetic algorithm, whereas FlexX uses an incremental construction method. There were some differences between the way in which each program was assessed, the most obvious approach being to calculate the RMS deviation between the theoretical and experimental structures, although this can sometimes be a rather simplistic and sometimes misleading metric. However, the best docking programs are able to get 'close' to the correct result for approximately 70% of the ligands.

Commercial 3D database systems for performing pharmacophore searches were available from the early 1990s, but it took several years for real applications to be reported in the literature, largely due to the confidential nature of many of the results. One example of a fairly typical study is that of Marriott and colleagues, who were looking for new lead molecules active against the muscarinic M₃ receptor [Marriott *et al.* 1999]. Antagonists of this particular receptor have potential therapeutic value in conditions such as irritable bowel syndrome, chronic obstructive airway disease and urinary incontinence. Three active molecules were used to define a series of 3D pharmacophores (using the clique detection method). The initial list of five pharmacophores was pruned following visual

examination to give two very similar pharmacophores containing a positively charged amine, a hydrogen-bond acceptor atom and two hydrogen-bond donor sites. Searching a 3D database and combining the selected molecules gave 172, which were tested. Three compounds were found to have significant activity in the assay, one of which proved to be of particular interest, being a simple molecule particularly amenable to lead optimisation.

12.8 Molecular Similarity and Similarity Searching

Substructure and 3D pharmacophore searching involve the specification of a precise query, which is then used to search a database in order to identify molecules for screening. In such an approach, either a molecule matches the query or it does not. Similarity searching offers a complementary approach, in that the query is typically an entire molecule. This query molecule is compared to all molecules in the database and a similarity coefficient calculated. The top-scoring database molecules (based on the similarity coefficient) are the 'hits' from the search. In a typical scenario the query molecule would be known to possess some desirable activity and the objective would be to identify molecules which will hopefully show the same activity. We therefore require some method for deciding how to compute the similarity between two molecules. In order to achieve this we need to choose a set of *molecular descriptors* for the compounds. These descriptors are then used to compute the similarity coefficient.

12.9 Molecular Descriptors

The descriptors that we will consider in this section are those which can be calculated readily from the molecular formula, the molecular graph or from one or more computed 3D conformations. It is also possible to include experimentally determined descriptors, but this is often not feasible due to the unavailability of experimental data or the expense in acquiring it, especially where there are many molecules to consider. Indeed, the molecules may not yet have been synthesised! Some descriptors can be calculated very rapidly, one obvious example being the molecular weight. Other descriptors may be time-consuming to calculate, such as those derived from quantum mechanics. Some descriptors have an obvious experimental counterpart with which the calculation can be compared, such as a partition coefficient. Others are purely computational, such as a binary fingerprint. Some descriptors refer to properties of the whole molecule; others refer to the properties of individual atoms. New descriptors are being invented continually, each purporting to provide novel insights into the relationship between a molecule's structure and its properties. Some commonly used descriptors are shown in Table 12.2.

12.9.1 Partition Coefficients

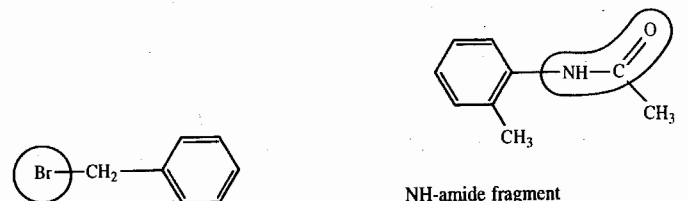
A very popular descriptor is $\log P$, the logarithm of the partition coefficient, most commonly for the partition between 1-octanol and water (the logarithm converts the

Descriptor	Information typically required for calculation	Comments
Molecular weight	Molecular formula	
Hashed fingerprints, structural keys	2D structure	See Section 12.2
Counts of specific atoms, rings or other features (e.g. hydrogen-bond donors, acceptors)	2D structure	Typically based on substructure searches
Octanol/water partition coefficient	2D structure	See Section 12.9.1
Calculated molar refractivity	2D structure	See Section 12.9.2
Molecular connectivity χ indices	2D structure	See Section 12.9.3
κ shape indices		See Section 12.9.3
Electrotopological indices		See Section 12.9.3
Atom pairs, topological torsions	2D structure	See Section 12.9.3
Dipole moment	3D structure	
Molecular volume, surface area, polar surface area	3D structure	Polar surface area is the amount of molecular surface due to polar atoms
Quantum mechanical descriptors (e.g. HOMO-LUMO energy gap)	3D structure	See Section 2.7.4
Partial atomic charge, polarisability	3D or 2D structure	See Section 4.9
Pharmacophore keys	Family of low-energy conformations	See Section 12.9.4
Geometric atom pairs, angles, torsions	3D structure	See Section 12.9.3

Table 12.2: A list of some of the more common descriptors. Details of some of these descriptors can be found elsewhere as indicated. This table is restricted to those descriptors which can be computed; it therefore excludes certain classes (such as the Hammett substituent constants) which are derived from experimental studies (see Section 12.12).

value onto a free-energy scale). Experimental determination of the partition coefficient can be difficult, particularly for zwitterionic and very lipophilic or polar compounds. The octanol/water system was selected by Hansch as a model system for hydrophobicity and has proved very successful, as we shall see in our discussion of quantitative structure-activity relationships in Section 12.12. Nevertheless, in some cases it would be more appropriate to measure the partition coefficient in an alternative system, and it is now possible to measure partition coefficients between lipid membranes and water directly.

Various theoretical methods can be used to calculate partition coefficients. The partition coefficient is an equilibrium constant and so is directly related to a free energy change. Partition coefficients can be calculated using free energy perturbation methods, as discussed in Section 11.3.2. These methods suffer from the limitations of force field parametrisation and the large amount of computer time that is required to perform such calculations. More widely used are fragment-based approaches, in which the partition coefficient is calculated as a sum of individual fragment contributions plus a set of correction factors. Such an approach clearly depends critically upon the definition of the fragments. The widely used CLOGP program of Hansch and Leo [Leo 1993] uses a small number of compounds to accurately define a set of fragment values. CLOGP breaks a molecule into fragments by identifying 'isolating carbons', which are carbon atoms that are not doubly



Bromide fragment	0.480	NH-amide fragment	-1.510
1 aliphatic isolating carbon	0.195	2 aliphatic isolating carbons	0.390
6 aromatic isolating carbons	0.780	6 aromatic isolating carbons	0.780
7 hydrogens on isolating carbons	1.589	10 hydrogens on isolating carbons	2.270
1 chain bond	-0.120	1 chain bond	-0.120
Total	2.924	1 benzyl bond	-0.150
		<i>ortho</i> substituent	-0.760
		Total	0.900

Fig. 12.24: CLOGP calculations on benzyl bromide and *o*-methyl acetanilide.

or triply bonded to a heteroatom. These carbon atoms and their attached hydrogens are considered hydrophobic fragments, with the remaining groups of atoms being the polar fragments. A partition coefficient is calculated by adding together appropriate values for the fragments and the isolating carbons, together with various correction factors. The process is illustrated in Figure 12.24 for two simple molecules, benzyl bromide and *o*-methyl acetanilide. Benzyl bromide contains one aliphatic isolating carbon and six isolating aromatic carbons, together with one bromide fragment. Each of these fragments contributes a characteristic score, to which are added values for the seven hydrogens on the isolating carbons and a contribution from one acyclic bond. *o*-Methyl acetanilide contains an amide fragment, two aliphatic isolating carbons and six isolating aromatic carbons. In addition, there are contributions from the hydrogen atoms, the acyclic bond, a benzyl bond (to the *o*-methyl group) and a factor due to the presence of an *ortho* substituent.

The CLOGP program remains the benchmark by which other methods for calculating octanol-water partition coefficients tend to be judged. One of its main drawbacks is the need for data for all the fragments in the molecule. Whilst the requisite data for a considerable number of fragments are included by default, the calculation will often not be correct for a fraction of the molecules in a typical pharmaceutical database. Of course, if the fragment is common to a molecular series of particular interest then it is usually straightforward to perform the experiment and add the necessary fragment value. An alternative is to use an atom-based approach to estimating the partition coefficient. This is very similar to the fragment-based method, but rather than checking for fragments the molecule is broken down into the atom types present. In the simplest case, the partition coefficient is given by a summation of the contributions from each atom type [Ghose and Crippen 1986; Ghose *et al.* 1998; Wang *et al.* 1997; Wildman and Crippen 1999]:

$$\log P = \sum_i n_i a_i \quad (12.4)$$

where n_i is the number of atoms of atom type i and a_i is the atomic contribution. These contributions are determined by regression analysis. The basic atomic contribution is in some cases moderated by correction factors to account for particular classes of molecules.

12.9.2 Molar Refractivity

The molar refractivity (MR) is given by:

$$MR = \frac{(n^2 - 1) MW}{(n^2 + 1) d} \quad (12.5)$$

In Equation (12.5) MW is the molecular weight, d is the density and n is the refractive index. The refractive index does not vary much from one organic compound to another and as the molecular weight divided by the density equals the volume, MR gives some indication of the steric bulk of a molecule. The presence of the refractive index term also provides a connection to the polarisability of the molecule. Molar refractivity can be calculated using atomic values with some correction factors for certain types of bonding (the CMR program [Leo and Weininger 1995]).

As we have seen, a common route to calculating both the partition coefficient and molar refractivity is by combining in some way the contributions from the fragments or atoms in the molecule. The fragment contributions are often determined using multiple linear regression, which will be discussed below (Section 12.12.2). Such an approach can be applied to many other properties, of which we shall mention only one other here, solubility. Klopman and colleagues were able to derive a regression model for predicting aqueous solubility based upon the presence of groups, most of which corresponded to a single atom in a specific hybridisation state but also included acid, ester and amide groups [Klopman *et al.* 1992]. This gave a reasonably general model that was able to predict the solubility of a test set within about 1.3 log units. A more specific model which contained more groups performed better but was of less generic applicability.

12.9.3 Topological Indices

Many of the descriptors which can be calculated from the 2D structure rely upon the molecular graph representation because of the need for rapid calculations. Kier and Hall have developed a large number of *topological indices*, each of which characterises the molecular structure as a single number [Hall and Kier 1991]. Every non-hydrogen atom in the molecule is characterised by two 'delta' values, the simple delta δ_i and the valence delta δ_i^v :

$$\delta_i = \sigma_i - h_i; \quad \delta_i^v = Z_i^v - h_i \quad (12.6)$$

where σ_i is the number of sigma electrons for atom i , h_i is the number of hydrogen atoms bonded to atom i and Z_i^v is the number of valence electrons for atom i . Thus the simple delta value will differentiate CH_3 from $-\text{CH}_2-$. CH_3 has the same simple delta value as NH_2 but a different valence delta. For elements beyond fluorine in the periodic table the

valence delta expression is modified as follows:

$$\delta_i^v = (Z_i^v - h_i)/(Z_i - Z_i^v - 1) \quad (12.7)$$

where Z_i is the atomic number. The *chi molecular connectivity indices* are obtained by summing functions of these delta values. Thus the chi index of order zero is defined as follows:

$${}^0\chi = \sum_{\text{atoms}} (\delta_i)^{-1/2}; \quad {}^0\chi^v = \sum_{\text{atoms}} (\delta_i^v)^{-1/2} \quad (12.8)$$

The summations are over the atoms in the molecule. This particular index does not encode much information about the structure. The first-order chi index involves a summation over bonds:

$${}^1\chi = \sum_{\text{bonds}} (\delta_i \delta_j)^{-1/2}; \quad {}^1\chi^v = \sum_{\text{bonds}} (\delta_i^v \delta_j^v)^{-1/2} \quad (12.9)$$

Higher-order chi indices involve summations over sequences of two, three, etc., bonds. To illustrate the difference between these indices for a series of related structures we show in Figure 12.25 the ${}^0\chi$, ${}^1\chi$ and ${}^2\chi$ indices for the isomers of hexane.

The *kappa shape indices* are generated by assessing how a molecular structure compares to molecular graphs with extreme shapes. As with the molecular connectivity indices, there


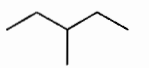
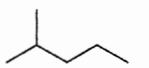
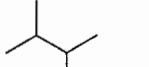
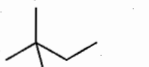
	Paths of length 2	Paths of length 3	Paths of length 4	Paths of length 5	${}^0\chi$	${}^1\chi$	${}^2\chi$
	4	3	2	1	4.828	2.914	1.707
	5	4	1	0	4.992	2.808	1.922
	5	3	2	0	4.992	2.770	2.183
	6	4	0	0	5.155	2.643	2.488
	7	3	0	0	5.207	2.561	2.914

Fig. 12.25: Chi indices for the various isomers of hexane. (Figure adapted in part from Hall L H and L B Kier 1991. *The Molecular Connectivity Chi Indexes and Kappa Shape Indexes in Structure-property Modeling*. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry Volume 2*. New York, VCH Publishers, pp. 367-422.)

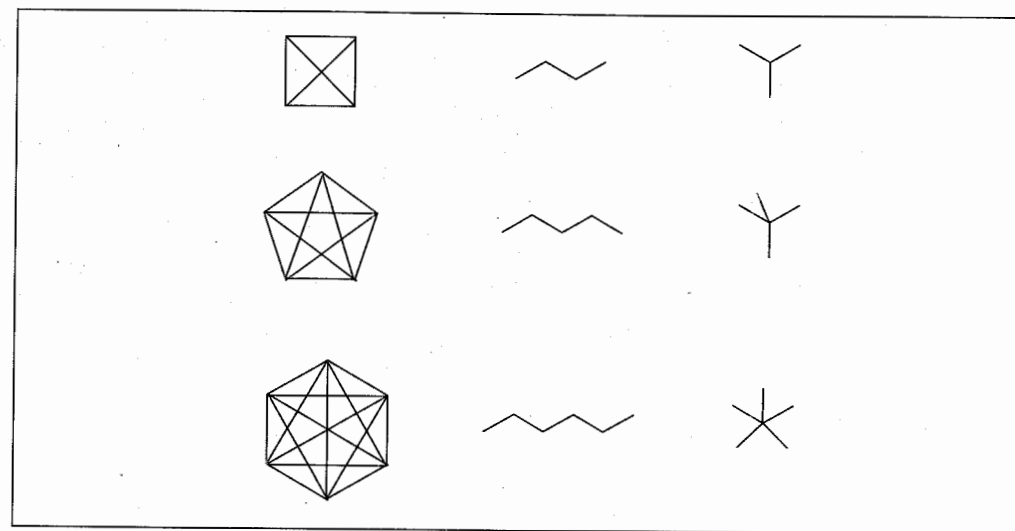


Fig. 12.26: First- and second-order extreme shapes for four-, five- and six-atom graphs (the linear molecule gives rise to the minimum in each case).

are shape indices of various order (first, second, etc.), with the first-order shape index involving a count over single-bond fragments, the second-order shape index involving a count of two-bond paths, and so on. In the first-order shape index the two extreme shapes are the linear molecule and the completely connected graph, where every atom is connected to every other atom (Figure 12.26). These two graphs contain $(A - 1)$ and $A(A - 1)/2$ bonds, respectively, where A is the number of atoms. If the number of bonds in our molecule is 1P then ${}^1P_{\max} \geq {}^1P \geq {}^1P_{\min}$, where ${}^1P_{\max}$ and ${}^1P_{\min}$ are the maximum and minimum number of bonds for that number of atoms. The first-order kappa index is written:

$${}^1\kappa = \frac{2{}^1P_{\max}{}^1P_{\min}}{({}^1P)^2} = \frac{A(A - 1)^2}{2({}^1P)^2} \quad (12.10)$$

The second-order kappa index is determined by the count of two-bond paths, written 2P . The maximum value is expressed by a 'star' shape (${}^2P_{\max} = (A - 1)(A - 2)/2$) and the minimum value again corresponds to the linear molecule (${}^2P_{\min} = A - 2$). The second-order shape index is then:

$${}^2\kappa = \frac{2{}^2P_{\max}{}^2P_{\min}}{({}^2P)^2} = \frac{(A - 1)(A - 2)^2}{2({}^2P)^2} \quad (12.11)$$

As with the molecular connectivity indices, higher-order shape indices have also been defined. The kappa indices themselves do not include any information about the identity of the atoms. This is the role of the 'kappa-alpha' indices. The alpha value for each atom is a measure of its size relative to some standard (chosen to be the sp^3 -hybridised carbon):

$$\alpha_x = \frac{r_x}{r_{Csp^3}} - 1 \quad (12.12)$$

An alpha value is calculated for the molecule by summing the individual atomic alphas and then incorporating them into the shape indices as follows:

$${}^1\kappa_\alpha = \frac{(A + \alpha)(A + \alpha - 1)^2}{2(1P + \alpha)^2} \quad (12.13)$$

$${}^2\kappa_\alpha = \frac{(A + \alpha - 1)(A + \alpha - 2)^2}{2(2P + \alpha)^2} \quad (12.14)$$

The final graph theoretical index popularised by Kier and Hall that we shall consider is the electrotopological state index [Hall *et al.* 1991]. Unlike the molecular connectivity and shape indices this is determined for each atom (including the hydrogen atoms, if so desired). This index depends upon the *intrinsic state* of an atom, which for an atom i (in the first row of the periodic table) is given by:

$$I_i = \frac{\delta_i^v + 1}{\delta_i} \quad (12.15)$$

This intrinsic state is a reflection of the electronic and topological characteristics of the atom i . The effects of interactions with the other atoms are incorporated by determining the number of bonds between the atom i and each of the other atoms, j . If this path length is r_{ij} then a perturbation is defined as:

$$\Delta I_i = \sum_j \frac{I_i - I_j}{r_{ij}^2} \quad (12.16)$$

The sum of ΔI_i and I_i gives the state of each atom (the E state). This descriptor is considered to encode the electronegativity of each atom (including the inductive effects of other atoms), together with its topological state. These atomic topological states can be combined into a whole-molecule descriptor by calculating the mean square value for the atoms. Vector or bit-string representations can also be produced. There is a finite number of possible I values, and so a bitstring representation can be obtained by setting the appropriate bit for each of the different I values present in a molecule. Alternatively, one can compute the sum or mean of the different E state values for each unique intrinsic state, to give a vector of real numbers.

Atom pairs [Carhart *et al.* 1985] and topological torsions [Nilakantan *et al.* 1987] are a related set of structural descriptors. Each atom pair descriptor codes the elemental type of a pair of atoms in the molecule together with the number of non-hydrogen atoms to which they are bonded, how many π -bonding electrons they have and the length of the shortest path between them. A topological torsion codes a sequence of four connected atoms together with their types, number of non-hydrogen connections and number of π electrons. Geometric atom pairs are a 3D equivalent, measuring the actual distance (in ångströms) between pairs of atoms, and likewise for other geometric descriptors.

12.9.4 Pharmacophore Keys

The development of 3D pharmacophore methods has spawned a new type of descriptor, the *pharmacophore key*. This is an extension of the binary keys used to facilitate 3D database

searching. When large sets of molecules are being considered then one is often restricted to properties that can be calculated relatively rapidly. This often precludes many of the properties dependent upon the 3D structure. Moreover, even when this 3D information is used it is often based upon a single conformation. The pharmacophore key can be computed relatively rapidly and it takes into account conformational flexibility and the pharmacophoric features in a molecule. In its simplest form, pharmacophores containing three features are represented. During the conformational analysis the pharmacophore features within each acceptable conformation are identified. All possible combinations of three features are enumerated, together with the distances between them (e.g. 'hydrogen-bond donor 6 Å from an aromatic ring centroid and 4 Å from a basic nitrogen with the third distance being 7 Å', Figure 12.27). Each distance is assigned to a distance bin, as described earlier. Every 3-point pharmacophore (distinguished by the features it contains and the distances)

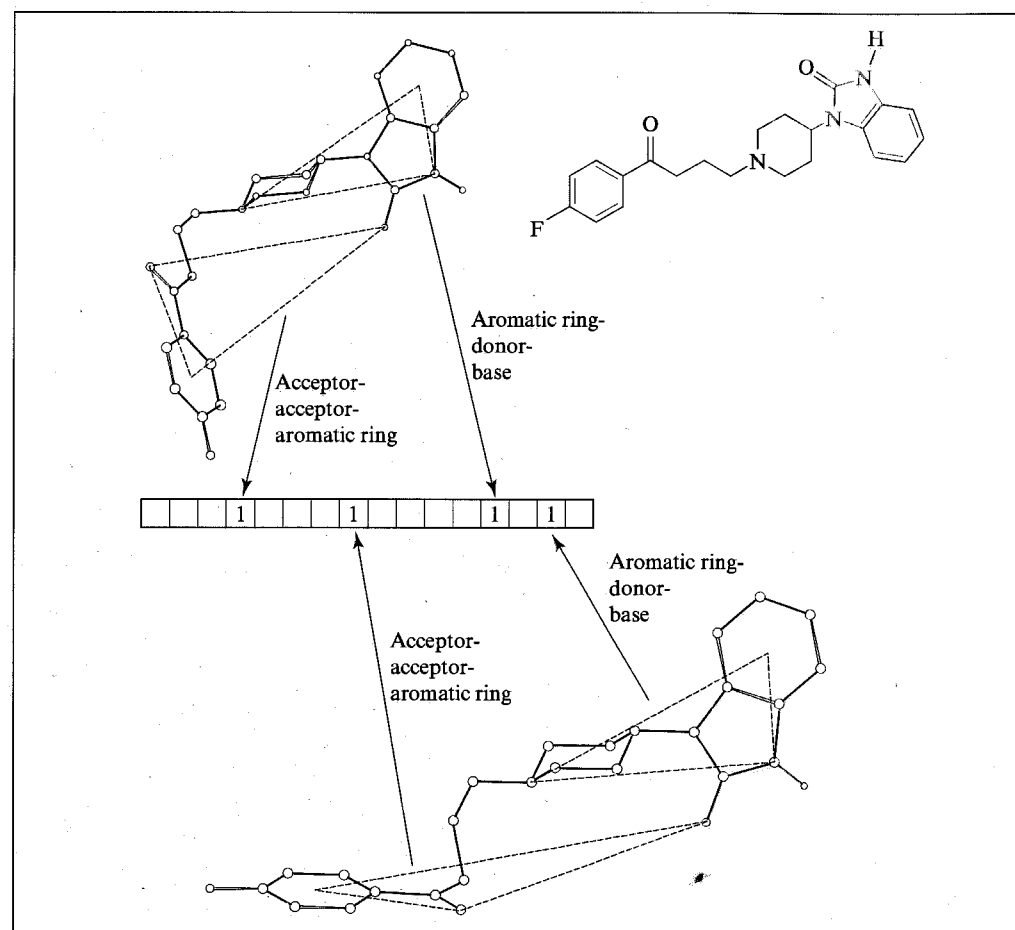


Fig. 12.27: The generation of 3-centre pharmacophore keys, illustrated using benperidol. Two different conformations are shown, together with two different combinations of three pharmacophore points.

is associated with a particular bit in the pharmacophore key bitstring. The pharmacophore key thus codes all possible 3-point pharmacophores that the molecule could express. These pharmacophore keys can be used in the same manner as any other binary descriptor. In addition to the use of 3-point pharmacophore keys [Good and Kuntz 1995; Pickett *et al.* 1996], 4-point pharmacophores are also possible [Mason *et al.* 1999]. These are claimed to contain more information and be more discriminatory than the 3-point keys (four points is required to differentiate stereoisomers, for example). However, the number of bits in a 4-point pharmacophore is considerably more than for the 3-point, and it is often not practical to store them simply as a single, large bitstring but as a sequence of integers, each of which identifies the bits that are set on.

12.9.5 Calculating the Similarity

The descriptors of a molecule can be considered a vector of attributes. These attributes may be real numbers or they may be binary in nature; in the case of the latter a value of 1 often indicates the presence of some feature and a value of 0 its absence. Having defined the descriptors, the next step is to compute a quantitative measure of the similarity [Willett *et al.* 1998]. Many similarity coefficients are in the range 0 to 1, with 1 indicating maximum similarity (note that this does not necessarily mean that the molecules are identical). Similarity is often considered to be complementary to distance, such that subtraction of the similarity coefficient from one gives the 'distance' between two molecules. Such distances may then be used in methods such as cluster analysis (see Section 9.13).

Here we will introduce three similarity coefficients that have been widely used for both real-valued (i.e. continuous) and binary (*dichotomous*) descriptors: the Tanimoto coefficient, the Dice coefficient and the Cosine coefficient. The formulae used to compute these coefficients are given in Table 12.3, where, for completeness, we have also provided the Euclidean and Hamming expressions that were introduced in Section 9.13. Different expressions are used for real-valued data (where the molecule is represented by a vector containing N real values x_i) and for binary data (where each molecule is represented by N binary values). For binary data, we additionally define a to be the number of bits 'on' in the bitstring for A, b to be the number of bits 'on' in the bitstring for B, and c to be the number of bits that are 'on' in both A and B (calculated using the AND operator).

Of the three similarity coefficients, the one most commonly used for binary molecular data (such as structural keys or hashed fingerprints) is the Tanimoto coefficient. It is important to recognise that there are some subtle differences between the way in which these metrics quantify the similarity between a series of compounds. These differences can be particularly marked for simple molecules which do not contain much functionality. Consider chlorpromazine and methoxypropazine, two phenothiazine neuroleptics (Figure 12.28). The Hamming 'distance' between these two molecules when calculated using the Daylight hashed fingerprints is 61 and the Soergel distance (equal to the complement of the Tanimoto coefficient for dichotomous data) is 0.28. These two molecules differ only by the substitution of a methoxy group for a chlorine atom. By contrast, the Hamming and Soergel distances between two smaller molecules that differ in the same way (methyl chloride and dimethyl

Name	Formula for continuous variables	Formula for binary (dichotomous) variables
Tanimoto similarity coefficient Also known as the Jaccard coefficient Complement equals the Soergel distance for dichotomous data	$S_{AB} = \frac{\sum_{i=1}^N x_{iA}x_{iB}}{\sum_{i=1}^N (x_{iA})^2 + \sum_{i=1}^N (x_{iB})^2 - \sum_{i=1}^N x_{iA}x_{iB}}$ Range: -0.333 to +1	$S_{AB} = \frac{c}{a + b - c}$ Range: 0 to 1
Dice similarity coefficient Also known as the Hodgkin index	$S_{AB} = \frac{2 \sum_{i=1}^N x_{iA}x_{iB}}{\sum_{i=1}^N (x_{iA})^2 + \sum_{i=1}^N (x_{iB})^2}$ Range: -1 to +1	$S_{AB} = \frac{2c}{a + b}$ Range: 0 to 1
Cosine similarity coefficient Also known as the Carbo index	$S_{AB} = \frac{\sum_{i=1}^N x_{iA}x_{iB}}{[\sum_{i=1}^N (x_{iA})^2 \sum_{i=1}^N (x_{iB})^2]^{1/2}}$ Range: -1 to +1	$S_{AB} = \frac{c}{(ab)^{1/2}}$ Range: 0 to 1
Euclidean distance	$D_{AB} = \left[\sum_{i=1}^N (x_{iA} - x_{iB})^2 \right]^{1/2}$ Range: 0 to ∞	$D_{AB} = [a + b - 2c]^{1/2}$ Range: 0 to N
Hamming distance Also known as the Manhattan or city-block distance	$D_{AB} = \sum_{i=1}^N x_{iA} - x_{iB} $ Range: 0 to ∞	$D_{AB} = a + b - 2c$ Range: 0 to N

Table 12.3: Formulae for various commonly used ways to compute the similarity or distance between molecules. For the binary data a is defined to be the number of bits 'on' in molecule A, b is the number of bits 'on' in molecule B and c is the number of bits that are 'on' in both A and B. Table based on [Willett *et al.* 1998].

thioether) are 16 and 0.80, respectively. The Hamming distance measure thus appears to suggest that the two smaller molecules are 'closer together' (more similar) than the pair of larger molecules. This contrasts with the Soergel/Tanimoto result. One reason for this difference is that the Hamming distance considers a common absence of features to indicate similarity, unlike the Soergel/Tanimoto measures. Moreover, the denominator in the Tanimoto coefficient has the effect of normalising the results according to the size of the molecule.

Another important feature of the Tanimoto coefficient when used with bitstring data is that small molecules, which tend to have fewer bits set, will have only a small number of bits in common and so can tend to give inherently low similarity values. This can be important when selecting 'dissimilar' compounds, as a bias towards small molecules can result.

A generalisation of the similarity formulae for binary data can be derived, based on the work of Tversky [Tversky 1977; Bradshaw 1997]. This takes the form:

$$S_{\text{Tversky}} = \frac{c}{\alpha(a - c) + \beta(b - c) + c} \quad (12.17)$$

where α and β are user-defined constants. The Tanimoto coefficient is recovered if $\alpha = \beta = 1$ and the Dice coefficient if $\alpha = \beta = \frac{1}{2}$. One interesting feature is that the Tversky coefficient is

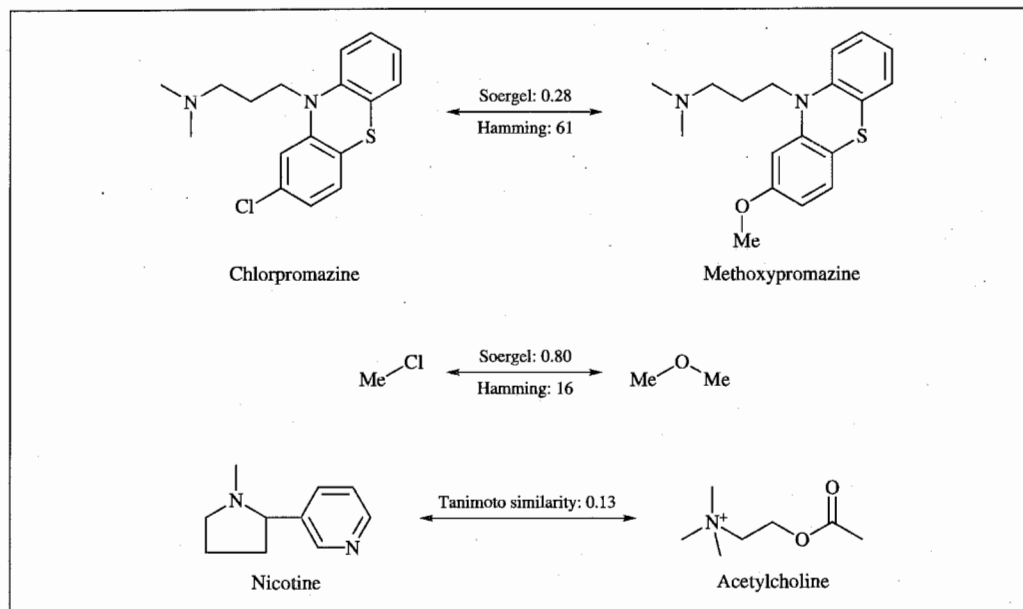


Fig. 12.28: The differences between the Soergel and Hamming distance measures for various molecules (see text).

asymmetric, such that $S_{Tversky}(A, B) \neq S_{Tversky}(B, A)$. If $\alpha = 1$ and $\beta = 0$ then the Tversky similarity value can be interpreted as being the fraction of the features in A which are also in B; a value of 1 with these parameters indicates that A is a 'substructure' of B.

12.9.6 Similarity Based on 3D Properties

Similarity methods based on rapidly calculated properties (particularly those derived from structural keys or hashed fingerprints) have become very popular, particularly for dealing with large numbers of molecules. Similarity measures derived from the 2D structure will tend to associate molecules with common substructures. However, molecular recognition depends on the three-dimensional structure and properties (e.g. electrostatics and shape) of a molecule rather than the underlying substructure. A simple illustration of this is provided by nicotine and acetylcholine (Figure 12.28), which have a very low similarity (0.13, as determined using the Tanimoto coefficient from the Daylight hashed fingerprints), despite the fact that they act at the same biological receptor. A simple 3D pharmacophore model to rationalise this comprises a positively charged or basic group an appropriate distance from an acceptor. For reasons such as this there has been much interest in similarity measures based upon three-dimensional properties.

Several measures of the shape and electronic similarity between pairs of molecules have been devised. The *Carbo index* was developed to enable the electron density of two molecules in some relative orientation to be compared [Carbo *et al.* 1980]. It is essentially the Cosine

coefficient (Table 12.3):

$$S_{AB} = \frac{\int \rho_A \rho_B d\nu}{(\int \rho_A^2 d\nu)^{1/2} (\int \rho_B^2 d\nu)^{1/2}} \quad (12.18)$$

The electron densities at each point are determined in the usual way from the square of the wavefunction. The value of the Carbo index runs from 0 (no similarity) to 1 (perfect similarity). Unfortunately, the electron density is not an ideal measure of similarity, because the density is strongest near the atomic nuclei and so the Carbo formula will be dominated by the extent to which the nuclei overlap. The electrostatic potential is a more appropriate property as it emphasises electronic effects away from the nuclei. Another drawback of the Carbo index is that it does not depend upon the magnitude of the property at a point but just its sign. This means that a location where the potential is positive from one molecule and equal but negative from the other would be weighted the same, irrespective of the magnitude of the potential. Hodgkin and Richards suggested the following alternative measure of similarity for use with the electrostatic potential [Hodgkin and Richards 1987]:

$$S_{AB} = \frac{2 \int \phi_A(\mathbf{r}) \phi_B(\mathbf{r}) d\mathbf{r}}{\int \phi_A^2(\mathbf{r}) d\mathbf{r} + \int \phi_B^2(\mathbf{r}) d\mathbf{r}} \quad (12.19)$$

A positive value of the Hodgkin-Richards index is obtained if large charges of the same sign are located in approximately the same regions of space; a negative value is obtained if large charges of opposite sign are located in the same regions of space. This index is effectively the Dice coefficient.

The integrals in the Hodgkin-Richards approach can be evaluated in a number of ways. One approach is to position the molecules within a rectangular grid and to evaluate the electrostatic potential due to each molecule at each grid point. The integrals in Equation (12.19) are then determined numerically by summing over the grid points. This can be rather slow, particularly if the molecules are allowed to vary their relative orientations and conformations in order to find the location of maximum similarity. An alternative is to represent the potential using an analytical function. For example, linear combinations of Gaussian functions can be fitted to the potential, enabling the similarity measure to be computed much more rapidly [Good *et al.* 1993].

3D similarity methods such as these provide the means to generate a structural alignment of molecules based upon some suitable property (electrostatics or shape). Methods such as pharmacophore mapping also provide a mechanism to align molecules, but in this case based upon their pharmacophore features. A structural alignment of a set of active molecules can be very useful in drug design, particularly when the structure of the target receptor is. Some of the varied techniques, such as 3D database searching and comparative molecular field analysis, which make use of structural alignments, are discussed in this chapter. The abundance of techniques for generating such alignments reflects the complex nature of the problem, in part a consequence of the need to consider the conformational flexibility of the molecules together with their relative orientations in space [Lemmen and Lengauer 2000].

12.10 Selecting 'Diverse' Sets of Compounds

Chemical diversity is a term which has been widely used and vigorously debated in the literature and at conferences since the advent of the era of high-throughput screening and combinatorial synthesis. But what is chemical diversity, and how can it be quantified? Given a finite number of screening slots or a finite number of compounds that could be synthesised, one could argue that it is desirable that the molecules are as diverse as possible, to maximise the chances of identifying one or more lead compounds. We therefore require methods which can be used to select diverse sets of compounds and techniques for comparing the diversity of one set against another. As there are $N!/k!(N-k)!$ possible ways to select a subset of k compounds from a total of N compounds it is obviously not feasible to examine all possible cases (for example, there are more than 10^{10} ways to select ten compounds out of just 50). Three popular ways to select 'diverse' sets of compounds are cluster analysis, dissimilarity-based methods and partition-based methods. Before applying any of these methods, however, it is usually recommended that the descriptors used to characterise the molecules are examined to determine whether some form of data manipulation is required.

12.10.1 Data Manipulation

Several tests and manipulations can be performed on a data set. For example, there is nothing gained by including a descriptor which shows no variation over the molecules in the set. It may also be useful to consider the distribution of values; some methods assume that the data are distributed according to a normal distribution, and significant deviations from this distribution will lead to invalid results. This is particularly the case with some of the methods used to derive quantitative structure-activity relationships. Two of the more common deviations from normality are skewness and kurtosis; the former indicates that the distribution is no longer symmetrical and the latter measures how 'peaked' the distribution is (Figure 12.29). Skewness and kurtosis are related to the third and fourth moments of the data (the mean and the variance corresponding to the first and second moments); these were encountered in our discussion of the moments theorem (see Section 4.23). Sometimes the distribution is bimodal, with two means. The *coefficient of variation* is a metric that can be used to identify descriptors which have a good spread over the range

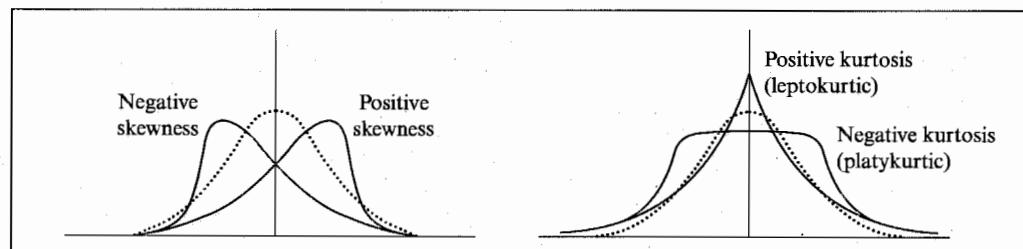


Fig. 12.29: Deviations from the normal distribution: skewness and kurtosis.

of the descriptor. This coefficient is equal to the standard deviation divided by the mean, and it automatically adjusts for different measurement scales. The larger the value the better the spread of values. Not all of the techniques described in this section need necessarily be used, but it certainly makes sense to perform some basic checks.

