

СООБЩЕНИЯ
ОБЪЕДИНЕННОГО
ИНСТИТУТА
ЯДЕРНЫХ
ИССЛЕДОВАНИЙ
ДУБНА



10519

ЭКЗ. ЧИТ. ЗАЛА

P11 - 10519

Д.Д.Арnaudов

МОДЕЛЬ ИПС ДЛЯ ОЦЕНКИ ПАРАМЕТРОВ
ИНФОРМАЦИОННЫХ СТРУКТУР

1977

P11 - 10519

Д.Д. Арнаудов

МОДЕЛЬ ИПС ДЛЯ ОЦЕНКИ ПАРАМЕТРОВ
ИНФОРМАЦИОННЫХ СТРУКТУР

Арнаулов Д.Д

P11 - 10519

Модель ИПС для оценки параметров информационных структур

Описана модель ИПС, предназначенная для оценки параметров инверсной, ассоциативно-адресной и частично инвертированной структуры. Параметры, относящиеся к запросам пользователей и банку данных, имеют статистический характер. Характеристики вычислительной машины относятся к параметрам устройств прямого доступа, типа магнитных дисков. Приводятся выражения для оценки среднего времени доступа в рассматриваемых структурах при работе информационно-поисковой системы, а также график функции, характеризующей запросы пользователей и частоту распределения дескрипторов в фонде ИНИС.

Работа выполнена в Лаборатории вычислительной техники и автоматизации ОНЯИ.

Сообщение Объединенного института ядерных исследований. Дубна 1977

Arnaudov D.D.

P11 - 10519

A Model for an Evaluation of Simulated Parameters of Information Structures

A model for an evaluation of simulated parameters of the multilist, inverted and partially inverted structure is described. Parameters, referring to the queries and the data-bank, have a statistical character. Parameters of the external devices refer to direct access devices, i.e. magnetic discs.

Expressions are derived for the evaluation of the average access time in these structures, during the work of an information retrieval system. The distributions of the queries and the frequencies of the descriptors in INIS data-bank are investigated.

The investigation has been performed at the Laboratory of Computing Techniques and Automation, JINR.

Communication of the Joint Institute for Nuclear Research. Dubna 1977

© 1977 Объединенный институт ядерных исследований Дубна

В информационно-поисковой системе имеется множество подсистем. Идеальная модель такой системы должна предоставлять возможность измерения и корректировки всех параметров данных подсистем. Однако этот вопрос в глобальном смысле пока еще не решен. Имеющиеся данные^{/1-4/} показывают, что большинство работ в области моделирования поисковых систем направлены на моделирование процедур поиска и на оценку стоимости работы этих процедур. Некоторые авторы^{/8/}, однако, описывают общую структурную модель поисковой системы, используя характеристики "реальных" массивов, а другие предлагают технику моделирования входного потока в информационно-поисковой системе^{/6/}.

Цель данной работы - построить модель, которая дала бы возможность исследовать основные параметры структуры информационно-поисковых систем (ИПС) на примере трех широко распространенных структур: инверсной, ассоциативно-адресной, и частично инвертированной^{/7-9/}. Как инверсная, так и ассоциативно-адресная организации массива представляют различные варианты структуры массива, реализация которых приводит к различным особенностям программирования. Однако обе организации могут представлять одну и ту же информационную структуру - иерархическую или ассоциативную. Они сходны в том смысле, что используют одни и те же средства логики разбиения, а именно адреса связи; в обоих случаях для разбиения массива необходимо указывать одно и то же количество адресов связи. Единственное отличие состоит в местонахождении этих адресов. Комбинируя инверсную и ассоциативно-адресную организации, можно получать различные варианты частично инвертированной структуры.

Основные параметры модели можно разбить на три группы: параметры, относящиеся к запросам пользователей, к банку данных (информационному массиву) и к характеристикам конкретной вычислительной машины, на которой реализуется система. Все параметры сведены в таблицу I, где массив ОМПД является основным информационным массивом, а массивы ОМД и МЗД представляют управляющую часть основного информационного массива [10]. Необходимо заметить, что большинство параметров, относящихся к запросам пользователей и к банку данных, имеют статистический характер. Кроме того, параметры вычислительной машины относятся к характеристикам устройств прямого доступа.

