

СООБЩЕНИЯ
ОБЪЕДИНЕННОГО
ИНСТИТУТА
ЯДЕРНЫХ
ИССЛЕДОВАНИЙ

ДУБНА



Ц840г
А-84

P10 - 9178

1/11

Д.Д. Арнаудов

4698/2-75

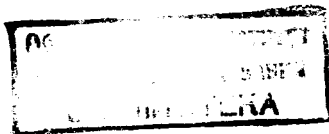
ОБ ОДНОМ СПОСОБЕ ОРГАНИЗАЦИИ
МНОГОУРОВНЕВОЙ АДАПТИРУЮЩЕЙСЯ
ИНФОРМАЦИОННО-ПОИСКОВОЙ СИСТЕМЫ

1975

P10 - 9178

Д.Д. Арнаудов

ОБ ОДНОМ СПОСОБЕ ОРГАНИЗАЦИИ
МНОГОУРОВНЕВОЙ АДАПТИРУЮЩЕЙСЯ
ИНФОРМАЦИОННО-ПОИСКОВОЙ СИСТЕМЫ



Проблема многоуровневой организации информационно-поисковой системы заключается прежде всего в организации основного информационного массива и системы поиска. Система поиска в иерархической ИПС подробно рассмотрена нами в^{1/}. Особенность этой системы заключается в необходимости определения последовательности прав сравнения. Под этим понимается последовательно усложняющийся анализ элемента с постепенным охватом все большего числа его деталей, при этом результаты, полученные на предыдущем этапе анализа, исследуются на последующем. Необходимость в использовании последовательности этапов анализа вызвана тем, что по мере того, как растет основной информационный массив, применять самый сложный анализ к каждому элементу оказывается слишком дорогостоящим делом. В связи с этим необходимо отметить, что, со своей стороны, основной информационный массив должен иметь подходящую структуру, обеспечивающую возможность использования менее сложных механизмов индексации и промежуточных мер для определения степени соответствия, чем окончательная мера соответствия. Именно этот аспект порождает значительные трудности при разработке системы информационного поиска, связанные с согласованием объема памяти поискового массива и времени для получения ответа. Эти трудности могут быть разрешены на основе создания самоорганизующейся многоуровневой информационно-поисковой системы, где число уровней иерархии основного информационного массива и величина сегмента

массива определенного уровня могла бы динамически изменяться в процессе нарастания информационного массива согласно определенному критерию и тем самым регулировать соответствующее отношение между временем, затрачиваемым на поиск элемента массива, и объемом расходуемой памяти.

Анализируя различные организационные структуры основного информационного массива, Лефковитц /2/ отмечает, что можно различить два полюса в организации массивов - инверсную и мультисписковую (ассоциативно-адресную) организацию. Все остальные варианты организации можно рассматривать как частично инвертированные мультисписковые схемы. Проведенный нами анализ инверсной и мультисписковой организации информации /3/ показал, что с нарастанием объема основного информационного массива ухудшаются поисковые характеристики как мультисписковой, так и инверсной организации. В данной работе мы покажем, что улучшение поисковых характеристик системы можно получить на базе комбинирования признаков инверсной и мультисписковой организации, создавая многоуровневую самоорганизующуюся поисковую систему на базе определенной стратегии поиска /1/. В основе методики разработки подобной системы заложен принцип модификации мультисписковой организации в сторону "инверсности".

На рис.1,2,3 показано возможное развитие структуры информационного массива. На рис.1 представлена классическая ассоциативно-адресная структура. Управляющая часть представляет набор дескрипторов словаря системы. Коды дескрипторов совпадают с номерами (индексами) элементов массива. В каждом дескрипторе имеется информация о первом члене узлового списка (АС- адрес связи) и о числе членов этого списка (ЧД - число членов). В массиве ОМПОД расположены соответствующие узловые списки.

На рис.2 показана уже частично-инвертированная структура. Все уровни списка в массиве ОМПОД сегментированы по зонам. Управляющая часть представляет сложную структуру, где хранятся заголовки дескрипторных списков. Сами эти заголовки соединены в цепочки при помощи адресов связи ($NZ1, AC1$). Расположение первого элемента данного списка определяется тоже с помощью адреса связи ($AC2, NZ2, AD2$). Более подробно эта структура рассмотрена в /4/. Ясно, что управляющая часть состоит из двух частей, т.е. имеет двухуровневую структуру.

Управляющая часть с трехуровневой структурой (массив ОМД из двух частей и массив МЗД) показана на рис.3. Здесь в массиве МЗД собраны заголовки дескрипторов узловых списков массива ОМПОД, а в массиве ОМД находятся начала цепных списков заголовков массива МЗД /4/.

Исходя из принятой стратегии поиска /1/, можно предложить следующие формулы для подсчета среднего времени обработки информационного массива с одноуровневой, двухуровневой и трехуровневой структурой.

$$T_{\text{одноур}} = K N_p T_z + L_s t_s = K N_p T_z + L_s (t_n + 1/2 t_c), \quad (1)$$

$$T_{\text{двухур}} = K N_p T_z + a_{\text{ОМД}} C_{\text{КМД}} T_{z_1} + a_{\text{ОМПОД}} C_{\text{МПОД}} T_{z_2}, \quad (2)$$

$$T_{\text{трехур}} = K N_p T_z + a_{\text{ОМД}} C_{\text{КМД}} T_{z_1} + a_{\text{МЗД}} C_{\text{МЗД}} T_{z_2} + a_{\text{ОМПОД}} C_{\text{КМПОД}} T_{z_3}. \quad (3)$$

Рассмотрим отдельные элементы предложенных формул. Первые одночлены всех формул одинаковы. Здесь учитывается время, которое расходуется на декодирование дескрипторов запроса. Необходимо отметить, что любой запрос представляет конъюнкцию дескрипторов /5/. В дальнейшем рассматриваются только положительные дескрипторы - N_p (дескрипторы без знака отрицания). Дескрипторы с отрицанием