If the descriptors are on different scales then those which naturally occupy a 'larger' scale may be given more weight in the subsequent analysis, simply because of their natural units. In *autoscaling* the descriptors are scaled to zero mean and a standard deviation of 1.

$$x'_i = \frac{x_i - \bar{x}}{\sigma} \quad (12.20)$$

An alternative to autoscaling is range scaling, where the denominator in Equation (12.20) equals the range (the difference between the maximum and minimum values). Range scaling gives a set of new values between -0.5 and $+0.5$.

It is also important to check for correlations between the descriptors. Highly correlated descriptors could lead to the information that they encode being over-represented. A straightforward way to determine the degree of correlation between two properties is to calculate a correlation coefficient. Pearson's correlation coefficient is given by:

$$r = \frac{\sum_{i=1}^N (x_i - \langle x \rangle)(y_i - \langle y \rangle)}{\sqrt{[\sum_{i=1}^N (x_i - \langle x \rangle)^2][\sum_{i=1}^N (y_i - \langle y \rangle)^2]}} \quad (12.21)$$

where $\langle x \rangle$ and $\langle y \rangle$ are the arithmetic means of x and y , and N is the number of compounds. A value of 1.0 indicates a perfect positive correlation such that the x , y coordinates lie on a straight line with a positive slope. A value of -1.0 indicates a perfect negative correlation with the x , y coordinates on a straight line with a negative slope. A value of 0.0 indicates either no correlation or that the x , y coordinates follow a non-linear scatter.

Another way to identify correlations is to plot the values of the parameters in graphical form; this can help to identify any correlations and the presence of 'outliers'. A *Craig plot* is a two-dimensional scatterplot of one parameter against another; ideally, the molecules should sample from all four quadrants of the plot.

One way to try to alleviate the problem of correlated descriptors is to perform a principal components analysis (see Section 9.13). Those principal components which explain (say) 90% of the variance may be retained for the subsequent calculations. Alternatively, those principal components for which the associated eigenvalue exceeds unity may be chosen, or the principal components may be selected using more complex approaches based on cross-validation (see Section 12.12.3). It may be important to scale the descriptors (e.g. using autoscaling) prior to calculating the principal components. However, unless each principal component is largely associated with any particular descriptor it can be difficult to interpret the physical meaning of any subsequent results.

An alternative to principal components analysis is *factor analysis*. This is a technique which can identify multicollinearities in the set - these are descriptors which are correlated with a linear combination of two or more other descriptors. Factor analysis is related to (and

often confused with) principal components analysis. Factor analysis seeks to express each descriptor as a linear combination of factors. For example, if we have some descriptor x_i ; then each of the N values of this descriptor ($x_{i,j}$ where j runs from 1 to N , corresponding to the N molecules in the set) can be written in terms of the factors as follows:

$$x_{i,j} = a_i^1 F_j^1 + a_i^2 F_j^2 + a_i^3 F_j^3 + \dots + a_i^d F_j^d + E_{i,j} \quad (12.22)$$

This can be expressed in matrix form $X = FA^T + E$. $F_j^1, F_j^2, \dots, F_j^d$ are the d common factors; for each factor there is a value for each of the N data items. Note that to achieve a data reduction d should be less than the total number of descriptors in the set. $E_{i,j}$ is the unique factor specific to the descriptor x_i for molecule j . The coefficients $a_i^1, a_i^2, \dots, a_i^d$, etc., are known as the *loadings* of the descriptor values onto the common factors. Recall that in principal components analysis each principal component is expressed as a linear combination of the variables; factor analysis uses a similar linear expression but one in which the variables themselves are expressed in terms of the factors. To retrieve the value of a descriptor for a particular molecule one uses Equation (12.22) with the appropriate factor values corresponding to that molecule.

The unique factors are usually removed as they are considered to represent irrelevant information specific to the particular descriptors alone (such as experimental error). This leaves the common factors, which like the principal components are orthogonal. It sometimes happens that several descriptors may be loaded onto one or more factors. The factors are often 'rotated' in order to try to arrange for each factor to be largely associated with a few variables as possible. Thus if any given factor is largely associated with only one or two variables then it can be much easier to interpret any subsequent analysis. Rotation can also be applied to a set of principal components.

12.10.2 Selection of Diverse Sets Using Cluster Analysis

Cluster analysis was considered in our discussion of conformational analysis (see Section 9.13); for compound selection one would typically want to select a representative molecule or molecules from each cluster. A practical consideration when deciding which cluster analysis method to use is that for large numbers of molecules some algorithms may not be feasible because they require an excessive amount of memory or may have a long execution time. Another consideration with cluster analysis (and with some of the other methods that we will discuss) is the need to calculate the distance between each pair of molecules from the vector of descriptors (or from their scaled derivatives or from a set of principal components, if these are being used). For binary descriptors such as molecular fingerprints this distance is often given by $1 - S$, where S is the similarity coefficient (Table 12.3).

Two examples of the use of cluster analysis in compound selection are the studies of Downs, Willett and Fisanick [Downs *et al.* 1994] and Brown and Martin [Brown and Martin 1996]. In the former study, each molecule was described by 13 properties. The objective was to determine how well each of the different clustering methods considered was able to predict these property values. The property value for each molecule was

predicted as the mean of the property values for the other molecules in the cluster. The predicted values for each molecule were compared with the actual values to score each cluster method. Both hierarchical and non-hierarchical methods were used and different numbers of clusters were formed. As would be expected, the more clusters that were formed the more accurate the predictions (because the molecules in each cluster were more alike). However, a balance is required because at least one other molecule needs to be in the cluster for a prediction to be made. Of the methods considered, the hierarchical algorithms performed significantly better than the non-hierarchical Jarvis-Patrick method in terms of their predictive ability.

The Brown and Martin study considered a variety of clustering methods, together with several structural descriptors such as structural keys, fingerprints and pharmacophore keys. The methods were evaluated according to their ability to separate a set of molecules so that active and inactive compounds were in different clusters. Four different data sets were considered, the results suggesting that a combination of 2D descriptors (particularly the structural keys) and hierarchical clustering methods were most successful. This was a rather surprising result which caused much subsequent debate, not least because the structural keys were not particularly designed for use in compound selection but rather for fast substructural searching. In particular, there was much discussion concerning the nature of the data sets used in this experiment, which contained a rather higher proportion of structurally similar, active molecules than might be the case in a typical high-throughput screening scenario.

12.10.3 Dissimilarity-based Selection Methods

In dissimilarity-based compound selection the required subset of molecules is identified directly, using an appropriate measure of dissimilarity (often taken to be the complement of the similarity). This contrasts with the two-stage procedure in cluster analysis, where it is first necessary to group together the molecules and then decide which to select. Most methods for dissimilarity-based selection fall into one of two categories: maximum dissimilarity algorithms and sphere exclusion algorithms [Snarey *et al.* 1997].

The maximum dissimilarity algorithm works in an iterative manner; at each step one compound is selected from the database and added to the subset [Kennard and Stone 1969]. The compound selected is chosen to be the one most dissimilar to the current subset. There are many variants on this basic algorithm which differ in the way in which the first compound is chosen and how the dissimilarity is measured. Three possible choices for the initial compound are (a) select it at random, (b) choose the molecule which is 'most representative' (e.g. has the largest sum of similarities to the other molecules) or (c) choose the molecule which is 'most dissimilar' (e.g. has the smallest sum of similarities to the other molecules).

To decide which molecule to add at each iteration requires the dissimilarity values between each molecule remaining in the database and those already placed into the subset to be calculated. Again, this can be achieved in several ways. Snarey *et al.* investigated two common definitions, MaxSum and MaxMin. If there are m molecules in the subset then

the scores for a molecule i using these two measures are given by:

$$\text{MaxSum: score}_i = \sum_{j=1}^m D_{i,j} \quad (12.23)$$

$$\text{MaxMin: score}_i = \min(D_{i,j}; j=1,m) \quad (12.24)$$

where $D_{i,j}$ is the dissimilarity between two individual molecules i and j . The molecule i that has the largest value of score _{i} is the one chosen. A useful modification of these two methods is to reject any compound that is too close to one already chosen, typically assessed using the Tanimoto coefficient.

At each iteration of the sphere-exclusion algorithm [Hudson *et al.* 1996], a compound is selected for inclusion in the subset and then all other molecules in the database which have a dissimilarity to this compound less than some threshold value are removed from further consideration. Variation is possible depending upon the way in which the first compound is selected, the threshold value, and the way in which the 'next' compound is selected at each stage. It is typical to try to select this next compound so that it is 'least dissimilar' to those already selected. Hudson *et al.* suggested the use of a 'MinMax' method, where the molecule with the smallest maximum dissimilarity with the current subset is selected. However, it is also possible to select this 'next' compound at random from those still remaining.

The behaviour of some of these methods is illustrated using a two-dimensional example in Figure 12.30. If the 'most dissimilar' compound is chosen as the first molecule in the maximum-dissimilarity cases then the MaxSum method tends to select compounds at the extremities of the distribution. This is also the initial behaviour of the MaxMin approach, but it then starts to sample from the middle. The sphere exclusion methods typically start somewhere in the middle of the distribution and work outwards.

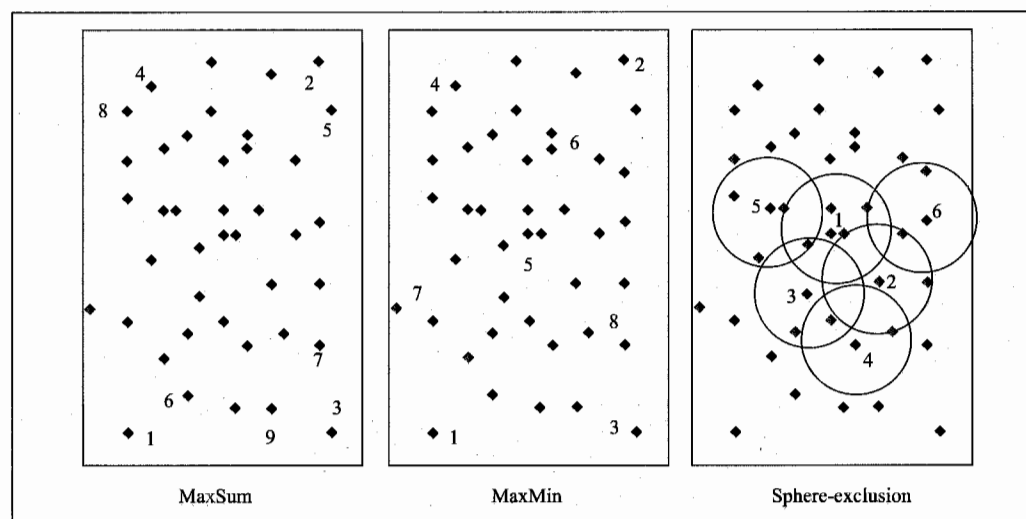


Fig. 12.30: Schematic comparison of various dissimilarity selection methods. The numbers indicate the order in which the molecules are selected.

To compare the behaviour of these various algorithms Snarey *et al.* performed a series of experiments using a set of compounds taken from the World Drug Index, a database that contains thousands of compounds, each with some biological activity. Every molecule in this database is assigned to one or more activity classes (e.g. antibiotics, antihistamines, analgesics). The objective of the exercise was to select subsets of compounds and determine the number of different activity classes. The more activity classes selected the 'better' the performance of the algorithm. In addition to investigating the various dissimilarity measures (MaxSum, MaxMin, etc.) there are a number of ways in which the dissimilarity between any pair of molecules can be quantified (D_{ij} in Equations (12.23) and (12.24)). The molecules were represented by hashed fingerprints or by a series of topological indices and physical properties, with the dissimilarity being given by the complement of either the Tanimoto or the Cosine coefficient. For the Cosine coefficient, a fast procedure is available which enables the dissimilarity between a molecule and a subset to be computed in a single operation rather than having to loop over all molecules currently in the subset [Holiday *et al.* 1995]. There was relatively little to choose between the best of these methods (though some proved significantly worse than a random selection), but the MaxMin maximum-dissimilarity algorithm was generally considered to be effective and efficient.

An alternative to the iterative procedures discussed thus far (where one compound at a time is added to the set) is to select the entire set as a single entity. A standard optimisation procedure such as a Monte Carlo search (often with simulated annealing) can be used to perform this [Hassan *et al.* 1996; Agrafiotis 1997]. A set is chosen at random, and an initial value for the diversity function (using an appropriate function, such as MaxMin) is calculated. At each iteration, a small number of compounds are replaced and a new diversity value computed. The change is accepted if the diversity has improved. If not, the Metropolis condition, $\exp[-\Delta E/k_B T]$, is applied (or, alternatively, the Felsenstein formula, $1/(1 + \exp[\Delta E/k_B T])$), which restricts the probability to less than 0.5 and so prevents the system just performing a random walk [Agrafiotis 1997].

12.10.4 Partition-based Methods for Compound Selection

A major potential drawback with cluster analysis and dissimilarity-based methods for selecting diverse compounds is that there is no easy way to quantify how 'completely' one has filled the available chemical space or to identify whether there are any 'holes'. This is a key advantage of the partition-based approaches (also known as cell-based methods). A number of axes are defined, each corresponding to a descriptor or some combination of descriptors. Each axis is divided into a number of 'bins'. If there are N axes and each is divided into b_i bins then the number of cells in the multidimensional space so created is:

$$\text{Number of cells} = \prod_{i=1}^N b_i \quad (12.25)$$

Each molecule is allocated to one cell according to its values along each axis. It is then a straightforward matter to select a 'representative' set of molecules; one just chooses one

(or more) molecules from each cell. The empty cells correspond to regions of the space not yet covered, which one might wish to target in order to increase the 'diversity' of the set.

The drawback with this is that a relatively low-dimensional space is required, as the number of cells increases exponentially with the number of dimensions, N . For this reason, it is not feasible to employ the binary descriptors that are commonly used to calculate intermolecular distances or dissimilarities (a 1024-long bitstring would contain 2^{1024} cells, an astronomically large number). It is therefore necessary to identify a reasonably low-dimensional space within which to work. This problem is nicely discussed by Lewis, Mason and McLay, who described the use of a partitioned space based upon a variety of molecular descriptors [Lewis *et al.* 1997]. The aim was to find a set of descriptors that would measure six key properties: hydrophobicity, polarity, shape, hydrogen-bonding properties and aromatic interactions. These properties were chosen because of their perceived importance in ligand-receptor interactions. A statistical analysis was performed to identify a set of six weakly correlated descriptors, each of which primarily measured one of the six properties. The partitions for each descriptor were chosen after plotting the distribution of each for approximately 47 000 molecules, after which two, three or four bins were chosen to give approximately equal areas of occupancies. Even such a relatively crude division gives nearly 600 partitions, which if the molecules were evenly distributed (which they are not) would provide about 80 molecules per cell. A representative set of compounds to act as a general-purpose screening set was then determined by selecting three representative molecules per cell, for a total of about 1000 compounds. One reason for the low occupancy of some cells is that they correspond to combinations of properties that are somewhat unlikely to exist (such as a very hydrophobic molecule with many hydrogen-bonding groups).

One approach to the problem of the exponential number of cells is to use principal components analysis or factor analysis to define a smaller set of orthogonal axes that are linear combinations of the original descriptors. This was the approach taken in a comparative study of various chemical databases [Cummins *et al.* 1996]. The initial descriptors comprised the computed free energy of solvation and a large set of topological indices. Descriptors with little variation, or which were highly correlated with another variable, were removed and then a factor analysis was performed. Four factors were able to explain 90% of the variation in the data. This four-dimensional space was partitioned into cells and the distribution of molecules from five databases was computed. Most of the molecules occupied a relatively small region of the space and so an iterative procedure was used to remove outliers, so enabling the resolution to be increased in the area populated by the majority of the molecules. Pairs of databases were compared by counting how many cells they had in common. Two of the databases contained only biologically active molecules and so it was of particular interest to identify the regions of space they occupied.

An alternative to the use of principal components or factor analysis is the BCUT method of Pearlman [Pearlman and Smith 1998]. In this method, three square matrices are constructed for each molecule. Each matrix is of a size equal to the number of atoms in the molecule and has as its elements various atomic and interatomic parameters. One matrix is intended to represent atomic charge properties, another represents atomic polarisabilities and the third hydrogen-bonding capabilities. These quantities can be computed with semi-empirical

quantum mechanical methods, but for large numbers of molecules more approximate methods are recommended. For each matrix, the highest and lowest eigenvalues are computed; for three matrices this gives six values per molecule, which form the axes for the partitioned space.

Finally, 3D pharmacophores can be used to provide a naturally partitioned space. By combining the pharmacophore keys of a set of molecules one can determine how many of the potential 3- or 4- point pharmacophores are accessible to the set and easily identify those which are not represented. This use of pharmacophores is the basis of a method named 'Pharmacophore-Derived Queries' (PDQ) [Pickett *et al.* 1996]. One feature of this particular method is that most molecules will occupy more than one 'cell' (as nearly all molecules will contain more than one 3-point pharmacophore due to the functionality present and conformational flexibility). This contrasts with the usual situation, wherein each molecule occupies just one cell.

12.11 Structure-based *De Novo* Ligand Design

Database searching is an attractive way to discover new lead compounds; in favourable cases the hits can be tested immediately or the molecule can be synthesised using a published method. However, database searching does not provide molecules that are structurally 'novel', although their new-found activity may be. Moreover, many databases are biased towards particular classes of compounds, so limiting the range of structures that can be obtained. In *de novo* design, the three-dimensional structure of the receptor or the 3D pharmacophore is used to design new molecules. There are two basic types of *de novo* design algorithm. The first class of methods have been described as 'outside-in' methods [Lewis and Leach 1994]. Here, the binding site is first analysed to determine where specific functional groups might bind tightly. These groups are connected together to give molecular skeletons, which are then converted into 'real' molecules. In the 'inside-out' approach, molecules are grown within the binding site under the control of an appropriate search algorithm, with each suggestion being evaluated using an energy function. These two approaches are compared in Figure 12.31.

12.11.1 Locating Favourable Positions of Molecular Fragments Within a Binding Site

One of the most widely used tools in structure-based ligand design is the GRID program [Goodford 1985]. A regular grid is superimposed upon the binding site. A probe group is then placed at the vertices of the grid and the interaction energy of the probe with the protein is determined using an empirical energy function. The result is a three-dimensional grid with an energy value at each vertex; this data can then be analysed to find those locations where it might be favourable to position a particular probe. An example of the output produced by GRID is shown in Figure 12.32 (colour plate section) for the binding site of neuraminidase. Parameters for many probes have been developed, covering a variety of small molecules and common functional groups. An alternative to the use of a grid is to

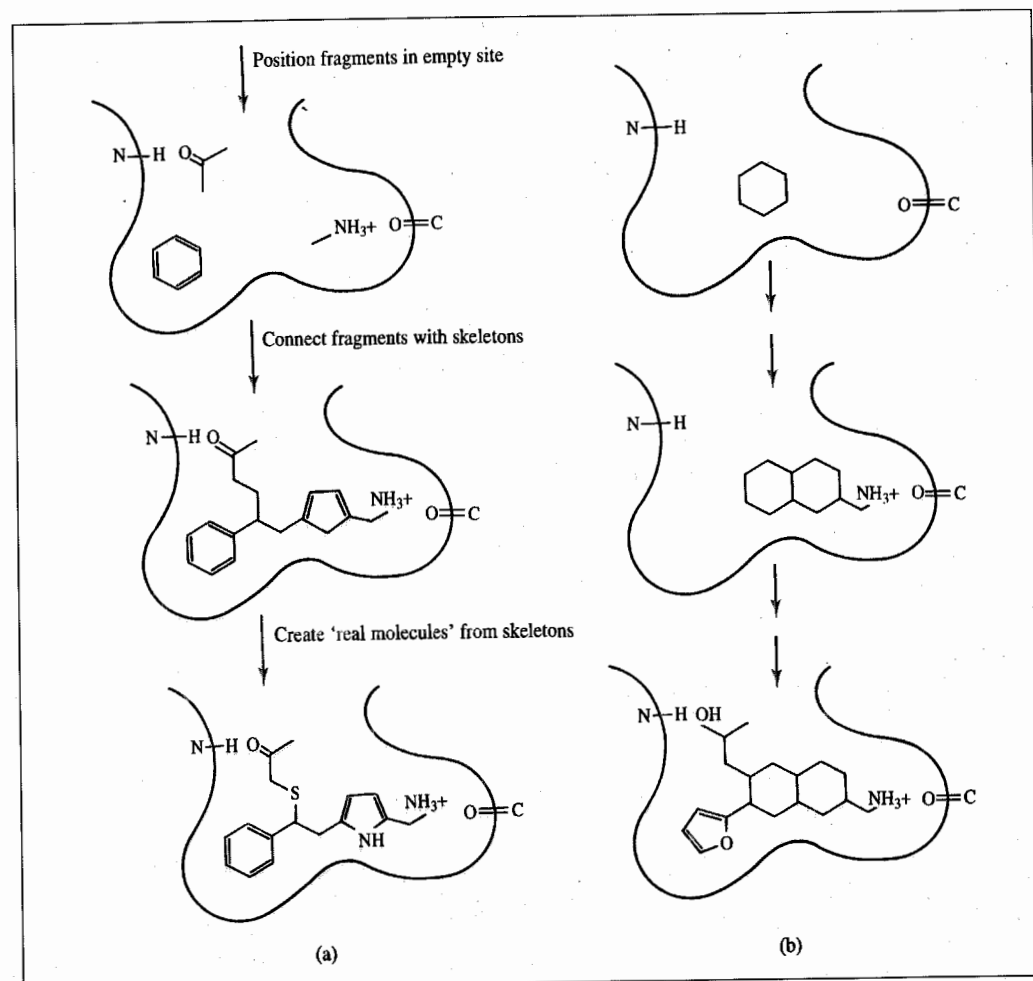


Fig. 12.31: Approaches to de novo design: (a) outside in, (b) inside out.

permit each fragment to explore the entire binding site using energy minimisation or some form of simulation method. In the multiple-copy simultaneous search (MCSS) approach [Miranker and Karplus 1991] the binding site is initially filled with many copies of the same fragment, distributed randomly. A molecular mechanics energy model is used in which the protein interacts simultaneously with all of the fragments, but there are no interactions between the individual fragments. Energy minimisation is then used to try to identify energetically favourable positions. The minimisation is done in a series of stages, with the orientations of the fragments being clustered at the end of each stage to remove duplicates.

An alternative to the energy-based methods such as GRID is to suggest possible binding positions using a knowledge-based approach. Analysis of experimentally determined

structures of ligand-receptor complexes reveals that they often contain certain types of interaction. For example, many ligands form hydrogen bonds with their receptors. The knowledge-based approaches generate binding modes that contain these commonly observed interactions, with the fragments being positioned to reproduce the most commonly observed geometries. For example, in most hydrogen bonds the distance between the donor hydrogen and its acceptor is close to 1.8 Å and the angle subtended at the hydrogen is rarely less than 120°. Information about the preferred geometries of such interactions can be obtained from analyses of X-ray crystallographic databases (described in Section 9.11). A program called LUDI has been widely used to dock small molecular fragments in protein binding sites using such an approach [Böhm 1992].

The knowledge-based docking approach to ligand design requires the receptor site to be surveyed to identify possible hydrogen-bonding donor and acceptor sites and regions where other groups might favourably be positioned. The results of such an analysis are often converted into a distribution of *site points*. A site point is a location within the binding site where an appropriate ligand atom or group could be placed. For example, a hydrogen-bonding analysis would typically result in a series of donor and acceptor site points. When generating the site points one should take account of any preferred geometries for that particular type of interaction. There is usually more than one site point associated with each donor or acceptor atom in the receptor to reflect the fact that a distribution of geometries is found in the crystal structure analyses. The range of preferred geometries can also be represented as a continuous region. Having surveyed the binding site, each molecular fragment is examined to determine which features it contains, and the fragment is positioned in the site by fitting the appropriate atoms to their corresponding site points.

The natural binding affinity of the small molecules typically used in a LUDI-type search is invariably rather low, and so it has been difficult to assess the accuracy of such computational techniques properly. However, it is now possible to use either X-ray crystallography or NMR to determine where such fragments bind to the protein. The technique of SAR-by-NMR in particular has attracted much attention [Shuker *et al.* 1996]. In SAR-by-NMR, a mixture of several highly soluble small molecules and a ¹⁵N-labelled protein is analysed using NMR. It is then possible to identify not only which small molecules bind but also where. Moreover, having identified one binder others can be found by repeating the assay with the protein together with the first 'hit', so identifying pairs of fragments that bind in a synergistic fashion.

12.11.2 Connecting Molecular Fragments in a Binding Site

Having positioned molecular fragments in favourable positions within the binding site using either an energy scheme or a knowledge-based approach, the next stage is to connect the fragments together into 'real' molecules. One way to tackle this problem is by searching a database of molecular connectors. Bartlett's CAVEAT program was one of the first methods for tackling this problem [Lauri and Bartlett 1994]. The relationship between each pair of fragments can be considered in terms of two vectors that represent the bonds they make

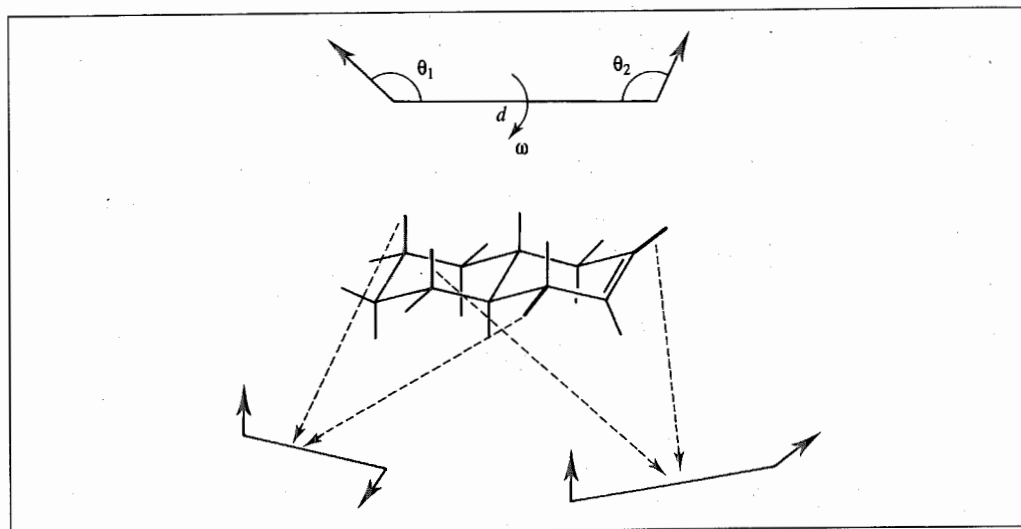


Fig. 12.33: The relationship between two bond vectors can be represented using a distance, two angles and a torsion angle as indicated (top). To derive the data for the database all possible pairs of exocyclic vectors are considered and these four geometric parameters calculated.

to the linker. The geometrical relationship between each pair of linking vectors can be described using a distance, two angles and a dihedral angle as shown in Figure 12.33. CAVEAT searches its database to find connectors that also contain two bond vectors with the same geometric parameters. The data is stored in an efficient form in the database, and so the search is very rapid. The connectors used by the first versions of CAVEAT consisted of ring systems extracted from the Cambridge Structural Database. The four geometric parameters were calculated between all pairs of bond vectors exocyclic to the ring, as shown in Figure 12.33. Connectors can also be generated using structure-generation programs.

Not all problems are amenable to the database searching approach exemplified by CAVEAT; the molecular fragments may be further apart than the longest connector in the database. Moreover, with such an approach one is inherently restricted to the fragments contained in the database. An alternative strategy is to create a skeleton that connects the fragments. The skeleton can be grown one atom at a time or by joining molecular templates. The templates typically comprise rings and acyclic fragments commonly found in drug molecules [Gillett *et al.* 1993].

The final stage in the outside-in approach to *de novo* design is to create 'real' molecules from the molecular skeletons that are generated. This is the most difficult part of the procedure, as there are so many ways in which even a small set of atom types can be assigned to a skeleton. The objective is to produce molecules that can be synthesised relatively easily and that also enhance the binding of the ligand to the binding site. It is very difficult to incorporate the concept of 'ease of synthesis' into a computer program, though some attempts have been made [Myatt 1995].

As can be seen from Figure 12.31, the outside-in procedure breaks the problem into a number of distinct stages. The inside-out method grows a ligand within the site in one step. Such an approach has been used successfully to generate peptides in protein binding sites [Moon and Howe 1991]. In this case, ligands are constructed from templates that are low-energy conformations of amino acids. Both systematic and random search algorithms can be used to explore the space of possible combinations of templates. In the systematic search, all of the amino acid building blocks are added to the growing ligand in turn. Each new structure is checked to ensure that it does not interact unfavourably with the protein and that it does not contain any high-energy intramolecular interactions. A molecular mechanics energy is then calculated for the structure. It is impractical to keep all of the structures from one stage to the next due to the combinatorial explosion and so only the lowest-energy structures are retained for the next iteration. Alternatively, a Monte Carlo simulated annealing search can be performed in which the Metropolis criterion is used at each stage to decide whether to accept or reject a given structure, based upon its energy and that of its predecessor. Genetic algorithm approaches have also been used to explore the search space [Glen and Payne 1995]. The advantage of applying this method to peptides is that there is a defined way of connecting the building blocks together, and the synthesis of peptides is straightforward. It is more difficult to generate general 'organic' molecules.

12.11.3 Structure-based Design Methods to Design HIV-1 Protease Inhibitors

An impressive example of the application of structure-based methods was the design of an inhibitor of the HIV protease by a group of scientists at DuPont Merck [Lam *et al.* 1994]. This enzyme is crucial to the replication of the HIV virus, and inhibitors have been shown to have therapeutic value as components of anti-AIDS treatment regimes. The starting point for their work was a series of X-ray crystal structures of the enzyme with a number of inhibitors bound. Their objective was to discover potent, novel leads which were orally available. Many of the previously reported inhibitors of this enzyme possessed substantial peptide character, and so were biologically unstable, poorly absorbed and rapidly metabolised.

The X-ray structures of the HIV protease revealed several key features that were subsequently incorporated into the designed inhibitor. The enzyme is a dimer with C_2 symmetry. It is a member of the aspartyl protease family, with the two aspartate residues lying at the bottom of the active site. Many of the crystal structures contained a tetra-coordinated water molecule that accepted two hydrogen bonds from the backbone amide hydrogens of two isoleucine residues in the 'flaps' of the enzyme and donated two hydrogen bonds to the carbonyl oxygens of the inhibitor (Figure 12.34, colour plate section).

A flow chart showing the various phases leading to the final compound is reproduced in Figure 12.35. The first step was a 3D database search of a subset of the Cambridge Structural Database. The pharmacophore for this search comprised two hydrophobic groups and a

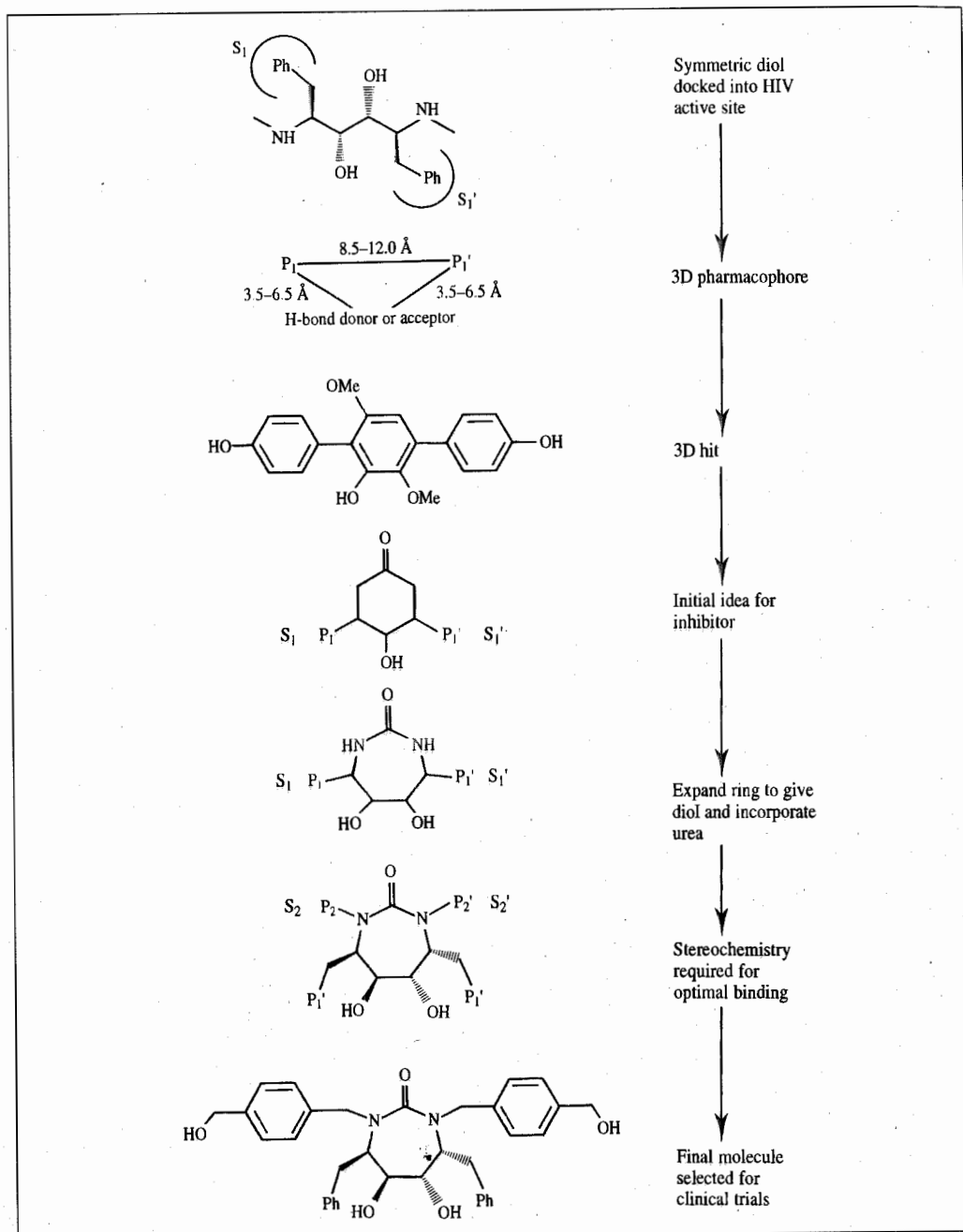


Fig. 12.35: Flow chart showing the design of novel orally active HIV-1 protease inhibitor. (Figure adapted from Lam P Y S, P K Jadhav, C E Eyermann, C N Hodge, Y Ru, L T Bachele, J L Meek, M J Otto, M M Rayner, Y N Wong, C-H Chang, P C Weber, D A Jackson, T R Sharpe and S Erickson-Viitanen 1994. *Rational Design of Potent, Bioavailable, Nonpeptide Cyclic Ureas as HIV Protease Inhibitors*. Science 263:380–384.)

hydrogen-bond donor or acceptor. The hydrophobic groups were intended to bind in two hydrophobic pockets (the S₁ and S₁' pockets) and the hydrogen-bond donor or acceptor to bind to the catalytic aspartate residues. The search yielded the hit shown in Figure 12.35. This molecule not only contained the desired elements of the pharmacophore but it also had an oxygen atom that could displace the bound water molecule. Displacement of the water was expected to be energetically favourable due to the increase in entropy. The benzene ring in the original compound was changed to a cyclohexanone, which was able to position the substituents in a more appropriate orientation.

The DuPont Merck group had previously explored a series of peptide-based diols that were potent inhibitors but with poor oral bioavailability. They were keen to retain the diol functionality, and so the next step was an expansion of the ring to a seven-membered diol. The ketone was then changed to a cyclic urea to strengthen the hydrogen bonds to the flaps and to aid the synthesis. Further modelling studies based upon the X-ray structure were performed to predict the optimal stereochemistry and the conformation required for optimal interaction with the enzyme. The results of these studies showed that the 4R, 5S, 6S, 7R configuration was most appropriate. Nitrogen substituents were predicted to bind to the S₂ and S₂' pockets of the enzyme, and so various analogues were synthesised in order to enhance the potency whilst maintaining the desired pharmacological properties. The compound eventually chosen for further studies, leading to clinical trials, was a *p*-hydroxymethylbenzyl derivative (Figure 12.35).

An increasing number of case histories of structure-based design have now been published in the literature, reflecting the widespread use of protein structures for drug design [Babine and Bender 1997; Kubinyi 1998]. All of the major pharmaceutical companies use structure-based design methods as part of their search for new drugs. Indeed, some smaller companies focus exclusively on structure-based design.

