

ハードウェア・ソートによるデータベース 基本演算の高速化について

BASIC OPERATIONS OF THE DATA BASE BASED ON THE HARDWARE SORTING

有澤 博

Hiroshi ARISAWA

横浜国立大学工学部

Yokohama National University

土肥 康孝

Yasunori DOHI

1. まえがき.

現在、広汎な分野で使われるようになった計算機システムの利用形態は、大きく分けて、
① 単発的に発生する数値計算的な仕事を処理、
② 外部システムの制御と直接データ処理、
③ 大規模データを恒常的に管理し、問合せ、
更新に応じる、の3種に分類できる。

今後、②のリアルタイム制御と、③のデータベース・マネジメントは電子技術の飛躍的発展と呼応して、ますます多くの分野に進出するであろう。ところで、②はハードウェアも含めた総合的なシステムとしてとらえられることが多いのに対し、③では多様な実務上の要求に対しても、長い間かなり現実的なレベルで、ソフトウェア的な対応がとられてきたにすぎなかった。最近になって、個別的なシステムを統合化し、新しい理論モデルを形成しようとする試みが活発になり、その一方で、モデルをその基本理念に即したハードウェアで効率よく実現させようとすう試み—データベース・マシンが考えられはじめている^{1)~3)}—データベース・マシンを特徴づけるものは、①ハードウェアを意識した記憶モデルと、②データ集合演算の並列処理である。しかし今までのところ、データベース・マシンに必要な機能と、それを実現させるためのアーキテクチャとの関係について、上記2点のもとに十分検討されているとは言えない。

本報告では、まず2章でこれらを検討と問題点の整理を行う。つづいて3章で有用な機能のひとつである3) ハードウェア・ソートについて、

特に並列性の点から議論する。4章では、実際にデータベースの基本演算子をいくつかとり上げ、ハードウェア・ソートによる演算の実現法を紹介する。

2. データベース・マシンのアーキテクチャと問題点.

本章ではデータベース・マシンのアーキテクチャと、それを実際に構成する上で特に考慮しなければならない点について、いくつかの側面から検討する。

2.1. データベース・マシンの位置づけ

データベース管理システム(DBMS)が担当する作業は、①データ操作言語(DML)で書かれた指令を、現実のデータベースに対するデータ処理の手順に変換する部分と、②データ処理をハードウェアの特性や、システム資源の利用状況を考慮して、最適効率で実行する部分とに分けられる。本稿では②の機能を、そのためには設計された特別のハードウェアで、効率よく実行させるものをデータベース・マシンとよぶ^{注1)}。

まず、②のデータ処理の手順を図2.1のように一般化して考えよう。ある情報の検索を行う

注1) ①の機能を中心に考え、ミニコンフミのディスクのような形のものもデータベース・マシンとよがることがある³⁾。

場合、はじめに、すべての情報を含むマスター・スペースの必要なファイルからデータ集合を取り出し(1)、システム・バッファ上で検索の対象になりうるデータを抽出して集合を小さくしたり、集合と集合を結びつけたりする(2)。そして結果を利用者のワーク・スペースに送り(3)、最終的なデータ操作を行って、出力する(4)。更新の場合には、ほぼこの反対の経路をとる(1, 2, 5, 6, 7)。

さて、データベース・マシンを構成するときに特に重要な分か水目になるのは、マスター・スペースをデータベース・マシン内にもつかどうかである。いま、マシン内にあくものを便宜的に出入力装置型、外部にある汎用計算機(GPC)のディスクなどを利用するものとチャネル型とよぶことにする。出入力装置型のデータベース・マシンは、メモリ自体をデータベース処理向きに特殊な機能とも言えることもできるし、システム・バッファと機能要素を多用して並列処理も行いやすい。たとえばトロント大学のRAP⁴⁾は、メモリはディスクで、各トラック毎にセルとヘッドがある。メモリとしては磁気バブルやCCDの利用も考えられ、新しい機能メモリの提案もされてる⁵⁾。しかし言うまでもなく、このようなメモリは主記憶用素子に比べ安いものの、従来のディスクに比べ高価であり、データベース本体を特殊メモリ上におくのは、将来の技術動向を考えてもなお疑問であろう。

チャネル型のデータベース・マシンのブロック図を図2.2に示した。この場合は、どうしても大量のデータ転送が必要になり、転送速度と必要なデータのみを通過させるフィルター的な機能が重要な問題になるであろう。チャネル型のデータベース・マシンの構想は、今のところ筆者らが考えてるものの⁶⁾の他はないが、データベース・マシンを実現性と云う点から考えると、何とかの形でチャネル型の形式をどうぞを得ないであろう。(RAPでも、チャネル型マシンとして考えようとする論文⁷⁾がでている。)