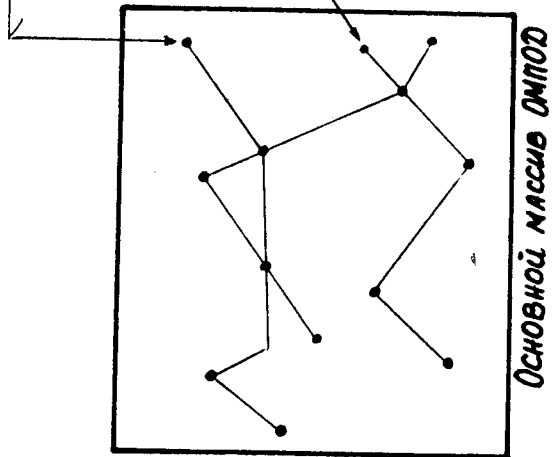
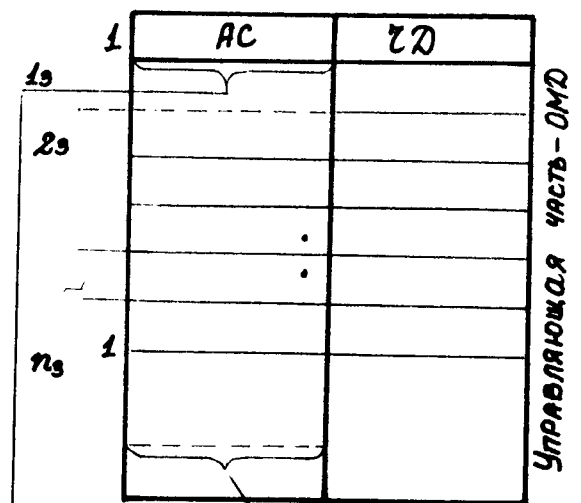


Рис. 1

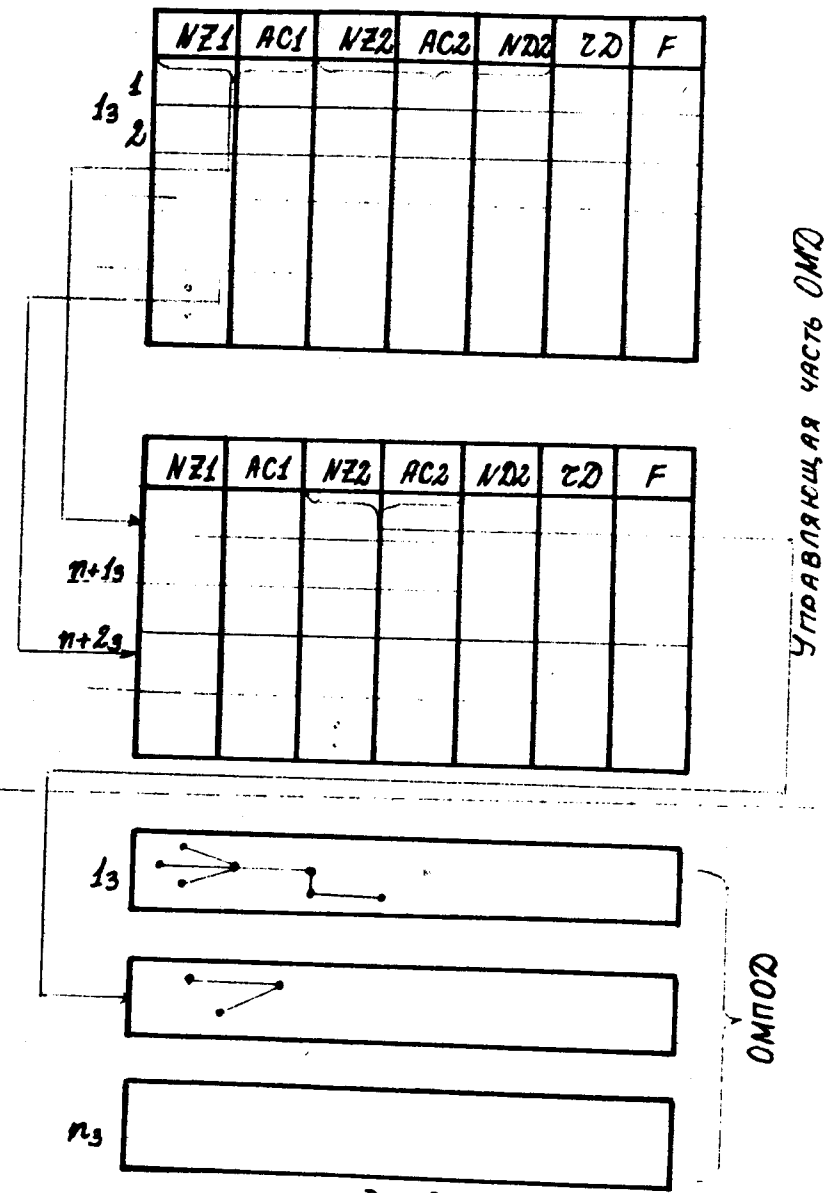


Рис. 2

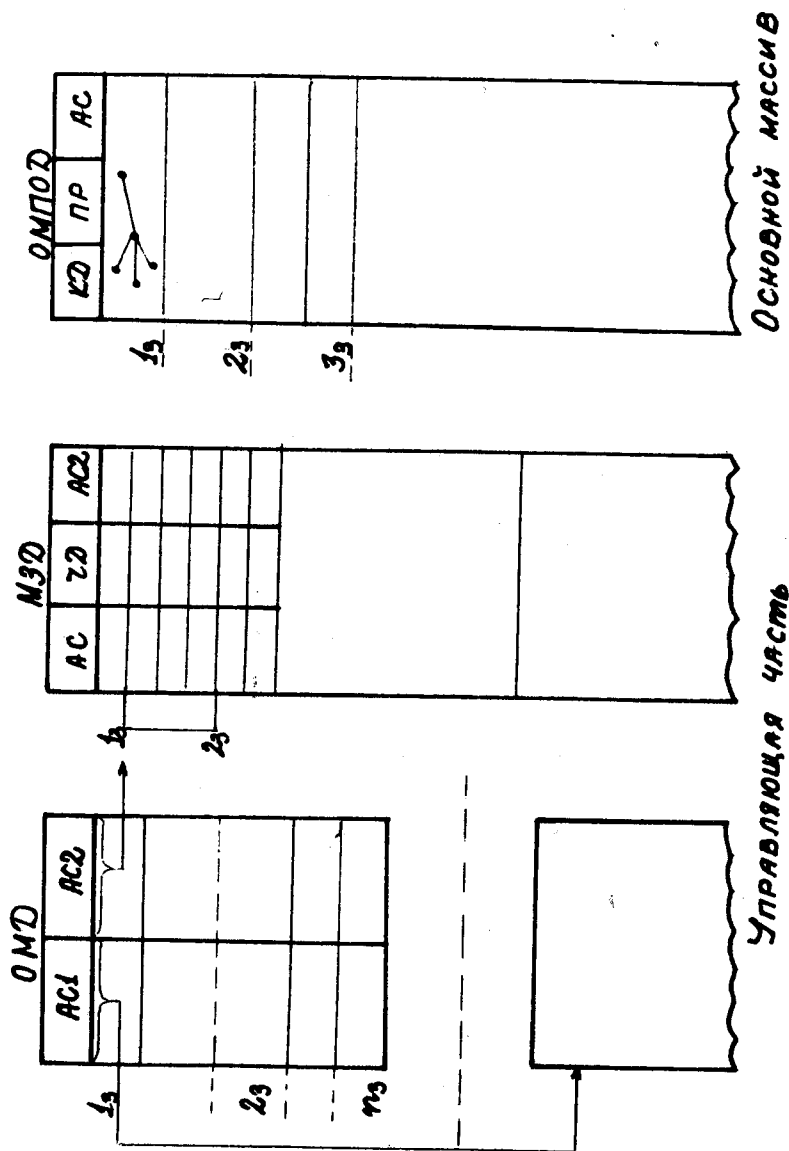


Рис. 3

используются только на последнем этапе поиска, в узлах ОМПОД, и практически не влияют на время ответа, вычисляемого по формулам 1-3. Если запрос состоит из суммы произведений (т.е. представлен в дизъюнктивной форме /5/), время обработки запроса меньше, чем суммы времени обработки отдельных конъюнкций благодаря принятой стратегии параллельной обработки запросов /1/.