12.11.4 Structure-based Design of Templates for Zeolite Synthesis

Although *de novo* design is most commonly associated with the design of biologically active molecules, it is not restricted to this area. One particularly interesting application is in the design of templating agents for the synthesis of microporous materials, such as zeolites. These materials are very important for processes such as catalysis, ion exchange and gas separation. Zeolites can be considered to have the general formula TO₂, where T is a tetrahedrally coordinated atom (silicon, aluminium or phosphorus, depending on the constitution). The variable constitution can result in the T sites having an average oxidation state less than +4, giving a net negative charge. The charge imbalance can be rectified by the presence of Na⁺ ions. These then exchange with NH₄⁺ ions, which donate protons to the framework oxygen atoms to give ammonia, leaving the zeolite as a solid Brønsted acid catalyst.

The synthesis of zeolites is traditionally performed by crystallisation from a sol-gel mixture comprising reagents such as silica, sodium aluminate, sodium hydroxide and water. Another key component of the sol-gel mixture is a base whose main role is to regulate the pH of the mixture. If an organic base is used then a *templating* effect may also be observed

whereby the base acts to control the shape and size of the zeolite pores. There are many factors involved in the templating process, but qualitatively at least it is necessary that there is a good fit between the base and the framework such that the template fills the empty space in the zeolite cavity.

Traditionally, the templates were chosen by trial and error or exhaustive enumeration. A computational method named ZEBEDDE (ZEolites By Evolutionary De novo DEsign) has been developed to try to introduce some rationale into the selection of templates [Lewis *et al.* 1996; Willock *et al.* 1997]. The templates are grown within the zeolite cavity by an iterative 'inside-out' approach, starting from a seed molecule. At each iteration an action is randomly selected from a list that includes the addition of new atoms (from a library of fragments), random translation or rotation, random bond rotation, ring formation or energy minimisation of the template. A cost function based on the overlap of van der Waals spheres is used to control the growth of the template molecule:

$$f = \frac{\sum_{i=1}^{N_t} d(i, \text{host})}{N_t} \quad (12.26)$$

where $d(i, \text{host})$ is the closest contact distance between the template atom i and its nearest host atom. The function is normalised by the number of atoms in the template, N_t . In addition, the template is checked to ensure that there are no unfavourable intramolecular contacts or conflicts with any of its symmetrically related images (it is possible to have more than one template molecule per unit cell). The procedure continues so long as the cost function decreases and stops when it has fallen below a predefined value. Different template molecules can be compared by their ability to minimise the cost function.

The molecules typically used for templating agents are relatively simple and this is reflected in the fragment library (in the published examples just nine fragments were used: methane, ethane, ammonia, benzene, ammonium, propane, pyrrole, adamantane and cyclohexane [Willock *et al.* 1997]). The cost function does not have an electrostatic component and so some restrictions were imposed on the number of nitrogen atoms (no more than two per molecule and no N–N bonds). The type of molecule that can be produced is obviously governed by the template library but can also be biased by applying different weighting schemes for the fragment addition step. Fragments are typically added by replacement of a hydrogen atom in the template molecule. Giving a higher weight to the hydrogen atoms in new fragments tends to encourage the growth of linear chains, whereas a lower weighting tends to produce highly substituted template molecules. The method was applied to zeolites for which templates were already known, to check that these could be 'discovered' and also to determine what types of structure were produced by the various weighting schemes. In general, the known templates were found (or at least close analogues) together with some new possibilities. However, unfavourable conformations were observed in some of the suggestions, but these could sometimes be relieved by forming a ring. Other suggested molecules were not commercially available, despite the limited fragment library. For these targets as well as the biological ones the synthetic accessibility of *de novo* designed molecules is a major issue.

12.12 Quantitative Structure–Activity Relationships

A quantitative structure–activity relationship (QSAR) relates numerical properties of the molecular structure to its activity by a mathematical model. The term 'quantitative structure–property relationship' (QSPR) is also used, particularly when some property other than biological activity is concerned. In drug design, QSAR methods have often been used to consider qualities beyond *in vitro* potency. The most potent enzyme inhibitor is of little use as a drug if it cannot reach its target. The *in vivo* activity of a molecule is often a composite of many factors. A structure–activity study can help to decide which features of a molecule give rise to its overall activity and help to make modified compounds with enhanced properties. The relationship between these numerical properties and the activity is often described by an equation of the general form:

$$v = f(p) \quad (12.27)$$

where v is the activity in question, p are structure-derived properties of the molecule (i.e. descriptors), and f is some function. An early example of a structure–activity relationship was the discovery by Meyer and Overton of a correlation between the potencies of narcotics and the partition coefficient of the compounds between oil and water. Overton's interpretation was that the narcotic effect is due to physical changes caused by the dissolution of the drug in the lipid component of cells.

The first use of QSARs to rationalise biological activity is usually attributed to Hansch [Hansch 1969]. He developed equations which related biological activity to a molecule's electronic characteristics and hydrophobicity. For example:

$$\log(1/C) = k_1 \log P - k_2 (\log P)^2 + k_3 \sigma + k_4 \quad (12.28)$$

where C is the concentration of the compound required to produce a standard response in a given time, $\log P$ is the logarithm of the partition coefficient of the compound between 1-octanol and water, which was chosen by Hansch as a suitable measure of relative hydrophobicity, σ is the Hammett substituent parameter and k_1 – k_4 are constants.

The hydrophobic component was considered to model the ability of the drug to pass through cell membranes. Hansch recognised that there is an optimal value of the hydrophobicity: too low and the drug would not partition into the cell membrane; too high and the compound would partition into the membrane but tend to remain there rather than proceeding to the actual target. This explains the parabolic dependence of the activity upon $\log P$. An alternative way to express the Hansch equation uses a parameter π . This is the logarithm of the partition coefficient of a compound with substituent X relative to a parent compound in which the substituent is hydrogen:

$$\pi = \log(P_X/P_H) \quad (12.29)$$

Thus

$$\log(1/C) = k_1 \pi - k_2 \pi^2 + k_3 \sigma + k_4 \quad (12.30)$$

The Hammett substituent parameter was used by Hansch as a concise measure of the electronic characteristics of the molecules. Hammett and others (such as Taft) showed that

the positions of equilibrium and the reaction rates of series of related compounds such as substituted benzoic acids could be expressed in the following way:

$$\log \left(\frac{k}{k_0} \right) = \rho\sigma \quad \text{or} \quad \log \left(\frac{K}{K_0} \right) = \rho\sigma \quad (12.31)$$

where k_0 and K_0 are the rate constant and equilibrium constant, respectively, for a 'reference' compound (usually a hydrogen-substituted compound). The substituent parameter σ depends only upon the nature of the substituent and whether it is *meta* or *para* to the carboxyl group. The reaction constant ρ is fixed for a given process under specified experimental conditions. The 'standard' reaction is the dissociation of benzoic acids, which have $\rho = 1$. A full discussion of linear free-energy relationships can be found in many physical organic chemistry textbooks. A reading of such material will reveal that many modifications to the original Hammett scale have been suggested. One important development was the introduction by Swain and Lupton of the field (F) and resonance (R) components [Swain and Lupton 1968; Swain *et al.* 1983]. It was suggested that any set of σ values could be expressed as a weighted linear combination of these two components. This greatly simplified the problem of selecting the 'correct' set of substituent values for any individual system.

An enormous number of QSAR equations have been reported in the literature, many having a functional form much more complicated than the original Hansch equation. Many different parameters have been used in QSAR equations, designed to represent the hydrophobic, electronic or steric characteristics of the molecule. The properties chosen for inclusion in the QSAR equation should be as uncorrelated with each other as possible. Many of the early QSARs were derived for sets of related series of compounds that differed in just one part of the molecule. These differences can often be characterised using appropriate substituent constants, which were available in published tables. There has been a trend in recent years towards the analysis of noncongeneric series of compounds, for which there is no 'parent' structure and the consequent use of whole-molecule descriptors that are calculated directly. Some of these are given in Table 12.2 but there are many others. An example would be molecular shape analysis, which includes descriptors that measure the relative shape of the compounds [Rhyu *et al.* 1995]. A conformational analysis of the compounds is performed to identify the minimum energy structures. These conformations are then overlaid on the reference structure (usually one of the most active compounds in the series). From these overlaid structures it is then possible to calculate the common overlap volume and the non-overlap volume, which can then be included in the QSAR equation together with other parameters.

Another type of 'parameter' that often appears in published QSAR equations is an *indicator variable*. Indicator variables are used to extend a QSAR equation over a variety of different types of molecule and so make the equation more generally applicable. For example, Hansch and colleagues derived the following equation for the binding constants of sulphonamides ($X\text{-C}_6\text{H}_4\text{-SO}_2\text{NH}_2$) to human carbonic anhydrase [Hansch *et al.* 1985]:

$$\log K = 1.55\sigma + 0.64 \log P - 2.07I_1 - 3.28I_2 + 6.94 \quad (12.32)$$

I_1 takes the value 1 for *meta* substituents (0 for others) and I_2 is 1 for *ortho* substituents (0 for others).

12.12.1 Selecting the Compounds for a QSAR Analysis

The derivation of a QSAR equation involves a number of distinct stages. First, it is obviously necessary to synthesise the compounds and determine their biological activities. When planning which compounds to synthesise, it is important to cover the range of properties that may affect the activity. This means applying the data-checking and -manipulation procedures discussed earlier. For example, it would be unwise to make a series of compounds with almost identical partition coefficients if this is believed to be an important property.

Experimental design techniques can be used to help to decide which compounds to synthesise to ensure that the most information can be extracted from the smallest number of molecules. A variety of experimental design methods have been devised, of which the most straightforward to understand is probably full factorial design. Suppose there are two variables (formally called *factors*) that might influence the outcome (often called the *response*) of an experiment. If the experiment were a chemical synthesis, the factors might be temperature and the pH, with the response being the product yield. In our case, the experiment might involve the inhibition of an enzyme where the factors could be the molecule's $\log P$ and the Hammett substituent parameter. The response could be the degree of inhibition, measured as the IC50 (a commonly used measure of binding affinity, being the concentration of inhibitor to reduce the binding of a ligand or the rate of reaction by one-half). Determination of an IC50 requires less data than to determine the dissociation equilibrium constant but, unlike equilibrium constants, IC50 values measured under different conditions or for different receptors cannot usually be compared. Suppose, moreover, that we are interested in just two different values of each factor. Four possible experiments could be performed; in the case of the chemical synthesis these would be $T_1\text{pH}_1$, $T_2\text{pH}_1$, $T_1\text{pH}_2$, $T_2\text{pH}_2$, where T_1 and T_2 are the two temperatures and pH_1 and pH_2 are the two pH values. The first three experiments would measure the effect of changing just one variable at a time, whereas the fourth experiment ($T_2\text{pH}_2$) measures the effect of changing both variables and could indicate possible interactions between the factors. If there were three factors (with two values for each) then such a *full factorial design* would involve $2^3 = 8$ experiments. With three variables there is the possibility of interactions between pairs of factors, or between all three factors. In general, it is found that single factors on their own are more important than pairwise interactions, which are more important than three-factor effects, and so on. A *fractional factorial design* involves fewer experiments than the full factorial design. A half-fractional factorial design involves half the number of experiments as are in the full factorial design; a quarter-fractional design involves one-quarter of the experiments. However, it may be less easy to determine unambiguously the most important factors or combinations of factors from a fractional factorial design.

Factorial design methods cannot always be applied to QSAR-type studies. For example, it may not be practically possible to make any compounds at all with certain combinations of factor values (in contrast to the situation where the factors are physical properties such as temperature or pH, which can be easily varied). Under these circumstances, one would like to know which compounds from those that *are* available should be chosen to give a well-balanced set with a wide spread of values in the variable space. *D-optimal design* is one technique that can be used for such a selection. This technique chooses subsets of

molecules from those that are possible in such a way as to maximise the determinant of the variance-covariance matrix (also referred to as the 'information matrix'). If A is a matrix with n rows (corresponding to n molecules) and p columns (corresponding to the p descriptors) then the variance-covariance matrix is AA^T and is of size $n \times n$. D-optimal design aims to find the subset of n molecules that optimises the determinant of this matrix. A maximal value of the determinant corresponds to maximum variance and minimum covariance. In other words, the molecules have a wide spread of values for the descriptors (large variance) but a small degree of correlation between them (small covariance).

12.12.2 Deriving the QSAR Equation

The most widely used technique for deriving QSAR equations is *linear regression*, which uses least-squares fitting to find the 'best' combination of coefficients in the QSAR equation (the technique is also referred to as ordinary least-squares). We can illustrate the least-squares technique using the simple case where the activity is a function of just one property (when the technique is known as simple linear regression). We therefore want to derive an equation of the form:

$$y = mx + c \quad (12.33)$$

where y is known as the dependent variable (the observations) and x is the independent variable (the parameters). For example, y might be the activity and x might be $\log P$. The objective of a regression analysis is to find the coefficients m and c that minimise the sum of the deviations of the observations from the fitted equation, as shown in Figure 12.36. The least-squares coefficients m and c in the linear regression equation (12.33) are given by:

$$m = \frac{\sum_{i=1}^n (x_i - \langle x \rangle)(y_i - \langle y \rangle)}{\sum_{i=1}^n (x_i - \langle x \rangle)^2}; \quad c = \langle y \rangle - m\langle x \rangle \quad (12.34)$$

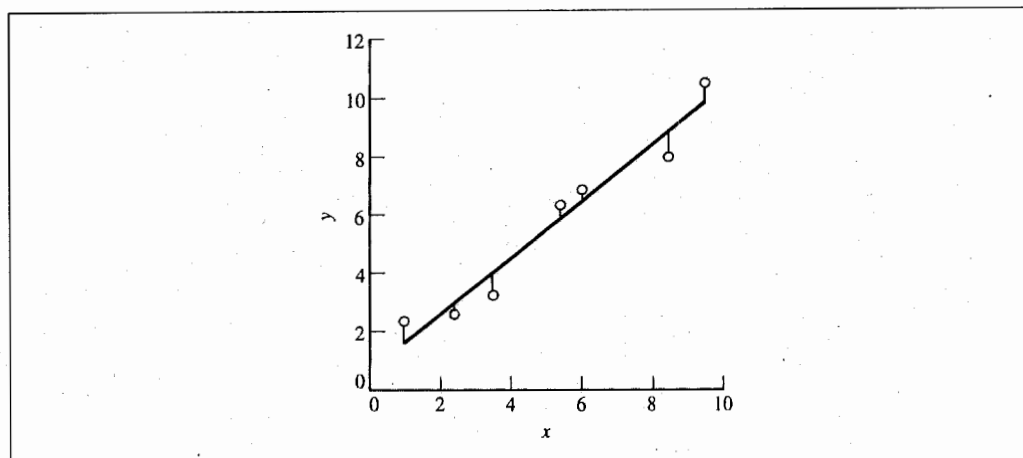


Fig. 12.36: The regression equation is the best-fit line through the data that minimises the sum of the deviations.

The regression equation passes through the point $(\langle x \rangle, \langle y \rangle)$, where $\langle x \rangle$ and $\langle y \rangle$ are the means of the dependent and independent variables, respectively. The 'quality' of a simple linear regression equation is often reported as the squared correlation coefficient, or r^2 value. This indicates the fraction of the total variation in the dependent variables that is explained by the regression equation. To determine r^2 , the total sum of squares (TSS) of the deviations of the observed y values from the mean $\langle y \rangle$ is calculated together with the explained sum of squares (ESS), which is the sum of squares of the deviations of the y values calculated from the model, $y_{\text{calc},i}$, from the mean:

$$\text{TSS} = \sum_{i=1}^N (y_i - \langle y \rangle)^2; \quad \text{ESS} = \sum_{i=1}^N (y_{\text{calc},i} - \langle y \rangle)^2; \quad \text{RSS} = \sum_{i=1}^N (y_i - y_{\text{calc},i})^2 \quad (12.35)$$

$y_{\text{calc},i}$ is obtained by feeding the appropriate x_i value into the regression equation. Another common squared term is the residual sum of squares (RSS), which is the sum of squares of the differences between the observed and calculated y values. TSS is equal to the sum of RSS and ESS. The r^2 is then given by:

$$r^2 = \frac{\text{ESS}}{\text{TSS}} \equiv \frac{\text{TSS} - \text{RSS}}{\text{TSS}} \equiv 1 - \frac{\text{RSS}}{\text{TSS}} \quad (12.36)$$

r^2 can adopt values between 0.0 and 1.0; a value of 0.0 indicates that none of the variation in the observations is explained by variation in the independent variables, whereas a value of 1.0 indicates that all of the variation in the observations can be explained. A disadvantage of the standard r^2 value is that it is dependent upon the number of independent variables, with higher r^2 values being obtained for larger data sets. More sophisticated statistical measures should ideally be used. These can, for example, help to determine whether the addition of a particular descriptor contributes significantly to the model.

It is straightforward to extend this analysis to more than one independent variable (known as *multiple linear regression*); such calculations are tedious to perform by hand but can be performed using a statistical package. Whereas simple linear regression involves fitting a straight line to the data, multiple linear regression corresponds to fitting a multidimensional surface. The quality of the regression is indicated by the multiple correlation coefficient, which we will write as R^2 (lowercase r^2 is also used). Another quantity that is commonly reported is the F statistic. This is the ratio of the explained mean square divided by the residual mean square. Values of F are available in statistical tables at different levels of confidence; if the calculated value is greater than the tabulated value then the equation is said to be significant at that particular level of confidence. It is important to note that the value of F depends upon the number of independent variables in the equation and the number of data points. As the number of data points increases and/or the number of independent variables falls so the value of F which corresponds to a particular confidence level also decreases. This is because we would like to be able to explain a large number of data points with an equation containing as few variables as necessary; such an equation would be expected to have greater predictive power. This is formally taken into account via the number of *degrees of freedom* associated with each parameter. A simple or multiple linear regression is associated with $N - 1$ degrees of freedom because the fitted line always passes through the means of the dependent and independent variables. The total

sum of squares is associated with $N - 1$ degrees of freedom. If there are p independent variables in the equation then there are $N - p - 1$ degrees of freedom associated with the residual sum of squares and p degrees of freedom associated with the explained sum of squares. Thus the explained mean square equals ESS divided by p , and the residual mean square equals RSS divided by $N - p - 1$, and so F is given by:

$$F = \frac{\text{ESS}}{p} \frac{N - p - 1}{\text{RSS}} \quad (12.37)$$

As to the significance of the individual terms in the equation, this can be assessed using the t statistic, which is obtained by dividing the relevant regression coefficient by the standard error of the coefficient. If k is the regression coefficient associated with a variable x , then the t statistic is obtained as follows:

$$t = \left| \frac{k}{s(k)} \right|; \quad s(k) = \sqrt{\frac{\text{RSS}}{N - p - 1} \frac{1}{\sum_{i=1}^N (x_i - \bar{x})^2}} \quad (12.38)$$

The value of t is compared with tabular values, which are listed according to the number of degrees of freedom associated with the residual sum of squares and for various significance levels. If the computed value is larger than the tabulated number then the coefficient can be considered significant.

There are some important criteria to consider when using multiple linear regression. To achieve statistically significant results there should be sufficient data; it is often considered that at least five compounds are required for each descriptor included in the regression analysis. The various checks on the data described in Section 12.10 should be performed to ensure that the selected compounds have a good spread of descriptor values, which should be as uncorrelated as possible. Compounds which have a value for some descriptor that is greatly different from the remainder (i.e. a significant outlier) should be examined very closely.

However, it is not sufficient to identify a set of non-correlated, well-distributed parameters and then simply 'press the button' and derive a multiple linear regression equation. To derive a QSAR equation properly requires a lot of care. Some of the descriptors may have little or no relevance to the property being modelled. Moreover, one generally wants to achieve a balance between an equation that captures the essence of the problem and yet is predictive. Fortunately, there are a number of procedures that can help with some of these problems.

One of the problems in deriving a QSAR equation is in deciding which descriptors to use. Two related procedures, forward-stepping regression and backward-stepping regression, can help in this. As the names imply, forward-stepping regression starts with an equation involving just one variable (the one that makes the most contribution, typically assessed using the t value). The second and subsequent terms are then added, again choosing the descriptor which makes the most contribution. Backward-stepping regression works in the reverse sense; initially, an equation is derived using all the descriptors, which are then removed (e.g. the one with the smallest t statistic). For both forward- and backward-stepping regression the equation finally chosen may be the one with the best fit to the data (as might be assessed using the F value).

Genetic algorithms can also be used to derive QSAR equations [Rogers and Hopfinger 1994]. The genetic algorithm is supplied with the compounds, their activities and information about their properties and other relevant descriptors. From this data, the genetic algorithm generates a population of linear regression models, each of which is then evaluated to give the fitness score. A new population of models is then derived using the usual genetic algorithm operators (see Section 9.9.1), with the parameters in the models being selected on the basis of the fitness. Unlike other methods, the genetic algorithm approach provides a family of models from which one can either select the model with the best score or generate an 'average' model.

12.12.3 Cross-validation

Cross-validation is a widely used and strongly recommended technique for checking the quality of a regression model. The technique (also known as jack-knifing) involves removing some of the values from the data set, deriving a regression model for the remainder and then predicting the values for the data left out. The most common form of cross-validation is 'leave-one-out', in which each data value is left out in turn and a model derived using the remainder of the data. A value can then be predicted for the data left out and compared with the true observed value. This is repeated for every data point in the set and permits the calculation of a 'cross-validated R^2 ' value (also written R_{cv}^2 or Q^2 , or their lowercase equivalents). Cross-validated R^2 values are typically lower than the normal R^2 values but are considered more indicative of the predictive ability of the equation. Indeed, Q^2 can have negative values (unlike R^2). Thus, whereas R^2 is a measure of goodness of fit, Q^2 is a measure of goodness of prediction. A more robust alternative to leave-one-out is to divide the data set into four or five groups, each of which is left out in turn to perform a cross-validation experiment. This process can be repeated many times (100 is typical) using different, randomly selected groups to obtain a mean Q^2 . The 'final' model (which might be used to predict the behaviour of as yet untested compounds) can then be derived in the normal way from all of the data; however, for a well-behaved system one would not expect the regression coefficients to change very much during the jack-knife procedure. Moreover, if the R^2 value from the whole data set is significantly larger than the mean Q^2 from the cross-validation experiment, then it is likely that the data has been over-fit. Another measure of predictive ability is the predictive residual sum of squares, PRESS, which is calculated in the same manner as the residual sum of squares except that, rather than using values $y_{calc,i}$, which are calculated from the model, we now use predicted y values $y_{pred,i}$, which are for data not used to derive the model. Q^2 is given by the following expression (compare with Equation (12.36)):

$$Q^2 = 1 - \frac{\text{PRESS}}{\sum_{i=1}^N (y_i - \langle y \rangle)^2}; \quad \text{PRESS} = \sum_{i=1}^N (y_{pred,i} - y_i)^2 \quad (12.39)$$

Strictly, the mean observed values, $\langle y \rangle$, which appear in Equation (12.39) should correspond to the mean of the values for each cross-validation group as appropriate rather than the overall mean value of the dependent variables, though often the mean of the entire data set will be used instead.

12.12.4 Interpreting a QSAR Equation

What does one do with a QSAR equation once it has been derived? An obvious use is for predicting the activities of as yet untested, and possibly not yet synthesised, molecules. The predictive ability of a QSAR is generally more accurate for interpolative predictions (i.e. for compounds that have parameter values within the range of those considered in the data set) than for extrapolative predictions (compounds that are outside the range). A QSAR equation may provide insights into the mechanism of the process being studied. As we have already noted, the presence of a parabolic relationship between the activity and the logarithm of the partition coefficient has been interpreted in terms of the transport of a compound to the receptor. An alternative model for transport is the *bilinear model*, in which the activity is related to the partition coefficient by an equation of the following form:

$$\log(1/C) = k_1 \log P - k_2(\log(\beta P + 1)) + k_3 \quad (12.40)$$

The bilinear model enables the ascending and descending parts of the function to have different slopes, whereas the parabolic equation is symmetrical. The parabolic model is generally most applicable to complex *in vivo* systems, where a drug must cross several barriers to reach its target, whereas the bilinear model often gives the best fit to the data for less complex *in vitro* systems.

Quantitative structure–activity relationships are often interpreted in terms of specific interactions with the macromolecular target. In a number of cases, the crystal structure of the ligand–receptor complex was subsequently determined and so it has been possible to use computer molecular graphics to discover whether the parameters in the QSAR equation have any real meaning [Hansch and Klein 1986]. For example, the presence of $\log P$ in the QSAR equation for the inhibition of carbonic anhydrase (Equation (12.32)) was interpreted as a hydrophobic interaction with the enzyme. The crystal structure of the enzyme revealed the presence of just such a hydrophobic surface along which a *para*-substituted group X could lie. The negative coefficients of the indicator variables for *meta* and *ortho* substitution also had a clear interpretation: such substituents would clash with the enzyme.

The absence of a correlation may also provide useful insights. For example, if one set of parameters gives a better correlation than another then this may indicate one particular mechanism is operating. If there is no correlation with a parameter (e.g. a steric measure) for a series of compounds then this may indicate that the associated property (i.e. steric volume) is of less importance.

12.12.5 Alternatives to Multiple Linear Regression: Discriminant Analysis, Neural Networks and Classification Methods

Whilst multiple linear regression is probably the most common technique used in QSAR and QSPR there are some other methods that have proved useful. One significant technique is partial least-squares, which is discussed separately in Section 12.13. Here we will describe a few of the other alternatives.

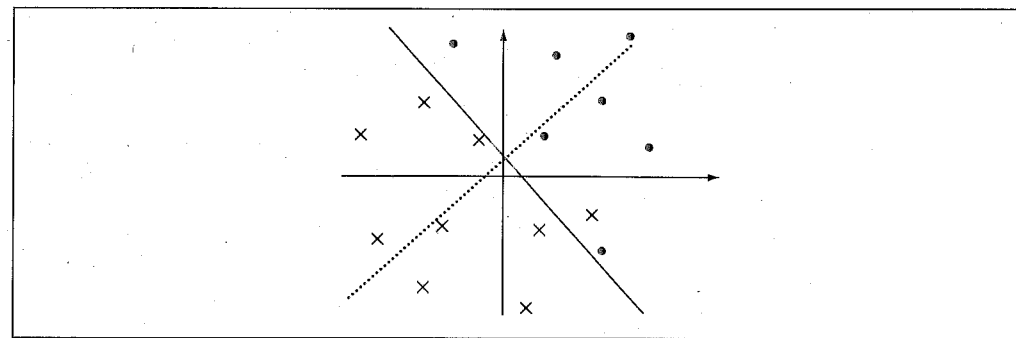


Fig. 12.37: Discriminant analysis defines a discriminant function (dotted line) and a discriminant surface (solid line).

Multiple linear regression is strictly a 'parametric supervised learning technique'. A parametric technique is one which assumes that the variables conform to some distribution (often the Gaussian distribution); the properties of the distribution are assumed in the underlying statistical method. A non-parametric technique does not rely upon the assumption of any particular distribution. A supervised learning method is one which uses information about the dependent variable to derive the model. An unsupervised learning method does not. Thus cluster analysis, principal components analysis and factor analysis are all examples of unsupervised learning techniques.

Discriminant analysis is a supervised learning technique which uses *classified dependent data*. Here, the dependent data (y values) are not on a continuous scale but are divided into distinct classes. There are often just two classes (e.g. active/inactive; soluble/not soluble; yes/no), but more than two is also possible (e.g. high/medium/low; 1/2/3/4). The simplest situation involves two variables and two classes, and the aim is to find a straight line that best separates the data into its classes (Figure 12.37). With more than two variables, the line becomes a hyperplane in the multidimensional variable space. Discriminant analysis is characterised by a *discriminant function*, which in the particular case of linear discriminant analysis (the most popular variant) is written as a linear combination of the independent variables:

$$W = c_1x_1 + c_2x_2 + \dots + c_Nx_N \quad (12.41)$$

The surface that actually separates the classes is orthogonal to this discriminant function, as shown in Figure 12.37, and is chosen to maximise the number of compounds correctly classified. To use the results of a discriminant analysis, one simply calculates the appropriate value of the discriminant function, from which the class can be determined.

Neural networks have been proposed as an alternative way to generate quantitative structure–activity relationships [Andrea and Kalayeh 1991]. A commonly used type of neural net contains layers of units with connections between all pairs of units in adjacent layers (Figure 12.38). Each unit is in a state represented by a real value between 0 and 1. The state of a unit is determined by the states of the units in the previous layer to which it is connected and the strengths of the weights on these connections. A neural net must first be trained to perform the desired task. To do this, the network is presented with a

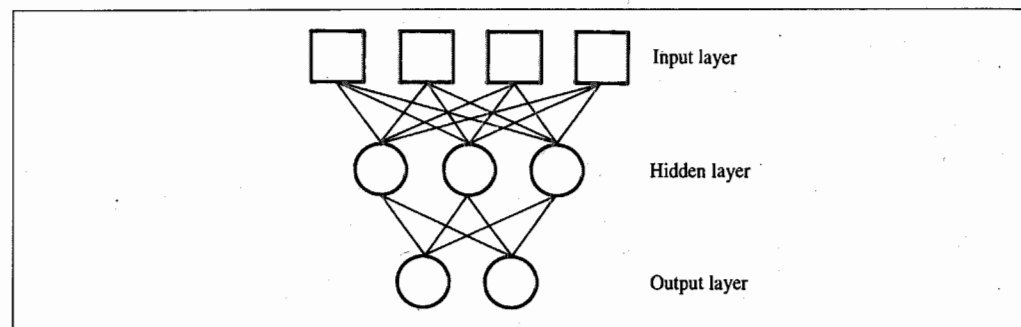


Fig. 12.38: Neural network with four input, three hidden and two output nodes.

set of sample inputs and outputs. Each input is fed along the connections to the nodes in the next layer, where they are operated upon and the results fed into the next layer, and so on. During the training period, the network adjusts the strengths of the connections using a method called back-propagation [Rumelhart *et al.* 1986] until it finds the set of values giving the best agreement between the input and output. Once trained, the net can then be used in a predictive fashion.

In QSAR, the inputs correspond to the value of the various parameters and the network is trained to reproduce the experimentally determined activities. Once trained, the activity of an unknown compound can be predicted by presenting the network with the relevant parameter values. Some encouraging results have been reported using neural networks, which have also been applied to a wide range of problems such as predicting the secondary structure of proteins and interpreting NMR spectra. One of their main advantages is an ability to incorporate non-linearity into the model. However, they do present some problems [Manallack *et al.* 1994]; for example, if there are too few data values then the network may simply 'memorise' the data and have no predictive capability. Moreover, it is difficult to assess the importance of the individual terms, and the networks can require a considerable time to train.

The output from a discriminant analysis or a neural network can often be difficult (if not impossible) to interpret in a manner that easily enables one to identify what features of a molecule give rise to the desired behaviour (or conversely, what features give rise to undesired behaviour!). This contrasts with a loosely associated group of methods that construct 'rules' which can be interpreted in terms of the association between specific molecular features and the activity. Various names are used to denote these methods, including classification trees, decision trees, regression trees, rule induction, machine learning and recursive partitioning. The output from one of these methods can be considered a tree-like structure. At each node there are usually two (but in some methods more than two) branches; which branch is followed for a particular molecule depends upon the rule associated with that node. In QSAR, each rule typically corresponds to the presence or absence of some structural feature or the value of some descriptor. An example is given in Figure 12.39 from a study on inotropic compounds (which increase the force of contraction of the heart without increasing its rate) [A-Razzak and Glen 1992]. Each molecule was

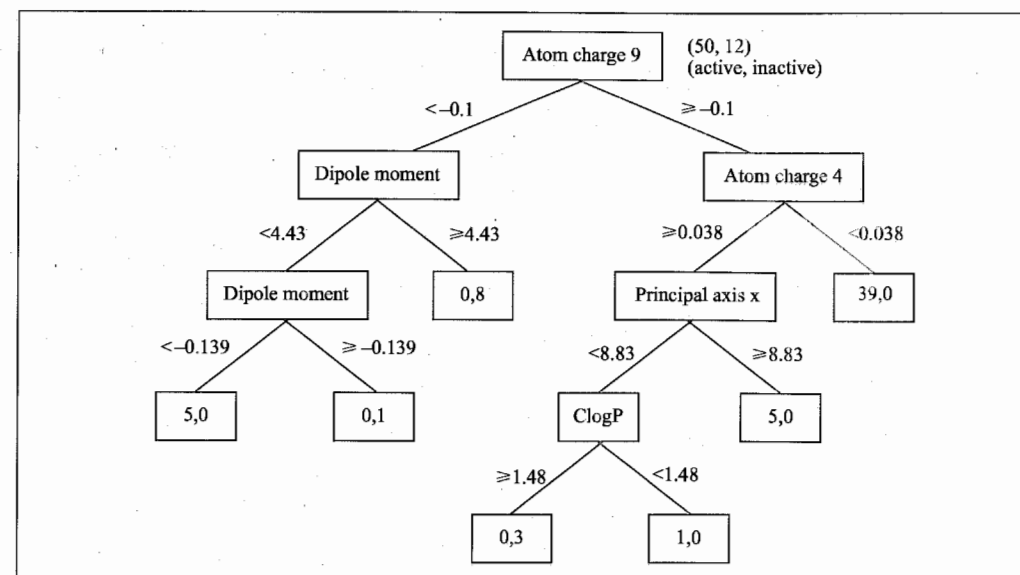


Fig. 12.39: Tree describing the 'rules' to differentiate active and inactive inotropic compounds. Each of the terminal nodes corresponds to the numbers of active and inactive molecules produced by the application of the preceding rules.

classified as active or inactive and was described by 44 descriptors. An algorithm called ID3 was used to decide how to construct the tree from the set of training data (42 compounds, in this case, with 20 compounds for testing). At each stage, the ID3 algorithm chooses to select the property that is 'most informative' from those not yet considered. This is placed on a mathematical footing using a powerful technique called information theory (effectively, it always tries to maximise the entropy gain). Other methods use statistical arguments to decide how to construct the tree. Recursive partitioning uses a statistical test (the *t* statistic) to identify the best descriptor to choose next. In this case, each terminal node in the tree is associated not just with a classification (e.g. active/inactive) but with an actual predicted activity. One particular implementation of recursive partitioning is able to handle large numbers of compounds, each of which is described by an extremely large number of descriptors [Rusinko *et al.* 1999]. This fast implementation is possible because of the binary nature of the descriptors (which are based upon the presence or absence of features such as atom pairs or topological torsions).

A third technique, which uses a somewhat different approach to the problem, is inductive logic programming (ILP) [King *et al.* 1992, 1996]. Initially, a large body of 'facts' about a series of both active and inactive molecules is created. At each stage of the subsequent procedure a machine learning algorithm then takes random pairs of molecules and determines what is common between them to produce a rule which is then evaluated to determine its effectiveness for predicting the remaining molecules. A typical rule can be expressed in a chemically meaningful form, such as 'molecule A is better than molecule B if B has no substitutions at positions 3 and 5 and A has no hydrogen-bond donors at position 3 and A has a π -donor at position 3 and A has a substituent at position 3 with fewer than three rotatable bonds'.

12.12.6 Principal Components Regression

Multiple linear regression cannot deal with data sets where the variables are highly correlated and/or where the number of variables exceeds the number of data values. Two methods are widely used to deal with such situations: principal components regression and partial least squares. In principal components regression, the variables are subjected to a principal components analysis (described in Section 9.13), and then regression analysis is performed using the first few principal components. When a principal components regression is performed using (say) forward-stepping regression then it will be found that the resulting equation is not necessarily expressed using just the lowest principal components. This is because the order of the principal components corresponds to their ability to explain the variance in the independent variables, whereas the regression analysis is concerned with explaining the dependent variable. A general rule of thumb is that only those principal components whose eigenvalues are greater than 1 should be considered for inclusion in a principal components regression. When an eigenvalue falls below 1 then one of the original variables in the set is more effective at explaining the variance than the principal component. Nevertheless, it is often the case that at least the first two principal components often give the best correlation with the dependent variable. Another interesting feature of principal components regression is that, as more principal components are incorporated, the regression coefficients of those already present do not change. This is due to the orthogonal nature of the principal components themselves and because the role of each new principal component is to explain variance not already covered.

12.13 Partial Least Squares

An alternative to principal components regression is to use the technique of *partial least squares* (PLS) [Wold 1982]. The PLS method expresses a dependent variable (y) in terms of linear combinations of the original independent (x) variables as follows:

$$y = b_1 t_1 + b_2 t_2 + b_3 t_3 + \dots + b_m t_m \quad (12.42)$$

where

$$t_1 = c_{11} x_1 + c_{12} x_2 + \dots + c_{1p} x_p \quad (12.43)$$

$$t_2 = c_{21} x_1 + c_{22} x_2 + \dots + c_{2p} x_p \quad (12.44)$$

$$t_m = c_{m1} x_1 + c_{m2} x_2 + \dots + c_{mp} x_p \quad (12.45)$$

t_1, t_2, \dots , are called *latent variables* (or components) and are constructed in such a way that they form an orthogonal set. The use of orthogonal linear combinations of the x values is very similar to principal components analysis. The major difference is that the latent variables in partial least squares are constructed to explain not only the variation in the independent variables x but also to simultaneously explain the variation in the observations, y .

We will illustrate the partial least squares method using a data set published by Dunn *et al.*, which provides the toxicity of a series of halogenated hydrocarbons together with eleven descriptor variables (see Table 12.4).