Все параметры модели хорошо описаны в работе [10]. Здесь мы приведем полученные выражения для вычисления среднего времени доступа к информации, расположенной на магнитном диске, для инверсной, ассоциативно-адресной и частично инвертированной структур. Вывод этих формул подробно рассмотрен в работах [9, 10].

$$T_{инв.} = N_p \frac{f(j)}{A} T_{z инв.} \quad (1)$$

$$T_{асс.} = L \cdot T_{z асс.} \quad (2)$$

$$T_{двуур.} = C_{к о м д} T_{z_1} + C_{о м п д} C_{к о м п д} T_{z_2} \quad (3)$$

В формулах (1) ÷ (3) $T_{z инв.}$, $T_{z асс.}$, T_{z_1} , T_{z_2} обозначают время, необходимое для чтения информации с соответствующих фрагментов дискового устройства. Если, например, величина узла меньше физической единицы записи, то время трансмиссии достаточно мало и тогда

$$T_{z асс.} = t_n + \frac{1}{2} t_0$$

Таблица I

к массиву	N — число дескрипторов словаря системы
	Z — величина зоны в физических записях
	V — количество различных дескрипторов, используемых в поисковых образах документов
	N_z — количество документов в массиве
	N_k — количество дескрипторов, индексирующих документ (среднее)
	$L = \frac{N_z N_k}{V}$ — количество документов на дескриптор (средняя длина списка)
	Q — количество символов на запись
	R_c — среднее количество записей в зоне
	$C_k = \left\{ \begin{array}{l} C_{к о м д} \\ C_{к м з д} \\ C_{к о м п д} \end{array} \right\}$ — среднее количество зон на дескриптор (соответственно в массивах ОМД, МЗД, ОМПД)
	$z_3 = \left\{ \begin{array}{l} z_{3 о м д} \\ z_{3 м з д} \\ z_{3 о м п д} \end{array} \right\}$ — число зон массивов ОМД, МЗД и ОМПД
к запросу	$z_3 = \left\{ \begin{array}{l} z_{3 з о м д} \\ z_{3 з м з д} \\ z_{3 з о м п д} \end{array} \right\}$ — среднее количество зон запроса (соответственно в массивах ОМД, МЗД, ОМПД)
	K — коэффициент, указывающий на наличие более чем одного дескриптора в зоне I части ОМД
	N_z — количество дескрипторов в запросе
	N_p — количество положительных дескрипторов в запросе
	L_s — длина самого короткого списка дескрипторов запроса
	ρ — отношение среднего числа ответов на запрос к L
	$\rho = \frac{z_{3 з о м п д}}{C_{к о м д}}$ — отношение среднего числа общих зон запроса в ОМПД к $C_{к о м д}$ (для двухуровневой схемы)
	$\delta = \frac{z_{3 з м з д}}{C_{к о м д}}$ — отношение среднего числа зон запроса в МЗД к $C_{к о м д}$ (для трехуровневой схемы)
	$\alpha = \left\{ \begin{array}{l} \alpha_{о м д} \\ \alpha_{м з д} \\ \alpha_{о м п д} \end{array} \right\}$ — отношение среднего количества общих зон запроса к C_k , где $\alpha_{о м д} = \frac{z_{3 з о м д}}{C_{к о м д}}$, $\alpha_{м з д} = \frac{z_{3 з м з д}}{C_{к м з д}}$, $\alpha_{о м п д} = \frac{z_{3 з о м п д}}{C_{к о м п д}}$
	γ — количество адресов записей массива на физическую запись
к ввб	$t_s = \left\{ \begin{array}{l} t_n \\ t_0 \\ t_{тр} \end{array} \right\}$ — время поиска цилиндра
	t_0 — время оборота
	$t_{тр}$ — время трансмиссии

В качестве частично инвертированной структуры рассматривается двухуровневая информационная структура $\langle 7, 10 \rangle$.