Коэффициент K учитывает одновременную выборку K/P дескрипторов запроса. Необходимо заметить, что первая часть массива ОМД (во всех рассматриваемых структурах) логически сегментирована по зонам и возможно, что для некоторых дескрипторов запроса необходимо читать только одну зону. Это намного сокращает время выборки и учитывается коэффициентом K . Параметр T_Z обозначает время чтения зоны (аналогично T_{Z1}, T_{Z2}, T_{Z3}). Параметр C_K обозначает среднее количество зон на дескриптор и имеет место в формулах 2 и 3. C_{KMND} показывает среднее количество зон на дескриптор в Π -й части массива ОМД. C_{KM3D} и C_{KMPPD} показывает среднее количество зон на дескриптор соответственно в массивах МЗД и ОМПОД. Параметр $\alpha_{MND}, \alpha_{M3D}, \alpha_{MPPD}$ показывает отношение среднего количества общих зон запроса к C_K . Следовательно:

$$\alpha_{MND} = \frac{Z_{3MND}}{C_{KMND}} ; \alpha_{M3D} = \frac{Z_{3M3D}}{C_{KM3D}} ;$$

$$\alpha_{MPPD} = \frac{Z_{3MPPD}}{C_{KMPPD}} .$$

Здесь $Z_{3MND}, Z_{3M3D}, Z_{3MPPD}$ обозначают среднее количество зон запроса соответственно в массивах ОМД, МЗД, ОМПОД.

Введем еще несколько параметров, которые встречаются в нашем изложении.

Параметры t_n , t_c , t_{mp} , t_s являются характеристиками операционной системы при работе с дисковым устройством, где

t_n - время выборки цилиндра;

t_c - время оборотов дискового пакета;

t_{mp} - время трансмиссии n символов;

t_s - время выборки информации одной физической записи с диска.

Они подробно рассмотрены нами в /1,4/.

Особенно важный параметр - это среднее количество документов на дескриптор, L , которым характеризуется средняя длина списка. Эта характеристика является среднестатистической величиной и связана с частотой встречаемости дескрипторов $-f(j)/|J|$. В данном случае ее можно вычислить по формуле

$$L = \frac{N_z N_k}{V}, \quad (4)$$

где N_z - количество документов основного информационного массива;

N_k - средняя глубина индексирования документа;

V - количество различных дескрипторов.

В формуле (1) через L_s обозначена длина самого короткого дескрипторного списка. Анализируя эту формулу, необходимо отметить, что в схеме на рис.1 поиск происходит при движении по самой короткой цепи L_s . Так как списки не сегментированы, то каждый раз происходит выбор отдельного узла с диска. Принимается, что максимальный размер узла не превышает одной физической единицы записи и поэтому

$$t_s = t_n + 1/2 t_c \quad (\text{более подробно см. /1,3,4/}).$$

Исходя из принятой стратегии поиска /1/ в двухуровневой и трехуровневой структуре необходимо просмотреть все элементы цепи дескрипторов в ОМД (т.е. в высшем уровне иерархии), поэтому в сред-

нем в этом массиве всегда просматривается среднее количество зон на дескриптор. Тогда $\alpha_{смп} = 1$ и формулы (2) и (3) принимают следующий вид:

$$T_{\text{обычур}} = K N_p T_z + C_{\text{кмпд}} T_{z1} + \alpha_{\text{смпд}} C_{\text{кмпд}} T_{z2}, \quad (2a)$$

$$T_{\text{трехур}} = K N_p T_z + C_{\text{кмпд}} T_{z1} + \alpha_{\text{мпд}} C_{\text{мпд}} T_{z2} + \alpha_{\text{смпд}} C_{\text{смпд}} T_{z3}. \quad (3a)$$

Основной вопрос, который необходимо решить на основе предложенных формул (1, 2a, 3a), - это когда и в каких случаях необходимо переходить от одной структуры к другой (см.рис.1,2,3) и какие величины зон сегментированных списков являются самыми подходящими при проведении поиска в соответствующих структурах.

Так как первое слагаемое одинаково во всех формулах (1+3), дальнейший сравнительный анализ проведем, используя следующие выражения:

$$(1') T_{\text{смпур}} = L (t_n + 1/2 t_c),$$

$$(2') T_{\text{обычур}} = C_{\text{кмпд}} T_{z1} + \alpha_{\text{смпд}} C_{\text{кмпд}} T_{z2},$$

$$(3') T_{\text{трехур}} = C_{\text{кмпд}} T_{z1} + \alpha_{\text{мпд}} C_{\text{мпд}} T_{z2} + \alpha_{\text{смпд}} C_{\text{смпд}} T_{z3}.$$

В формуле (1') вместо L_s берем L - среднюю длину списка в ОМПОД, т.к.

$$\frac{\sum_{s=1}^N L_s}{N} \leq L,$$

где N - число всех дескрипторов. Средняя длина списка L вычисляется по формуле (4) и, используя формулу $L_{\text{смпд}} |J|$, можно вычислить, при каком количестве документов ожидается, что все N дескрипторов словаря системы будут использованы.

$$V = \frac{N_z N_k}{L n N + \gamma}. \quad (5)$$

Для информационного массива принимаем $N_k = 9$.

Тогда

$$10000 = \frac{N_z \cdot 9}{\ln 10000 + 0,57} ; N_z = 10720.$$

Это означает, что для информационного массива из документов ИНИСа^{16/} при объеме 10720 документов ожидается, что все N дескрипторов словаря системы будут встречаться в документах информационного массива.

В дальнейшем в качестве характеристик массива будут использоваться характеристики информационного массива, получаемого на базе магнитной ленты ИНИСа^{16/}.

На рис.4 приведена зависимость $L = f(N_z)$.