Compound	y LD ₂₅	x_1 MR	x_2 log P	x_3 BP	x_4 H_{vap}	x_5 MW	x_6 d_{20}	x_7 n_{20}^{Na}	x_8 q_{C}	x_9 q_{Cl}	x_{10} $E \text{ In C}$	x_{11} $E \text{ In Cl}$
1 CH ₂ Cl ₂	0.96	16.56	1.25	40.0	7.57	85	1.326	1.424	0.097	-0.1083	8.88	9.96
2 CF ₂ CHBrCl	1.31	23.54	2.30	50.0	7.11	197	1.484	1.448	0.1883	-0.1001	9.72	10.04
3 CHCl ₂	1.45	21.43	1.97	61.7	7.50	119	1.483	1.370	0.1805	-0.0870	9.69	10.16
4 CCl ₄	1.53	26.30	2.83	76.5	8.27	154	1.589	1.461	0.2662	-0.0666	10.55	10.36
5 Cl ₂ C=CHCl	2.26	26.05	2.29	86.5	8.01	131	1.465	1.456	0.1175	-0.0696	9.90	10.33
6 Cl ₂ C=CCl ₂	2.26	30.45	2.60	121.0	9.24	166	1.623	1.506	0.1360	-0.0680	10.08	10.34
7 CHCl ₂ CHCl ₂	2.42	30.92	2.66	146.0	9.92	168	1.587	1.494	0.1370	-0.1018	9.27	10.02

Table 12.4: Data on halogenated hydrocarbons [Dunn *et al.* 1984].

The parameters are as follows: LD₂₅: total toxicity measure; MR: molar refractivity; log P : logarithm of the partition coefficient; BP: boiling point; H_{vap} : latent enthalpy of vaporisation; MW: molecular weight; d_{20} : density at 20°C; n_{20}^{Na} : refractive index at 20°C measured using sodium light; $q_{\text{C}}, q_{\text{Cl}}$: charges on the chlorine-bearing carbon and the chlorine atom, respectively; $E \text{ In C}, E \text{ In Cl}$: orbital electronegativities of C, Cl.

The first seven variables (x_1, \dots, x_7) are standard global molecular descriptors, and the final four variables (x_8, \dots, x_{11}) are calculated measures of the electronic character of each molecule. Each of the seven compounds is thus described in terms of eleven variables. However, many of these variables are highly correlated with each other. For example, the charge on the chlorine is perfectly correlated with the electronegativity of the chlorine (a consequence of the way in which these two parameters were calculated, using the Gasteiger and Marsili method). Another strong correlation is that between the molar refractivity and the boiling point (correlation coefficient = 0.92).

A partial least-squares analysis [Malpass 1994] provides the weightings of the original variables in the latent variables. For example, the weightings for the first three latent variables are given in Table 12.5. These results suggest that all of the variables contribute to the first component, with the higher weightings being due to the molar refractivity (x_1), log P (x_2), boiling point (x_3), latent heat of vaporisation (x_4), d_{20} (x_6) and n_{20}^{Na} (x_7). The first latent variable thus represents a combination of steric, hydrophobic and electronic factors. The highest weightings in the second component are for the charge on the carbon (x_8) and the electronegativity of the carbon (x_{10}) and so this component has a higher contribution from electronic effects. Note that because partial least squares attempts to explain not only the variation in x but also in y , these weightings will differ from those obtained from a principal components analysis on the x variables alone. It is also possible to calculate the

Component	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}
1	0.4320	0.3197	0.4428	0.3875	0.1896	0.3265	0.3271	-0.1038	0.2133	0.1219	0.2105
2	-0.0850	0.2172	-0.2863	-0.1765	0.2833	0.2322	0.0479	0.7111	0.0936	0.4147	0.0982
3	0.1273	0.0985	0.0346	-0.3307	-0.1061	-0.1416	-0.5048	-0.2248	0.4841	0.2352	0.4853

Table 12.5: Weightings of the various parameters in the first three latent variables.

Component	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}	Total
1	97.2	81.5	78.8	67.7	39.8	86.1	66.6	3.6	32.6	30.2	32.1	74.1
2	0.2	14.6	16.4	21.2	12.8	6.8	5.4	84.5	18.4	57.3	19.2	12.8
3	0.3	1.1	0.5	2.2	22.5	0.9	5.7	3.4	45.9	11.4	45.5	4.7

Table 12.6: The degree to which each component explains the variance in each variable, and in the dependent variable (last column).

degree to which each component explains the variance in each variable, and how far each component explains the variation in the dependent variable (Table 12.6).

The results in Table 12.6 reinforce our earlier conclusion that the first component explains most of the steric and hydrophobic effects, with the second component explaining electronic effects. The first component explains 74.1% of the variation in the observed activity, with the first two components explaining a total of 86.9% of the variation.

In this illustration, we have only considered one dependent (y) variable, LD_{25} . In fact, partial least squares can deal with *multivariate* problems, where there is more than one dependent variable, such as different measures of biological activity or different properties. Indeed, for the above set of compounds five different measures of biological activity were reported in the original paper and a partial least-squares analysis performed on the entire data set [Dunn *et al.* 1984]. The algorithm effectively finds pairs of vectors through both the x data and the y data such that the vector pairs are maximally correlated with each other whilst simultaneously explaining as much of the variance in their individual data blocks as possible.

12.13.1 Partial Least Squares and Molecular Field Analysis

One of the most popular uses of the partial least-squares method in molecular modelling and drug design is comparative molecular field analysis (CoMFA), first described by Cramer and co-workers [Cramer *et al.* 1988]. The starting point for a CoMFA analysis is a set of conformations, one for each molecule in the set. Each conformation should be the presumed active structure of the molecule. The conformations must be overlaid in the proposed binding mode. The molecular fields surrounding each molecule are then calculated by placing appropriate probe groups at points on a regular lattice that encompasses the molecule, in a manner analogous to that used by the GRID program. The results of this analysis can be represented as a matrix, S , in which each row corresponds to one of the molecules and the columns are the energy values at the grid points (Figure 12.40). If there are N points in the grid and P probe groups are used, then there will be $N \times P$ such columns. The table is completed by adding a group additional column that contains the activity of the molecule. A correlation between the biological activity and the field values is then determined. The general form of the equation that we desire is:

$$\text{activity} = C + \sum_{i=1}^N \sum_{j=1}^P c_{ij} S_{ij} \quad (12.46)$$

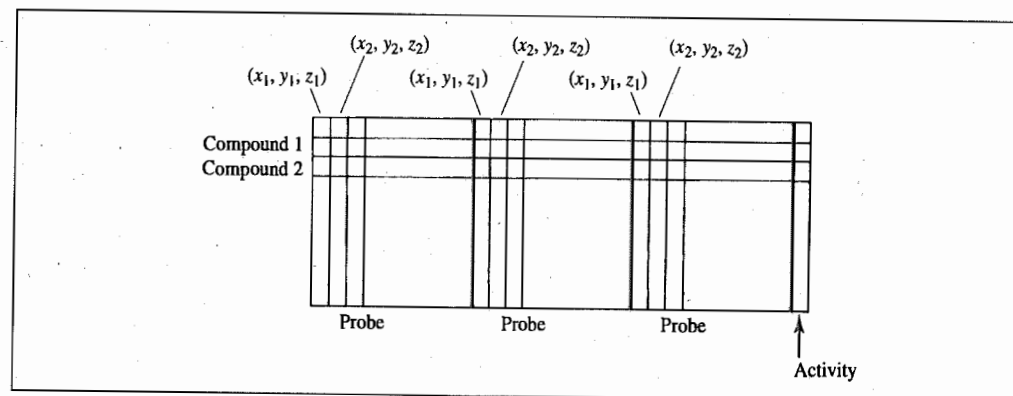


Fig. 12.40: The data structure used in a CoMFA analysis.

where c_{ij} is the coefficient for the column in the matrix that corresponds to placing probe group j at grid point i . As such, the problem is massively overdetermined as there may be thousands of grid points but often fewer than 30 compounds. Nevertheless, a successful analysis may often be performed using partial least squares.

The maximum number of latent variables is the smaller of the number of x values or the number of molecules. However, there is an optimum number of latent variables in the model beyond which the predictive ability of the model does not increase. A number of methods have been proposed to decide how many latent variables to use. One approach is to use a cross-validation method, which involves adding successive latent variables. Both leave-one-out and the group-based methods can be applied. As the number of latent variables increases, the cross-validated R^2 will first increase and then either reach a plateau or even decrease. Another parameter that can be used to choose the appropriate number of latent variables is the standard deviation of the error of the predictions, s_{PRESS} :

$$s_{\text{PRESS}} = \sqrt{\frac{\text{PRESS}}{N - c - 1}} \quad (12.47)$$

where c is the number of components in the current model. One would generally like to find the smallest number of latent variables which gives a reasonably high Q^2 such that each latent variable gives a fall in the value of s_{PRESS} of at least 5% [Wold *et al.* 1993]. Another measure of the predictive ability is SDEP, which is favoured by some practitioners:

$$\text{SDEP} = \sqrt{\frac{\text{PRESS}}{N}} \quad (12.48)$$

However, SDEP does not penalise an increase in the number of components; should the inclusion of an additional component slightly increase PRESS then the SDEP metric will select the model with more components, whereas s_{PRESS} will not. Bootstrapping is another procedure for assessing the stability of a PLS model, which attempts to overcome the all-too-familiar scenario of having fewer data values than might be ideal. In bootstrapping, N random selections are made from the original set several times in order to simulate different samplings from a larger set. Thus in each bootstrapping run some of the data would be

included more than once. This enables one to assess the variation in the different terms in the PLS model and so its stability.

An alternative way to assess the significance of a model is to randomly reassign the activities, thereby associating the 'wrong' activity with each set of grid values. When this is done, the predictive ability of the model should be significantly better for the true data set than for any of the randomised sets. This is a useful technique to check for random correlations when using descriptors that are not easy to interpret.

The CoMFA approach generates a coefficient for each column in the data table. This coefficient indicates the significance of each grid point in explaining the activity. Such data can usefully be represented as a three-dimensional surface that connects points having the same coefficients. These diagrams have been used to identify regions where (for example) changing the steric bulk would increase or decrease binding. An example is shown in Figure 12.41 (colour plate section). These contours can also be very useful for checking that a sensible model has been generated.

Since its introduction, partial least squares has been widely used to calculate such so-called '3D' QSARs. These studies have demonstrated its validity and usefulness but have also highlighted the sensitivity of the approach to several factors [Thibaut *et al.* 1993]. These factors include the selection of the active compounds, the different types of probe group that can be employed, the force-field models to describe the interactions between the probe and each compound, the size and spacing of points in the grid, and indeed the way in which the PLS analysis is performed. One of the main requirements (and indeed limitations) of the CoMFA technique is that it requires the structures of the molecules to be correctly overlaid in what is assumed to be the bioactive conformation (this in turn implies that the compounds have a common binding mode). The first application of CoMFA was to a series of steroid molecules binding to two different targets: human corticosteroid globulins and testosterone-binding globulins. In this case, the steroid nucleus of each molecule was least-squares fitted to the nucleus of the most active steroid. It can be more difficult to determine the appropriate binding mode in other cases, though the pharmacophore identification programs discussed in Section 12.4 may help with this problem. CoMFA can be particularly useful in the design of compounds that are selective for one target over another (related) target; a comparison of the contour maps can highlight regions where the two receptors have different requirements, which can be used to guide subsequent synthesis. It is also worth noting that although it is by far the most well-known approach, CoMFA is by no means the only 3D QSAR technique available [Greco *et al.* 1997].

The vast number of grid-field variables in a typical CoMFA analysis are obviously closely coupled; even the smallest structural change in the compounds will cause changes in not just one variable but in a group of variables that are connected in space. It is thus possible to envisage groups of spatially contiguous grid-field variables which are affected in the same way by structural variations in the compounds. Within such a group, all of the variables contain the same information. The use of such groups in the PLS analysis should give rise to 'better' models (i.e. greater predictive ability and enhanced interpretability). This is the basis of a method termed Smart Region Definition (SRD) [Pastor *et al.* 1997]. The same research group had previously developed automated procedures for selecting

subsets of variables in order to enhance the quality of PLS models [Cruciani *et al.* 1993]. In this latter GOLPE approach (GOLPE stands for Generating Optimal Linear PLS Estimations) multiple combinations of variables are selected using fractional factorial design. For each combination a PLS model is derived and only those variables which significantly influence the predictive ability of the model are retained.

There are three steps in the SRD approach. The first stage involves the selection of a set of grid nodes which have a high importance for the model (i.e. high weights in a 'traditional' CoMFA calculation). Each of these nodes is characterised by a point in 3D space around the molecules and by the particular field (e.g. electrostatic, van der Waals) to which they belong. These nodes act as seeds, with each of the remaining variables in the data set being assigned to its nearest seed (in a distance sense) or, if the distance to the nearest seed is greater than some cutoff, then the variable is removed from the analysis. It is important to note that the seeds are not uniformly distributed throughout the space; information-rich areas will have many seeds, whilst areas with less information will have fewer seeds. The spatial extent of the regions around the seeds in the information-rich areas will be correspondingly smaller. Those variables which are removed from the analysis usually correspond to areas far away from the compounds or where there is no chemical variation in the compounds. In the third step, the algorithm attempts to merge together neighbouring regions which contain the same information (i.e. are correlated), thus leading to an even greater information reduction. The SRD method was evaluated using a series of glycogen phosphorylase inhibitors, which are particularly relevant for such studies because the structure of each inhibitor bound to the enzyme had been determined by X-ray crystallography. As such, the problem of identifying the active conformation of each ligand and producing a molecular overlay disappeared. It was also possible to try to interpret the results of the 3D QSAR analysis in terms of specific interactions with the enzyme. In this case, the energy values were determined using the GRID program with a phenolic hydroxyl probe group (OH). Comparisons with the standard PLS method together with other procedures for grouping variables showed that the SRD algorithm did improve the fit, predictive ability and interpretability of the analysis.

The ability of partial least squares to cope with data sets containing very many x values is considered by its proponents to make it particularly suited to modern-day problems, where it is very easy to compute an extremely large number of descriptors for each compound (as in CoMFA). This contrasts with the traditional situation in QSAR, where it could be time-consuming to measure the required properties or where the analysis was restricted to traditional substituent constants.

12.14 Combinatorial Libraries

Combinatorial chemistry has significantly increased the numbers of molecules that can be synthesised in a modern chemical laboratory. The 'classic' approach to combinatorial synthesis involves the use of a solid support (e.g. polystyrene beads) together with a scheme called 'split-mix'. Solid-phase chemistry is particularly appealing because it permits excess reagent to be used, so ensuring that the reaction proceeds to completion. The excess

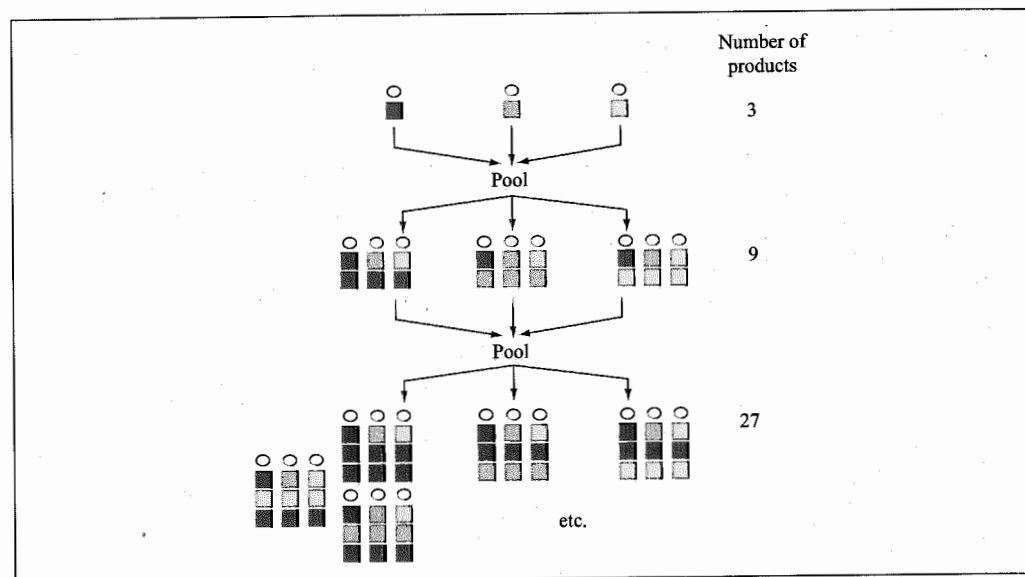


Fig. 12.42: Illustration of the split-mix approach to combinatorial synthesis, using sets containing three monomers.

reagent can then be simply washed away. The split-mix approach is illustrated schematically in Figure 12.42; initially, we start with separate pots of the first solid-supported reagent, A (there are n_A of these). These are mixed together and then divided into n_B equal amounts for reaction with the second set of reagents, B. The beads are now mixed together again and divided out for reaction with the third set of reagents, C. At this stage there should be $n_A \times n_B \times n_C$ products. The number of products grows exponentially with the number of reagents (hence the term 'combinatorial'). Due to the earlier mixing steps one only knows the identity of the final reagent for a given bead. However, it is important to realise that each bead contains just one discrete compound. The traditional split-mix approach has become less popular because there can still be much work involved in ascertaining the precise identity of any compounds that show activity. The use of 'tags' to encode information about the reagents (henceforth referred to as 'monomers') used to prepare the compound on any particular bead is one alternative.

The initial excitement about combinatorial chemistry was undoubtedly due to the number of molecules that could be synthesised, the assumption being that more molecules would surely lead to more hits in biological assays. However, on the whole this was not observed in practice, with combinatorial libraries often giving rise to fewer leads than historical sets of compounds. In part, this can be ascribed to the fact that as all the molecules in a library will have been synthesised using the same reaction scheme there is an inherent constraint on the amount of structural 'diversity' possible. However, it is also a reflection of the fact that many of the early combinatorial libraries did not (with the benefit of hindsight) contain molecules that were likely to have biological activity let alone be of interest as leads. Two general trends that have emerged recently are the move towards the synthesis of libraries that have been designed for activity against a particular biological target or group of related targets, such

as an enzyme family (focused library design) and towards the synthesis of libraries that contain more 'drug-like' molecules.

As was alluded to above, the early emphasis in combinatorial synthesis was on the generation of large numbers of 'diverse' molecules. Despite the general shift in emphasis to more 'focused' libraries, combinatorial methods do still offer a very powerful way to explore the range of chemical diversity. A useful way to reconcile this apparent conflict is to align the degree of diversity with the amount of knowledge about the molecular target for which the library is being synthesised. Thus, when one has a lot of knowledge about the target (for example, when an X-ray structure is available) then rather less diversity would be required than is the case when (for example) one only knows what general class the target belongs to on the basis of a sequence analysis. The difficulty is in quantifying these factors and the balance between them. An early attempt was described by Martin and colleagues, who used experimental design techniques to select monomers for peptoid libraries (peptoids are synthetic oligomers with a peptide backbone but with the side chain attached to the nitrogen atom rather than the alpha-carbon atom) [Martin *et al.* 1995]. A variety of descriptors were selected to represent features such as lipophilicity, shape and chemical functionality. These properties could then be displayed graphically for each monomer in a very intuitive fashion that greatly facilitated comparison.

12.14.1 The Design of 'Drug-like' Libraries

Solid-phase combinatorial synthesis has its roots in the work of Merrifield on the synthesis of peptides and so it is not surprising that many of the early libraries were made using this type of chemistry. However, peptides do not generally make very good drugs; not only are they readily broken down *in vivo* but they also tend to be rather large, high molecular weight compounds with many rotatable bonds when compared with typical drug molecules. Figure 12.43 shows the distribution of molecular weight, the number of rotatable bonds (a simple measure of conformational flexibility) and the calculated $\log P$ for a number of combinatorial libraries. As can be seen, for some libraries the distributions are quite close to those for the drug molecules, but for others there is a significant difference [Leach and Hann 2000]. What makes a molecule 'drug-like'? This is clearly an almost impossible question to answer, given the vast range of drug structures. Nevertheless, several attempts have been made to try to quantify 'drug-likeness'. The expectation is that libraries that are more drug-like will be more likely to contain molecules with biological activity and that any hits from such a library will represent more attractive starting points for lead optimisation. Of course, the need for drug-likeness is not restricted to combinatorial libraries but indeed can include any molecule that might be put through an assay, and so the techniques can also be applied to the selection of in-house or external compounds for screening.

Most practical implementations of drug-likeness use a computational model which takes as input the molecular structure, together with various properties, and predicts whether the molecule is drug-like or not. Some of these models may be very simple, such as a series of substructural filters. Only those molecules which pass all of these filters are output. Such filters can be used to eliminate molecules that contain inappropriate functionality.

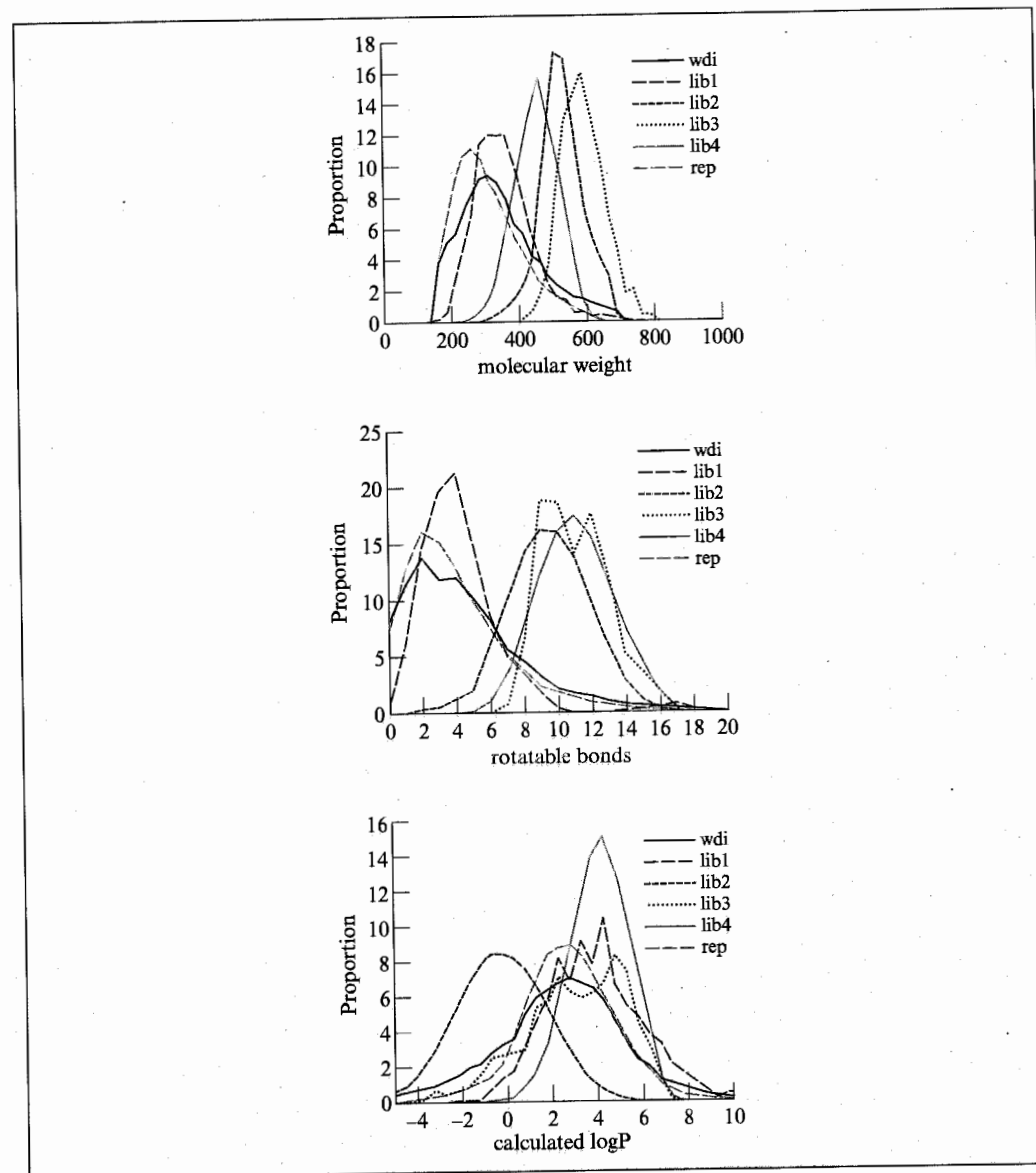


Fig. 12.43: Distribution of molecular weight, number of rotatable bonds and calculated partition coefficient for a number of early libraries synthesised at Glaxo Wellcome (lib1–lib4), together with the corresponding distributions for the World Drug Index (wdi) and a set of representative 'historical' compounds (rep).

For example, certain types of reactive group (e.g. alkyl halides, acid chlorides) will almost always give a positive reading in any biological assay. Similar filters can be used to eliminate molecules that have a high molecular weight or are very flexible. The 'rule of 5' is a set of empirical filters that can be used to suggest whether or not a molecule is likely to be

poorly absorbed [Lipinski *et al.* 1997]. Any molecule with a molecular weight greater than 500, a calculated $\log P$ greater than 5, more than 5 hydrogen-bond donors (defined as the sum of OH and NH) or more than 10 hydrogen-bond acceptors (defined as the sum of nitrogen and oxygen atoms in the molecule) is predicted to be poorly absorbed.

More sophisticated models are also possible. Such models may use neural networks [Sadowski and Kubinyi 1998; Ajay *et al.* 1998] or a regression-type equation with coefficients derived using a genetic algorithm [Gillet *et al.* 1998] to predict drug-likeness from a set of molecular properties. To train such models it is usually necessary to have a set of molecules that are considered drug-like and a set that is considered non-drug-like. The model is then optimised to achieve the optimal discrimination between the drug and non-drug sets.

Filters and drug-likeness models can be extremely valuable in library design and compound selection. They typically require just a 2D representation (such as a SMILES string) and so can be used to rapidly eliminate molecules of no interest and to score or rank the remainder (an example of *virtual screening*). However, it should always be remembered that such approaches are usually very general in nature and that any specific target may require molecules that violate one or more of these more general criteria.

12.14.2 Library Enumeration

The term 'enumeration' when applied to a combinatorial library refers to the process by which the connection tables for the product structures in a real or virtual library are produced. It should be noted that a single compound can be considered as a library of one and so enumeration can equally well be applied in this case. However, whereas it is considered reasonable for a chemist to draw the structure of a single compound manually (which may have taken days, if not months or years, to synthesise), it is clearly not practical to do so even for small combinatorial libraries. Hence the need for automated tools to perform this procedure.

Generally speaking, there are two different approaches to the enumeration problem. The first of these is often referred to as 'fragment marking'. In this method, a central core template, common to all product structures, is identified. The template will contain one or more points of variation where different substituents (often termed R groups) can be placed. By varying the R groups at the points of substitution, different product structures can be generated. In order to enumerate a combinatorial library, it is first necessary to construct sets of R group substituents from the relevant monomer sets. In the simplest cases, this is done by replacing the reactive functional group in the monomer with a 'free valence'. By creating a bond between the template and the required R groups the connection table for the product molecule can be generated. Enumeration of the full library corresponds to systematically generating all possible combinations of R-group substituents at the different points of variation.

The alternative approach is to use the computational equivalent of a chemical reaction, or *reaction transform*. Here, one does not need to define a common template or to generate sets of 'clipped' reagents. Rather, the library can be enumerated using as input the initial reagent structures and the chemical transforms required to operate upon them. In this

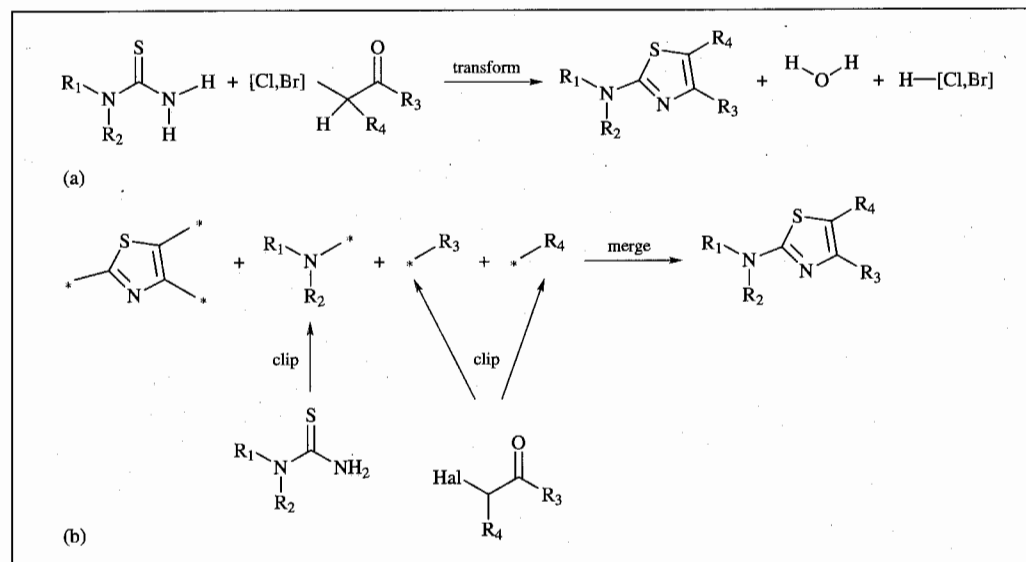


Fig. 12.44: Comparison of the reaction transform (a) and fragment-marking (b) approaches to the enumeration of aminothiazoles.

way, it more closely replicates the stages involved in the actual synthesis, wherein reagents react together according to the rules of synthetic chemistry (at least, when the chemistry works as planned!).

The key elements of the fragment-marking and transform approach can be illustrated using as an example the synthesis of aminothiazoles from thioureas and alpha halo ketones (Figure 12.44). With the reaction transform approach (Figure 12.44(a)), one would simply define an appropriate transform. The enumeration engine then applies this transform to the initial starting materials (i.e. the thiourea and the alpha halo ketone) to produce the aminothiazole (with water and the hydrogen halide as byproducts). Using the fragment-marking approach (Figure 12.44(b)), one would construct three sets of 'clipped' fragments (two from each alpha halo ketone and one from each thiourea), which would then be grafted on to give the central thiazole core to give the appropriate products.

Both the marking-up and reaction transform approach have advantages and disadvantages. In favour of the marking-up approach is the fact that for some kinds of library (i.e. those that most obviously fit the 'core plus R group' definition) it can be the fastest way to enumerate the library. This is because the fragment-marking approach involves only some rather elementary connection table operations once the R groups have been generated. Although most systems offer automated ways to generate the R groups (i.e. 'clipping algorithms') problems almost invariably arise which need to be corrected by hand. This can make the fragment-marking approach time-consuming to perform for a non-expert unless sets of pre-defined R groups are already available. In addition, there are certain reactions which are not properly handled by the fragment-marking approach, one well-known example being the Diels-Alder reaction, where a simple fragment-marking approach would generate a

number of extraneous and incorrect products. Moreover, in some cases there is no clear core structure. The advantages of the reaction method include the ability to enumerate directly from the reagents without having to perform any pre-processing and the ability to reuse the same transforms many times (once they have been defined). However, this method requires more computational steps and so is typically slower. Perhaps the key advantage, however, is that this approach models the actual chemical steps involved in the experiment, so bringing the experimental and computational systems closer together (and thus easier for the bench scientist to appreciate and to do themselves).

Enumeration of a library by either the fragment-marking or the reaction transform approach typically involves just one molecule or reaction scheme being considered at a time. However, the nature of most combinatorial libraries is that there is often much in common between the molecules. A common 'core' can usually be identified (else the fragment-marking method would not work) but in addition some subsets of the products may also have parts in common. For example, a fraction of the products may have a phenyl ring at a particular position. By recognising these relationships (using Markush structures, which were originally developed for the computer representation of patents) enumeration can be performed much more efficiently [Downs and Barnard 1997].

12.14.3 Combinatorial Subset Selection

For any synthetic scheme, the key issue in combinatorial library design is monomer selection, the objective of which is to identify those monomers which when combined together provide the 'optimal' combinatorial library. By 'optimal' we mean that library which best meets the prescribed objectives: it might be the most diverse, have the maximum number of molecules that could fit a 3D pharmacophore or a protein binding site, best match a particular distribution of some physicochemical property, or some combination of these or other criteria. An important consideration when designing a combinatorial library is the *subset selection constraint*. In a 'true' combinatorial library of the form $A \times B \times C$, every molecule from the set of reagents A reacts with every molecule from B and every molecule from C to generate $n_A \times n_B \times n_C$ product structures, where n_A , n_B , n_C are the numbers of reagent molecules A, B and C. Typically, there will be many more possible reagents A, B, C available to us than we can actually incorporate into the library, hence the need to select the subset of monomers which give rise to the 'optimal' library. Suppose the number of possible reagents A is N_A , etc. The size of the so-called *virtual library* is thus $N_A \times N_B \times N_C$. The number of ways of selecting n objects from N is ${}^N C_n$, and so the number of different combinatorial libraries of size $n_A \times n_B \times n_C$ that could be made for this three-component library is ${}^{N_A} C_{n_A} \times {}^{N_B} C_{n_B} \times {}^{N_C} C_{n_C}$. If we have available 100 reagents for each of A, B and C and we wish to make a $10 \times 10 \times 10$ library then the number of possible libraries that we could make is approximately 10^{40} . Identifying the one 'optimal' library from this extremely large number of possible libraries is clearly a difficult problem, which cannot be solved by a systematic examination of every possible solution.

As might be expected, established optimisation techniques such as simulated annealing and genetic algorithms have been used to tackle the subset selection problem. These methods

gradually evolve possible solutions until either no better solution can be found or until the predetermined number of iterations is exceeded. In the genetic algorithm approach, the chromosome encodes for a particular set of monomers, which, when combined together in a combinatorial fashion, would give a particular set of products. As we alluded earlier, the initial efforts were directed towards 'diverse' libraries and so the optimisation functions were formulated accordingly. The subsequent emphasis on more focused libraries has required alternative functions that aim to optimise the number of molecules with a particular property. More generic are those approaches which are able to simultaneously optimise for both diversity and target constraints [Gillet *et al.* 1999].

This type of library design is often known as *product-based monomer selection* as it is the properties of the product molecules that determine the ultimate monomer selections. Enumeration is clearly key to this approach, as it requires product structures to be generated. The alternative is monomer-based selection, where one only considers the properties of the individual monomers and not the properties of the product molecules. The main advantage to monomer-based selection is that the size of the search space is much smaller; in product-based selection one has to directly or indirectly consider the $N_A \times N_B \times N_C$ potential product molecules, whereas in monomer-based selection one only need consider the $N_A + N_B + N_C$ monomers. It has been shown that product-based selection gives superior results (as one would expect) [Gillet *et al.* 1997], but it is also clear that in some cases the virtual libraries will be so large that full enumeration may be impossible, and some combination of the two approaches may be required. Moreover, the synthetic strategy may not require the library to be fully combinatorial so enabling this constraint to be relaxed somewhat.

12.14.4 The Future

The techniques of combinatorial chemistry and high-throughput screening have developed extremely rapidly since their potential application to drug discovery was recognised [Leach and Hann 2000]. Most large pharmaceutical companies have a significant investment in these areas, and many smaller companies have been founded specifically to exploit these new techniques. In addition, combinatorial techniques are starting to be applied to other areas such as materials science. Nevertheless, there remain some important issues in the way that combinatorial chemistry is practically applied in drug discovery and the role of computational methods in supporting that process. Many of these issues are more concerned with practice than theory, such as the need for adequate supplies of appropriate monomers, the development of new solid-phase chemistries and the alignment of synthesis and screening resources. On the theoretical side, we have seen a gradual shift in emphasis from 'diversity' as the sole factor in library design towards 'biased' or 'focused' libraries that also try to take into account relevant knowledge about the biological target(s) for which the library is intended [Hann and Green 1999]. In addition, some of the notions concerning 'drug-likeness' and the type of molecule that should be in a library are being scrutinised. For example, as the chance of finding the 'drug molecule' in one step is so unlikely it is perhaps more appropriate to expect that molecules with only modest affinity will be found in the early stages. Thus exploratory libraries should contain smaller, less complex molecules than typical drugs [Teague *et al.* 1999]. As with so many other areas of molecular

modelling, progress requires not only good algorithms and faster computers but also a close integration with experiment and a better understanding of the underlying chemical and physical principals involved.