2.2. 並列処理のアーキテクチャ

データベースの処理が、本質的に集合に対する、あるいは集合間の演算である以上、並列処理は効率向上のための有力な手段である。また、

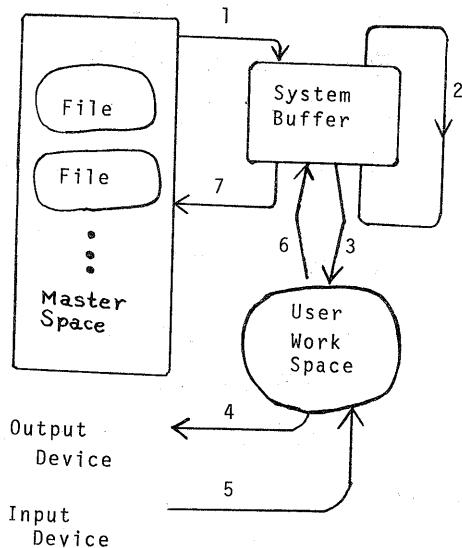


Fig. 2.1 Data Flow in DBMS

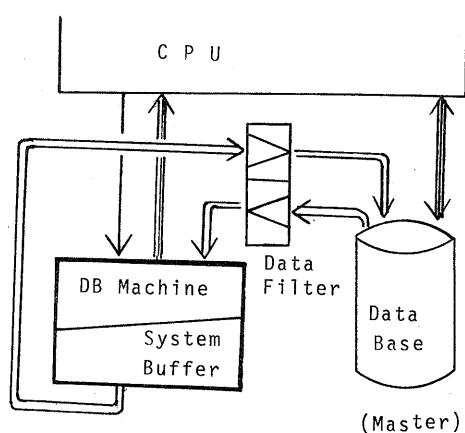


Fig. 2.2 Block Diagram of DB Machine

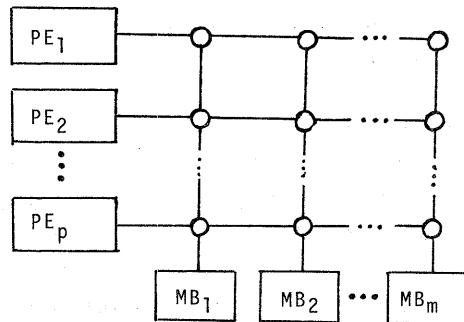
最近の技術的進歩から、多数のプロセッサによる並列処理方式は十分に実現性がある⁸⁾。ここでは図2.2のマシン形態をとることとして、DBマシンとシステム・バッファ上で、並列処理などのように実現させるかを考える。

ところで、あるひとつの仕事(JOB)を複数のプロセッサで並列処理(parallel processing)させるとには、つきの2つの形態が考えられる。やーは、複数のプロセッサ間で共有するメモリをもち、各プロセッサが共有メモリ上にあるデータあるいは処理プログラムの一部を、任意に占有したり解放したりする、いわゆる多重処理(multi-processing)の形態である。や2は、各プロセッサが独自のメモリをもち、1まとまりの処理に必要なデータと処理プログラムをとり込んで、各プロセッサ毎に全体の仕事の1部を請け負う形の処理形態、いわゆる分散処理(distributed processing)である。^{注1)}

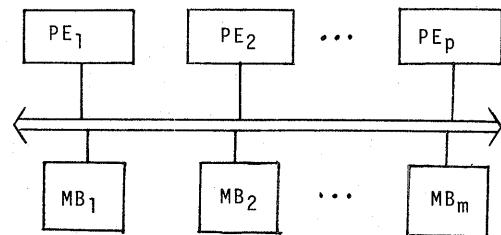
や1の形態では、1つの仕事を合計いくつの単位(プロセス)に分解するかは、あまり考慮する必要がない。すなはちプロセッサとプロセスは動的に任意に組合せることができるために、プロセスの個数は問題にならない。またこの意味で、仕事に対する各プロセッサの寄与は均等であると言える。これに反してや2の形態では各プロセッサに割当てられたプロセスは、始めから終りまでそのプロセッサ上で実行されるのが普通なので、処理効率向上のために、データの配分や処理内容を各プロセッサに対して均質にし、負荷をそろえる工夫が必要である。

このように、一般的には分散処理よりも多重処理の形態の方が有利であるが、プロセッサの数が多くなると多重処理には問題がある。すなはち、メモリ共有型並列処理の最大のネックは、メモリ・アクセスの衝突が起こることである。このたび通常はメモリをいくつかのメモリ・バスク(MB)に分けており、各プロセッサ(PE)に対してマトリックス状のスイッチを経由して割当てる(図2.3(a))か、共通バス線と各PEが時分割を利用して(図2.3(b))方法となる。しかし、前者では p 台のPEと m 台のMBに対して $p \times m$ 個のバス・スイッチが必要であり、これは高価で制御も複雑である上にPEやMBの追加が困難である。一方後者は共通バス線を1

注1) や2の形態で、各プロセッサが同じ処理プログラムを同期して実行する、という制限をつけたものがILLIAC-IVにみられるようなアレイ・プロセッサである。ここでは、もっと広く、むしろ緊密に結合されたコンピュータ・ネットワークに近いものを考えている。