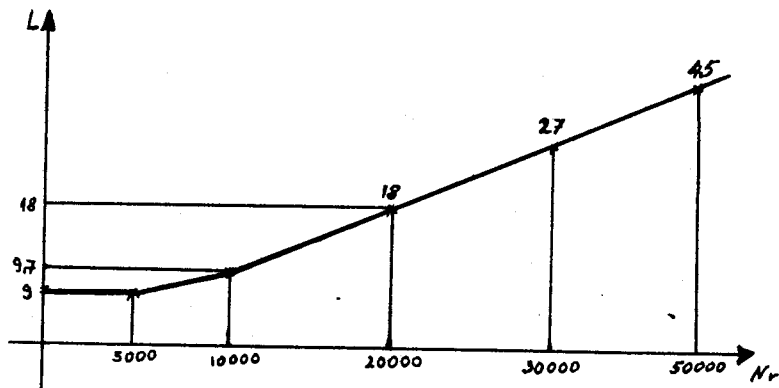


Рис.4

Величина L является основной при решении вопроса о переходе от одной структуры к другой.

Первый вопрос, который необходимо решить, - это при каком L_{min}

$$T_{авчур} < T_{одноур} , \text{ т.е. } T_{авчур} < L(t_n + 1/2 t_0).$$

С целью определения L_{min} проведен машинный эксперимент.

В качестве экспериментального массива моделировался массив по закону Цифа с поправкой L_{OWE} ^{13/}. Получено, что $T_{авчур} < T_{одноур}$ для массива с параметрами $L = 9, N_z = 3220$. Оптимальная вели-

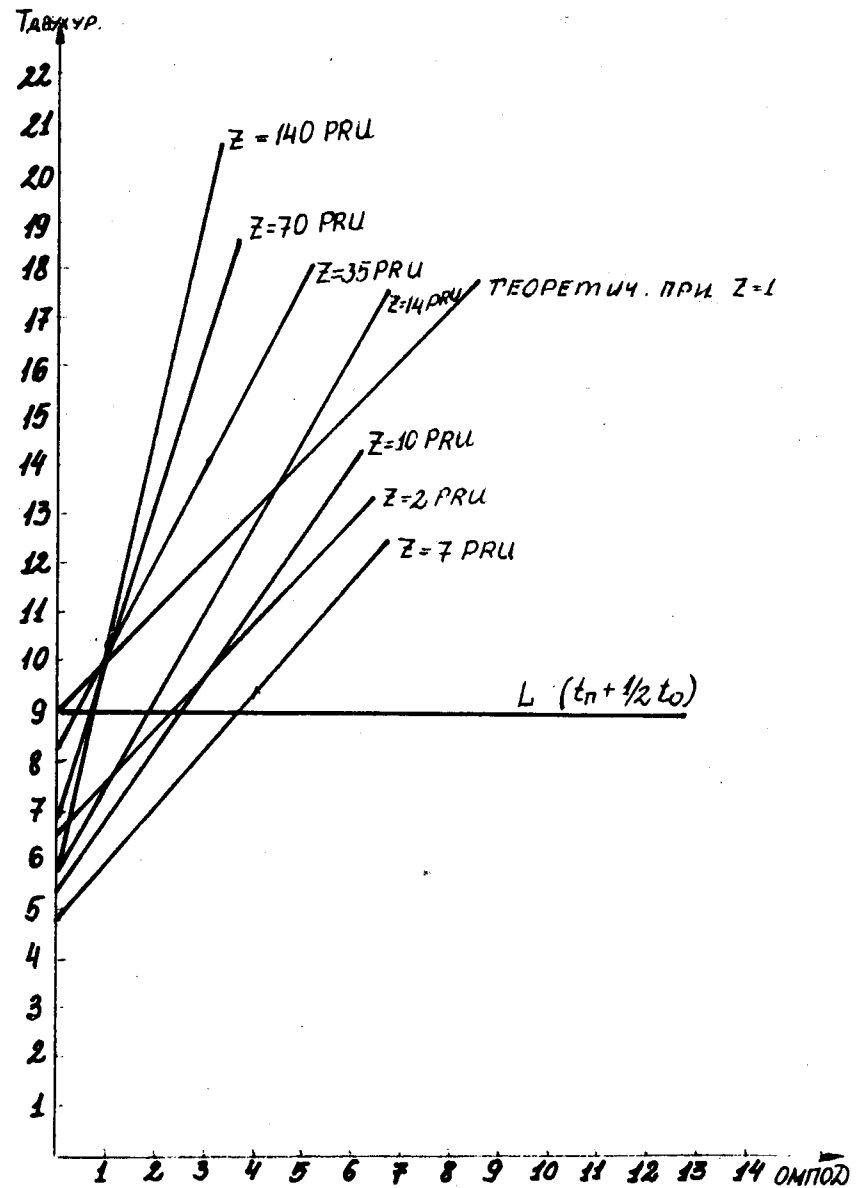


Рис. 5

чина зоны (Z), в этом случае $Z_1 = Z_2 = 7PRU$ (о вычислении величины зоны в физических единицах см.¹⁷¹). На рис.5 показан график $T_{\text{выкуп}} = f(Z_{\text{ОМПОД}})$ при $Z = \text{CONST}$. В этом эксперименте Z меняется от $2PRU$ до $140PRU$. Самая хорошая прямая времени - при $Z = 7PRU$. Как видно из рисунка, при обработке $Z_{\text{ОМПОД}} \leq 3$ T двухуровн. $< T$ одноуровн. Обработка трех зон основного массива является вполне достаточной, и в этом случае $\alpha_{\text{ОМПОД}} = -0,43$. Исследования Лефковитца¹²¹ показали, что при запросах, состоящих не меньше чем из трех дескрипторов, $\alpha_{\text{ОМПОД}}$ варьируется в пределах $0,1 \pm 0,2$. В наших экспериментах (среднее число дескрипторов запроса равно четырем) $\alpha_{\text{ОМПОД}} \leq 0,1$. Значения параметров двухуровневой структуры при объеме массива 3200, 10000, 20000 и 50000 документов экспериментально получены и сведены в таблицу 1А. На базе данного эксперимента исследованы и получены коэффициенты K_1 и K_2 , которые являются основными в расчете основных параметров массива. Об этих коэффициентах речь пойдет дальше, а сейчас приведем формулы для расчета основных параметров массива. Эти формулы будут использованы дальше для эмпирического расчета параметров в выражении (1, 2, 3) при различных объемах информационного массива.