Further Reading

- Agrafiotis D K, J C Myslik and F R Salemme 1999. Advances in Diversity Profiling and Combinatorial Series Design. *Molecular Diversity* 4:1-22.
- Charifson P S (Editor) 1997. *Practical Application of Computer-Aided Drug Design*. New York, Dekker.
- Clark D E, C W Murray and J Li 1997. Current Issues in *De Novo* Molecular Design. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 11. New York, VCH Publishers, pp. 67-125.
- Dean P M (Editor) 1995. *Molecular Similarity in Drug Design*. London, Blackie Academic and Professional.
- Downs G M and Peter Willett 1995. Similarity Searching in Databases of Chemical Structures. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 7. New York, VCH Publishers, pp. 1-66.
- Drewry D H and S S Young 1999. Approaches to the Design of Combinatorial Libraries. *Chemometrics in Intelligent Laboratory Systems* 48:1-20.
- Good A C and J S Mason 1995. Three-Dimensional Structure Database Searches. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 7. New York, VCH Publishers, pp. 67-117.
- Graham R C 1993. *Data Analysis for the Chemical Sciences. A Guide to Statistical Techniques*. New York, VCH Publishers.
- Guner O F (Editor) 2000. *Pharmacophore Perception, Development, and Use in Drug Design*. International University Line Biotechnology Series, 2.
- Jurs P C 1990. Chemometrics and Multivariate Analysis in Analytical Chemistry. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 1. New York, VCH Publishers, pp. 169-212.
- Kubinyi H (Editor) 1993. *3D QSAR in Drug Design, Theory, Methods and Applications*. Leiden, ESCOM.
- Kubinyi H 1995. The Quantitative Analysis of Structure-Activity Relationships. In Wolff M E (Editor) *Burger's Medicinal Chemistry and Drug Discovery*, 5th Edition, Volume 1. New York, John Wiley & Sons, pp. 497-571.
- Livingstone, D 1995. *Data Analysis for Chemists*. Oxford, Oxford University Press.
- Livingstone, D 2000. The Characterisation of Chemical Structures Using Molecular Properties. A Survey. *Journal of Chemical Information and Computer Science* 40:195-209.
- Marshall G R 1995. Molecular Modeling in Drug Design. In Wolff M E (Editor) *Burger's Medicinal Chemistry and Drug Discovery*, 5th Edition, Volume 1. New York, John Wiley & Sons, pp. 573-659.
- Martin E J, D C Spellmeyer, R E Critchlow Jr and J M. Blaney 1997. Does Combinatorial Chemistry Obviate Computer-Aided Drug Design? In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 10. New York, VCH Publishers, pp. 75-100.
- Martin Y C 1978. *Quantitative Drug Design. A Critical Introduction*. New York, Marcel Dekker.
- Martin Y C, M G Bures and P Willett 1990. Searching Databases of Three-Dimensional Structures. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 1. New York, VCH Publishers, pp. 213-263.
- Montgomery D C and A A Peck 1992. *Introduction to Linear Regression Analysis*. New York, John Wiley & Sons.

- Murcko M A 1997. Recent Advances in Ligand Design Methods. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 11. New York, VCH Publishers, pp. 1–66.
- Oprea T I and C L Waller 1997. Theoretical and Practical Aspects of Three-Dimensional Quantitative Structure–Activity Relationships. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 11. New York, VCH Publishers, pp. 127–182.
- Otto M. *Chemometrics. Statistics and Computer Application in Analytical Chemistry*. New York, Wiley–VCH.
- Spellmeyer D C and P D J Grootenhuys 1999. Recent Developments in Molecular Diversity: Computational Approaches to Combinatorial Chemistry. *Annual Reports in Medicinal Chemistry* 34:287–296.
- Tute M S 1990. History and Objectives of Quantitative Drug Design. In Hansch C, P G Sammes and J B Taylor (Editors) *Comprehensive Medicinal Chemistry* Volume 4. Oxford, Pergamon Press, pp. 1–31.
- Waterbeemd H van de 1995. *Chemometric Methods in Molecular Design*. Weinheim, VCH Publishers.
- Willett P (Editor) 1997. *Computational Methods for the Analysis of Molecular Diversity. Perspectives in Drug Discovery and Design* Volumes 7/8. Dordrecht, Kluwer.

References

- Agrafiotis D K 1997. Stochastic Algorithms for Maximising Molecular Diversity. *Journal of Chemical Information and Computer Science* 37:841–851.
- Ajay A and M A Murcko 1995. Computational Methods to Predict Binding Free Energy in Ligand–Receptor Complexes. *Journal of Medicinal Chemistry* 38:4951–4967.
- Ajay A, W P Walters and M A Murcko 1998. Can We Learn to Distinguish Between ‘Drug-like’ and ‘Non-drug-like’ Molecules? *Journal of Medicinal Chemistry* 41:3314–3324.
- Andrea T A and H Kalayeh 1991. Applications of Neural Networks in Quantitative Structure–Activity Relationships of Dihydrofolate Reductase Inhibitors. *Journal of Medicinal Chemistry* 34:2824–2836.
- A-Razzak M and R C Glen 1992. Applications of Rule-induction in the Derivation of Quantitative Structure–Activity Relationships. *Journal of Computer-Aided Molecular Design* 6:349–383.
- Babine R E and S L Bender 1997. Recognition of Protein–Ligand Complexes: Applications to Drug Design. *Chemical Reviews* 97:1359–1472.
- Barnum D, J Greene, A Smellie and P Sprague 1996. Identification of Common Functional Configurations among Molecules. *Journal of Chemical Information and Computer Science* 36:563–571.
- Baxter C A, C W Murray, D E Clark, D R Westhead and M D Eldridge 1998. Flexible Docking using Tabu Search and an Empirical Estimate of Binding Affinity. *Proteins: Structure, Function and Genetics* 33:367–382.
- Blaney J M and J S Dixon 1993. A Good Ligand Is Hard to Find: Automated Docking Methods. *Perspectives in Drug Discovery and Design* 1:301–319.
- Böhm H-J 1992. LUDI – Rule-Based Automatic Design of New Substituents for Enzyme Inhibitor Leads. *Journal of Computer-Aided Molecular Design* 6:593–606.
- Böhm H-J 1994. The Development of a Simple Empirical Scoring Function to Estimate the Binding Constant for a Protein–ligand Complex of Known Three-Dimensional Structure. *Journal of Computer-Aided Molecular Design* 8:243–256.
- Böhm H-J 1998. Prediction of Binding Constants of Protein Ligands: A Fast Method for the Prioritisation of Hits Obtained from *De Novo* Design or 3D Database Search Programs. *Journal of Computer-Aided Molecular Design* 12:309–323.
- Böhm H-J and G Klebe 1996. What Can We Learn From Molecular Recognition in Protein–Ligand Complexes for the Design of New Drugs? *Angewandte Chemie International Edition in English* 35:2588–2614.

- Boström J, P-O Norrby and T Liljefors 1998. Conformational Energy Penalties of Protein-bound Ligands. *Journal of Computer-Aided Molecular Design* 12:383–396.
- Bradshaw J 1997. Introduction to Tversky Similarity Measure. At http://www.daylight.com/meetings/mug97/Bradshaw/MUG97/tv_tversky.html.
- Bron C and J Kerbosch 1973. Algorithm 475. Finding All Cliques of an Undirected Graph. *Communications of the ACM* 16:575–577.
- Brown R D and Y C Martin 1996. Use of Structure–Activity Data to Compare Structure-Based Clustering Methods and Descriptors for Use in Compound Selection. *Journal of Chemical Information and Computer Science* 36:572–583.
- Carbo R, L Leyda and M Arnau 1980. An Electron Density Measure of the Similarity Between Two Compounds. *International Journal of Quantum Chemistry* 17:1185–1189.
- Carhart R E, D H Smith and R Venkataraghavan 1985. Atom Pairs as Molecular Features in Structure–Activity Studies: Definition and Applications. *Journal of Chemical Information and Computer Science* 25:64–73.
- Charifson P S, J J Corkery, M A Murcko and W P Walters 1999. Consensus Scoring: A Method for Obtaining Improved Hit Rates from Docking Databases of Three-Dimensional Structures into Proteins. *Journal of Medicinal Chemistry* 42:5100–5109.
- Cramer R D III, D E Patterson and J D Bunce 1988. Comparative Molecular Field Analysis (CoMFA). 1. Effect of Shape on Binding of Steroids to Carrier Proteins. *Journal of the American Chemical Society* 110:5959–5967.
- Cruciani G, S Clementi and M Baroni 1993. Variable Selection in PLS Analysis. In Kubinyi H (Editor) *3D QSAR in Drug Design*. Leiden, ESCOM, pp. 551–564.
- Cummins D J, C W Andrews, J A Bentley and M Cory 1996. Molecular Diversity in Chemical Databases: Comparison of Medicinal Chemistry Knowledge Bases and Databases of Commercially Available Compounds. *Journal of Chemical Information and Computer Science* 36:750–763.
- Dalby A, J G Nourse, W D Hounshell, A K I Gushurst, D L Grier, B A Leland and J Laufer 1992. Description of Several Chemical Structure File Formats Used by Computer Programs Developed at Molecular Design Limited. *Journal of Chemical Information and Computer Science* 32:244–255.
- Dammkoehler R A, S F Karasek, E F B Shands and G R Marshall 1989. Constrained Search of Conformational Hyperspace. *Journal of Computer-Aided Molecular Design* 3:3–21.
- Desjarlais R L, R P Sheridan, G L Seibel, J S Dixon, I D Kuntz and R Venkataraghavan 1988. Using Shape Complementarity as an Initial Screen in Designing Ligands for a Receptor Binding Site of Known Three-Dimensional Structure. *Journal of Medicinal Chemistry* 31:722–729.
- Downs G M and J M Barnard 1997. Techniques for Generating Descriptive Fingerprints in Combinatorial Libraries. *Journal of Chemical Information and Computer Science* 37:59–61.
- Downs G M, P Willett and W Fisanick 1994. Similarity Searching and Clustering of Chemical-Structure Databases Using Molecular Property Data. *Journal of Chemical Information and Computer Science* 34:1094–1102.
- Dunn W J III, S Wold, U Edlund, S Hellberg and J Gasteiger 1984. Multivariate Structure–Activity Relationships Between Data from a Battery of Biological Tests and an Ensemble of Structure Descriptors: The PLS Method. *Quantitative Structure–Activity Relationships* 3:131–137.
- Eldridge M D, C W Murray, T R Auton, G V Paolini and R P Mee 1997. Empirical Scoring Functions: I. The Development of a Fast Empirical Scoring Function to Estimate the Binding Affinity of Ligands in Receptor Complexes. *Journal of Computer-Aided Molecular Design* 11:425–445.
- Gasteiger J, C Rudolph and J Sadowski 1990. Automatic Generation of 3D Atomic Coordinates for Organic Molecules. *Tetrahedron Computer Methodology* 3:537–547.
- Gelhaar D K, G M Verkhivker, P A Rejto, C J Sherman, D B Fogel, L J Fogel and S T Freer 1995. Molecular Recognition of the Inhibitor AG-1343 by HIV-1 Protease: Conformationally Flexible Docking by Evolutionary Programming. *Chemistry and Biology* 2:317–324.

- Ghose A K and G M Crippen 1986. Atomic Physicochemical Parameters for Three-dimensional Structure-directed Quantitative Structure-Activity Relationships. I. Partition Coefficients as a Measure of Hydrophobicity. *Journal of Computational Chemistry* 7:565-577.
- Ghose A K, V N Viswanadhan and J J Wendoloski 1998. Prediction of Hydrophobic (Lipophilic) Properties of Small Organic Molecules Using Fragmental Methods: An Analysis of ALOGP and CLOGP Methods. *Journal of Physical Chemistry* 102:3762-3772.
- Gillet V J, A P Johnson, P Mata, S Sik and P Williams 1993. SPROUT - A Program for Structure Generation. *Journal of Computer-Aided Molecular Design* 7:127-153.
- Gillet V J, P Willett and J Bradshaw 1997. The Effectiveness of Reactant Pools for Generating Structurally Diverse Combinatorial Libraries. *Journal of Chemical Information and Computer Science* 37:731-740.
- Gillet V J, P Willett and J Bradshaw 1998. Identification Of Biological Activity Profiles Using Substructural Analysis And Genetic Algorithms. *Journal of Chemical Information and Computer Science* 38:165-179.
- Gillet V J, P Willett, J Bradshaw and D V S Green 1999. Selecting Combinatorial Libraries to Optimize Diversity and Physical Properties. *Journal of Chemical Information and Computer Science* 39:169-177.
- Glen R C and A W R Payne 1995. A Genetic Algorithm for the Automated Generation of Molecules within Constraints. *Journal of Computer-Aided Molecular Design* 9:181-202.
- Good A C, E E Hodgkin and Richards W G 1993. The Utilisation of Gaussian Functions for the Rapid Evaluation of Molecular Similarity. *Journal of Chemical Information and Computer Science* 32:188-192.
- Good A C and I D Kuntz 1995. Investigating the Extension of Pairwise Distance Pharmacophore Measures to Triplet-based Descriptors. *Journal of Computer-Aided Molecular Design* 9:373-379.
- Goodford P J 1985. A Computational Procedure for Determining Energetically Favorable Binding Sites on Biologically Important Macromolecules. *Journal of Medicinal Chemistry* 28:849-857.
- Goodsell D S and A J Olson 1990. Automated Docking of Substrates to Proteins by Simulated Annealing. *Proteins: Structure, Function and Genetics* 8:195-202.
- Greco G, E Novellino and Y C Martin 1997. Approaches to Three-dimensional Quantitative Structure-Activity Relationships. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 11. New York, VCH Publishers, pp. 183-240.
- Greene J, S Kahn, H Savoj, P Sprague and S Teig 1994. Chemical Function Queries for 3D Database Search. *Journal of Chemical Information and Computer Science* 34:1297-1308.
- Hall L H and L B Kier 1991. The Molecular Connectivity Chi Indexes and Kappa Shape Indexes in Structure-Property Modeling. In Lipkowitz K B and D B Boyd (Editors) *Reviews in Computational Chemistry* Volume 2. New York, VCH Publishers, pp. 367-422.
- Hall L H, B Mohny and L B Kier 1991. The Electrotopological State: An Atom Index for QSAR. *Quantitative Structure-Activity Relationships* 10:43-51.
- Hann, M and R Green 1999. Chemoinformatics - A New Name for an Old Problem? *Current Opinion in Chemistry and Biology* 3:379-383.
- Hansch C 1969. A Quantitative Approach to Biochemical Structure-Activity Relationships. *Accounts of Chemical Research* 2:232-239.
- Hansch C and T E Klein 1986. Molecular Graphics and QSAR in the Study of Enzyme-Ligand Interactions. On the Definition of Bioreceptors. *Accounts of Chemical Research* 19:392-400.
- Hansch C, J McClarin, T Klein and R Langridge 1985. A Quantitative Structure-Activity Relationship and Molecular Graphics Study of Carbonic Anhydrase Inhibitors. *Molecular Pharmacology* 27:493-498.
- Hassan M, J P Bielawski, J C Hempel and M Waldman 1996. Optimisation and Visualisation of Molecular Diversity of Combinatorial Libraries. *Molecular Diversity* 2:64-74.
- Head R D, M L Smythe, T I Oprea, C L Waller, S M Green and G R Marshall 1996. VALIDATE: A New Method for the Receptor-Based Prediction of Binding Affinities of Novel Ligands. *Journal of the American Chemical Society* 118:3959-3969.

- Hodgkin E E and W G Richards 1987. Molecular Similarity Based on Electrostatic Potential and Electric Field. *International Journal of Quantum Chemistry. Quantum Biology Symposia* 14:105-110.
- Holiday J D, S R Ranade and P Willett 1995. A Fast Algorithm For Selecting Sets Of Dissimilar Molecules From Large Chemical Databases. *Quantitative Structure-Activity Relationships* 14:501-506.
- Holloway M K, J M Wai, T A Halgren, P M D Fitzgerald, J P Vacca, B D Dorsey, R B Levin, W J Thompson, J Chen, S J deSolms, N Gaffin, A K Ghosh, E A Giuliani, S L Graham, J P Guare, R W Hungate, T A Lyle, W M Sanders, T J Tucker, M Wiggins, C M Wiscount, O W Woltersdorf, S D Young, P L Darke and J A Zugay 1995. A Priori Prediction of Activity for HIV-1 Protease Inhibitors Employing Energy Minimisation in the Active Site. *Journal of Medicinal Chemistry* 38:305-317.
- Hudson B D, R M Hyde, E Rahr, J Wood and J Osman 1996. Parameter Based Methods for Compound Selection from Chemical Databases. *Quantitative Structure-Activity Relationships* 15:285-289.
- Jones G, P Willett and R C Glen 1995a. A Genetic Algorithm for Flexible Molecular Overlay and Pharmacophore Elucidation. *Journal of Computer-Aided Molecular Design* 9:532-549.
- Jones G, P Willett and R C Glen 1995b. Molecular Recognition of Receptor Sites Using a Genetic Algorithm with a Description of Desolvation. *Journal of Molecular Biology* 245:43-53.
- Jones G, P Willett, R C Glen, A R Leach and R Taylor 1997. Development and Validation of a Genetic Algorithm for Flexible Docking. *Journal of Molecular Biology* 267:727-748.
- Judson R S, E P Jaeger and A M Treasurywala 1994. A Genetic Algorithm-Based Method for Docking Flexible Molecules. *Journal of Molecular Structure: Theochem* 114:191-206.
- Kennard R W and L A Stone 1969. Computer Aided Design of Experiments. *Technometrics* 11:137-148.
- King R D, S Muggleton, R A Lewis and M J E Sternberg 1992. Drug Design by Machine Learning: The Use of Inductive Logic Programming to Model the Structure-Activity Relationships of Trimethoprim Analogues Binding to Dihydrofolate Reductase. *Proceedings of the National Academy of Sciences USA* 89:11322-11326.
- King R D, S H Muggleton, A Srinivasan and M J E Sternberg 1996. Structure-Activity Relationships Derived by Machine Learning: The Use of Atoms and Their Bond Connectivities to Predict Mutagenicity by Inductive Logic Programming. *Proceedings of the National Academy of Sciences USA* 93:438-442.
- Klopman G, S Wang and D M Balthasar 1992. Estimation of Aqueous Solubility of Organic Molecules by the Group Contribution Approach. Application to the Study of Biodegradation. *Journal of Chemical Information and Computer Science* 32:474-482.
- Kramer B, M Rarey and T Lengauer 1999. Evaluation of the FLEXX Incremental Construction Algorithm for Protein-Ligand Docking. *Proteins: Structure, Function and Genetics* 37:228-241.
- Kubinyi H 1998. Structure-based Design of Enzyme Inhibitors and Receptor Ligands. *Current Opinion in Drug Discovery and Development* 1:5-15.
- Kuntz I D 1992. Structure-Based Strategies for Drug Design and Discovery. *Science* 257:1078-1082.
- Kuntz I D, J M Blaney, S J Oatley, R Langridge and T E Ferrin 1982. A Geometric Approach to Macromolecule-Ligand Interactions. *Journal of Molecular Biology* 161:269-288.
- Kuntz I D, E C Meng and B K Shoichet 1994. Structure-Based Molecular Design. *Accounts of Chemical Research* 27:117-123.
- Lam P Y S, P K Jadhav, C E Eyermann, C N Hodge, Y Ru, L T Bachelor, J L Meek, M J Otto, M M Rayner, Y N Wong, C-H Chang, P C Weber, D A Jackson, T R Sharpe and S Erickson-Viitanen 1994. Rational Design of Potent, Bioavailable, Nonpeptide Cyclic Ureas as HIV Protease Inhibitors. *Science* 263:380-384.
- Lauri G and P A Bartlett 1994. CAVEAT - A Program to Facilitate the Design of Organic Molecules. *Journal of Computer-Aided Molecular Design* 8:51-66.
- Leach A R 1994. Ligand Docking to Proteins With Discrete Side-chain Flexibility. *Journal Of Molecular Biology* 235:345-356.

- Leach A R and M M Hann 2000. The In Silico World of Virtual Libraries. *Drug Discovery Today* 5:326-336.
- Leach A R and I D Kuntz 1990. Conformational Analysis of Flexible Ligands in Macromolecular Receptor Sites. *Journal of Computational Chemistry* 13:730-748.
- Lemmen C and T Lengauer 2000. Computational Methods for the Structural Alignment of Molecules. *Journal of Computer-Aided Molecular Design* 14:215-232.
- Leo A and Weininger A 1995. CMR3 Reference Manual. At <http://www.daylight.com/dayhtml/doc/cmr/cmrref.html>.
- Leo A J 1993. Calculating log P_{oct} from Structures. *Chemical Reviews* 93:1281-1306.
- Lewis D W, D J Willock, C R A Catlow, J M Thomas and G J Hutchings 1996. De Novo Design of Structure-directing Agents for the Synthesis of Microporous Solids. *Nature* 382:604-606.
- Lewis R A, J S Mason and I M McLay 1997. Similarity Measures for Rational Set Selection and Analysis of Combinatorial Libraries: The Diverse Property-Derived (DPD) Approach. *Journal of Chemical Information and Computer Science* 37:599-614.
- Lewis R M and A R Leach 1994. Current Methods for Site-Directed Structure Generation. *Journal of Computer-Aided Molecular Design* 8:467-475.
- Lipinski C A, F Lombardo, B W Dominy and P J Feeney 1997. Experimental and Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and Development Settings. *Advanced Drug Delivery Reviews* 23:3-25.
- Malpass J A 1994. Continuum Regression: Optimised Prediction of Biological Activity. PhD thesis, University of Portsmouth, UK.
- Manallack D T, D D Ellis and D J Livingstone 1994. Analysis of Linear and Nonlinear QSAR Data Using Neural Networks. *Journal of Computer-Aided Molecular Design* 37:3758-3767.
- Marriott D P, I G Dougall, P Meghani, Y-J Liu and D R Flower 1999. Lead Generation Using Pharmacophore Mapping and Three-Dimensional Database Searching: Application to Muscarinic M₃ Receptor Antagonists. *Journal of Medicinal Chemistry* 42:3210-3216.
- Martin E J, J M Blaney, M A Siani, D C Spellmeyer, A K Wong and W H Moos 1995. Measuring Diversity: Experimental Design of Combinatorial Libraries for Drug Discovery. *Journal of Medicinal Chemistry* 38:1431-1436.
- Martin Y C, M G Bures, A A Danaher, J DeLazzer, I Lico and P A Pavlik 1993. A Fast New Approach to Pharmacophore Mapping and its Application to Dopaminergic and Benzodiazepine Agonists. *Journal of Computer-Aided Molecular Design* 7:83-102.
- Mason J S, I Morize, P R Menard, D L Cheney, C Hulme and R F Labaudiniere 1999. New 4-Point Pharmacophore Method for Molecular Similarity and Diversity Applications: Overview of the Method and Applications, Including a Novel Approach to the Design of Combinatorial Libraries Containing Privileged Substructures. *Journal of Medicinal Chemistry* 42:3251-3264.
- Meng E C, B K Shoichet and I D Kuntz 1992. Automated Docking with Grid-Based Energy Evaluation. *Journal of Computational Chemistry* 13:505-524.
- Miranker A and M Karplus 1991. Functionality Maps of Binding Sites - A Multiple Copy Simultaneous Search Method. *Proteins: Structure, Function and Genetics* 11:29-34.
- Moon J B and W J Howe 1991. Computer Design of Bioactive Molecules - A Method for Receptor-Based De Novo Ligand Design. *Proteins: Structure, Function and Genetics* 11:314-328.
- Morgan H L 1965. The Generation of a Unique Machine Description for Chemical Structures - A Technique Developed at Chemical Abstracts Service. *Journal of Chemical Documentation* 5:107-113.
- Myatt G 1995. Computer-aided Estimation of Synthetic Accessibility. PhD thesis, University of Leeds.
- Nilakantan R, N Bauman, J S Dixon and R Venkataraghavan 1987. Topological Torsion: A New Molecular Descriptor for SAR Applications. Comparison with Other Descriptors. *Journal of Chemical Information and Computer Science* 27:82-85.
- Oshiro C M, I D Kuntz and J S Dixon 1995. Flexible Ligand Docking Using a Genetic Algorithm. *Journal of Computer-Aided Molecular Design* 9:113-130.

- Pastor M, G Cruciani and S Clementi 1997. Smart Region Definition: A New Way to Improve the Predictive Ability and Interpretability of Three-Dimensional Quantitative Structure-Activity Relationships. *Journal of Medicinal Chemistry* 40:1455-1464.
- Patani G A and E J LaVoie 1996. Bioisosterism: A Rational Approach in Drug Design. *Chemical Reviews* 96:3147-3176.
- Pearlman R S and K M Smith 1998. Novel Software Tools for Chemical Diversity. *Perspectives in Drug Discovery and Design* vols 9/10/11(3D QSAR in Drug Design: Ligand/Protein Interactions and Molecular Similarity), pp. 339-353.
- Pickett S D, J S Mason and I M McLay 1996. Diversity Profiling and Design Using 3D Pharmacophores: Pharmacophore-Derived Queries (PDQ). *Journal of Chemical Information and Computer Science* 36:1214-1223.
- Poso A, R Juvonen and J Gynther 1995. Comparative Molecular Field Analysis of Compounds with CYP2A5 Binding Affinity. *Quantitative Structure-Activity Relationships* 14:507-511.
- Priestle J P, A Fassler, J Rosel, M Tintelnog-Blomley, P Strop and M G Gruetter 1995. Comparative Analysis of The X-Ray Structures of HIV-1 and HIV-2 Proteases in Complex with a Novel Pseudosymmetric Inhibitor. *Structure (London)* 3:381-389.
- Rarey M, B Kramer, T Lengauer and G Klebe 1996. A Fast Flexible Docking Method Using an Incremental Construction Algorithm. *Journal of Molecular Biology* 261:470-489.
- Rhyu K-B, H C Patel and A J Hopfinger 1995. A 3D-QSAR Study of Anticocccidal Triazines Using Molecular Shape Analysis. *Journal of Chemical Information and Computer Science* 35:771-778.
- Rogers D and A J Hopfinger 1994. Application of Genetic Function Approximation to Quantitative Structure-Activity Relationships and Quantitative Structure-Property Relationships. *Journal of Chemical Information and Computer Science* 34:854-866.
- Rumelhart D E, G W Hinton and R J Williams 1986. Learning Representations by Back-propagating Errors. *Nature* 323:533-536.
- Rusinko A III, M W Farnen, C G Lambert, P L Brown and S S Young 1999. Analysis of a Large Structure/Biological Activity Data Set Using Recursive Partitioning. *Journal of Chemical Information and Computer Science* 39:1017-1026.
- Rusinko A III, J M Skell, R Balducci, C M McGarity and R S Pearlman 1988. *CONCORD: A Program for the Rapid Generation of High Quality 3D Molecular Structures*. St Louis, Missouri, The University of Texas at Austin and Tripos Associates.
- Sadowski J and H Kubinyi 1998. A Scoring Scheme for Discriminating Between Drugs and Nondrugs. *Journal of Medicinal Chemistry* 41:3325-3329.
- Sheridan R P, R Nilakantan, J S Dixon and R Venkataraghavan 1986. The Ensemble Approach to Distance Geometry: Application to the Nicotinic Pharmacophore. *Journal of Medicinal Chemistry* 29:899-906.
- Shuker S B, P J Hadjuk, R P Meadows and R P Fesik 1996. Discovering High-affinity Ligands for Proteins: SAR by NMR. *Science* 274:1531-1534.
- Snarey M, N K Terrett, P Willett and D J Wilton 1997. Comparison of Algorithms for Dissimilarity-based Compound Selection. *Journal of Molecular Graphics and Modelling* 15:372-385.
- Swain C G and E C Lupton 1968. Field and Resonance Components of Substituent Effects. *Journal of the American Chemical Society* 90:4328-4337.
- Swain C G, S H Unger, N R Rosenquist and M S Swain 1983. Substituent Effects on Chemical Reactivity. Improved Evaluation of Field and Resonance Components. *Journal of the American Chemical Society* 105:492-502.
- Teague S J, A M Davis, P D Leeson and T Oprea 1999. The Design of Leadlike Combinatorial Libraries. *Angewandte Chemie International Edition in English* 38:3743-3748.
- Thibaut U, G Folkers, G Klebe, H Kubinyi, A Merz and D Rognan 1993. Recommendations for CoMFA Studies and 3D QSAR Publications. In Kubinyi H (Editor) *3D QSAR in Drug Design*. Leiden, ESCOM, pp. 711-728.

- Thornber C W 1979. Isosterism and Molecular Modification in Drug Design. *Chemical Society Reviews* 8:563-580.
- Tversky A. 1977. Features of Similarity. *Psychological Reviews* 84:327-352.
- Ullmann J R 1976. An Algorithm for Subgraph Isomorphism. *Journal of the Association for Computing Machinery* 23:31-42.
- Von Itzstein M, W Y Wu, G B Kok, M S Pegg, J C Dyason, B Jin, T V Phan, M L Smythe, H F Whites, S W Oliver, P M Colman, J N Varghese, D M Ryan, J M Woods, R C Bethell, V J Hotham, J M Cameron and C R Penn 1993. Rational Design of Potent Sialidase-Based Inhibitors of Influenza Virus Replication. *Nature* 363:418-423.
- Wang R, Y Fu and L Lai 1997. A New Atom-Additive Method for Calculating Partition Coefficients. *Journal of Chemical Information and Computer Science* 37:615-621.
- Weininger D 1988. SMILES, A Chemical Language and Information System. 1. Introduction to Methodology and Encoding Rules. *Journal of Chemical Information and Computer Science* 28:31-36.
- Weininger D, A Weininger and J L Weininger 1989. SMILES. 2. Algorithm for Generation of Unique SMILES Notation. *Journal of Chemical Information and Computer Science* 29:97-101.
- Welch W, J Ruppert and A N Jain 1996. Hammerhead: Fast, Fully Automated Docking of Flexible Ligands to Protein Binding Sites. *Chemistry and Biology* 3:449-462.
- Wildman S A and G M Crippen 1999. Prediction of Physicochemical Parameters by Atomic Contributions. *Journal of Chemical Information and Computer Science* 39:868-873.
- Willett P, J M Barnard and G M Downs 1998. Chemical Similarity Searching. *Journal of Chemical Information and Computer Science* 38:983-996.
- Willock D J, D W Lewis, C R A Catlow, G J Hutchings and J M Thomas 1997. Designing Templates for the Synthesis of Microporous Solids Using *De Novo* Molecular Design Methods. *Journal of Molecular Catalysis A: Chemical* 119:415-424.
- Wiswesser W J 1954. *A Line-Formula Chemical Notation*. New York, Crowell Co.
- Wold H 1982. Soft Modeling. The Basic Design and Some Extensions. In Joreskog K-G and H Wold (Editors) *Systems under Indirect Observation* Volume II. Amsterdam, North-Holland.
- Wold S, E Johansson and M Cocchi 1993. PLS - Partial Least-squares Projections to Latent Structures. In Kubinyi H (Editor) *3D QSAR in Drug Design*. Leiden, ESCOM, pp. 523-550.

Index

Note: **boldened** page references indicate chapters.

- ab initio* defined 65
- ab initio* molecular dynamics 616-22
- ab initio* potentials for water 216-18
- ab initio* quantum mechanics, calculating properties using 74-86
see also advanced *ab initio* methods
- absolute free energies 573-4
- accessible surface 7
- ACE (angiotension converting enzyme) 649-51
- acetaldehyde 180, 578
- acetamide 573
- acetic acid 504-5, 573, 643
SMILES notation 644, 645
- 4-acetamido benzoic acid 661
- acetonitrile 597
- acetylcholine 678
- acronyms and abbreviations 104-5, 553-4
- adenine 227
- adiabatic mapping 286
- adjacency matrix 647
- adjoint matrix 15
- adsorption processes, Monte Carlo simulations of 441-2
- advanced *ab initio* methods **108-64**
density functional theory 126-37
electron correlation 110-17
energy component analysis 122-4
open-shell systems 108-10
practical considerations 117-22
solid state quantum mechanics 138-60
valence bond theories 124-6
- agglomerative cluster analysis methods 493-4
- agonists 640
- AINT function 350
- alanines 169, 459, 511, 525, 542, 546, 556-7
energy minimisation methods 277, 280, 286
free energy calculations 583-4
- aldehyde 610-11
- aldol reactions 610-12
- aliovalent substitution 623
- alkaline earth oxides 147
- alkanes 449-50
- α -helix 513-15
- AM1 86, 97-8, 102-3, 230
- AMBER force field 169-70, 175-6, 191, 211, 230-2
- amino acids 511, 549, 602
computer simulation 329-30
conformational analysis 459, 487
energy minimisation methods 277, 280, 286
force fields 169-70, 221
free energy calculations 572, 583-4
motifs 522
PAM matrices 524-6, 531, 556-7
'threading' 546
torsion angles 515
see also peptides; proteins
- aminothiazoles 716
- AMPAC program 8, 99
- amphiphiles, molecular dynamics simulation of 394-404
- angiotension converting enzyme 649-51
- angle bending 166, 173
- annealing, simulated 483-9, 504-5, 519, 691
- annotations 513
- antagonists 640
- antisymmetry principle 35
- arbitrary step energy minimisation 264
- Argand diagram 17
- arginine 329-30, 510, 525, 556-7
- argon 253, 323
force fields 205, 214
J-walking 434-5
time-steps 361-2
velocity autocorrelation 377
- arithmetic mean 20
- aromatic systems and charge schemes 197-9
- asparagine 330, 510, 525, 546, 556-7
- aspartic acid 510, 525, 556-7
- atoms/atomic
charges 157-9, 181, 192-5
marker 329-30
one-electron 30-4
orbitals 41-2, 56, 100, 241
polyelectronic 34-41
type 169
units 29
- atoms in molecules theory 80-1
- Aufbau principle 35

- Austin Model 1 (AM1) 86, 97–8, 102–3
autocorrelation function 376–8
automated protein modelling 548–9
autoscaling 681
availability 300
Axilrod–Teller term (triple-dipole) 213–15, 239
Azimuthal quantum number 31
- backtracking 462
backward sampling 567
band theory 141–2
orbital-based approach 142–6
Barker–Fisher–Watts potential 214
Basic Local Alignment Search Tool *see* BLAST
basis sets/functions 56, 85, 123
computational quantum mechanics 65–74
Gaussian functions 65–73 *passim*, 120, 137, 195
superposition error 121–2
see also STOs
BCUT method 686–7
bead model of polymers 428
Becke *see* BLYP
Beeman's algorithm 357
bending 166, 173, 176–8
benperidol 675
benzamidine 586–8
benzene 123
force fields 170, 174, 178, 186, 197–8
Hückel theory 99, 100
ring 81, 170, 178
SMILES notation 644
spin-coupled valence bond theory 126
benzyl bromide 670
beryllium 40, 113
 β -strand structures 513–14
 β -turns 513
BFCS (Broyden–Fletcher–Goldfarb–Shanno)
method 269–70
bilinear model 702
binding site 662, 689–91
binomial expansion 11
bioinformatics 513
see also amino acids; DNA; proteins
bioisosteres 648
biotin 576, 641
bitstring 645–7
BLAST (Basic Local Alignment Search Tool)
521, 524, 531–4, 548
Bloch's theorem/function 142–5, 146, 148, 161
block-diagonal Newton–Raphson minimisation
268
BLOSUM matrices 526
BLYP (Becke gradient-exchange correction and
Lee–Yang–Parr correlation functional) 135,
136, 137
B3LYP density function 615
- Bohr radius 31
Boltzmann distribution 192, 214, 274, 483, 611
computer simulation 306–7, 347
conformational analysis 457
Monte Carlo simulation 415–16, 433, 435–6,
445–6
Boltzmann factor 306–7, 413
Boltzmann weighted average 576, 581
bond/bonding 644
inorganic molecules 234–5
lack of *see* non-bonded interactions
orders 81–3
stretching 166, 170–3
valence 124–6
see also under carbon; hydrogen
bond fluctuation model 424–5
Born equation/model 238, 593–4, 598–601, 609
Born–Oppenheimer approximation 4, 35–6, 50
boundary
computer simulation 317–21
element method 598
Bravais lattices 138–9
Brillouin's theorem 112–13, 115
Brillouin zone 140–1, 145–6, 150–1, 157–8,
298–9
bromine 571
Broyden–Fletcher–Goldfarb–Shanno method
269–70
BSSE (basis set superposition error) 121–2
Buckingham potential 209, 238
build-up approach 517
butadiene 233, 294–5
butane 85–6, 186, 449–50, 582
butanone 611–12
- cage structure of liquids 377
calcium 589, 626
calix[4]arene 291–2
Cambridge Structural Database *see* CSD
canonical ensemble 563, 569
canonical genetic algorithms 479–80
canonical representation 644
canonical structures 541
captopril 649
Carbo index 678–9
carbon
bonds 4–5, 98, 362, 612, 652
energy minimisation methods 253, 280–1
force fields 167, 180, 211, 233, 236
five-carbon fragment 472
force fields 167, 180, 211, 233, 236, 244
valence electron density 160
carbon dioxide 181, 316, 616–17
carbonic anhydrase 616
carboxylic acid 660
Car–Parrinello scheme 610, 617–19

