

## Л и т е р а т у р а

1. Евдокимов В.П., Покрасс В.М.  
Методы обработки данных в научных космических  
экспериментах. М., "Наука", 1977.

И.С. Ризаев

### СЖАТИЕ БАЗЫ ДАННЫХ

Запоминающие устройства с прямой записью (ЗУПД) позволяют создавать информационные системы, базирующиеся на единой интегрированной базе данных. Для эффективной работы с базой данных необходимо, чтобы прямой доступ был обеспечен в любое время ко всей базе. Но ввиду ограниченности дисковой памяти и больших объемов информации такой доступ, как правило, не обеспечивается. Решением этой проблемы может быть сокращение избыточности хранимой информации в базе данных путем уплотнения или сжатия этих данных.

Пусть  $M = \{m_i\}$  - множество массивов информации, составляющих базу данных до сжатия;  $M' = \{m'_i\}$  - множество полученных массивов информации после сжатия данных;  $A = \{a_i\}$  - алгоритмы преобразования,  $i = \overline{1, n}$ . Видимо, сжатие данных можно проводить только в том случае, если будет выполняться соотношение

$$M' + A < \bar{M},$$

где  $\bar{M}$ ,  $M'$ ,  $A$  - объемы соответствующих массивов и алгоритмов.

Таким образом, задача сводится к созданию или подбору соответствующих алгоритмов сжатия. Существуют различные методы сжатия: машинные, немашинные, методы оптимального размещения, комбинированные. Хотя методы сжатия разрабатываются давно, применительно к базам данных они стали рассматриваться сравнительно недавно [1]. Большинство известных методов сводятся к рассмотрению сжатия информации, расположенной непосредственно в памяти ЭВМ. Чаще используют лексическое кодирование [2], суть которого заключается в том, что более длинному сообщению ставится в соответствие более короткое, а именно, порядковый номер данного сообщения. В справочниках или в памяти ЭВМ хранятся два словаря:

словарь слов или сообщений и словарь их номеров, представленных в виде двоичных кодов. При уплотнении записи с символьной информацией заменяются их двойничными номерами.

В работе рассматривалось сжатие информации применительно к базам данных, используемым в системах интегрированной обработки, типа СИОД-1, СИОД-2 [3]. База данных, создаваемая ППП СИОД-1 и СИОД-2, состоит из массивов двух типов: главных и связующих. Каждый массив состоит из набора записей фиксированной длины. Формат записей  $Z$  представляет собой набор полей, часть из которых образуют системные поля  $\{C\}$ , а часть представляют поля пользователей  $\{P\}$ :

$$Z = \langle C_1, C_2, \dots, C_m, P_1, P_2, \dots, P_n \rangle.$$

Поля записей могут быть обязательными и необязательными. Обязательные поля в любом входном массиве могут быть расположены произвольно, но их длина и местоположение задаются. Поля пользователей могут быть произвольной длины и произвольного содержания. Поскольку поля пользователей имеют значительно большую длину, чем системные поля, то уплотнение необходимо проводить именно по отношению к этим полям. Выбор тех или иных алгоритмов сжатия зависит от вида полей записей. В СИОД-1, СИОД-2 встречаются поля, включающие только символьную информацию, например наименование предмета. Встречаются поля, включающие только цифровую информацию хранящуюся в упакованном формате, поля с двоичной информацией и смешанные поля, включающие как символьную так и цифровую информацию. Естественно, что для различных видов полей должны применяться различные виды сжатия. Например, для малосимвольных полей, длиной в один, два, три байта, могут подойти обычные таблицы перекодировок. Для полей длиной в четыре и более символов, таблицы перекодировок могут не подойти в связи с большими затратами времени на поиск информации. Для таких полей, как "наименование", выгоднее в памяти хранить набор дескрипторов (ключевых слов). В некоторых случаях для ускорения процесса преобразования исходного слова в сжатый код можно использовать методы однозначного преобразования без промежуточного обращения к таблицам перекодировок. Цифровую десятичную информацию выгоднее хранить в виде двоичных чисел путем непосредственного преобразования десятичных чисел в двоичные.