Сначала приведем формулы расчета параметров двухуровневой структуры.

$$C_{\text{КОМПОД}} = \frac{K_1 Z_{\text{ОМПОД}} \cdot 64 Z_{\text{ОМПОД}}}{2V} \quad (6)$$

Вывод формулы покажем на базе анализа физической сущности $C_{\text{КОМПОД}}$. Так как $C_{\text{КОМПОД}}$ является средней характеристикой числа зон (основного массива ОМПОД) на дескриптор, то

$$C_{\text{КОМПОД}} = \frac{\text{число всех заголовков}}{V} \quad (7)$$

b_1	1,42	1,63	1,71	1,78	2	2,32
b_2	1,1	1,1	1,12	1,19	1,24	1,32
K_1	0,70	0,67	0,64	0,56	0,50	0,43
K_2	0,90	0,90	0,89	0,84	0,80	0,75
$Z (PRU)$	7PRU	10PRU	14PRU	35PRU	70PRU	140PRU
$C_{\text{КОМПОД}}$	6,33	5,8	5,3	4,0	3,2	2,6
$C_{\text{К ОМД}}$	4,46	3,85	3,3	2,1	1,4	0,6
$Z_{\text{ОМПОД}}$	132	93	66	27	13	6,5
$Z_{\text{ОМД}}$	76	48	32	9	3,2	1,2
$C_{\text{КОМПОД}}$	6,72	6,43	6,13	5,34	4,73	4,09
$C_{\text{К ОМД}}$	5,2	4,9	4,57	3,66	2,97	2,27
$Y_{\text{ОМПОД}}$	401	283	201	81	41	21
$Z_{\text{ОМД}}$	242	162	109	37	16	1
$C_{\text{К ОМПОД}}$	12,3	12,24	11,68	10,18	9,01	7,79
$C_{\text{К ОМД}}$	10,98	10,36	9,72	8,00	6,62	5,16
$Z_{\text{ОМПОД}}$	801	565	401	162	82	41
$Z_{\text{ОМД}}$	523	351	237	82	36	16
$C_{\text{КОМПОД}}$	18,59	18,01	17,47	15,23	13,5	11,0
$C_{\text{К ОМД}}$	16,08	16,05	15,09	12,53	10,0	8,1
$Z_{\text{ОМПОД}}$	1201	847	601	242	120	60
$Y_{\text{ОМД}}$	782	543	367	128	55	23
L						
	N_c					
	V					
9	3200					
9,49	10000					
18,06	20000					
27,01	30000					

Таблица 1а

Ясно, что число всех зон массива ОМПОД равняется числу заголовков узловых списков, сегментированных в конкретных зонах. Априори ясно, что число всех заголовков меньше числа всех элементов массива ОМПОД, т.к. предполагается, что в узловых списках имеется более одного члена (в худшем случае, если есть в списках по одному члену, то число всех заголовков будет равно числу элементов массива ОМПОД). Под "элементом" массива ОМПОД (а также ОМД и МЗД) понимаем физическую величину, состоящую из двух машинных слов (более подробно см. /4/). Тогда

$$\text{Число всех заголовков} = \frac{K_1 Z_{\text{ОМПОД}} 64 Z_{\text{ОМПОД}}}{2} \quad (8)$$

В формуле (8) K_1 является коэффициентом встречаемости и связан с числом различных дескрипторов в данной зоне (каждый дескриптор является основной частью одного элемента). Ясно, что $K_1 \leq I$. Величины K_1 вычислены экспериментально и сведены в таблицу 1А. Видно, что величина этого коэффициента не зависит от длины массива, а зависит только от величины зоны. Следовательно, физическая природа этого коэффициента показывает, что он узко связан с характеристиками входного потока и частотой встречаемости дескрипторов - $f(i)$

Заметим, что в дальнейшем используется и обратный коэффициент $G_1 = 1/K_1$, который называется коэффициентом связанности и обозначает число членов в списке данного дескриптора в данной зоне ОМПОД. Этот коэффициент тоже вычислен экспериментально и показан в таблице 1А. Видно, что для распределения, генерируемого по закону Ципфа (с поправкой L_{OWF}), видно, формируемые списки являются недостаточно плотными (в среднем по 2,32 члена в списке при $Z_{\text{ОМПОД}} = 140 \text{ PRU}$).

В используемых формулах $Z_{\text{ОМПОД}}$ обозначает величину зоны в физических единицах (PRU) ($Z_{\text{ОМПОД}}$ меняется от 7 до 140).

Коэффициент 64 обозначает, что в одной физической единице (PRU) находятся 64 машинных слова.

Параметры $Z_{\text{ОМПОД}}$ и V были рассмотрены нами раньше. Имея в виду все вышеописанное, подставляя (8) в (7), получим формулу для вычисления $C_{\text{КОМПОД}}$ (6).

С другой стороны, можно показать взаимосвязь этого параметра с другими параметрами информационного массива. Для этого будем использовать формулу для подсчета $Z_{\text{ОМПОД}}$.

$$Z_{\text{ОМПОД}} = \frac{2 N_z N_k}{64 Z_{\text{ОМПОД}}} \quad (9)$$

Подставляя (9) в (6), получим:

$$C_{\text{КОМПОД}} = \frac{K_1 N_z N_k}{V} \quad (10)$$

Подставляя (4) в (10), получим:

$$C_{\text{КОМПОД}} = K_1 L \quad (11)$$

Формулы (10) и (11) выражают $C_{\text{КОМПОД}}$ при помощи числа документов массива - N_z , среднего числа дескрипторов на документ - N_k и средней длины узлового списка - L .

На рис.6 показан график изменения $C_{\text{КОМПОД}}$ при изменении величины зоны, т.е. $C_{\text{КОМПОД}} = f(Z_{\text{ОМПОД}})$. Здесь использованы экспериментальные данные массива из 3200 документов.

Перейдем к анализу формулы для вычисления $C_{\text{КОМД}}$.

$$C_{\text{КОМД}} = \frac{K_2 Z_{\text{ОМД}} Z_{\text{ОМД}} 64}{2V} \quad (12)$$

Физическая сущность этого параметра в том, что он обозначает среднее число зон на дескриптор в массиве ОМД. Необходимо, однако, отметить, что речь идет о числе зон во второй части массива (см. рис.2), поскольку в первой части совершается автоматическая выборка по кодам дескрипторов и она влияет только на первое сла-