- Cartesian coordinates 2–4
conformational analysis 466–8
energy minimisation methods 255, 257–8, 275,
290
molecular dynamics simulation 370, 372, 379,
393
Monte Carlo simulation 417, 420–1, 423
vectors 11–12
CASP3 548
CASSCF (complete active-space SCF) 113, 295
CAVEAT program 689
CBMC (configurational bias Monte Carlo)
simulation 443–50
cell
cubic 315–19
index method 326
multipole method 341–3
unit 138
Wigner–Seitz 140, 350
Central Dogma 509, 512
central multipole expansion 181–7
CFF (consistent force field) 231
chain amphiphiles and molecular dynamics
394–404
charge 187
atomic 157–9, 181, 191–5
density matrix 58–9
image simulation 340–1
oscillating 201–2
schemes 197–9
CHELP procedure 191
chemical potential calculation in Monte Carlo
simulation 442–3
chemical reactions 610–22
molecular dynamics 616–22
potential of mean force of 612–14
quantum and molecular mechanics combined
614–16
simulation empirically 610–12
chemokine 475
chi molecular connectivity indices 672
chi-squared test 344–5
chiral constraints 473–4
chlorides 238, 612–13, 614, 676–7
chlorine 181, 231, 280–1, 571, 612, 620–1
chloroform 227, 573
chlorpromazine 678
chymosin 545
chymotrypsin 522–3
CI (configuration interaction) 111–13, 120
CID (configuration interaction doubles) 112, 113
CISD (configuration interaction singles and
doubles) 112, 113
city block *see* Hamming
Clausius, virial theorem of 309
Clausius–Mosotti relationship 238–9
- clique detection 653–6
CLOGP program 669–70
closed-shell systems 51, 56–9, 86–8, 109
cluster analysis 494, 534, 682–3
clustering algorithms 491–7
CMR program 671
CNDO (complete neglect of differential
overlap) 86, 89–92, 93, 94, 95
coefficients 148–9, 374
new molecules 676–8, 680–1, 685
partition 572–3, 668–71
cofactor 14–15
combination generator 420, 453–4
combinatorial explosion 460–1
combinatorial libraries 711–19
CoMFA (comparative molecular field analysis)
679, 708–11
comparative modelling of proteins 539–45
comparative molecular field analysis 679, 708–11
complete active-space SCF 113, 295
complete neglect of differential overlap 86,
89–92, 93, 94, 95
complete-linkage (furthest-neighbour) cluster
algorithm 493–4
complex numbers 16–18
computational quantum mechanics 26–107
acronyms used in 104–5
approximate orbital theories 86
atomic units 29
basis sets 65–74
calculating properties 74–86
Hückel theory 99–102
one-electron atoms 30–4
operators 28–9
polyelectronic atoms and molecules 34–41
semi-empirical methods 65, 86–99, 102–3
see also calculations *under* orbital;
Hartree–Fock equations
computer simulation 303–52
boundaries 317–21
equilibration, monitoring 321–3
free energy calculation difficulties 563–4
long-range forces 334–43
molecular dynamics 305–6, 307
phase space 312–15
practical aspects 315–16
real gas contribution to virial 309, 349–50
results and errors 343–7
statistical mechanics 347–8
thermodynamic properties, simple 307–12
time and ensemble averages 303–5
translating particle back into control box 350
truncating potential and minimum
image convention 324–34
see also molecular dynamics simulation,
Monte Carlo simulation

- computers
 hardware 8-9
 Internet and World Wide Web 9-10, 548, 553
 software 8-9, 99
see also databases
- concepts 1-25
see also coordinates; mathematical concepts
- CONCORD program 659
- conditionally convergent series 336
- conduction band 142
- conductor-like screening model 597
- configuration interaction *see* CI; CID; CISD
- configurational bias *see* CBMC
- conformational analysis 457-508
 choice of method 476-7
 clustering algorithms and pattern recognition 491-7
 conformational search 457, 662
 random 465-7, 476
 systematic 458-64, 476, 505
 crystal structures predicted 501-5
 dimensionality of data set reduced 497-9
 distance geometry 467-75, 476
 fitting, molecular 490-1
 global energy minimum 458, 479-83
 model-building 464-5 and NMR/x-ray crystallography 468, 474-5, 483-9
 poling 499-501
 structural databases 482, 489-90, 493-4, 499
 variations on standard methods 477-9
- conformational changes in molecular dynamics simulation 392-3
- conformationally flexible docking 662-3
- CONGEN program 542
- conjugate gradients 262, 264-7, 473
- conjugate peak refinement 290-1
- connection table 643
- consistent force field (CFF) 231
- constraints
 chiral 473-4
 constraints, holonomic versus non-holonomic 370
 in molecular dynamics 368-74
 simulation 368-74
 and restraints, difference between 369-70
 subset selection 717
 systematic search 649-51
- contact surface 7
- 'continuous' models of polymers 428-31
- continuum models and solvation, free energy of 592-3, 598-601
- contraction, basis set 69
- convergence sphere 186
- coordinates 2-4
 internal 2-4, 257
- intrinsic reaction 288-9
 mass-weighted 274-5
 scaled 438-9
see also Cartesian coordinates
- Corey-Pauling-Koltun (CPK) models 5-6
- CORINA program 659
- correlation
 BLYP 135, 136, 137
 coefficients 374, 681
 electron 110-17
 exchange-correlation functional 129-34
 functions and molecular dynamics simulation 374-80
 spectroscopy (COSY) 474-5, 486
- Cosine coefficient 676, 685
- COSMO (conductor-like screening model) 597
- COSY (correlated spectroscopy) 474-5, 486
- Coulomb interaction/integral 598-9
 advanced *ab initio* methods 122, 127, 128, 132, 133-4, 146-7
 computational quantum mechanics 30, 42, 45, 49-53, 58, 60, 85, 100
 force fields 184-5, 187, 194, 202, 238
- Coulomb potential 167, 244, 338, 341-2
- Coulomb's law 95, 194, 202, 212, 596, 603-4, 607
- counterpoise correction 121-2
- coupling parameter 567
- Craig plot 681
- cross-correlation function 376
- crossover operator 480-1
- crystal momentum 148
- crystal structures, predicted 501-5
- CSD (Cambridge Structural Database)
 conformational analysis 482, 489, 493-4, 499
 new molecules 659, 691, 693
- Cu-Zn superoxide dismutase 607
- cut-offs in computer simulation 324-7
 group-based 327-30
 problems with 330-4
- cyclic urea, HIV protease inhibitor 691-3
- cyclobutane 176-7
- cyclobutanone 176
- cyclobutene 117
- cycloheptadecane 476
- cyclohexane 286, 463, 465, 497, 597, 644
- cyclopropene 117
- cyclosporin 196-7, 391
- cysteine 511, 525, 556-7
- cytosine 84, 227
- databases 489, 537, 539
 3D 659-61, 679
see also CSD; structural databases
- Davidon-Fletcher-Powell method 269-70
- de Broglie thermal wavelength 411, 440-1
- de novo* ligand design 687-94

- degrees of freedom 699-700
- delocalised π -systems, force fields for 233-4
- DelPhi program 604-5
- density
 charge density matrix 58-9
 electron 77-9, 80-2, 160
 functional theory 126-37, 156, 619
 of levels 154
 spin 129-31
 of states and Fermi surface 153-5
- depth-first search 462, 663
- derivative, energy
 calculating 120-1
 function 225-6
 minimisation 257-8, 261-2, 268-9
- descriptors and new molecules 668-79
- determinant of matrix 13-14
- DFP (Davidon-Fletcher-Powell) method 269-70
- DFT (density functional theory) 126-37, 156, 619
- DHFR (dihydrofolate reductase) 278-9, 320, 460
- diagonalisation of matrix 16
- diatomic overlap *see* MNDO
- 1,2-dichloroethane 387-8
- Dick-Overhauser shell model 239
- dielectric constant 297
- dielectric models 202-4
- Diels-Alder reaction 294-5, 615, 716-17
- differences, free energy 564-74
 applications of methods 569-74
 formula for 568, 630-1
 methods for calculating 564-9
- differential overlap
 neglect of 86, 89-96
 zero (ZDU) 88-9, 91
- dihedral angle, definition 4
- dihydrofolate reductase 278-9, 320, 460
- DIIS (direct inversion of iterative subspace) 118
- dimensionality reduction 497-9
- dimethyl formamide 613-14
- dimethyl thioether 676-7
- N,N*-dimethyl-ketopropanamide 230
- dipole 75-7
 force fields 181, 182-5, 189, 199-201, 219, 246
 models of solvation, free energy of 593-5, 601-3
 moment, net 378-9
 triple 213-15, 239
- dipole correlation time 378-8
- direct inversion of iterative subspace 118
- direct SCF method 118-20
- director 396
- discriminant analysis and QSAR 703-5
- dispersion curve 298-9
- dispersive interactions 204-6
- displacements 322-3, 624
- dissimilarity-based methods 683-5
- dissipative particle dynamics 402-4
- distance
 bounds 468
 geometry 467-75, 476, 651-3, 663
 Hamming (city block) 492, 676-8
 map 651
 matrix 652
 Soergel 676-8
- distance-dependent dielectric 203
- distributed multipole analysis 195-7
- diverse sets of compounds, selecting 680-7
- DMA (distributed multipole analysis) 195-7
- DMF (dimethyl formamide) 613-14
- DNA 197, 227, 452, 489, 604
 computer simulation 319, 338-9
 Human Genome Project 512, 548-9
 inhibitor 270-1
 new molecules 662
 proteins 509, 512, 549
- DOCK program/algorithm 662-3, 665, 667
- docking 661-8, 689
- domain 515
- D-optimal design 697-8
- double dynamic programming 537-9
- double zeta basis sets 70
- double-wide sampling 567-8
- DPD (dissipative particle dynamics) 402-4
- Dreiding models 5-6
- Drude molecules 205-6
 interaction between 246-7
- drugs, new *see* new molecules
- dual topology 578
- dummy atoms, in Z-matrix 271-2
- Dunning basis sets 73
- dynamic/dynamics 353-409
- dynamically modified windows 578
 programming and protein prediction 526-9
 and statics in energy minimisation 295-300
see also molecular dynamics
- EA (evolutionary algorithms) 479-83
- edges of search trees 461
- effective medium theory 243-4
- effective pair potentials 214-15
- eigenvalues and eigenvectors 15-16, 114
 conformational analysis 469, 471, 479, 498
 energy minimisation methods 272, 282-5
- Einstein relationships 381, 627-8
- Eisenberg's 3D profiles 543-4
- elastic constants 240, 296-7
- electric multipoles, calculation of 75-7
- electron
 affinity (EA) 194
 correlation 110-17
 density 77-9, 80-2, 160

- electron (*cont.*)
 gas theory 240
 integrals, one- and two- 50-1
 nearly free-electron approximation 147-53
 polyelectronic atoms and molecules 34-41
 spin 34-5, 38
 electronegativity 192-3
 electrostatic interactions
 force fields 166, 181-204, 221, 237
 free energy calculations 566-7, 576, 580, 588, 613
 potentials 83-5, 188, 189-91
 solvation free energy calculations 593-608
 electrotological state index 674
 embedded-atom model 241, 243-4
 embedding 469
 empirical bond-order potential *see* Tersoff
 endothiapepsin 589-91
 energy
 calculation from wavefunction 41-6
 of closed-shell system 51
 component analysis 122-4
 computer simulation 308, 348-9
 conservation in molecular dynamics
 simulation 359, 405-6
 derivatives, calculating 120-1
 force field 240
 function, derivatives of 225-6
 of general polyelectronic system 46-50
 global minimum 253, 458, 479-83, 551-2
 Koopman's theorem and ionisation
 potentials 74-5
 lower-energy regions 564
 minimum, global 253, 458, 479-83, 551-2
 potential 4-5, 238, 253
 strain 226-7, 627
 surface (hypersurface) 4-5, 253, 475
 units of 9
see also derivative; energy minimisation; free energy; quantum mechanics
 energy minimisation methods 253-302, 623
 applications of 273-9
 choice of 270-3
 derivative 257-8, 261-2, 268-9
 first-order 262-7
 Newton-Raphson 267-8, 270, 288
 non-derivative 258-61
 quasi-Newton 268-9
 solid-state systems 295-300
 statement of problem 255-7
 transition structures and reaction pathways 279-95
 enol borate/aldehyde reaction 610-11
 ensemble
 averages 303-5
 distance geometry 651-3
 molecular dynamics 653
 Monte Carlo simulation 438-42, 450-1
 enthalpy 159, 574
 entropy 574
 enumeration of libraries 715-17
 equilibration monitoring and computer simulation 321-3
 equilibria phases in computer simulation 315, 450-1
 ergodicity 304, 313
 quasi ergodicity 433-8
 error, estimating in a simulation 343-7
 ESS (explained sum of squares) 699-700
 ethane
 carbon-carbon bond 4-5
 force fields 174
 Monte Carlo simulation 441-2
 SMILES notation 644
 thiol 564-9
 torsion angles 286-7
 Z-matrix 2-3, 9
 ethanol 564-9, 611-12
 ethene 83, 236, 293-5
 ethylene 621
 ethyne 83
 Euclidean distance measure 492
 Euler angles 421-2
 even-tempered basis set 71-2
 evolutionary algorithms 479-83
 evolutionary design 694
 evolutionary planning (EP) and strategies (ES) 479-80, 482
 Ewald summation 238, 334-9, 342, 402, 625-6
 exchange
 -correlation functional 129-34
 forces *see* repulsive forces
 gradient 135, 136, 137
 integral 50, 52-3, 58, 60
 interaction 46
 exclusion spheres 658-9
 exons 512
 explained sum of squares 699-700
 extended Hückel theory (EHT) 101-2
 extended system method 384
 extreme value distribution 532
 fabric softeners 401-2
 face-centred cubic lattice 139, 316
 factor analysis 681-2, 686
 factors (variables) 697
 family (proteins) 539
 fast Fourier transform 24, 338-9, 342
 fast multipole method 341-3, 364
 FASTA 524, 531
 FDPB (finite difference Poisson-Boltzmann method) 604-8

- Fermi surface and energy 153-5
 ferrocene 234
 FFT (fast Fourier transform) 24, 338-9, 342
 Fick's laws 380-1
 finite difference methods 355-8, 604-8
 Finnis-Sinclair potential 241-5 *passim*
 first principles method for predicting proteins 517-22
 first-order energy minimisation 262-7
 fitting, molecular 490-1
 flexible fitting 491
 flexible molecules 423, 582-5
 FlexX program 667
 Fock matrix 100
 Hartree-Fock equations 57-9, 61, 63-4
 open-shell systems 108-9
 semi-empirical methods 89-90, 94, 95, 96
 solid state quantum mechanics 146
 Fock operator 53, 57, 114
 focusing 606
 folding *see under* proteins
 force field models, empirical 165-252, 610
 angle bending 166, 173
 bond stretching 166, 170-3
 Class 1, 2 and 3 178-80
 cross terms 178-80
 derivatives of molecular mechanics energy function 225-6
 Drude molecules, interaction between 246-7
 effective pair potentials 214-15
 general features 168-70
 hydrogen bonding 215-16
 improper torsions and out-of-plane bending 176-8
 inorganic molecules 234-6
 many-body effects in empirical potentials 212-14
 metals and semiconductors 240-5
 parameters 221, 224-5, 228-32
 π systems, delocalised 233-4
 simple 165-6
 solid-state systems 236-40
 thermodynamic properties calculated using 226-8
 torsional terms 173-6
 united atom, reduced representations and 221-5
 water simulation 216-20
see also non-bonded interactions
 force-bias Monte Carlo method 432-3
 formaldehyde 76-7, 236
 formamide
 electron density around 78-9, 81-2
 gradient vector path 81
 HOMO and LUMO for 79
 forward sampling 567
 Fourier
 analysis 379, 392-3
 coefficient 148-9
 series 21-3, 155, 235, 237, 392
 transform 22-4, 392
 fractional factorial design 697
 fragments 472
 binding 689-91
 conformational analysis 464-5, 472
 locating 687-9
 marking 715-17
 free energy calculations 563-639
 chemical reactions 610-22
 computer difficulties 563-4
 enthalpy and entropy differences 574
 linear response method 631-2
 partitioning 574-6
 pitfalls 577-70
 potentials of mean force 580-5
 rapid methods, approximate 585-92
 solid-state defects 622-30
see also differences, free energy; Helmholtz free energy; solvation
 freely rotating chain model 428-9
 Frenkel defect 623, 626
 friction coefficient 388-9
 frontier orbitals 293
 full configuration interaction 112
 full factorial design 697
 fullerenes 101
 functional genomics 512
 future 160-1, 718-19
 GA *see* genetic algorithms
 gap penalties 526-8
 Gasteiger-Marsili approach 192-3
 Gaussian functions/distribution 20-1
 basis sets 65-73 *passim*, 120, 137, 195
 computer simulation 336-7, 339
 conformational analysis 481-2
 density functional theory 131-2
 force fields 195-6
 Gaussian-3 (G3) theory 116-17
 many-body perturbations 116-17
 molecular dynamics simulation 365, 381, 384, 389-90
 new molecules 679, 703
 proteins 551
 SCF 119
 semi-empirical methods 92-8
 solid state quantum mechanics 146
 Gay-Berne potential 222-5
 GB (generalised Born equation) 598-601
 surface area model (GB/SA) 609
 Gear algorithm 358-9
 general polyelectronic systems 38-41, 46-50

- generalised coordination 370
 generalised valence bond 125
 Generating Optimal Linear PLS Estimations 711
 generator matrices 429
 genetic algorithms
 conformational analysis 479–82
 new molecules 653, 663, 691, 701
 genomics 512, 548–9
 geometry 658–9
 distance 467–75, 476, 651–3, 663
 germanium 159–60, 244
 Gibbs ensemble Monte Carlo method 439, 450–1
 Gibbs free energy 563, 569
 global energy minimum 253, 458, 479–83, 551–2
 D-glucose 575
 glutamic acid 510, 525, 556–7
 glutamine 510, 525, 556–7
 glycine 221, 459, 511, 525, 556–7
 Go–Scheraga chain closure algorithm 541–2
 goal nodes 461
 GOLD program 667
 GOLPE (Generating Optimal Linear PLS Estimations) 711
 gradient
 -corrected functional 134–5
 exchange 135, 136, 137
 vector path 80–1
 grand canonical Monte Carlo simulations 440–2
 graphics, molecular 5–6
 graphite, adsorption 441–2
 graphs 642–3, 654
 Green–Kubo formula 382
 GRID program 215, 687–8, 708, 711
 grid search 459, 505
 GROMOS program 330
 Grothuss mechanism 620
 Group 14 elements 244
 solid state quantum mechanics applied to 158–60
 see also carbon; germanium; silicon
 group average 493–4
 group-based cut-offs 327–30
 G3 theory 116–17
 guanine 227
 GVB (generalised valence bond) 125
- haemagglutinin 667
 haemerythrin 544
 halides 237, 239, 504, 571
 halogenated hydrocarbons 707
 Hamiltonian operator
 advanced *ab initio* methods 114–15, 120–1
 computational quantum mechanics 27–9, 30, 32, 36, 42, 46, 53, 90–2
 computer simulation 312, 313, 410
 force fields 246
 free energy calculations 565, 567–9, 574, 577–9, 586, 595–6, 614
 Hammett substituent parameter 695–7
 Hamming (city block) distance 492, 676–8
 hard-sphere model 353–4
 harmonic approximation 278
 harmonic potential *see* Hooke's law
 Hartree atomic unit 29
 Hartree product 38–9
 Hartree–Fock equations/theory 51–65, 85
 application 65
 closed-shell 109
 configuration interaction 112–13
 density functional theory 126, 128, 129, 135–7
 free energy calculations 615–16
 LCAO 56
 many-body perturbation 114–15, 116
 RHF 108–10
 SCF 75, 87, 119
 Slater's Rules 54–6
 solid state quantum mechanics 146–7
 two-electron integrals 19
 see also Fock; Roothaan–Hall equations; UHF
 hashed fingerprint 645–6, 677, 678, 685
 heat bath, molecular dynamics 384
 heat capacity and computer simulation 308–9, 348–9
 Heitler–London model of hydrogen 124–5
 helium 36–8, 39
 hydrogen molecular ion (HeH^+) 62–5
 Slater determinant for 41
 helix 513–15, 583, 584
 Hellmann–Feynman theorem 121
 Helmholtz free energy 299–300, 411, 563–9
 computer simulation 307, 313–14
 Hessian matrix 267–9, 274–5, 280, 282–5, 288
 heterovalent substituent 623
 heuristic searches and protein prediction 531–4
 hexane 449, 462, 463, 672
 hexapeptide 482
 HF 123, 189, 196
 hierarchical cluster analysis 494, 534
 high throughput screening (HTS) 641
 high- T_c superconductor $\text{YBa}_2\text{Cu}_3\text{O}_{7-x}$ 628–30
 Hill potential 209
 histidine 169–70, 510, 525, 556–7
 HIV-1 protease 666, 667, 691–3
 HMMs (Hidden Markov Models) 536–7, 548
 Hodgkin–Richards index 679
 Hohnberg–Kohn theorem 128
 holonomic constraints 370
 HOMO (highest occupied molecular orbit) 79, 112, 293–4
 Hooke's law 172–3, 275, 486
 HP model 518–19
 HTS (high throughput screening) 641

- Hückel theory 99–102
 Human Genome Project 512, 548–9
 Hund's rules 35
 Hunter–Saunders approach 197–8
 hybrid Monte Carlo/molecular dynamics methods 452–3
 hydrodynamic vortex 377–8
 hydrogen 70
 bonding 122–3, 291–2, 391, 578, 689
 bond order 83
 C–H bonds/interactions 98, 167, 180, 211, 233, 236, 362
 conformational analysis 490, 504–5
 force fields 196, 215–16, 221, 227–8
 new molecules 655, 658
 O–H 98, 620
 configuration interaction 112
 dissociation 109–10
 electron correlation 110–11
 energy minimisation methods 282–3, 291–2
 fluoride 81–2
 Heitler–London model of 124–5
 molecule 41–6
 -suppressed notation 644
 hydrogen fluoride 620
 hydrophobic effect 515–17, 518–19, 669
 hysteresis 577
- iceberg model 516
 ID3 algorithm 705
 ILP (inductive logic programming) 705
 image charge computer simulation 340–1
 immunoglobulin 544
 immunosuppressant FK506 229–30
 importance sampling 410
 independent (random) samples 345
 indicator variable 696
 INDO (intermediate neglect of differential overlap) 86, 92–3, 94–6
 inductive logic programming 705
 initial configuration, prior to simulation 315–16
 inorganic molecules, force fields for 234–6
 Inorganic Structural Database 489
 inside-out ligand design 687–8
 integration
 algorithms for molecular dynamics simulation 359–60
 calculating properties by 412–14
 thermodynamic 568–9, 574, 577, 630–1
 intermediate neglect of differential overlap 86, 92–3, 94–6
 intermolecular processes and energy minimisation 278–9
 internal coordinates 2–4, 257
 Internet and World Wide Web 9–10, 548, 553
 interstitials 622–3, 627
 intrinsic reaction coordinate 288–9
 inverse agonists 640
 inverse of matrix 15–16
 ionic solids, force fields for 238–40
 ionisation potentials 74–5
 IRC (intrinsic reaction coordinate) 288–9
 iron, liquid 621–2
 Isis system 645
 island model 481
 isodemic reactions 116
 isoleucine 511, 525, 556–7
 isomerism, subgraph 645
 isothermal–isobaric ensemble, definition 307, 563, 569
- Jahn–Teller effect 234
 Jarvis–Patrick algorithm 496–7, 683
 JBW (jumping between wells) 435
 jellium 244
 jump frequency 627–8
 jumping between wells (JBW) 435
 J-walking 533–5
- kappa shape/kappa-alpha indices 672–4
 keys, pharmacophore 674–6
 Kohn–Sham scheme/orbitals 128–9, 131, 132, 134, 135–7, 156, 157, 616
 Koopman's theorem 74–5
 Kronecker delta 30, 396
 kurtosis 680
- lag, *ab initio* molecular dynamics 618
 Lagrange multiplier 18, 52, 127, 191, 371–2
 lags 453
 Laguerre polynomials 31, 55
 lambda dynamics 585–8
 Langevin dipole method 601–3
 Langevin equation 388–9, 391, 580, 601–3, 616
 Langmuir–Blodgett films/layers 395, 400–2
 large structures, reaction path for 289–92
 large systems, deriving charge models for 191–2
 latent variables 706
 lattices
 models of polymers 424–8
 models for proteins 518–20
 solid state quantum mechanics 138–60 *passim*
 statics and dynamics in energy minimisation 295–300
 LCAO (linear combination of atomic orbitals) 41–2, 56, 100, 241
 LDA/LSDA (local (spin) density approximation) 130–1
 leap-frog algorithm 356–7
 least-squares approach 230–1
 leave-one-out 701

- Lee-Yang-Parr *see* BLYP
 Legendre polynomials 32
 length, units of 9
 Lennard-Jones potential 253
 computer simulation 305, 319, 324, 327, 331-3,
 341-2
 force fields 167, 207-10, 212, 214-16, 225-6,
 237
 free energy calculations 579, 586, 613
 molecular dynamics simulation 361, 368, 402
 Monte Carlo simulation 428, 439, 441, 448, 450
 LES (locally enhanced sampling) 575-6
 leucine 487, 511, 525, 556-7
 Levinthal paradox 550
 libraries, combinatorial 711-19
 LIE (linear interaction energy) 588-9
 line search in one direction 262-3
 linear combination of atomic orbitals 41-2, 56,
 100, 241
 linear congruential method 418-19
 linear interaction energy 588-9
 linear potential, piecewise 665
 linear regression 666, 698-9, 702
 linear response 588-9, 591, 631-2
 linkage methods 493
 lipids 338
 simulation of 397-400
 liquid crystals 222-3
 literature 9
 lithium 111, 238, 323, 626
 loadings 682
 local density approximation 130-1
 local spin DFT 129, 135
 locally enhanced sampling 575-6
 logP 668-70
 London force 204-5
 long-range correction 327
 long-range forces and computer simulation
 334-43
 long time-tails, molecular dynamics 377
 loop conformations 541-2
 Lorentz-Berthelot mixing rules 210
 low-mode search 478-9
 Löwdin population analysis 80
 lower-energy regions 564
 lowest unoccupied molecular orbit 79, 112, 293-4
 LR (linear response) 588-9, 591, 631-2
 LSDFT (local spin density functional theory)
 129, 135
 LUDI program 689
 LUMO (lowest unoccupied molecular orbit) 79,
 112, 293-4
 lysine 510, 525, 556-7
 MACCS system 645
 Maclaurin series 11
 magnesium 238, 623, 626
 many-body
 effects in empirical potentials 212-14
 perturbation theory 114-17
 potentials 241
 mapping
 adiabatic 286
 distance 651
 pharmacophore 648
 Ramachandran 459-60, 514, 543, 547
 marker atom 329-30
 Markov chain 414-15
 Markov models, hidden 536-7, 538
 Marsaglia random number generator 420, 453-4
 mass-weighted coordinates 274-5
 mathematical concepts 10-24
 complex numbers 16-18
 multiple integrals 19-20
 series expansions 10-11
 statistics 20-1
 see also eigenvalues; Fourier; Lagrange;
 matrices; vectors
 matrices 2-3, 9, 12-16, 415
 adjacency 647
 charge density 58-9
 distance 652
 elastic constant 296-7
 PAM 524-6, 531, 556-7
 positive definite 16, 258
 statistical weight 430
 stochastic 415
 see also Fock matrix; Hessian; Z-matrix
 maxima 273
 maximal segment pair 531-3
 maximum dissimilarity algorithms 683-4
 maximum likelihood method 657-8
 MaxSum and MaxMin 683-4, 685
 Maxwell-Boltzmann distribution 365, 367, 384
 Mayer bond order 83
 MC *see* Monte Carlo
 MCSCF (multiconfiguration SCF) 113
 MCSS (multiple-copy simultaneous search) 688
 MDL (Molecular Design mol) format 643-4
 mean field approach 307-9
 mean squared displacement 322-3
 mean square end-to-end distance, polymers 426
 mechanics, molecular *see* force field
 mesoscale modelling 402-4
 messenger RNA (mRNA) 509
 metals 147, 589, 607, 626, 649-50, 693
 force field potentials for 240-5
 met-enkephalin 517
 methane
 bond order 83
 force fields 189
 Monte Carlo simulation 441-2

- methane (*cont.*)
 octopole moment 76
 population analysis 79
 SMILES notation 644
 methanol 573
 methionine 511, 525, 556-7
 methoxypropazine 678
 methyl chloride 612-13, 614, 676-7
 2-methyl propane 644
 o-methylacetanilide 670
 methylalanine 583-4
 methylene group 160, 162, 330, 396, 448
 energy minimisation methods 280, 291
 force fields 181, 221
 4-methyl-2-oxetanone 137
 metric matrix 469
 metrisation 472
 Metropolis Monte Carlo simulation 306, 433,
 436, 437, 447
 conformational analysis 467, 505
 implementation 417-20
 new molecules 663, 685, 691
 proteins 518
 theoretical background 414-16
 microcanonical ensemble, definition 307
 MINDO/3 86, 94-6, 102-3
 minima 272-3
 minimal basis set 69-70
 minimisation *see* energy minimisation
 minimum image convention and computer
 simulation 324-34
 mixing rules 210
 MM2/MM3/MM4 programs 8, 615
 force fields 169-71, 173, 176, 179, 187, 211,
 233-4
 MNDO (modified neglect of diatomic overlap)
 86, 96-7, 98-9, 102-3, 192
 MOD function 418-19
 Modeller program 541, 549
 modified INDO (MINDO/3) 86, 94-6, 102-3
 modified neglect *see* MNDO
 molar refractivity 671
 molecular dynamics simulation 354-409, 623
 of chain amphiphiles 394-404
 computer simulation 305-6, 307
 conformational analysis 457, 475-6, 483-9
 conformational changes from 392-3
 constant pressure dynamics 385-7
 constant temperature dynamics 382-5
 constraint dynamics 368-74
 continuous methods 355-64
 energy conservation in 405-6
 ensemble 653
 free energy calculations 564, 572, 577, 579,
 581, 588, 616-22, 628
 Monte Carlo compared with 307, 387, 452-3
 new molecules 664
 proteins 552
 setting up and running 364-8
 simple models 353-4
 solvent effects 387-90
 time-dependent properties 374-82
 see also computer simulation
 molecular field analysis 708-11
 molecular fitting 490-1
 molecular fragments *see* fragments
 molecular modelling *see* advanced *ab initio*;
 computer simulation; concepts;
 conformational analysis; energy
 minimisation; force field; free energy;
 molecular dynamics; Monte Carlo; new
 molecules; proteins; quantum mechanics
 molecular orbital theories, semi-empirical 86,
 89-96, 102-3
 molecular surface *see* surface
 Möller-Plesset *see* MP
 moments theorem 241-2
 monomers 289-90, 423, 550
 new molecules 712-13, 717-18
 Monte Carlo
 configurational bias 443-50
 force-bias 432
 Grand canonical 440-2
 smart 432
 Monte Carlo simulation 410-56
 bias 432-3, 443-50
 chemical potential, calculating 442-3
 computer 306-7
 conformational analysis 457, 475-6, 479, 483,
 504-5
 density functional theory 130
 different ensembles, sampling from 438-42
 force fields 189
 free energy calculations 564, 577, 579, 588
 chemical reactions 613, 616
 PMF 581-2, 584
 solid-state defects 623, 628
 thermodynamic perturbation 572-3
 Gibbs ensemble 450-1
 integration, calculating properties by 412-14
 molecular dynamics compared with 307, 387,
 452-3
 molecules 420-3
 new 662-3, 685, 691
 polymers 423-31
 proteins 517-19, 551
 quasi ergodicity 433-8
 random number generators 418-20, 453-4
 see also computer simulation; Metropolis
 MOPAC program 8, 99
 Morgan algorithm 644
 Morokuma analysis 122-4

- Morse potential/curve 170–2, 210
 motifs 522
 Mott–Littleton method 623–4, 625–7
 MP (Möller–Plesset) perturbation theory 114, 115–16, 119
 MR (molar refractivity) 671
 MS (Murtaugh–Sargent) method 269–70
 MSP (maximal segment pair) 531–3
 Mulliken population analysis 79–80, 189
 multicanonical Monte Carlo simulation 435–8
 multiconfiguration SCF 113
 multiple integrals 19–20
 multiple linear regression 666, 699, 702
 multiple sequence alignment 534–7
 multiple-copy simultaneous search 688
 multipole
 electric, calculation of 75–7
 fast 341–3, 364
 models 195–7, 219
 multivariate problems 708
 Murtaugh–Sargent method 269–70
 mutation
 operator 480
 probability matrices for proteins 556–7
- naphthalene 233
 NCC (Nieser–Corongiu–Clementi) model 219–20
 NDDO (neglect of diatomic differential overlap) 86, 93–4, 95, 96
 nearly free-electron approximation 142, 147–53
 Needleman–Wunsch algorithm 526–9, 534
 neglect of differential overlap 86, 89–96
 neighbour lists 325–7, 493–6
 net dipole moment 378–9
 net (partial) atomic charges 157–9, 181
 netropsin 270–1
 neural networks and QASR 703–5
 new molecules **640–726**
 combinatorial libraries 711–19
 computer representations 642–7
 de novo structure based ligand design 687–94
 descriptors 668–79
 discovery of drugs 640–1
 diverse sets of compounds, selecting 680–7
 docking 661–8, 689
 partial least squares 702, 706–11
 similarity 668, 676–9
 3D 674–5, 687
 databases 659–61, 679
 pharmacophores 648–59, 674–5, 687
 searching 645, 647, 667–8
 similarity 678–9
 see also QSAR
 Newton–Raphson energy minimisation 267–8, 270, 288, 625
- Newton's laws 304, 309, 353, 366, 371
 niching 481
 nickel oxide 147
 nicotine/nicotinic pharmacophore 653, 678
 nitrogen 490
 amide 660
 basis sets 73
 bond order 83
 charge models 187–8
 distributed multipole model 196
 electrostatic potentials 188
 force fields 181
 substituents 693
 NM23 547
 NMR and X-ray crystallography 316
 conformational analysis 468, 474–5, 483–9, 490
 molecular dynamics simulation 379, 383, 395
 new molecules 647, 659, 661, 667, 689, 691, 693, 704, 713
 proteins 516, 522, 546–7, 552, 512.514
 nodes
 on graphs 642–3
 on search trees 461
 NOESY (nuclear Overhauser enhancement spectroscopy) 474–5, 486, 488
 non-bonded cutoffs 324–34
 non-bonded interactions 166, 181–212, 324
 cell multipole method for 341–3
 electrostatic 166, 181–204
 neighbour lists 325–7
 Van der Waals 166, 204–12
 non-derivative energy minimisation 258–61
 non-electrostatic contributions to solvation free energy calculations 608–9
 non-holonomic constraints 370
 non-periodic boundary methods 320–1
 normal distribution *see* Gaussian functions
 normal mode analysis and energy minimisation 273–8
 normal vibrational modes 274
 nuclear Overhauser *see* NOESY
 nucleic acids 196–7
- 1-octanol and water, partition between 668–9
 octopole 76, 181
 one-electron
 atoms 30–4
 integrals 50–1
 ONIOM approach 615
 Onsager dipole model 593–5
 open-shell systems 108–10
 operators 28–9, 53, 57, 114, 480–1
 see also Hamiltonian
 OPLS (optimised parameters for liquid simulations) 210, 228, 599

- orbital
 -based approach to band theory 142–6
 calculations, molecular 26, 41–51
 approximate theories 86
 energy of closed-shell system 51
 energy of general polyelectronic system 46–50
 hydrogen 41–6
 one- and two-electron integrals 50–1
 semi-empirical 86, 89–96, 102–3
 total electron density 77–9
 see also STOs
 electronegativity 192–3
 linear combination of atomic 41–2, 56, 100, 241
 virtual 61
see also Kohn–Sham
- order
 bond 81–3
 order, of integration algorithm 358
 parameters 321–2
 orientational correlation 379–80
 orthogonalisation, symmetric 60
 orthonormal wavefunctions 30
 oscillating charge 201–2
 out-of-plane bending 176–8
 outside-in ligand design 687–8
 overlap
 differential 86, 88–96
 forces *see* repulsive forces
 integral 52
 oxides 147, 238
 oxygen bonds/interactions 98, 237, 328, 620, 652
- pairwise potential models 240–1
 PAM matrices 524–6, 531, 556–7
 parameters 567, 599
 force field 221, 224–5, 228–32
 substituent 695–7
see also Verlet
- partial equalisation of orbital electronegativity 192–3
 partial least squares (PLS) 702, 706–11
 partial (net) atomic charges 157–9, 181
 partition/partitioning 683–5
 coefficients 572–3, 668–71
 electron density 80–1
 free energy 574–6
 pattern recognition 491–7
 Pauli principle 206
 PCA (principal components analysis) 497–9, 681, 686
 PCM (polarisable continuum method) 596–7, 598
 PDB (Protein Databank) 489–90, 539
 pdf (probability density function) 304, 541
- Pearson correlation coefficient 681
 penalty functions 483
 pentane 253, 430, 462, 582
 pepsin 545
 peptides/polypeptides 277, 423, 509, 515, 517–18, 520, 691
 conformational analysis 459, 482
 dynamic programming 527–8
 folding 552
 force fields 196–7, 221, 231
 free energy calculations 571, 583–4
 loop conformations 541–2
 peptoids 713
 'threading' 546
 see also amino acids; proteins
 percentage sequence identity 524
 pericyclic reactions transition structures 292–5
 periodic boundary conditions 317–19
 perturbation and free energy 566–8, 573–4, 582, 584, 595
 thermodynamic 564–6, 569–73, 577, 592
 perturbation theories 36, 114–17, 119
 pharmacophore
 mapping 648
 keys 674–6
 pharmacophores 647
 see also new molecules
 phase equilibria, simulation of 450–1
 phase problem, in X-ray crystallography 484
 phase space and computer simulation 312–15
 phenylalanine 169, 286, 511, 525, 542, 546, 556–7
- phonons and dispersion curve 298–9
 π systems 197–9
 benzene 126
 delocalised 233–4
 pivot algorithm 423
 plane waves 155–6
 PLS (partial least squares) 702, 706–11
 PM3 98–9, 102
 PMF (potentials of mean force) 387–90, 546, 580–5, 612–14
 point defect 622
 point-charge electrostatic models 187
 Poisson equation 133
 Poisson–Boltzmann equation 603–8
 polarisation/polarisable basis functions 71
 continuum method 596–7, 598
 electrostatic non-bonded interactions 199–202, 203
 energy component analysis 122
 force field models for simulation of water 218–19
 poling and conformational analysis 499–501
 polyatomic systems 210–12
 polyelectronic atoms and molecules 34–41