Как уже заметили, формулы 1 ÷ 5 и таблица 1 содержат параметры, которые имеют ярко выраженный статистический характер. Так, например, N_p — указывает на среднее число положительных дескрипторов запроса, L — является средней длиной дескрипторного списка запросов, $f(j)$ — указывает на число случаев, когда j — й дескриптор используется для индексирования в поисковых образах документов. В этом случае самый общий вид рассматриваемых формул имеет вероятностный характер, и они могут быть выражены следующим образом (в параметрах частично инвертированной структуры $S_{\text{инв. под. код.}} \cdot S_{\text{код.}} \cdot S_{\text{инв. под. код.}}$ уже учтен вероятностный характер рассматриваемых распределений $\langle 10 \rangle$):

$$T_{\text{инв.}} = T_{z \text{ инв.}} \sum_{j=1}^N \frac{f(j)}{A} \cdot P(j), \quad (4)$$

где A — число записей в одной физической записи.

$$T_{\text{acc.}} = T_{z \text{ acc.}} \sum_{j=1}^N f(j) P(j). \quad (5)$$

В приведенных формулах принимается, что в рассматриваемых структурах имеется $f(j)$ документов в списке, соответствующем j — му дескриптору; следовательно, j — й список имеет длину $f(j)$. Тогда время, необходимое для обработки, например, всех узлов в списке ассоциативно-адресной структуры, равно $f(j) \cdot T_{z \text{ acc.}}$. Это есть и среднее, и максимальное время, необходимое для обработки всех узлов. Кроме того, имеется в виду, что каждая запись охватывает не больше чем одну физическую единицу записи.

Из предложенных формул видно, что выбор соответствующих дескрипторов запроса осуществляется на основе распределения функции $P(j)$.

а распределение дескрипторов информационного массива описывается функцией $f(j)$. Следовательно, для моделирования входного потока запросов пользователей и формирования банка данных необходимо иметь математическое выражение для функций $P(j)$ и $f(j)$.

Раскроем более подробно физическую сущность этих функций. Для этого необходимо напомнить, что под банком данных ИПС, или, иначе говоря, под понятием "информационного массива" системы, будем подразумевать набор записей поисковых образов документов, каждый из которых характеризуется некоторым количеством дескрипторов.

Пусть общее количество дескрипторов в дескрипторном словаре равно N . Код каждого дескриптора характеризуется некоторым числовым индексом j , где $1 \leq j \leq N$.

Пусть $f(j)$ — это общее число случаев, когда j — й дескриптор используется для индексирования в поисковых образах документов (иначе говоря, это частота использования данного дескриптора в различных документах массива).

Пусть S — это общее число случаев, в которых все дескрипторы встречаются в различных документах массива. Следовательно,

$$S = \sum_{j=1}^N f(j).$$

данная формула характеризует суммарное использование дескрипторов в различных документах массива и является, как уже заметили, характеристикой самого массива.

В качестве функции, характеризующей запросы пользователей, можно выбрать $P(j)$ — вероятность того, что выбранный в запросе дескриптор есть j — й дескриптор.

Функции $f(j)$ и $P(j)$ могут быть неизвестными или очень трудно определяемыми. Вид этих функций может зависеть как от конкретной информации в массиве, так и от метода индексирования этой информации и т.п.

Поэтому особый интерес представляет расчет отдельных параметров модели при различных распределениях функций $f(j)$ и $P(j)$, и их сравнение с параметрами модели при использовании реального информационного массива и реальных запросов пользователей. Тогда можно показать, какое распределение этих функций является самым близким к реальному, и использовать это при моделировании указанных структур на вычислительной машине.

Теоретически показано, что возможные распределения $f(j)$ и $P(j)$ лежат в области, определяемой равномерным распределением и распределением по закону Ципфа^{/9/}. Тогда возможны четыре случая изменения функций $f(j)$ и $P(j)$.

1. Равномерное распределение предполагает, что функции равномерно распределены и постоянны для всех j , т.е.

$$f(j) = \frac{S}{N} ; P(j) = \frac{1}{N}$$

2. Функции имеют распределение по закону Ципфа, т.е.

$$f(j) = \frac{S}{j(\ln N + \gamma)} ; P(j) = \frac{1}{j(\ln N + \gamma)}$$

В этом случае используется формула Ципфа с поправкой Lowe^{/9, II/}. Здесь $\gamma = 0,5772$ и употребляется в качестве константы Эйлера.