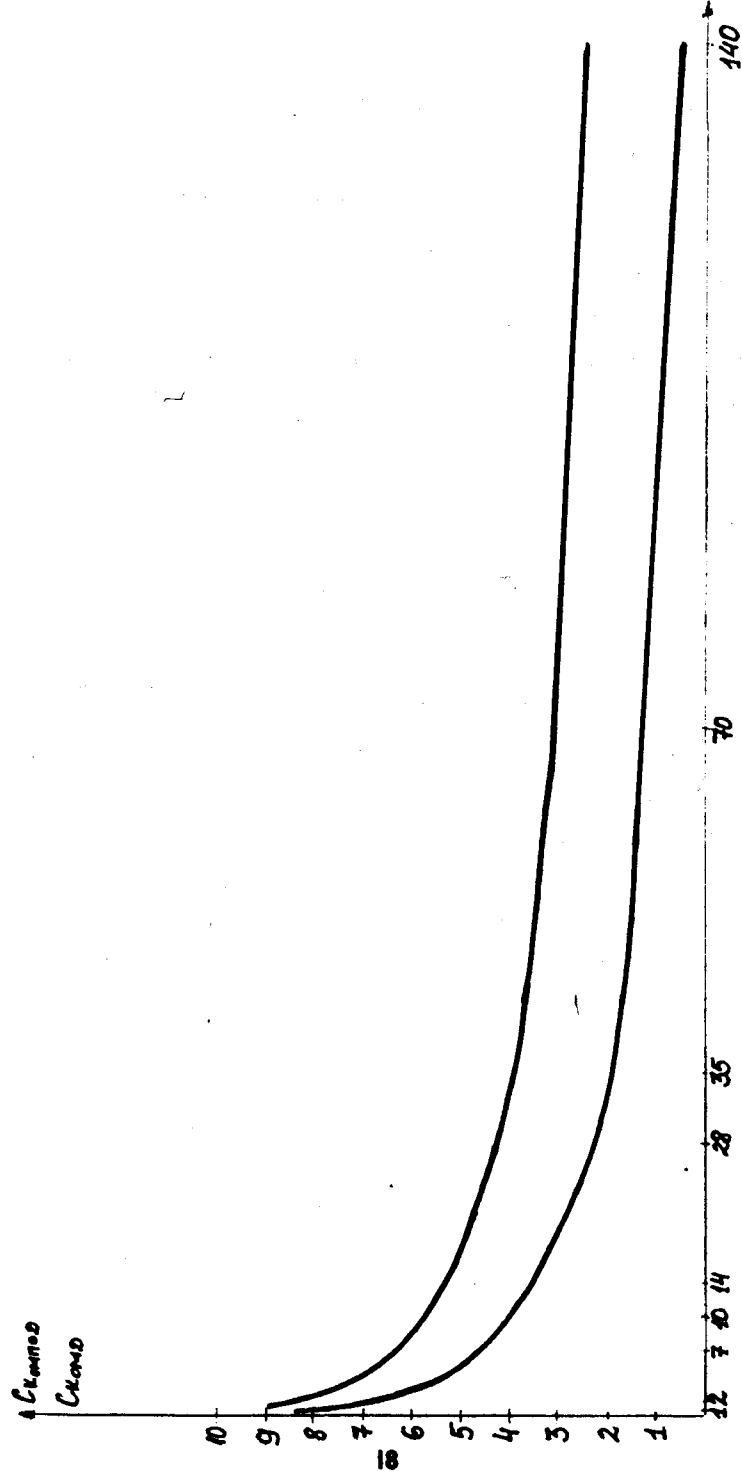


Рис. 6

гаемое формулы (2), которое является общим для всех трех формул и в сравнительном анализе не рассматривается. Так как к первой части массива ОМД относится столько элементов массива (в данной структуре элементами являются заголовки цепных списков массива ОМПОД), сколько имеется различных дескрипторов (V), то

$$C_{КОМПОД} = \frac{K_2 (Z_{ОМД} - V)}{V}. \quad (13)$$

Как уже заметили,

$$Z_{ОМД} = \text{числу всех заголовков}. \quad (14)$$

Здесь вводится впервые коэффициент на K_2 , выполняющий такую же функцию, как и коэффициент K_1 по отношению к параметру $C_{КОМПОД}$. В данном случае K_2 связан с той частью элементов массива ОМД, принадлежащих данной зоне, которые являются представителями различных дескрипторных заголовков. Ясно, что $K_2 < 1$. Обратный коэффициент $G_2 = 1/K_2$ является коэффициентом связанности и показывает наличие больше чем одного заголовка одного и того же дескриптора в одной зоне ОМД. Коэффициенты K_2 и G_2 получены экспериментальным путем и сведены в таблицу IА. Из таблицы видно, что даже если зона равняется 140 PRU, то там встречается не больше чем 1,32 членов в цепном списке дескриптора. А для зоны $Z = 7PRU$ - только по 1,1 члену.

Исходя из приведенных рассуждений, нужно отметить, что

$$Z_{ОМД II} = Z_{ОМД} - V, \quad (15)$$

где $Z_{ОМД II}$ - это число элементов второй части массива ОМД.

Для вычисления по формуле (12) необходимо выразить $Z_{ОМД II}$ через $Z_{ОМД II}$ (число зон второй части массива ОМД).

$$Z_{ОМД II} = \frac{Z_{ОМД} \cdot Z_{ОМД} \cdot 64}{2}. \quad (16)$$

Подставляя (16) в (15), а (15) в (13), получим формулу (12).

Можно установить некоторую взаимосвязь между $C_{КОМД}$ и $C_{КОМПОД}$. Если выразить $Z_{КОМД}$ через $Z_{КОМПОД}$, то получим:

$$Z_{КОМД} = \frac{2 Z_{КОМПОД}}{Z_{КОМД} \cdot 64} = \frac{2 (\text{число всех заголовков} - V)}{Z_{КОМД} \cdot 64} \quad (17)$$

число всех заголовков = $\frac{K_1 Z_{КОМД} \cdot 64 \cdot Z_{КОМПОД}}{2}$

$$\begin{aligned} \text{тогда } Z_{КОМПОД} &= \frac{2 \left(\frac{K_1 Z_{КОМПОД} \cdot 64 \cdot Z_{КОМПОД}}{2} - V \right)}{Z_{КОМД} \cdot 64} = \\ &= \frac{K_1 Z_{КОМПОД} \cdot 32 \cdot Z_{КОМПОД} - V}{Z_{КОМД} \cdot 32} \end{aligned}$$

Подставляя результат в (12), получим:

$$C_{КОМД} = K_2 (C_{КОМПОД} - 1) \quad (18)$$

График зависимости $K_{КОМД}$ от величины зоны при $L = const$ показан на рис. 6. Необходимо отметить, что с изменением величины объема массива $C_{КОМД}$ изменяется по линейному закону.

Используя приведенные уже формулы (1,2), а также формулы (1+18), покажем изменение оптимальных зон ($Z_{ОПТИМ.ОМД}$, $Z_{ОПТИМ.ОМПОД}$) при изменении объема массива (соответственно изменяется и L).

В формулах (1+18) в качестве коэффициентов K_1 и K_2 использованы полученные экспериментальным путем данные для этих коэффициентов. Они сведены в таблицу 1А. Как уже заметили, эти коэффициенты являются постоянными для определенной величины зоны и не изменяются (очень слабо изменяются) с увеличением объема массива. Вычислены $C_{КОМД}$ и $C_{КОМПОД}$ для массива из 100000 и 500000 документов. Данные приведены в таблице 2а.

Используя приведенные данные и формулы (1*) и (2*), получаем следующие результаты.