- polymers
 energy minimisation methods 289–90
 free energy calculations 621, 622
 molecular dynamics simulation 391, 404, 550, 551
 Monte Carlo simulation of 423–31
see also amino acids; peptides; proteins
 population analysis 79–80, 189
 porphyrins 197–8
 positive definite matrix 16, 268
 potential 156–7, 275, 486, 546
 computational quantum mechanics 74–5, 83–5
 computer simulation 305, 319, 324–34, 338, 341–2
 electrostatic 83–5
 energy 4–5, 238, 253
 force fields 167, 170–3, 188–92, 207–10, 212–17, 222–6, 237–8, 240–5
 free energy calculations 549, 579, 580–5, 586, 612–14
 ionisation 74–5
 of mean force *see* PMF
 models, pairwise 240–1
 molecular dynamics simulation 387–90
 Monte Carlo simulation 442–3
 new molecules 665, 666
 prediction of crystal structures 501–5
 predictive residual sum of squares 701
 predictor-corrector methods of molecular dynamics simulation 358–9
see also under proteins
 preferential sampling 432
 PRESS (predictive residual sum of squares) 701
 pressure 309, 385–7
 principal components analysis *see* PCA
 principal components regression 706
 probability density function 304, 541
 probability matrices for proteins 556–7
 product-based monomer selection 718
 production phase in simulation 315
 profile 535
 proline 221, 511, 525, 556–7
 PROMET 502–3
 propane 167, 644
 proteins 6, 423
 computer simulation 329–30, 338–9
 conformational analysis 475, 489–90
 force fields 192, 221
 free energy calculations 571
 predicting structure of 509–62
 acronyms and abbreviations 553–4
 basic principles 513–17
 comparative model 539–45
 comparison of methods 547–9
 databases, list of 555
 first principles methods 517–22
 folding and unfolding 512, 516–17, 539, 545–7, 549–53
 mutation probability matrices 556–7
 sequence alignment 522–39
 threading 545–7
 Protein Databank 489–90, 539
see also amino acids; peptides
 pseudo-acyclic molecules 463–4
 pseudopotentials 156–7
 pyrazine/pyridine 573
 2-pyridone 597–8
 Q² (cross-validated R²) 701
 QCISD (quadratic CISD) 113, 117, 119
 QSAR (quantitative structure-activity relationships) 695–706, 710, 711
 cross-validation 701
 deriving equation 698–70
 discriminant analysis 703–5
 interpreting equation 702
 neural networks 703–5
 principal components regression 706
 -property relationship 695, 702
 selecting compounds for analysis 697–8
 QSPR (quantitative structure-property relationship) 695, 702
 quadratic region 283–4
 quadrupole 76, 181, 183, 185–6, 196
 quantitative structure-activity *see* QSAR
 quantum mechanics
 future role 160–1
 and molecular mechanics combined in
 chemical reactions 614–16
 solvation, free energy of 594–8
see also ab initio quantum mechanics;
 advanced *ab initio*; computational
 quantum mechanics
 quasi ergodicity and Monte Carlo
 simulation 433–8
 quasi-Newton energy minimisation 268–9
 quaternions 422
 R² 699
 R groups 716–17
 radial distribution functions and computer
 simulation 310–12
 Ramachandran map 459–60, 514, 543, 547
 random number generators 418–20, 453–4
 random sampling 345–7
 random search 465–7, 476, 517–18
 random tweak 542
 range scaling 681
 ranitidine 489, 644
 RANTES 475
 rapid free energy calculations, approximate 585–92

- Rappé-Goddard method 193–4
 RATTLE method 373–4
 Rayleigh-Schrödinger perturbation theory 114
 reaction
 field 339–40, 595–6, 597
 isodemic 116
 pathways 279–95
 transform 715–16
 zone 320–1
see also chemical reactions
 real gas contribution to virial 309, 349–50
 reciprocal lattice 139–40
 recombination operator 480–1
 reduced units, in non-bonded interactions 212
 re-entrant surface 7
 refractivity, molar 671
 regression 706
 equation 698–9
 linear 666, 698–9, 702
 relative energies 226
 relaxation time 376
 reptation 427
 repulsive forces 206
see also Coulomb attraction/repulsion
 residual sum of squares 699–700
 RESP (restrained electrostatic potential fit) 191–2
 response 697
 restraints/restrained
 and constraints, difference between 369–70
 electrostatic potential fit 191–2
 molecular dynamics 483–4
 spatial, satisfaction of 540–1
 reversible reference system 363–4
 RHF (spin-restricted Hartree-Fock theory) 108–10
 ribose phosphate 493–4
 rigid molecules, simulation of 420–2
 rigid-body method 540
 ring critical point 81
 RIS (rotational isomeric state) model 429–31
 RMS (root-mean-square) 273, 359–60, 552, 667
 RMSD (root-mean-square-distance) 491–3
 RNA 509, 512
 root nodes 461
 root-mean-square 273, 359–60, 552, 667
 Rootaan-Hall equations
 closed-shell systems 56–9, 86–8
 density functional theory 132
 illustrated 62–5
 solving 59–62
 Rosenbluth weight 444–7
 rotational isomeric state 429–31
 rotational order 322
 roulette wheel selection 480
 r-RESPA (reversible reference system
 propagation algorithm) 363–4
 RSS (residual sum of squares) 699–700
 rule-based approaches to protein
 prediction 520–2
 saddle points 253, 272–3, 280, 282–3, 291
 location 285–8, 478
 quadratic region 283–4
 SAM1 (Semi-Ab-initio Model 1) 99, 102
 sampling 345, 346–7, 410, 432, 438–42, 567–8, 575–6
 SC24/halide system 571
 scaling/scaled
 autoscaling 681
 coordinates 438–9
 mesoscale modelling 402–4
 particle theory 609
 range 681
 scalar product and triple product 12, 14
 SCF (self-consistent field) 117, 280
 complete active-space 113, 295
 computational quantum mechanics 54, 64, 73, 75, 87
 direct method 118–20
 energy component analysis 122
 free energy calculations 595–6, 597
 Hartree-Fock 75, 87, 119
 multiconfiguration 113
 Schottky defect 622–3, 626
 Schrödinger equations and solutions to
 computational quantum mechanics 27–8, 29–30, 32, 34–7 *passim*, 128
 computer simulation 347–8
 density functional theory 127, 128
 for Drude molecules 205, 246–7
 solid state quantum mechanics 147, 148
 SCOP (Structural Classification of Proteins) 539
 scoring functions for docking 664–7
 SCRf (self-consistent reaction field) 595–6, 597
 SCRs (structurally conserved regions) 539–40
 SDEP measure 709
 search
 depth-first 462, 663
 grid 459, 505
 heuristic 531–4
 line 262–3
 low-mode 478–9
 multiple-copy 688
 new molecules (3D) 645, 647, 667–8
 random 465–7, 476, 517–18
 systematic 458–64, 476, 505
 trees 461–5
see also under conformational analysis
 second-moment approximation 242
 secondary structure of proteins 513
 segment matching 540
 self-consistent field *see* SCF; SCRf
 self-penalty walk (SPW) 289–90, 584

- Semi-Ab-initio Model 1 (SAM1) 99, 102
 semi-empirical methods of computational quantum mechanics 65, 86-99, 102-3
 semi-empirical molecular orbital theories 86, 89-96, 102-3
 semiconductors, force field potentials for 244-5
 separation of variables 36-7
 sequence alignment of proteins 522-39
 sequence identity 546-7
 sequential univariate minimisation 260-1
 series expansions 10-11
 serine 511, 525, 556-7
 SHAKE procedure 369-74, 582, 618
 shape anisotropy parameter 224-5
 SHAPES force field 235-7
 shear viscosity 381
 shielding constant 55-6
 shifted potential 330-1
 shorthand representation of electron integrals 50
 sigmoidal dielectric model 202-4
 silica 297-8
 silicalite 449-50
 silicon 483, 693
 -O bond 237
 and chlorine 620-1
 force fields 237, 245
 phases of 159-60
 shielding constant 55-6
 valence electron density 160
 similarity
 calculating 676-8
 searching 668
 and 3D properties 678-9
 simple force field models for simulation of water 216-18
 simplex method of non-derivative energy minimisation 258-60
 Simpson's rule 412-13
 simulated annealing
 in conformational analysis 483
 in *ab initio* molecular dynamics 616-18
 in X-ray refinement 484-6
 simulations *see* computer simulation; conformational analysis; molecular dynamics; Monte Carlo
 SINDO1 program 99
 single-linkage cluster algorithm 493-6
 site points 689
 skewness 680
 Slater determinants
 density functional theory 135, 136
 general polyelectronic systems 38-41
 many-body perturbation 115
 orbitals *see* STOs
 Slater functions and basis sets 67-9
 Slater's Rules and Hartree-Fock equations 54-6
 slow growth free energy calculations 568-9, 577, 631
 smart Monte Carlo method 432-3
 Smart Region Definition 710-11
 SMILES notation 643-5, 715
 Smith-Waterman algorithm 529-30
 S_N2 reaction
 potential of mean force 612-14
 transition state of 280-2
 sodium 181, 589, 626
 sodium chloride 238
 Soergel distance 676-8
 solid-state
 defects and free energy calculations 622-30
 energy minimisation methods 295-300
 force fields for 236-40
 quantum mechanical methods for studying 138-60
 solvation/solvents 320
 dielectric models of electrostatic non-bonded interactions 202-4
 free energy of 576, 592-610
 continuum models 592-3, 598-601
 electrostatic contributions 593-608
 non-electrostatic contributions 608-9
 simple models 609-10
 molecular dynamics simulation 387-90
 Monte Carlo simulation 432, 452
 space group 138
 spatial restraints, satisfaction of 540-1
 SPC model 216-18
 sphere-exclusion algorithm 684
 spherical cut-off 324
 spherical harmonic 30-1
 spin
 -coupled valence bond theory 125-6
 density 109
 local 129, 135
 orbitals 35
 -polarised density functional theory 129
 -restricted Hartree-Fock 108-10
 -unrestricted Hartree-Fock 108-10
 split valence double zeta basis sets 70
 SPV (self-penalty walk) 289-90, 584
 squares
 least-squares approach 230-1
 partial least 702, 706-11
 root-mean 273, 359-60, 552, 667
 square-well potentials 354
 sum of 699-701
 SRD (Smart Region Definition) 710-11
 ST2 potential 217
 standard deviation 20
 statistical inefficiency 346

- statistics/statistical 20-1
 mechanics and computer simulation 347-8
 weight matrix 430
 steady state genetic algorithm 481
 stearic acid 394, 400-1
 steepest descent energy minimisation 262
 step size 264
 steric energy 226
 Stillinger-Weber model 241, 244-5
 stochastic boundary conditions 320-1
 stochastic collisions method 384
 stochastic dynamics simulations 390-2
 stochastic matrix 415
 Stokes law 388
 STOs (Slater-type orbitals) 46-8, 55, 72
 density functional theory 131-2
 force fields 194
 STO- π G basis sets 62, 69, 85, 123
 strain 296
 energy 226-7, 627
 stratified sampling 346
 streptavidin 641
 stress 296
 Structural Classification of Proteins 539
 structural databases
 conformational analysis 482, 489-90, 493-4, 499
 proteins 537-9, 555
 structural genomics 512
 structural key 645
 structural properties, calculating 85-6
 structurally conserved regions 539-40
 structurally variable regions 539-40
 structure factor 484
 subgraph 642-3, 645
 subset selection 717-18
 substitutionals 623
 substructure search 465, 642-7
 sulphur dioxide 117
 sum of squares 699-701
 SUMM (systematic unbounded multiple minimum) 477-8
 superfamily (proteins) 539
 superoxide dismutase 607
 surface 6-8
 area model 609
 energy 4-5, 253, 475
 Fermi 153-5
 van der Waals 7, 600
 Sutton-Chen potential 241, 243
 SVRs (structurally variable regions) 539-40
 SWISSPROT database 537
 switching function 331-4
 symmetric matrix 13
 symmetric orthogonalisation 60
 systematic sampling 346-7
 systematic search 649-51
 conformational 458-64, 476, 505
 systematic unbounded multiple minimum 477-8
 Tanimoto coefficient 676-8, 685
 Taylor series 10-11, 230, 267, 342, 355, 358, 439, 592
 temperature 240
 computer simulation 309-10
 molecular dynamics simulation 368, 382-5
 template forcing 491
 templating effect 694
 tensor properties 183
 terminal nodes 461
 Tersoff model 241, 244-5
 thermodynamic(s)
 computer simulation 307-12
 cycles 569-70
 force fields 226-8
 integration 568-9, 574, 577, 630-1
 perturbation 564-6, 569-73, 577, 592
 properties 85-6, 307-12
 thermolysin 571-2
 thiazole 490
 Thomas-Fermi model 127
 threading, predicting proteins by 545-7
 three-body problem/effects 34, 212-14, 244
 3D
 profiles 543-4, 547
see also under new molecules
 3₁₀ 513, 583-5
 threonine 511, 525, 546, 556-7
 thrombin 522-3
 thymidylate synthase 667
 thymine 227
 'TIM barrel' 522-3
 time
 -averaged NMR 487-9
 averages 303-5
 correlation coefficients 374
 -dependent properties and molecular dynamics simulation 374-82
 step and molecular dynamics simulation 360-4
 TiN 155
 TIP3P/TIP4P models 216-17, 219, 327-8
 topological indices 671-4
 torsion/torsional 166
 angle/bend 3-4, 179-80, 254, 515
 driving 286-8
 improper 176-8
 parameters 229
 terms 173-6
 total electron density and molecular orbitals 77-9
 total sum of squares 699-700
 Toxvaerd anisotropic model 221-2
 transfer RNA (tRNA) 509

- transferability and force fields 168
 transition structures 255, 279-95
 transport and molecular dynamics 380-2
 transpose of matrix 15
 trapezium rule 412-13
 tree representation 461-2
 trial and error and parametrisation 228-9
 triangle smoothing 468-9, 660-1
sym-triazine 198-9
 1,3,5-trifluorobenzene 198-9
 trimethoprim 278-9
 truncating potential 324-34
 trypsin 522-3, 587, 607
 tryptophan 169, 511, 525, 556-7
 TSS (total sum of squares) 699-700
 Tversky similarity 677-8
 twin-range method 327
 2D substructure searching 642-7
 two-electron integrals 50-1
 tyrosine 286, 510, 525, 542, 556-7
- UFF (Universal Force Field) 193-4, 232, 235, 237
 UHF (spin-unrestricted Hartree-Fock theory) 108-10
 Ullmann algorithm 646-7
 umbrella sampling 581-2, 584
 underlying matrix 415
 unit cell 138
 united atom force fields and reduced representations 221-5
 Universal Force Field *see* UFF
 uracil-2,6-diaminopyridine (DAP) 227-8
 Urey-Bradley force field 179, 235
- vacancy 622, 627
 formation energy 240
 valence
 band 142
 bond theories 124-6
 electron density 160
 split 70
 valine 511, 525, 546, 556-7
 van der Waals
 energy 253
 interactions 320, 486
 force fields 166-7, 204-12, 231, 237
 free energy calculations 566-7, 576, 586, 588-9
 parameters 229
 potentials 666
 radii 84, 470, 592, 596, 660
 surface 7, 600
 variable metric *see* quasi-Newton
 variables (factors) 697
 variance 20
 variance-covariance matrix 498
- variation
 coefficient of 680-1
 theorem 51-2
 vectors 11-12
 path, gradient 80-1
 product 12, 14
 Veillard-Baron order parameter 322
 velocity autocorrelation 376-8
 velocity Verlet algorithm 357
 Verdier Stockmayer algorithm 427
 Verlet algorithm/parameters 321-2, 325-6, 355-9, 403-4
 see also SHAKE
 vibrational modes 274
 virial 309, 349-50
 virial theorem of Clausius 309
 virtual molecules 642
 virtual orbitals 61
 virtual screening 715
 VWN (Vosko-Wilk-Nusair) standard local correlation 136
- water
 bond order 83
 and carbon dioxide 616-17
 computer simulation 317, 327-8
 dimer analysis 123-4, 327-8
 force field models 201, 216-20
 free energy calculations 573
 differences 569-70
 see also solvation
 infrared spectra 379
 normal modes of 274
 and 1-octanol, partition between 668-9
 wavefunction 41-6, 161
 Wigner-Seitz cells 140, 350
 Wiswesser line notation 643-4
 Woodward-Hoffmann rules 102, 292-3, 295
 World Drug Index 685
 world wide web (WWW) 9-10
- X-ray crystallography *see* NMR
- YETI force field 215-16
- ZDO (zero-differential overlap) 88-9, 91
 ZEBEDDE (ZEolites By Evolutionary De novo DEsign) 694
 zeolites 298, 449-50
 force fields for 236-7
 synthesis 693-4
 zero-differential overlap 88-9, 91
 zero-point energy 274
 zinc 607, 649-50
 ZINDO program 99
 Z-matrix 2-3, 9, 74, 255, 271-2
 Zwanzig expression 588, 631-2

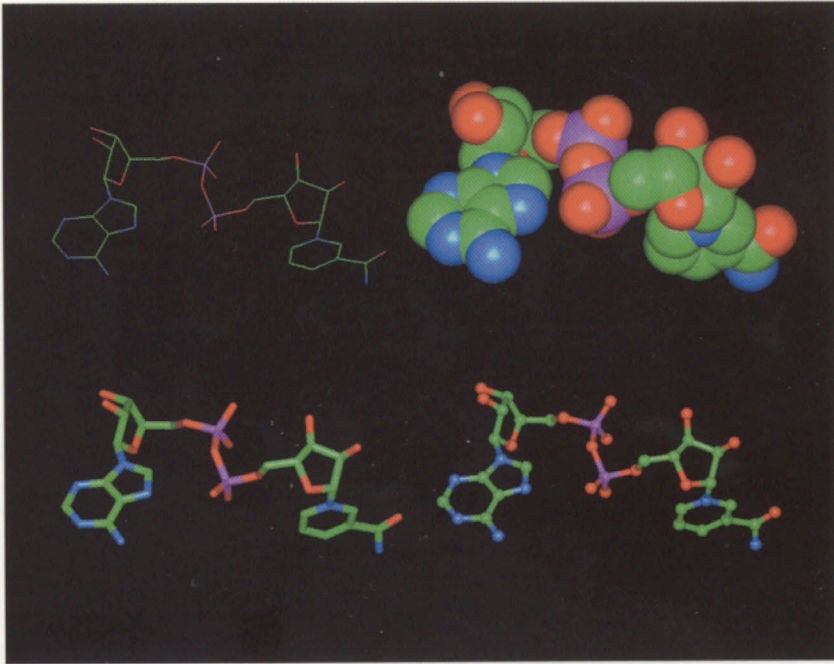


Fig. 1.4: Some of the common molecular graphics representations of molecules, illustrated using the crystal structure of nicotinamide adenine dinucleotide phosphate (NADPH) [Reddy et al. 1981]. Clockwise, from top left: stick, CPK/space filling, 'balls and stick' and 'tube'.

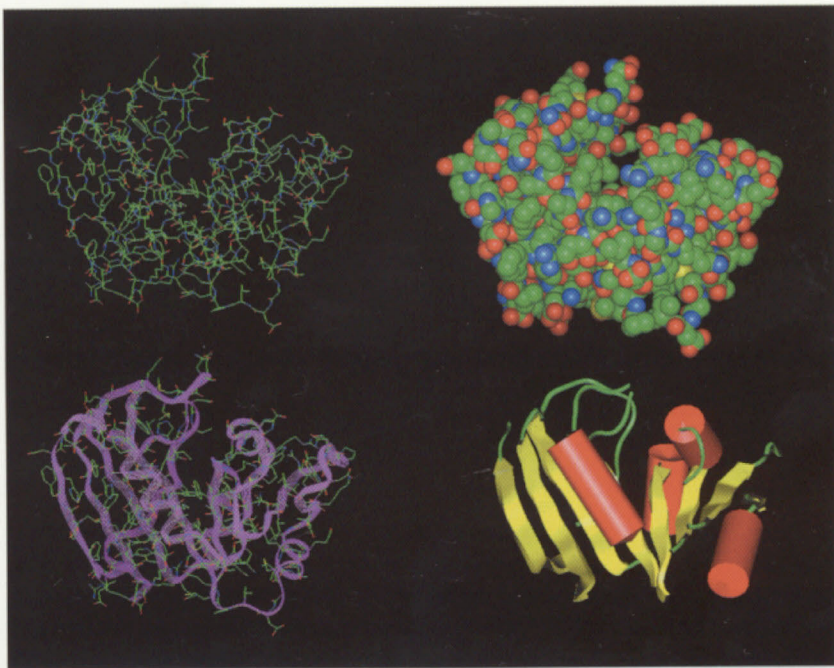


Fig. 1.5: Graphical representations of proteins illustrated using the enzyme dihydrofolate reductase [Bolin et al. 1982]. Clockwise from top left: stick, CPK, 'cartoon' and 'ribbon'.

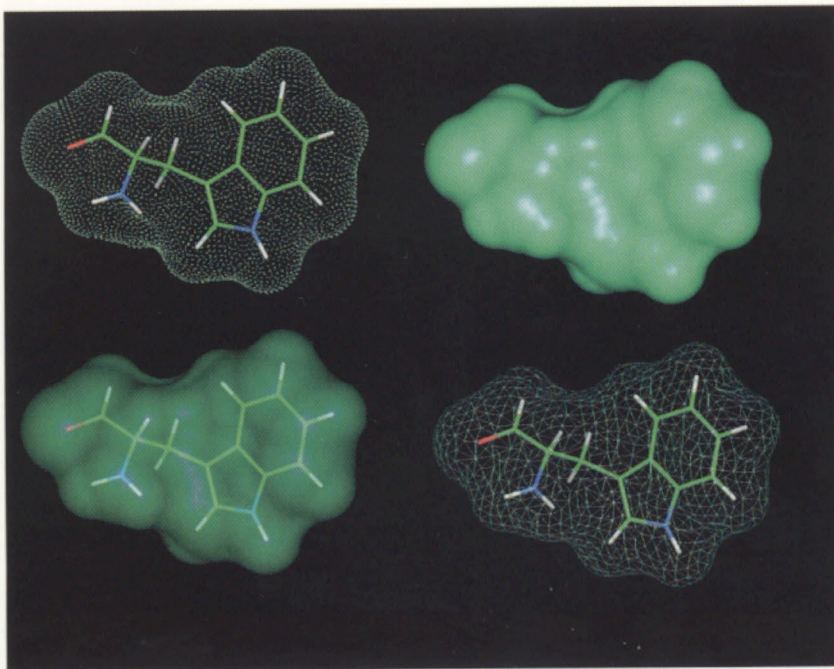


Fig. 1.7: Graphical representations of the molecular surface of tryptophan. Clockwise from top left: dots, opaque solid, mesh, translucent solid.

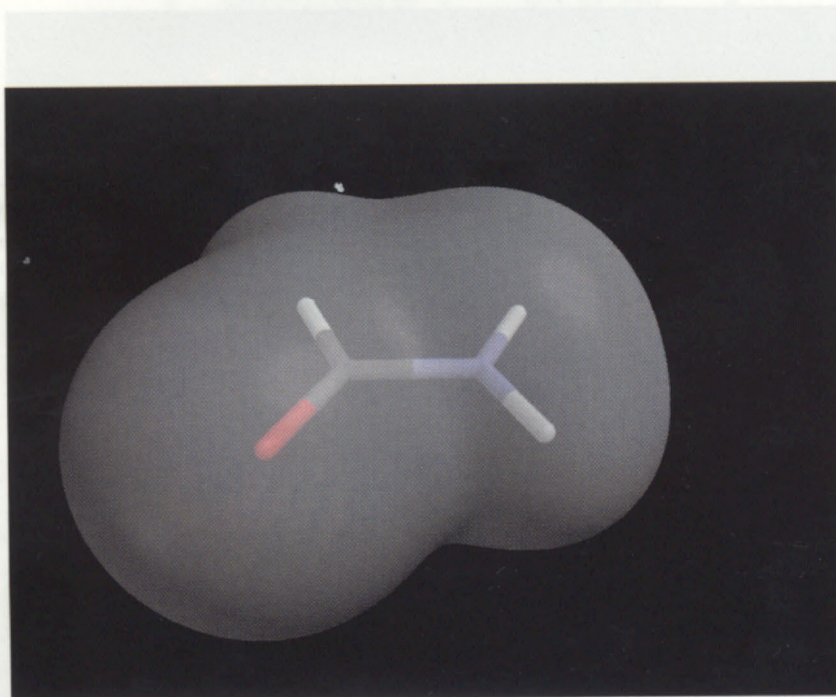


Fig. 2.11: Surface representation of electron density around formamide at a contour of 0.0001 au (electrons/bohr³).

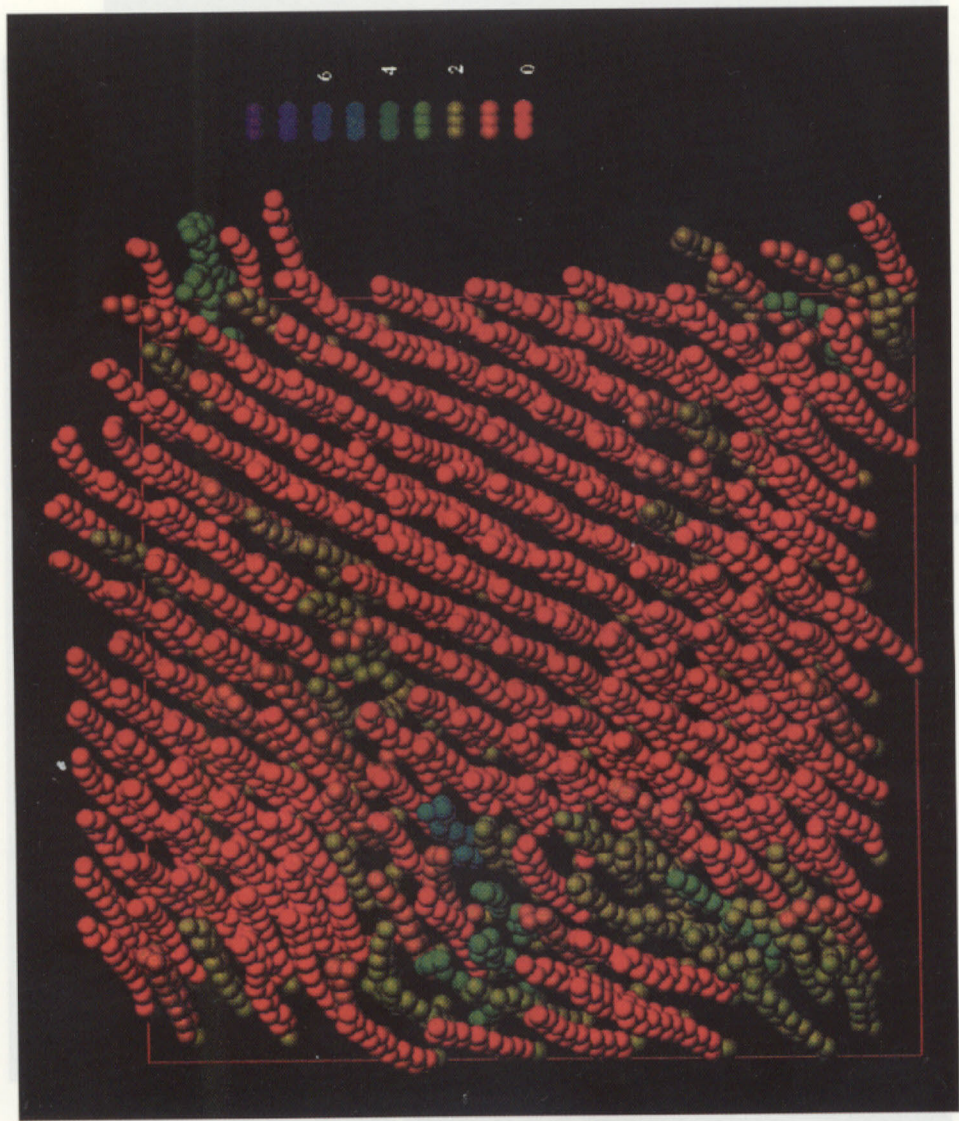
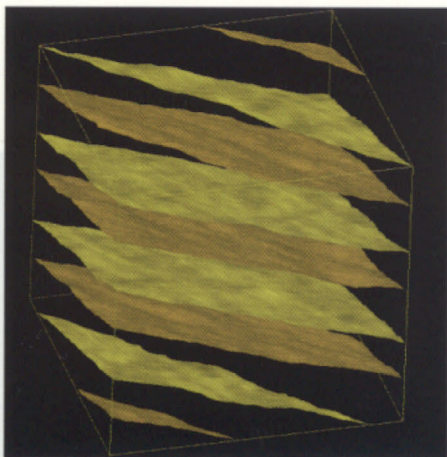
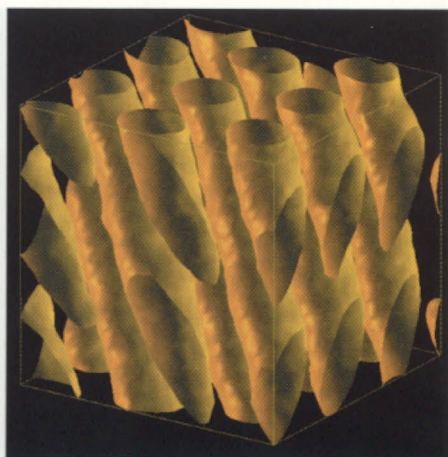


Fig. 8.21: Final configuration obtained from a Configurational Bias Monte Carlo simulation of thioalkanes absorbed on a gold surface [Siepmann and MacDonald 1993a]. The system contains 224 molecules which are colour coded according to the number of gauche defects, with red chains being all trans, yellow chain containing three gauche bonds and green chains containing five gauche bonds.

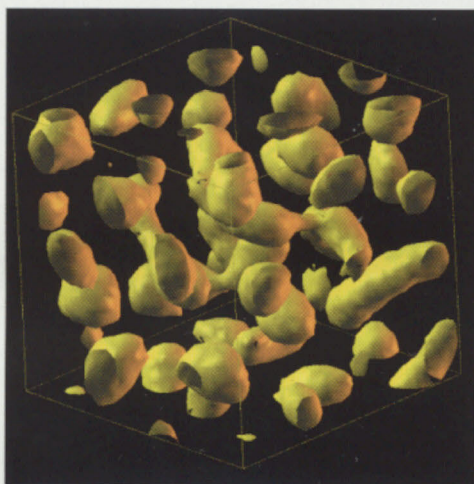
Fig. 7.24: Graphical representation of final configurations obtained from dissipative particle dynamics simulations on block copolymers. (a) shows the lamellar phase obtained for the A_5B_5 system, (b) the hexagonal phase from A_3B_7 and (c) the body-centred-cubic phase obtained for A_2B_8 . Figure redrawn from Groot, R. D. and Madden, T. J. 1998. Dynamic simulation of diblock copolymer microphase separation. *The Journal of Chemical Physics*, 108: 8713–8724.



(a)



(b)



(c)

Fig. 7.24: Graphical representation of final configurations obtained from dissipative particle dynamics simulations on block copolymers. (a) shows the lamellar phase obtained for the A_5B_5 system, (b) the hexagonal phase from A_3B_7 and (c) the body-centred-cubic phase obtained for A_2B_8 . Figure redrawn from Groot, R. D. and Madden, T. J. 1998. Dynamic simulation of diblock copolymer microphase separation. *The Journal of Chemical Physics*, 108: 8713–8724.

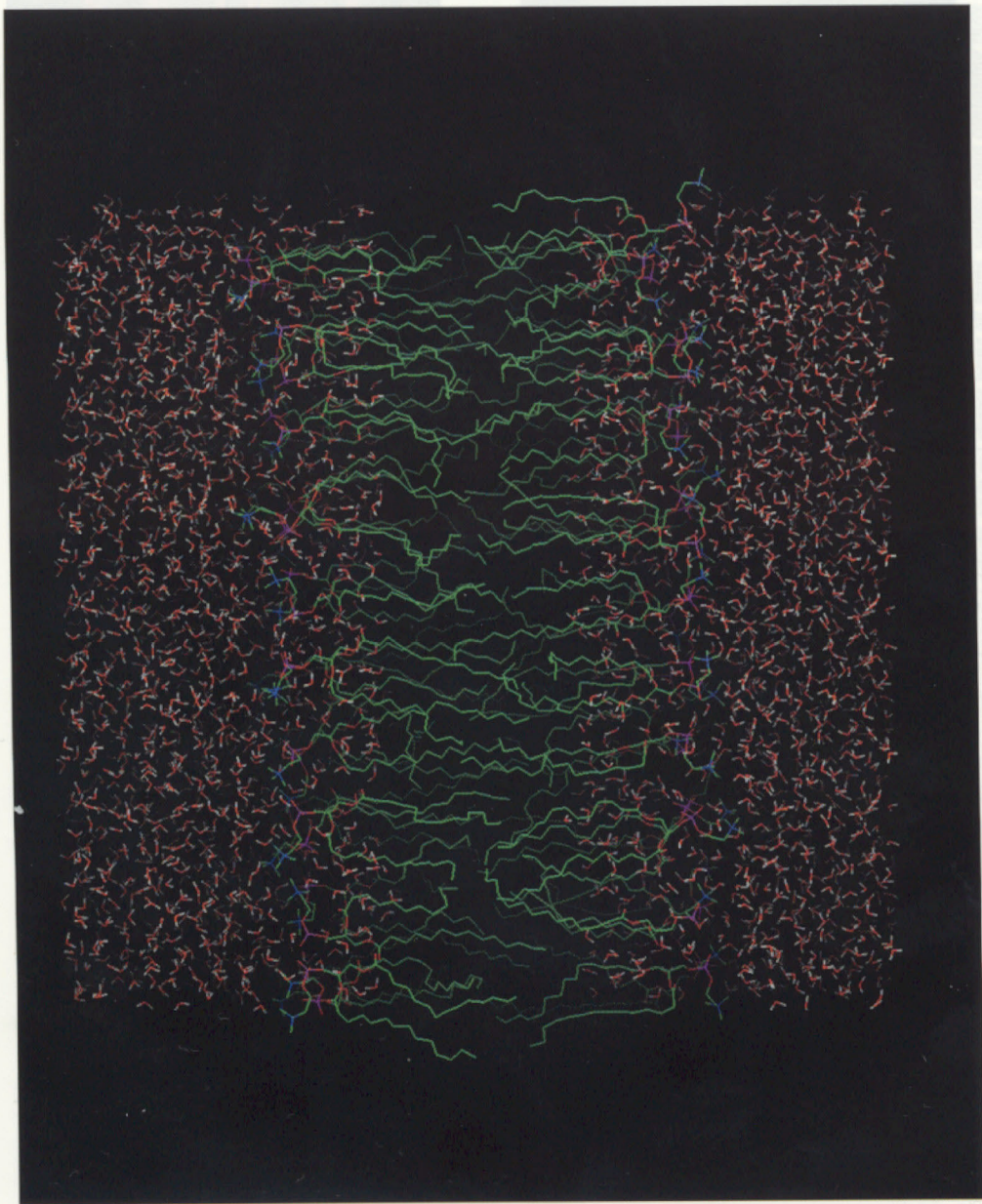


Fig. 7.21

Fig. 7.21: Snapshot from a molecular dynamics simulation of a solvated lipid bilayer [Robinson et al. 1994]. The disorder of the alkyl chains can be clearly seen.

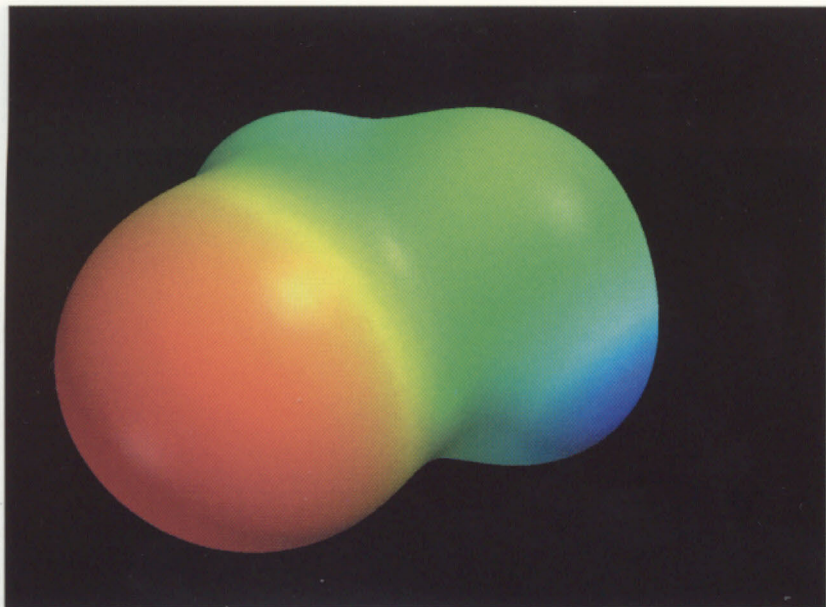


Fig. 2.18: Electrostatic potential mapped onto the electron density surface for formamide. The orientation of the molecule is as in Fig. 2.11. Red indicates negative electrostatic potential and blue is positive potential.