3. Функция $f(j)$ распределена по закону Ципфа, а $P(j)$ имеет равномерное распределение, т.е.

$$f(j) = \frac{S}{j(\ln N + \gamma)} ; P(j) = \frac{1}{N}$$

4. Функция $P(j)$ распределена по закону Ципфа, а $f(j)$ имеет равномерное распределение, т.е.

$$f(j) = \frac{S}{N} ; P(j) = \frac{1}{j(\ln N + \gamma)}$$

Основные параметры модели для массива из 3200 и 10000 документов показаны в таблицах 2,3. В таблице 4 показаны параметры модели реального массива ИНИСа^{/13/} для 3200 и 10000 документов. Эксперименты проводились для двух значений величины зоны информационного массива: $Z = 7$ физических единиц (PRU) и $Z = 140$ физических единиц.

В качестве запросов использовались 300 различных дескрипторов, генерированных в соответствии с различными случаями ($I \div 4$), а также пакет реальных запросов пользователей (300 запросов). Запросы были представлены в конъюнктивной форме. В каждом запросе имелось в среднем по три дескриптора.

На рис.1 показаны графики функции $f(j)$ для реального массива из 3200 документов (штрихпунктирная линия), для массива из 10000 документов (пунктирная линия) и сплошной линией изображен график функции Ципфа. На рис.2 показаны график распределения дескрипторов в реальном запросе пользователей (пунктирная линия) и график функции Ципфа.

Необходимо заметить, что для моделирования различных распределений ($I \div 4$) использовался хорошо известный в теории вероятности метод воспроизведения событий с заданной вероятностью. Этот метод состоит в следующем. Разобьем интервал $/0,1/$ на N подинтервалов $D_1, D_2 - D_N$. Правый конец δ_i интервала D_i вычисляется по формуле:

$$\delta_i = \sum_{j=1}^i P(j),$$

$$P(j) = \frac{1}{j(\ln N + \gamma)}$$

Для каждого документа моделируемого массива генерируются N_k случайных точек / где N_k - среднее количество дескрипторов на документ / с равномерным распределением в интервале $/0,1/$. Если случайная точка

Таблица 2

$$f(j) = \frac{S}{j(\theta \pi N + \gamma)}$$

	N_z	V	L	N_z	V	L
	3200	6552	4,3	10000	9459	9,5
	$S = N_z N_k = 28800$			$S = N_z N_k = 90000$		
$Z (PRU)$	7	140	7	140		
C_k омпод	3,14	1,97	6,81	4,26		
C_k омпд	1,52	0,47	4,6	1,9		
z_3 омпод	134	7	417	21		
z_3 омпд	54	1,2	246	11		
число всех заголовков омпд	20582	12912	64420	40373		
z_3 омпод $P(j) = \frac{1}{j(\theta \pi N + \gamma)}$	1	0,7	1,4	1,05		
α омпод при $P(j) = \frac{1}{j(\theta \pi N + \gamma)}$	0,3	0,3	0,2	0,25		
z_3 омпод $P(j) = \frac{1}{N}$	0,3	0,6	0,68	0,45		
α омпод при $P(j) = \frac{1}{N}$	0,09	0,3	0,1	0,1		

Таблица 3

$$f(j) = \frac{1}{N}$$

	N_z	V	L	N_z	V	L
	3200	9419	3	10000	9997	9
	$S = N_z N_k = 28800$			$S = N_z N_k = 90000$		
$Z (PRU)$	7	140	7	140		
C_k омпод	3	2,54	8,9	7,4		
C_k омпд	2,8	2	8,8	6		
z_3 омпод	134	7	417	21		
z_3 омпд	86	4	353	15		
число всех заголовков омпд	28481	23972	89022	74510		
z_3 омпод $P(j) = \frac{1}{j(\theta \pi N + \gamma)}$	1,1	0,9	2,7	2,4		
α омпод при $P(j) = \frac{1}{j(\theta \pi N + \gamma)}$	0,3	0,36	0,3	0,32		
z_3 омпод $P(j) = \frac{1}{N}$	1,24	1,1	3,8	3,6		
α омпод при $P(j) = \frac{1}{N}$	0,4	0,44	0,42	0,48		