Таблица 2А

Титуляр. (одна зона)	L	V	Mz	Z (PRU)		7	10	14	35	70	140
				$C_{КОМД}$	$C_{КОМПОД}$						
71 (7n + 1/2t)	90	10000	100000	$C_{КОМД}$	43	63	53	50	41	34,5	27
				$C_{КОМПОД}$	2900	4000	252	289	45	38,5	
356 (7n + 1/2t)	450	10000	500000	$C_{КОМД}$	220	315	-	-	-	-	-
				$C_{КОМПОД}$	13000	20000	-	-	-	-	-

Если $Z_{\text{ОМД}} = Z_{\text{ОМПОД}} = 7PRU$, то как при $N_z = 100000$, так и при $N_z = 500000$ Т двухур. < Т одноур. при $\alpha_{\text{ОМПОД}} \leq 0,43$. Это означает, что $Z = 7PRU$ является оптимальной при исследовании не больше чем 43% массива ОМПОД. Однако при этом получается довольно большая затрата памяти на массив ОМД - около 62% от массива ОМПОД. Начиная примерно с 30000 документов массив ОМД занимает большое место на диске.

Необходимо отметить, что в реальных случаях исследуется не больше, чем 20% ОМПОД (по данным Лефковитца^[2]), а в наших исследованиях (для данного экспериментального массива) - не больше 10%, т.е. $\alpha_{\text{ОМПОД}} \leq 0,1$. В таком случае при данной оптимальной зоне $Z_{\text{ОМД}} = 7PRU$ имеется лишний запас, т.к. $\alpha_{\text{ОМПОД}} \leq 0,43$. Память, расходуемая на ОМД, можно уменьшить, не увеличивая Т двухур., но уменьшая границы изменения коэффициента $\alpha_{\text{ОМПОД}}$. Этого можно достичь путем увеличения размера зоны массива ОМПОД. Хорошие результаты получаются при увеличении зоны, начиная с массива величиной в 30000 документов (если массив меньше этой длины, то при увеличении зоны увеличивается Т двухур). Значение основных характеристик для массивов из 30000, 100000, 150000, 500000 документов сведены в таблицу 3А. Имеется в виду, что $Z_{\text{ОМПОД}} = 140PRU$, а $Z_{\text{ОМД}} = 7PRU$. В данном случае характеристики лучше, чем в случае $Z_{\text{ОМД}} = Z_{\text{ОМПОД}} = 7PRU$ (см. табл. 2А и 3А). Кроме того, в данном случае за счет изменения границы $\alpha_{\text{ОМПОД}}$ достигнуто соотношение между массивами $\frac{\text{ОМД}}{\text{ОМПОД}} = 0,37$, при котором ОМД уже составляет 37% от основного массива. Это лучшие возможные результаты, т.к. $Z_{\text{ОМПОД}}$ принимает максимально допустимое значение - $140PRU$. Конечно, в зависимости от конкретных требований к поисковой системе в других случаях можно увеличить границы $\alpha_{\text{ОМПОД}}$ (например, до 0,25), но тогда необходимо, чтобы $Z_{\text{ОМПОД}} = 70PRU$.

Таблица 3А

	30000		100000		150000		500000		1000000	
	L	V	L	V	L	V	L	V	L	V
$\alpha_{\text{ОМД}}$	27	10000	8,9	2,6	42,5	170	347			
$\alpha_{\text{ОМПОД}}$	II	39	58	194						
$Z_{\text{ОМД}}$	450	1290	200	300	2200	7500	7300			
$Z_{\text{ОМПОД}}$	60	200	300	300	2200	7500	2000			
$\alpha_{\text{ОМПОД}}$	$\leq 0,1$	$\leq 0,1$	$\leq 0,1$	$\leq 0,1$	$\leq 0,1$	$\leq 0,1$	$\leq 0,1$	$\leq 0,1$	$\leq 0,1$	$\leq 0,1$
Т двухур.	104 ($t_n \cdot 1/2 t_o$)	62,5 ($t_n \cdot 1/2 t_o$)	9 ($t_n \cdot 1/2 t_o$)	351 ($t_n \cdot 1/2 t_o$)						

Увеличение зоны ОМД больше чем на $7PRU$ ведет к ухудшению временных характеристик. Понятно, что это происходит из-за медленного уменьшения $C_{КОМД}$ (см.рис.6) с увеличением величины зоны, а в этом случае T двухур. резко повышается.

Исходя из проведенных экспериментов и вышеописанного анализа, можно сделать следующие выводы.

Для основного информационного массива величиной до 3000 документов необходимо, чтобы структура системы была одноуровневой (см.рис.1).

Переход к двухуровневой структуре необходимо осуществить при массиве выше 3000 документов. Тогда при увеличении объема массива до 30000 документов оптимальной является зона величиной $7PRU$, т.е. $Z_{ОМД} = Z_{ОМПОД} = 7PRU$.

При массиве выше 30000 документов необходимо переключить систему на следующие оптимальные зоны: $Z_{ОМД} = 7PRU$, $Z_{ОМПОД} = 140PRU$.

График изменения $Z_{ОПТ.ОМД}$ и $Z_{ОПТ.ОМПОД}$ с изменением величины массива показан на рис.7.

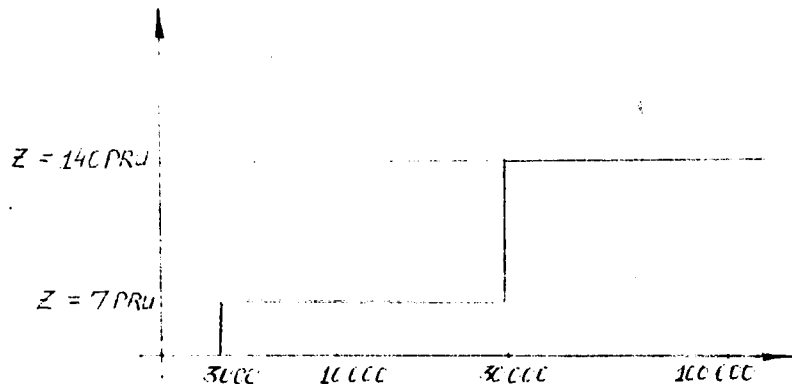


Рис.7

Здесь сплошной линией показано изменение $Z_{ОПТ.ОМПОД}$, штриховой линией — $Z_{ОПТ.ОМД}$.

Все до сих пор сказанное подчеркивает тот факт, что величина оптимальной зоны после достижения определенного объема фактически не зависит от величины информационного массива. В качестве оптимальных зон $Z_{ОМПОД}$ и $Z_{ОМД}$ рассматриваются такие зоны, при которых T двухур. $< T$ одноур. и при этом имеется самый маленький возможный расход памяти на управляющую часть массива (т.е. в данном случае на массив ОМД). Время обработки для одноуровневой структуры зависит прежде всего от средней длины дескрипторного списка $-L$. Чем длиннее L , тем больше времени уходит на поиск (см.формулу I). Однако с увеличением L увеличивается (после $N_z = 10000$) прямо пропорционально коэффициенты $C_{КОМД}$ и $C_{КОМПОД}$ при уже выбранной оптимальной зоне. Напомним, что $C_{КОМПОД} = C_{КОМД} = K_2 (C_{КОМПОД} - 1)$.