В общем случае задача скатия должна решаться как оптимизационная. Критериями оптимизации являются заданный объем оперативной памяти и время перекодировок.

Пусть  $\{a_j\}$  - множество алгоритмов скатия информации  $p_i$  поля,  $j = \overline{1, n}$ ,  $i = \overline{1, m}$ ;

$b_j$  - число ячеек памяти, требуемое для алгоритма  $a_j$ ;

$t_j$  - время скатия, затрачиваемое алгоритмом  $a_j$ ;

$K_j$  - степень скатия информации алгоритмом  $a_j$ ;

$m$  - число полей (реквизитов) в  $Z$  записи.

Таким образом, запись  $Z$  включает набор полей, предназначенных для скатия  $p_1, p_2, \dots, p_i, \dots, p_m$ .

Соответственно запишем матрицу алгоритмов скатия:

$$\begin{array}{cccc} a_{11} & a_{21} & \dots & a_{i1} & \dots & a_{m1} \\ a_{1j} & a_{2j} & \dots & a_{ij} & \dots & a_{mj} \\ a_{1n} & a_{2n} & \dots & a_{in} & \dots & a_{mn} \end{array} \parallel$$

$A = \{a_{ij}\}$  - множества алгоритмов скатия.

Задача выбора вариантов скатия будет состоять в том, что необходимо выбрать такой набор алгоритмов  $a_{ij}$ , чтобы их объем не превосходил  $S$  ячеек оперативной памяти, а время счета  $T$  было бы минимальным.

Введем переменные

$$x_{ij} = \begin{cases} 1, & \text{если } a_{ij} \text{ входит в } A \\ 0, & \text{в противном случае} \end{cases}$$

Тогда

$$\sum_{i=1}^m \sum_{j=1}^n b_{ij} x_{ij} \leq S;$$

$$\sum_{i=1}^m \sum_{j=1}^n K_{ij} x_{ij} \rightarrow \max;$$

$$\sum_{i=1}^m \sum_{j=1}^n t_{ij} x_{ij} \rightarrow \min.$$

Элементы  $x_{ij}$  должны удовлетворять следующему условию:

$$x_{ij} = 1, \quad j = \overline{1, n}.$$

Если ввести ограничение на время  $T$ , то данную задачу можно свести к задаче линейного программирования.

При скатии важно провести предварительную оценку временных соотношений. Такую оценку можно провести путем сравнения временных затрат на обращение к записям в ЗУПД и временных затрат на перекодировки.

Рассмотрим время доступа к записям, расположенным на магнитных дисках. Время, требуемое на доступ к данным в ЗУПД, и время на пересылку данных состоит из четырех частей: времени перемещения механизма доступа, время на выбор головки, задержки от вращения носителя информации и время на пересылку данных [4].

Время перемещения механизма доступа к цилиндру  $t'$ , содержащему нужную запись, является функцией от числа цилиндров и может быть определено по следующей приближенной формуле:

$$t' = (40 + 0,5d) 10^{-3} \text{ с} .$$

Здесь  $d$  - количество пересекаемых дорожек ( $1 \ll d \ll 200$ ). При этом  $t_{min} = 25 \text{ мс}$ ,  $t_{max} = 135 \text{ мс}$ . Среднее время доступа составит  $t'_{cp} = 90 \text{ мс}$ .

Время на выбор головки очень мало и им пренебрегают. Задержка от вращения носителя составит  $t_{max}^* = 25 \text{ мс}$ ,  $t_{cp}'' = 12,5 \text{ мс}$ . Время на передачу данных  $t_{cp}'' \approx 6 \text{ мкс} \approx 0,006 \text{ мс}$ . Таким образом, общее среднее время доступа к записи составит:

$$T_{cp} = t'_{cp} + t_{cp}'' + t_{cp}''' .$$

К  $N$  записям  $T = NT_{cp}$ .