Таблица 4

Реальный массив ИНИСА

	N_z	V	L	N_z	V	L
	3200	2220	9	10000	7780	11
	$S = N_z N_k = 28800$			$S = N_z N_k = 89851$		
Z (PPI)	7	140	7	140		
Ск омпод	6,33	2,6	8,03	3,9		
Ск омпд	4,46	0,6	5,58	1,37		
τ_{30} омпод	132	6,5	410	21		
τ_{30} омпд	78	1,2	245	6		
число всех заголовков ОМД	20000	12100	62540	31076		
τ_{32} омпод	0,6	1,2	1,5	0,9		
α омпод	0,1	0,2	0,18	0,23		

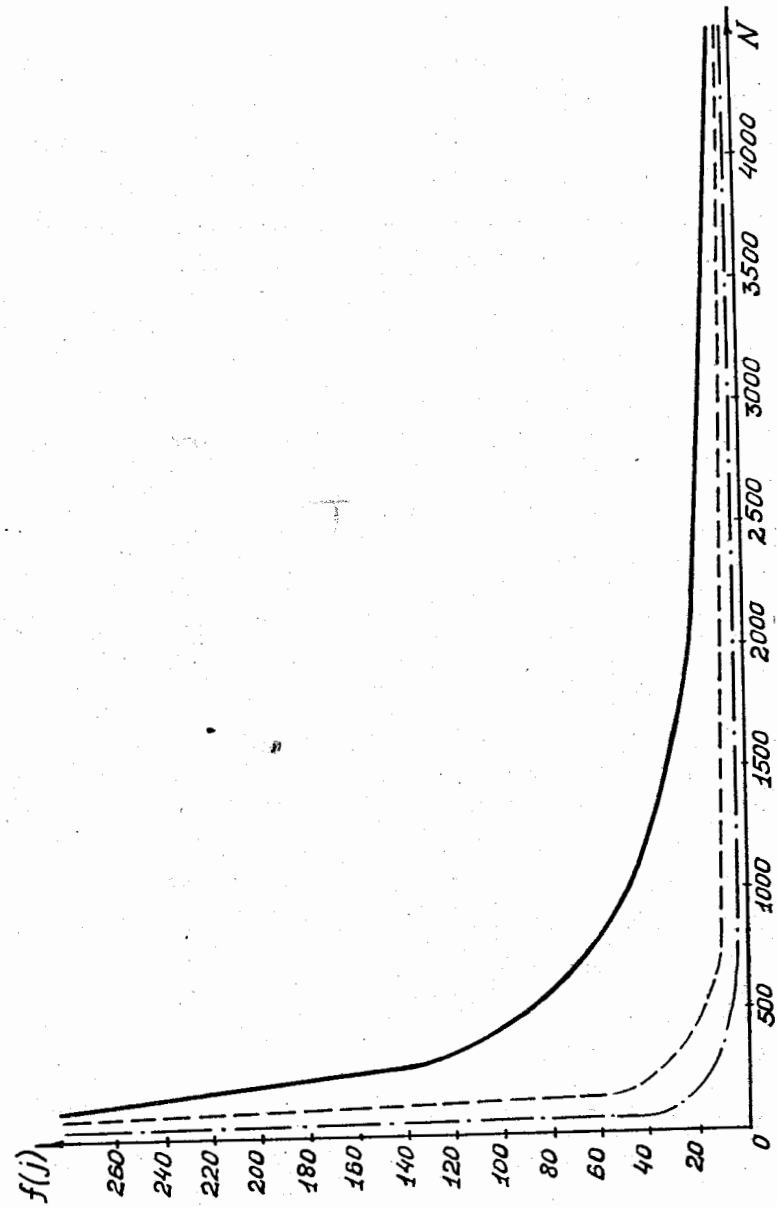


Рис. I

попадает в i -й подинтервал, считаем, что в данном документе встречается i -й дескриптор. Поскольку длина i -го подинтервала равняется $\delta_i - \delta_{i-1}$, то данная процедура приводит /при S обращениях к генератору случайных чисел/ к появлению i -го дескриптора в среднем $P(i) \cdot S$ раз, т.е. в данном случае по закону Цифа частота использования дескриптора получается равной $\frac{S}{i(\theta n N + \gamma)}$.