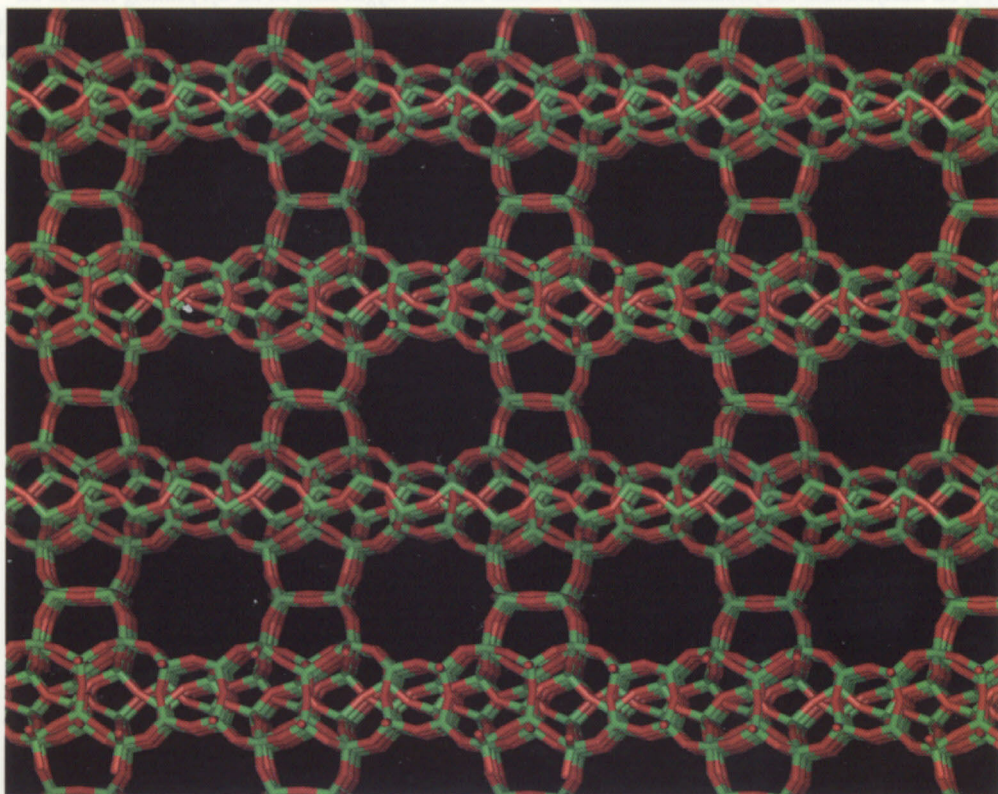


Fig. 5.36: The zeolite NU-87.

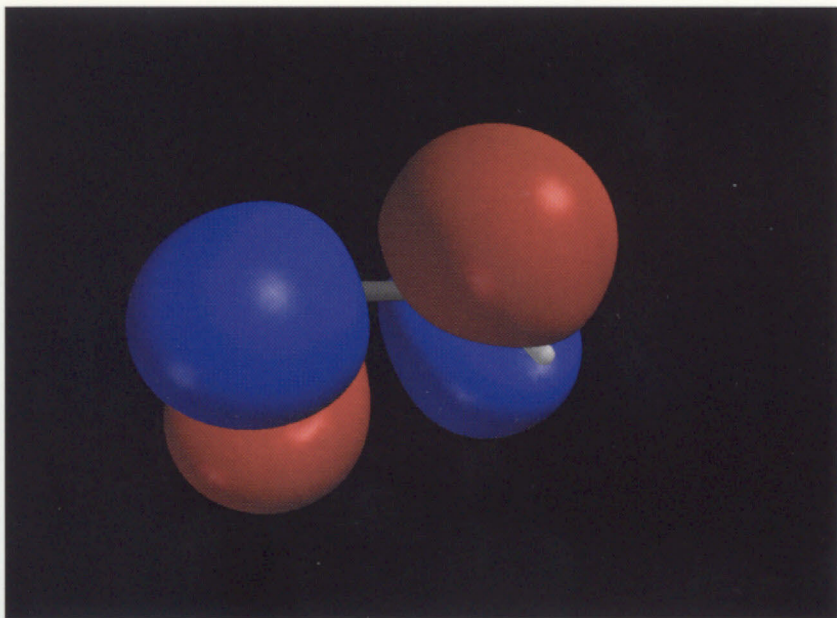


Fig. 2.12: HOMO of formamide. The red contour indicates the negative part of the wavefunction and blue the positive part of the wavefunction. The formamide molecule is oriented with the oxygen atom on the left pointing towards the viewer, as in Fig. 2.11.

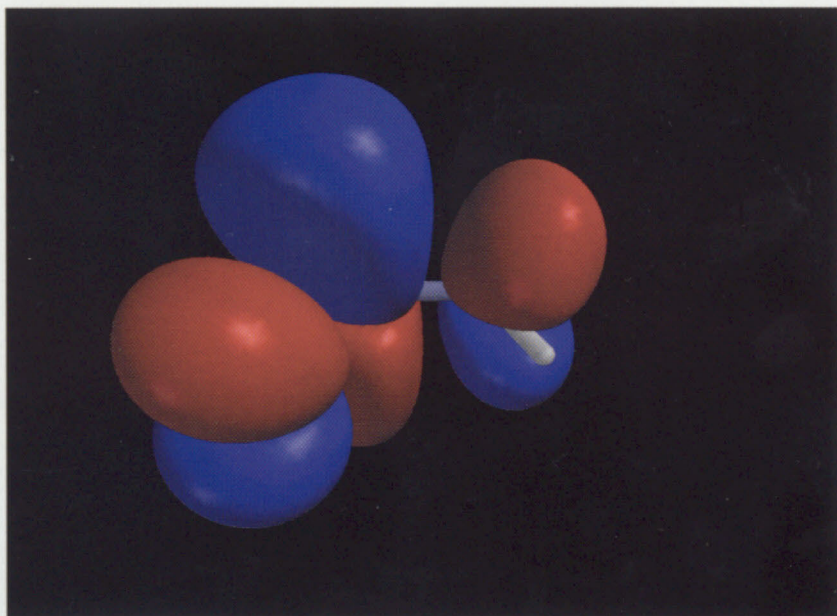


Fig. 2.13: LUMO of formamide.

2.11. Snapshot from a molecular dynamics simulation of a simulated liquid nitrogen. The number of the cell of chains can be clearly seen.

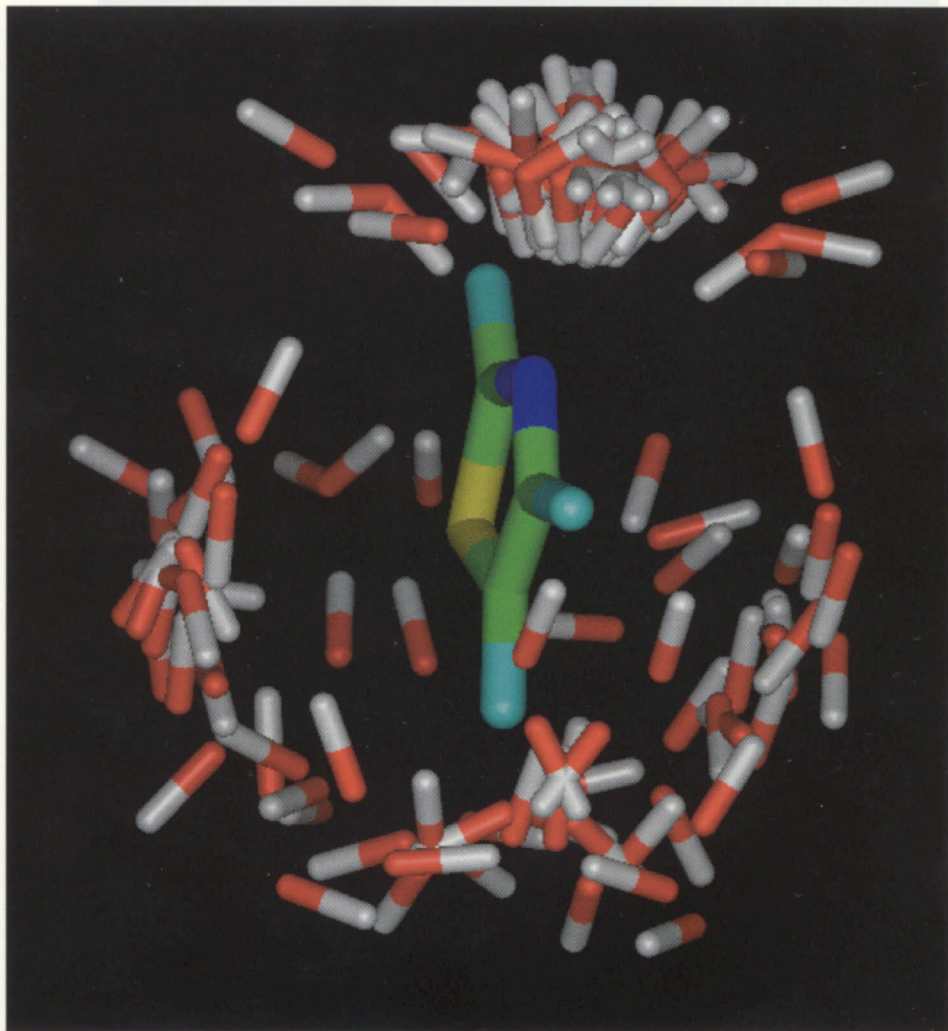


Fig. 9.27: Distribution of hydroxyl groups around thiazole ring systems as extracted from the Cambridge Structural Database [Bruno et al. 1997], illustrating the greater propensity of the nitrogen atom to act as a hydrogen-bond acceptor.

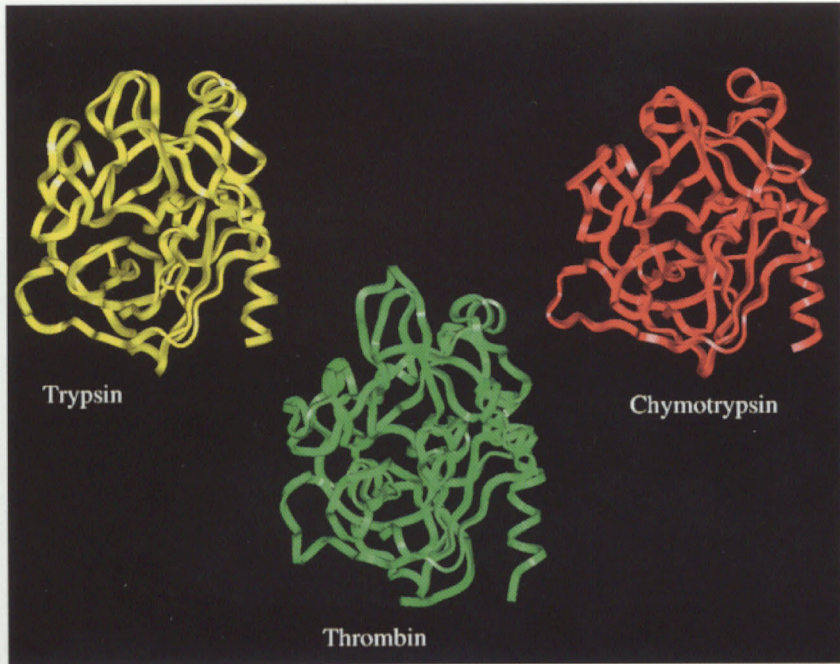


Fig. 10.9: Trypsin (top left) [Turk et al. 1991], chymotrypsin (top right) [Birktoft and Blow 1972] and thrombin (bottom) [Turk et al. 1992] have similar three-dimensional structures.

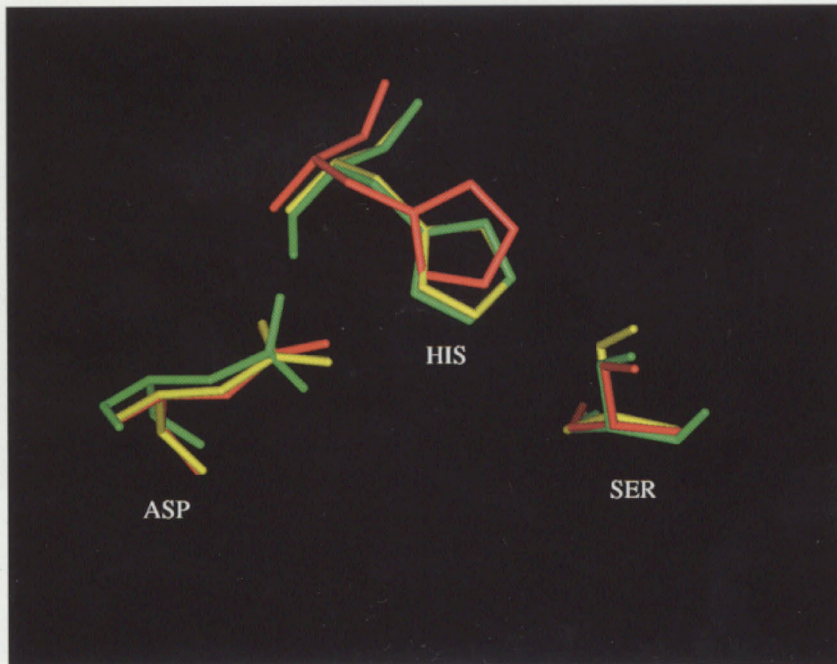


Fig. 10.11: A superposition of the aspartic acid, histidine and serine amino acids in the active sites of trypsin (yellow), chymotrypsin (red) and thrombin (green).

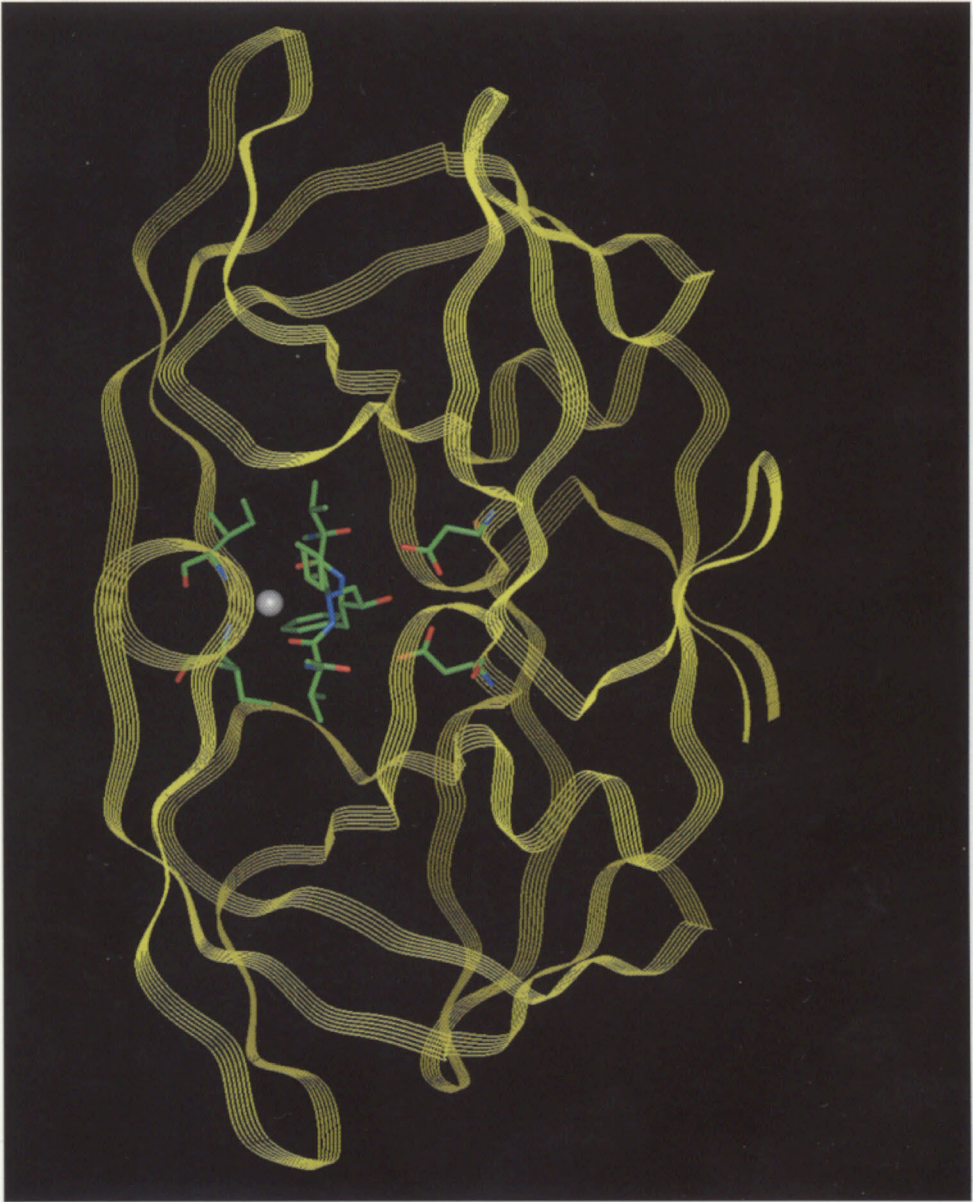


Fig. 12.34: The HIV-1 protease with the inhibitor CGP53820 bound [Priestle et al. 1995]. The water molecule that forms hydrogen bonds both to the inhibitor and to the 'flaps' of the protein is drawn as a white sphere and the catalytic aspartate groups are also represented.

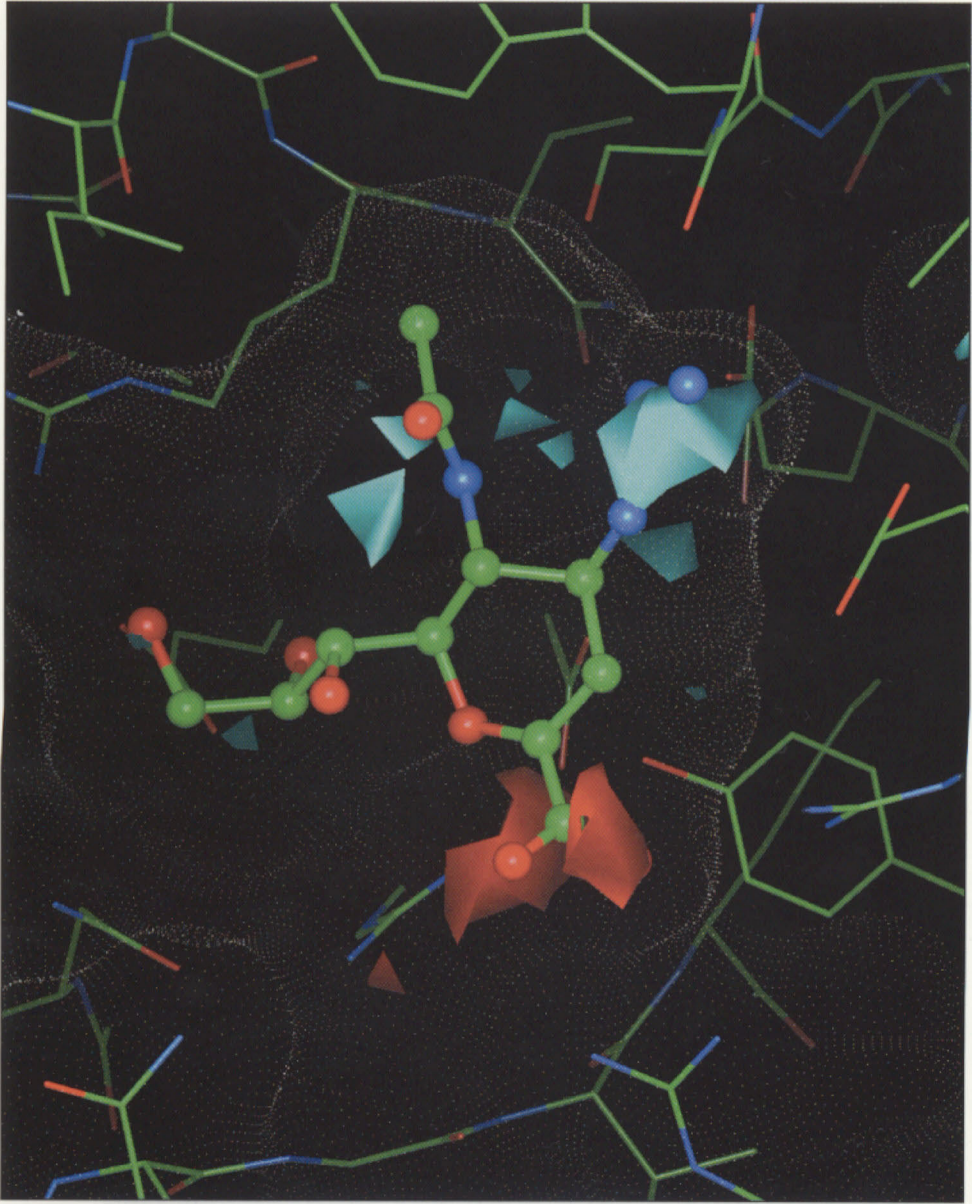


Fig. 12.32: The result of a GRID calculation using carboxylate and amidine probes in the binding site of neuraminidase. The regions of minimum energy are contoured (carboxylate red; amidine blue). Also shown is the inhibitor 4-guanidino-Neu5Ac2en which contains two such functional groups [von Itzstein et al. 1993].

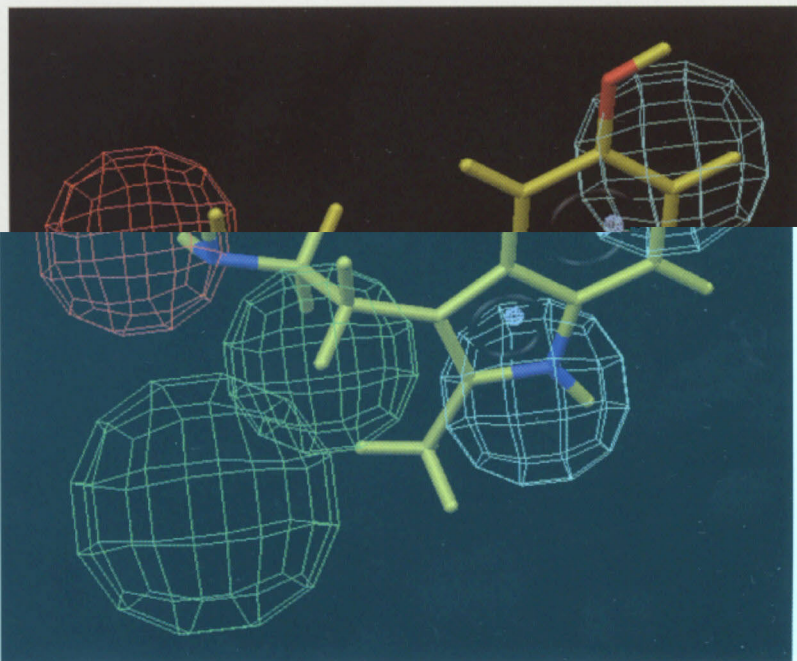
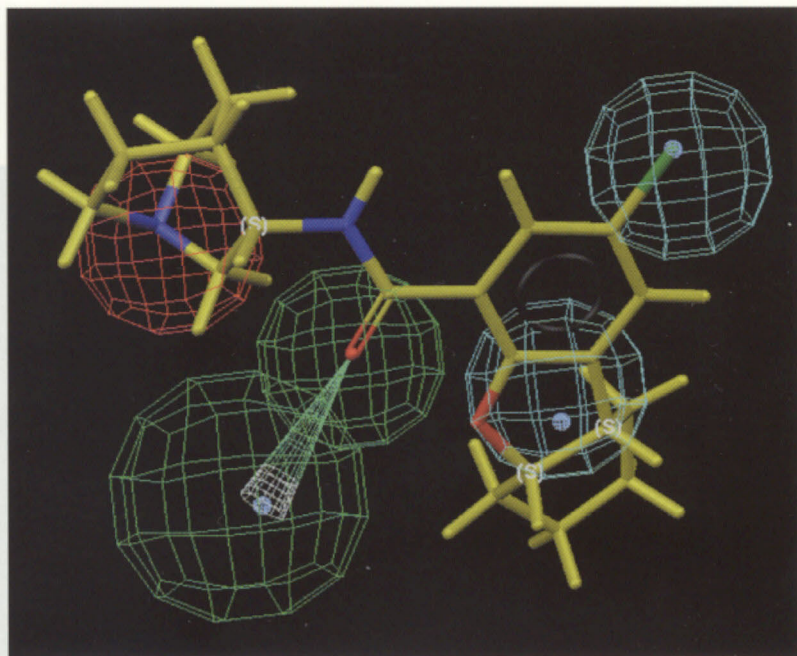


Fig. 12.16: 3D pharmacophore derived for a series of molecules with activity at the 5HT₃ receptor. The spheres indicate location constraints where an appropriate pharmacophore group should be located (red: positively ionisable, green: hydrogen-bond acceptor, blue: hydrophobic region). The figure shows a very active molecule, JMC-35-903-10 superimposed on the pharmacophore (top) and a much less potent molecule, 2-Me-5HT (bottom). The inactive molecule is not able to match all of the points in the pharmacophore in a low-energy conformation.

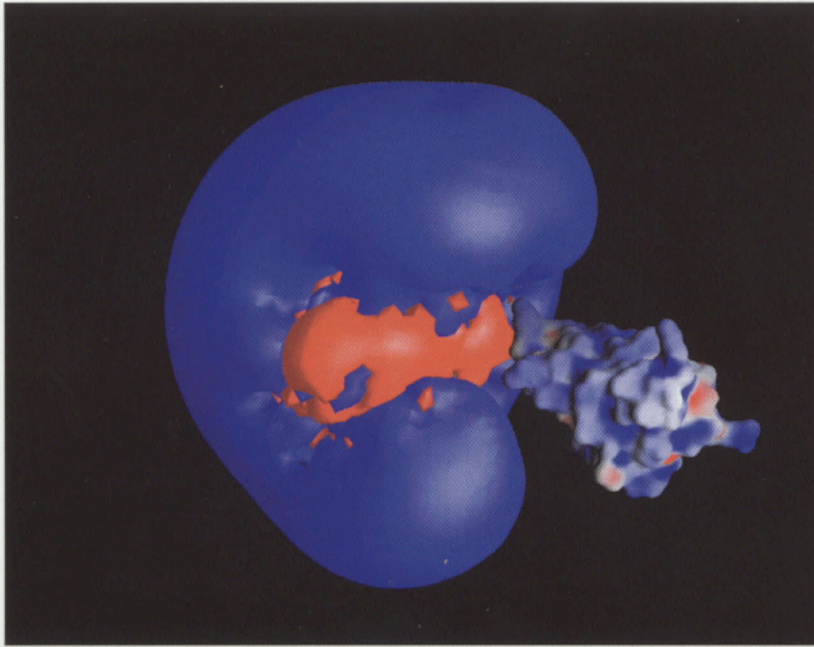


Fig. 11.29: 3D Electrostatic isopotential contours around trypsin [Marquart et al. 1983]. Contours are drawn at $-1k_B T$ (red) and $+1k_B T$ (blue). The trypsin inhibitor is also represented with its electrostatic potential mapped onto its molecular surface.

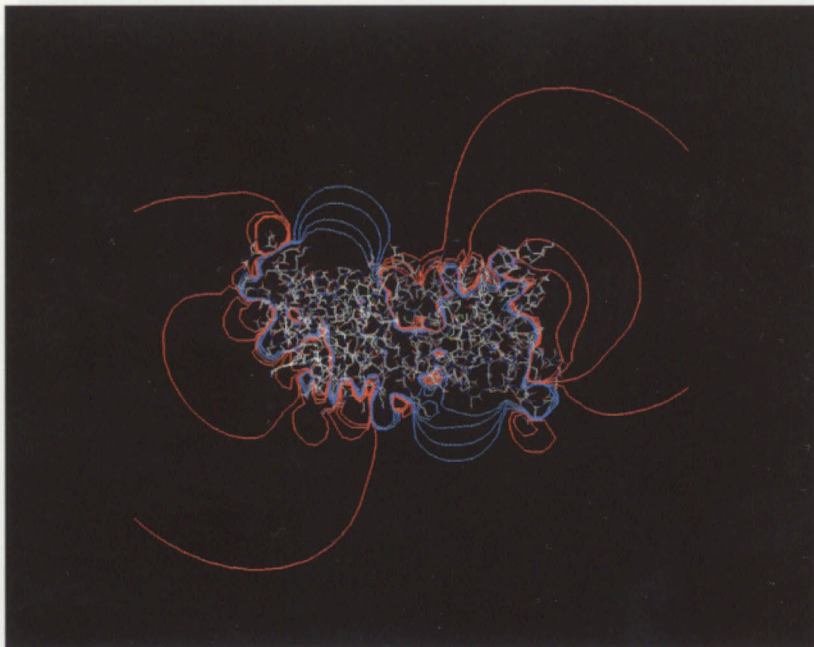


Fig. 11.30: Electrostatic potential around Cu-Zn superoxide dismutase [McRee et al. 1990]. Red contours indicate negative electrostatic potential and blue contours indicate positive electrostatic potential. Two active sites are present in each dimer, at the top left and bottom right of the figure where there is a significant concentration of positive electrostatic potential.

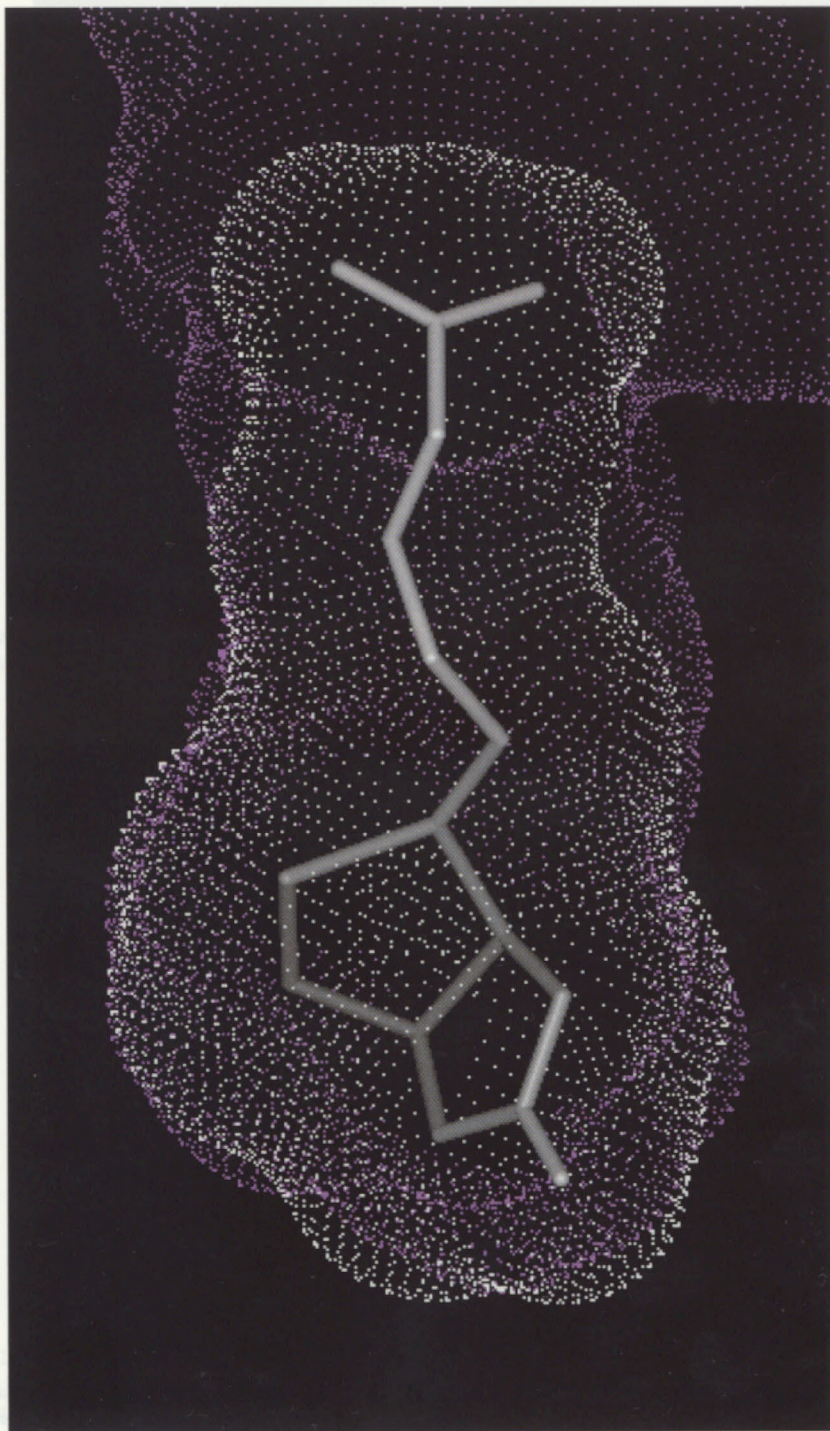


Fig. 11.12: Surface complementarity of the protein streptavidin (purple) and the ligand biotin (white) [Freitag et al. 1997].

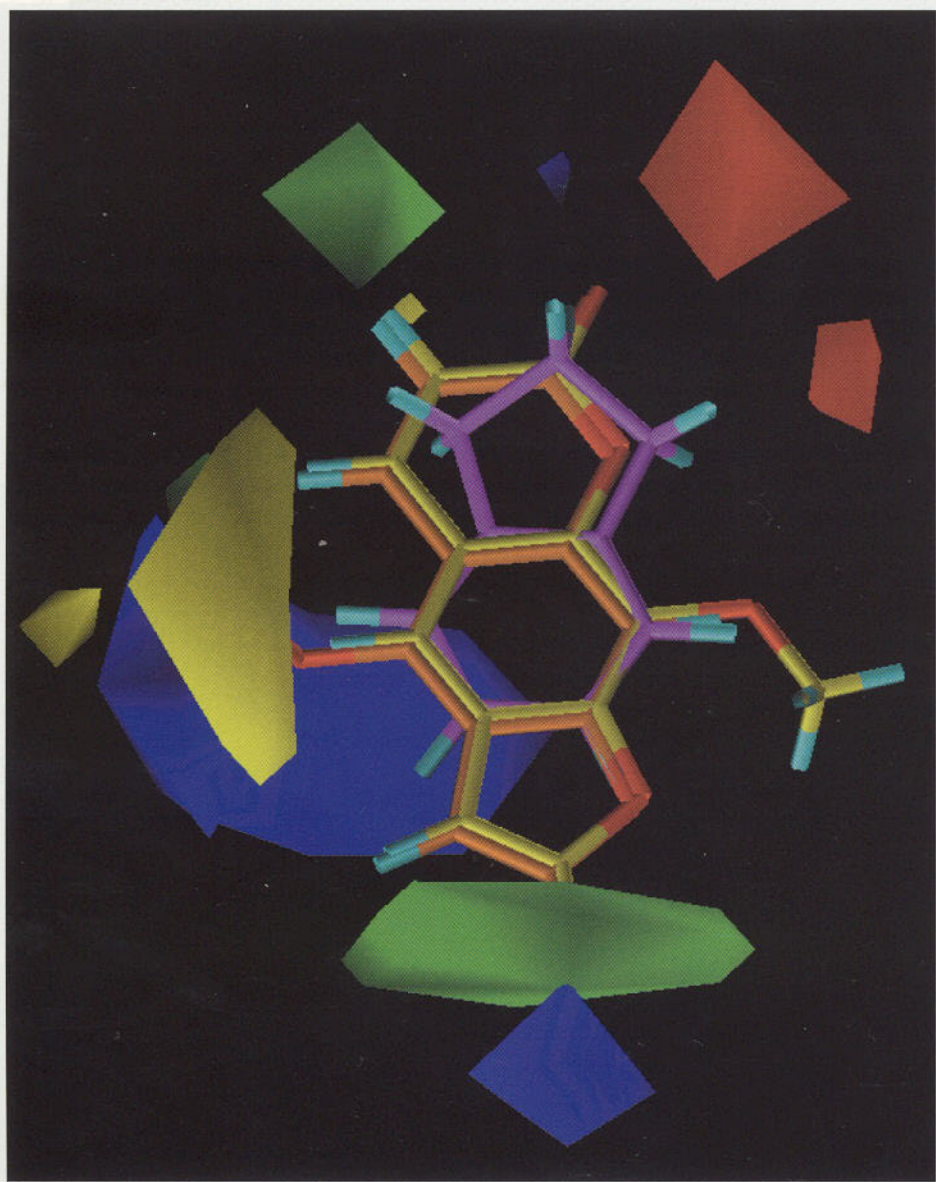


Fig. 12.41: Contour representation of key features from a CoMFA analysis of a series of coumarin substrates and inhibitors of cytochrome P₄₅₀2A5 [Poso et al. 1995]. The red and blue regions indicate positions where it would be favourable and unfavourable respectively to place a negative charge and the green/yellow regions where it would be favourable/unfavourable to locate steric bulk.