Сжатие информации будет в основном влиять на время  $t'_{cp}$ , так как при этом будет изменяться число пересекаемых дорожек. При скатии в  $K$  раз получим:

$$t'_{cp}^{сж} = (40 + 0,5 \frac{d}{K}) 10^{-3} ;$$

$$T_{cp}^{сж} = t'_{cp}^{сж} + t_{cp}'' + t_{cp}''' ;$$

$$T_{пол}^{сж} = NT_{cp}^{сж} .$$

При оценке временных затрат на перекодирование необходимо исходить из того, что наибольшее время будет затрачиваться на поиски информации в таблицах уплотнения. Пусть система скатии включает  $n$  таблиц с длинами  $L_1, L_2, \dots, L_n$ . Средняя длина таблицы составит  $L_{cp} = \frac{1}{n} \sum_{i=1}^n L_i$ .

Число операций сравнения составит величину  $Q_{cp} = L_{cp} q$ , где  $q$  - цикл, выраженный количеством операций, затрачиваемых на выборку кода слова из словаря в стандартные ячейки и на сравнение с кодом исходного слова,  $q \approx 20$  операций [2]. Время, затрачиваемое на один поиск, составит

$$t = \frac{1}{2} Q_{cp} t,$$

где  $t$  - время одной операции сравнения. При обращении к  $N$  записям  $E_{пол} = NE$ .

Время доступа к записям, расположенных на магнитных дисках, после сжатия составит величину

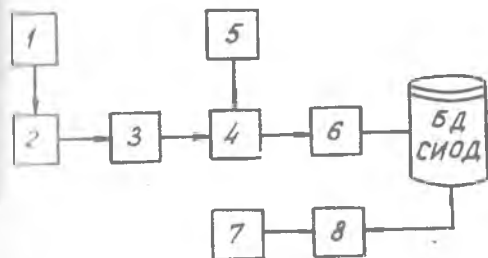
$$T^{сж} = T_{пол}^{сж} + E_{пол}.$$

Целесообразность применения алгоритмов сжатия определится из сравнения величин  $T^{сж}$  и  $T$ .

Рассмотрим схему сжатия базы данных. Сжатие базы данных можно проводить тремя методами: ручным, полуавтоматическим и автоматическим. При ручном методе заранее составляются таблицы и методы уплотнения, которые хранятся в специальных альбомах. При вводе исходной информации производится поиск соответствующих сжатых кодов в альбоме и их ввод в ЭВМ. Декодирование производится аналогичным образом.

При полуавтоматическом методе таблицы и схемы уплотнения, оставленные предварительно вручную, размещаются в памяти ЭВМ. При вводе исходной информации производится автоматическое обращение к таблицам и замена исходных кодов слов сжатыми кодами.

При использовании автоматической системы из исходной необработанной информации автоматически создаются таблицы и схемы уплотнения и производится автоматическое обращение к ним (рис. 1).



Р и с. 1.

Из исходных неуплотненных данных (блок 1), находящихся на перфокартах или на МЛ, создается сводная ведомость (блок 2), в которой по каждому реквизиту записей будут представлены все значащие характеристики. Из сводной ведомости

мости создаются таблицы перекодировок (блок 3) или различные схемы уплотнения. Блок оптимизации (5) предназначен для настройки программ уплотнения (4) на выбор тех или иных алгоритмов сжатия в соответствии с заданными критериями оптимизации. Уплотненные данные (блок 6) размещаются на любом из внешних носителей, после чего, используя стандартную программу генерации, можно сформировать базу данных СИОД с уплотненной информацией.

Функциональные пакеты (блок 7) смогут работать с уплотненной базой данных СИОД через программы перекодировок (блок 8).

Проведенная подобным образом работа с системой интегрированной обработки СИОД позволила получить сжатие базы данных более чем в два раза без существенных временных затрат.

#### Л и т е р а т у р а

1. *Alsberg Peter A. Space and time savings through large data base compression and dynamic restructuring. Proc. IEEE, 1975, 63, №.*
2. Курбаков К.Н. Кодирование и поиск информации в автоматическом словаре. М., "Советское радио", 1968
3. Келехсаев А.А., Беляев А.П. Системы интеграции и обработки данных СИОД-1, СИОД-2. М., "Статистика", 1977.
4. Система ИБМ/360. Введение в запоминающие устройства прямого доступа и методы организации данных. М., "Статистика", 1974.