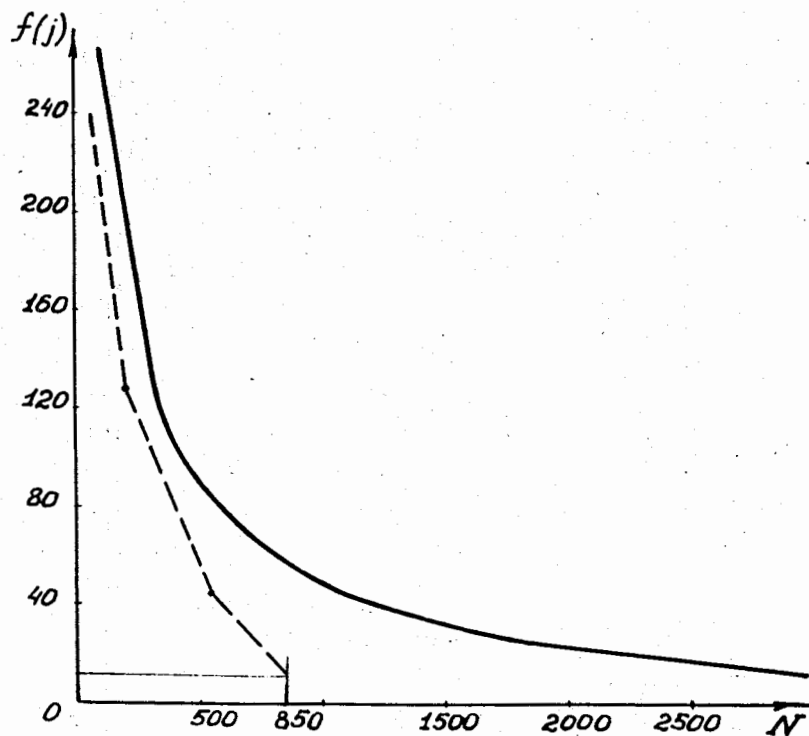


Рис. 2

Расчет среднего времени доступа для инверсной, ассоциативно-адресной и частично инвертированной структур на основе банка данных ИНИСа (из 10000 документов) на ЭВМ CDC 65000.

Значения основных параметров модели сведены в таблицу 4. Кроме указанных параметров для вычисления времени доступа используются и характеристики ЭВМ CDC, которые подробно описаны в работах^{/9, 12/}. С помощью этих характеристик вычисляются величины $T_{2 \text{ инв}}, T_{2 \text{ асс}}, T_2, T_2$.

В связи с этим необходимо иметь в виду, что для чтения с магнитного диска информации величиной в одну физическую единицу записи в среднем затрачивается время $t_n + \frac{1}{2} t_0$. Следовательно, для чтения информации в семь физических единиц (т.е. 7PRU) необходимо время $T_7 = t_n + \frac{1}{2} t_0 + t_0$, где $t_n = 75$ мс, $t_0 = 25$ мс. Следовательно, $T_7 = 1,3(t_n + \frac{1}{2} t_0)$. Для чтения информации величиной в 140 физических единиц записи необходимо $T_{140} = 6,7(t_n + \frac{1}{2} t_0)$. Более подробно методика расчета времени доступа описана в работе^{/12/}. Необходимо заметить, что время доступа в двухуровневой частично инвертированной структуре рассчитано для величин зон $Z_7 = 7PRU$ и $Z = 140PRU$.

$T_{\text{инв}}$ для запроса, состоящего из пяти дескрипторов, можно рассчитать по формуле (I). Для этого, если на условный номер документа уходят в среднем два слова CDC, то формула (I) преобразуется в формулу (Ia), где

$$T_{\text{инв}} = 5(t_n + \frac{1}{2} t_0 + \frac{2L}{448} t_0) = 5t_n + \frac{5}{2} t_0 + \frac{5L}{224} t_0, \\ \text{при } t_n = 75 \text{ мс, } t_0 = 25 \text{ мс, } L = 11 \quad T_{\text{инв}} = 5,1(t_n + \frac{1}{2} t_0).$$

В данном случае считается, что к каждому дескриптору в среднем относится L документов.