При достижении некоторого L^* , в данном случае это 27 (при 30000 документах), как уже говорилось, за счет изменения границы $\alpha_{ОМПОД}$ имеется предельная возможность (при постоянных в данном случае K_1 и K_2) изменения зоны ОМПОД, т.е. достигается ее максимальное значение. И в дальнейшем, конечно, двухуровневая структура всегда будет работать лучше одноуровневой, но она не будет являться оптимальной структурой по отношению ко времени ответа для некоторого достаточно большого информационного массива. В связи с этим отметим, что оптимизация двухуровневой структуры связана с уменьшением коэффициентов K_1 и K_2 , что реализуется путем "уплотнения" дескрипторных списков. Но даже и после применения соответствующих оптимизационных методов, после того как достигнута максимальная величина оптимальной зоны, если информационный массив продолжает расти, то приходит момент, когда возможно большое уменьшение коэффициента P ($P = \frac{Z_{ОМПОД}}{C_{КОМД}}$).

Это означает ($C_{КОМД}$ достаточно большое), что "впустую" просматривается большое число зон в массиве ОМД, а в результате только маленькое количество из них является релевантным запросу. Это те зоны, которые в дальнейшем должны быть просмотрены в массиве ОМПОД - $Z_{ОМПОД}$. В этом случае улучшение времени поиска может быть достигнуто за счет некоторого увеличения объема управляющей части вследствие перехода к новой иерархии, в данном случае к трехуровневой структуре. Рассмотрим этот вопрос более подробно. Напомним, что время трехуровневой структуры вычисляется по формуле

$$3. T_{трехур} = C_{КОМД} T_{Z_1} + Z_{МЗД} T_{Z_2} + Z_{СМПСО} T_{Z_3}.$$

Рассмотрим отдельные параметры этой формулы. Заметим, что $Z_{ОМПОД}$ и $C_{КОМПОД}$ вычисляются по уже сформулированным выражениям 6+II.

Число зон массива МЗД - $Z_{МЗД}$ можно вычислить по формуле

$$Z_{МЗД} = \frac{K_1 \cdot Z_{СМПСО} \cdot Z_{СМПСО}}{Z_{МЗД}} \quad (19)$$

Физическая сущность этого параметра состоит в том, что он отражает число зон массива заголовков. Известно, что число всех заголовков, выходящих из зон массива ОМПОД, вычисляется следующим образом:

$$\text{число всех заголовков} = \frac{K_1 Z_{СМПСО} \cdot 64 \cdot Z_{СМПСО}}{2} \quad (20)$$

С другой стороны,

$$Z_{МЗД} = \frac{\text{число всех заголовков} \cdot 2}{Z_{МЗД} \cdot 64} \quad (21)$$

Подставляя (20) в (21), получим формулу (19).

В /4/ рассмотрена подробно структура массива МЗД. Там же отмечено, что в пределах данной зоны все одинаковые заголовки соединяются в цепь. В данном случае в ОМД накапливаются только заголовки цепных списков зон массива МЗД. Тогда все элементы массива ОМД можно вычислить следующим образом:

$$Z_{ОМД} = K_2 \text{ число всех заголовков} = \frac{K_2 K_1 Z_{СМПСО} \cdot 64 \cdot Z_{СМПСО}}{2} \quad (22)$$

Формула (22) может быть использована для вычисления $C_{К МЗД}$, которое обозначает среднее число зон на дескриптор в массиве МЗД.

$$C_{К МЗД} = \frac{Z_{ОМД}}{V} \quad (23)$$

Подставляя (22) в (23), получим

$$C_{К МЗД} = \frac{K_2 K_1 Z_{СМПСО} \cdot 64 \cdot Z_{СМПСО}}{2V} \quad (24)$$

По аналогии с формулой (17) можно привести формулу для вычисления числа зон второй части массива ОМД трехуровневой структуры.

$$Z_{СМПСО II} = \frac{K_2 K_1 Z_{СМПСО} \cdot 32 \cdot Z_{СМПСО} \cdot V}{32 \cdot Z_{СМПСО}} \quad (25)$$

Здесь же по аналогии с коэффициентами K_1 и K_2 вводится коэффициент K_3 , который связан с числом различных заголовков в данной зоне массива ОМД и который служит для вычисления параметра $C_{КОМД}$.

$$C_{КОМД} = \frac{K_3 Z_{СМПСО II} \cdot Z_{СМПСО} \cdot 64}{2V} \quad (26)$$

Подставляя (25) в (26), получим:

$$C_{КОМД} = \frac{K_3 K_2 K_1 Z_{СМПСО} \cdot 32 \cdot Z_{СМПСО} \cdot V}{V} \quad (27)$$

$$C_{КОМД} = K_3 (C_{К МЗД} - 1) \quad (28)$$

Применяя далее формулы к различным объемам массива при $Z_{ОМД} = Z_{МЗД} = 7PRU$ и $Z_{ОМПОД} = 140$ и при значениях коэффициентов $K_1 = 0,43$ и $K_2 = 0,9$, $K_3 = 0,9$ (полученных экспериментальным путем), получаем, что T трехур. $>$ T двухур. Различные комбинации разных величин зон не изменяют это неравенство. Дальнейший анализ пока-

зывает, что для лучшей работы трехуровневой структуры необходимо уменьшить коэффициент K_3 . Однако этот коэффициент зависит прежде всего от частоты встречаемости дескрипторов $f(j)$, т.е. от характеристик входного потока, т.к. он непосредственно связан с K_1 и K_2 . Поэтому если исследуемый информационный массив такой, что при реализации иерархической многоуровневой структуры K_1 , K_2 и K_3 получаются близкими к единице, то необходимо после нахождения оптимальных зон двухуровневой структуры оптимизировать эту структуру (путем уплотнения списков методом автоматической сортировки^{/2/}), а потом принимать решение о целесообразности перехода к трехуровневой структуре.

ЛИТЕРАТУРА

1. Д.Д.Арнаулов. Сообщение ОИАИ, РЮ-8622, 1975.
2. Д.Левковитц. Структуры информационных массивов оперативных систем, М., 1973.
3. Д.Д.Арнаулов. Препринт ОИАИ, Ю-7953, Дубна, 1974.
4. Д.Д.Арнаулов. Сообщение ОИАИ, РЮ-8621, Дубна, 1975.
5. Д.Д.Арнаулов, З.И.Коженкова. Сообщение ОИАИ, Ю-8799, Дубна, 1975.
6. Д.Д.Арнаулов, Н.И.Янев. Сообщение ОИАИ, П-8555, Дубна, 1975.
7. Д.Д.Арнаулов. Сообщение ОИАИ, Ю-7949, Дубна, 1974.

Рукопись поступила в издательский отдел
15 октября 1975 года.