Из данного примера видно, что время доступа в инверсном массиве зависит как от средней длины списка, так и, в еще большей степени, от числа положительных дескрипторов запроса.

T_{acc} вычисляется по формуле (2). Для массива из 10000 документов $L = 11$ (см. табл.4). При расчете T_{zacc} принимается, что величина узла не превосходит одной физической единицы записи. Тогда

$$T_{zacc} = t_n + \frac{1}{2}t_0 \quad \text{и} \quad T_{acc} = 11 \left(t_n + \frac{1}{2}t_0 \right).$$

Из данного примера видно, что среднее время обработки ассоциативно-адресного массива зависит прежде всего от средней длины списка.

$T_{двухур.}$ можно рассчитать по формуле (3). При $T_{z1} = 1,3 \left(t_n + \frac{1}{2}t_0 \right)$ и $T_{z2} = 6,7 \left(t_n + \frac{1}{2}t_0 \right)$, $L_{сма} = 2,61$. Тогда $T_{двухур.} = 9,5 \left(t_n + \frac{1}{2}t_0 \right)$.

Если обе зоны имеют одинаковую величину в семь физических единиц, то $T_{двухур.} = 5,58 \cdot 1,3 \left(t_n + \frac{1}{2}t_0 \right) + 0,18 \cdot 8,03 \cdot 1,3 \left(t_n + \frac{1}{2}t_0 \right) = 9,1 \left(t_n + \frac{1}{2}t_0 \right)$.

данный пример показывает, что, выбирая соответствующим образом параметры частично инвертированной структуры, можно улучшить среднее время поиска в этой структуре.

Выводы.

1. На основе полученных результатов можно утверждать, что самым близким к реальному распределению является распределение функций $f(j)$ и $P(j)$ по закону Ципфа.
2. Используя формулы (1-3) и параметры предложенной модели, можно рассчитать среднее время доступа в указанных структурах. Тем самым для конкретной ИИС можно выбрать самую подходящую структуру.
3. Данная модель дает возможность судить об эффективности разрабатываемой ИИС на конкретной ЭВМ в зависимости от характера $f(j)$ и $P(j)$. Следовательно, на данной ЭВМ можно построить ИИС в среднем достаточно эффективную для различного рода банков данных при условии, что все они имеют одинаковый закон распределения.
4. Описанная модель дает возможность выбирать параметры рассмотренных структур в зависимости от конкретной ЭВМ и природы банка данных.

Литература

1. IEM Corporation: Bibliography on Simulation, White Plains, C20-I649, 1969.
2. Cooper M.D., J. Information Storage and Retrieval, v 9, pp.13-32, 1973.
3. Senko M.E. oth., Proceedings of IFIP, pp.514-519, 1969.
4. Griffiths J.M. oth., ASLIB Conference, pp.141-152, Amsterdam, 1976.
5. Cardens A.F., Comm. ACM, v 16, No.9, pp.540-548, 1973.
6. Griffiths J.M., J. Documentation, v. 31, No.3, pp.185-190, 1975.
7. Lefkowitz D., "File Structures for On-Line Systems", Spartan Books, 1969.
8. Prywes N.S., Proceedings of IFIP, pp.273-278, 1962.
9. Арнаудов Д.Д. Препринт ОИЯИ, IO-7953, 1974, Дубна.
10. Арнаудов Д.Д. Сообщение ОИЯИ, P10-9178, 1975, Дубна.
11. Lowe T.C. J. ACM, v 15, No.4, pp.534-548, 1968.
12. Арнаудов Д.Д., Сообщение ОИЯИ, IO-7949, 1974, Дубна.
13. Арнаудов Д.Д. и др. Депонированная публикация ОИЯИ, Б1-11-3553, 1975, Дубна.

Рукопись поступила в издательский отдел

22 марта 1977г.