(a) Matrix Switch



(b) Common Bus

Fig. 2.3 Shared Memory

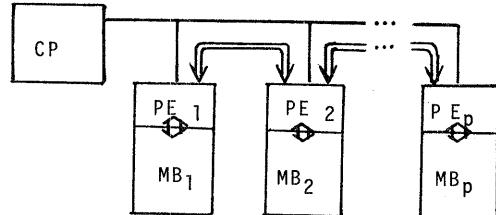


Fig. 2.4 Distributed Memory

台のPEあたり $1/p$ の時間しか使はず、データのやりとりが多いためベース処理には不向きである。データベース・マシンで並列処理の利点を生かすためには、比較的簡単なプロセッサでより個数は100以上は必要と思われ、多重処理方式では問題が多い。

図2.4に分散処理の1形態を示した。ここでは各PEに一定の大きさのMBを管理させている。ただし各PEは隣接したPEとデータのやりとりはできるものとする。また全体の制御用に別のプロセッサ(CP)があり、各PEに同時に命令を与えることができる。この方式は、多重処理に比べて制約が大きいが、後章で示すように、十分並列処理の特長が生かせる。

2.3. データベース・マシン向きのデータ・モデル

データベース・マシンは、論理的に定義された集合演算を、ハードウェアによって高速に実行せることが目的であるから、データの物理構造と論理構造が近いものであることが望ましい。またハードウェアで並列的にアクセスするという点で言えば、レコード内構造(intra record)のように可変長で複雑な構造や、陽に定義されたポインタによる逐次的なアクセス手段もできるだけ避けたい。

このような点と、多くのデータ・モデルの基本的性質を含んでいるという点から、E. F. Codd のリレーションナル・データモデル¹⁰⁾が有用である。しかし Codd モデルにはいくつか問題点もあり¹¹⁾、またデータの依存性(dependency)をはっきりさせることと、レコードの長さを小さくするために、ドメイン単位に「たて割り」の形でデータ蓄積するのがよい¹²⁾。ここでは、図 2.5. に簡単に例を示すにとどめる。

2.4. ハードウェア・ソートの有用性

データベースの集合演算をハードウェアで実現する場合に、もっとも必要な機能として連想記憶(associative memory)がある。たとえば図 2.5 で、R1 と R2 の egz-join¹³⁾

$$R1[B = B] \text{ } R2$$

をとる場合を考えてみると、R1 上の B-domain の各具体値に対して、R2 の B-domain で連想サートを行なう必要がある。(連想サートを行なう、毎回全データをスキャンすることは、この join 演算に必要なデータ参照回数は、 $(R1 \text{ の tuple 数}) \times (R2 \text{ の tuple 数})$ となり、実際的でない。)

ところで、連想サートを行なうアーキテクチャとして次の 3 種が考えられる。

- ① メモリ自体に連想機能をもたせる。
 - ② ハッシュ法などのリフト的連想記憶アルゴリズムをハード化する。
 - ③ 着目してくるフィールドについて、ハードウェアで高速ソートできるようにする。
- ① は機能メモリ¹⁴⁾の提案などが以前からあるが、LSI 化したときの集積度やピン数などに問題があり、実際的でないと思われる。② は同

R1 (A B C)			R2 (B D)	
a ₁	b ₂	c ₁	b ₁	d ₂
a ₂	b ₁	c ₃	b ₂	d ₁
a ₃	b ₁	c ₂	b ₃	d ₂
a ₄	b ₃	c ₁		
a ₅	b ₂	c ₂		

Fig. 2.5 Sample Data

じ具体値が多数現れる場合問題があり、また大小関係の join (G.T.-join, L.E.-joinなど) がとれないとなどの欠点があるが、使われ方によれば高効率を期待できる^{15), 16)}。③では、ソートが大容量のメモリに対しても、比較的簡単なハードウェアで高速に実行できるように工夫されていれば、効率のよい処理ができる。すなはち、ソートされた N 個の要素中をバイナリ・サーチすればよいから、 $\log_2 N$ 回のサーチで必ず終了する。また大小関係の join や具体値の重複にも問題がない。ハードウェア・ソートは、この他、データベース上の他の集合演算や、DBMS の報告書作成機能でも有用である。本稿の次章からは、ハードウェア・ソートとその応用について述べられていく。

3. ハードウェア・ソートの方式

本章では、チャネル型のデータベース・マシンで、流れこんでくるデータをソートしながら蓄積するためのアーキテクチャについて述べる。本稿のデータベース・マシンでは、少くともソートすべき領域(domain) の全具体値が、マシン上のシステム・バッファに収納できるものとしている。また、いったん蓄積したデータ集合を昇順または降順に高速に掃き出した後、データ集合の一部を任意に消去したりでき子機能も必要である。(消去しても、全体はソートされたままでなければならぬ。)

3.1. 特殊メモリによる方法

流れこんでくるデータをソートするためには、必ずデータの相対的な位置関係の変更を伴う。

この変更が、いかに並列的にうまく処理されるかが、ソートの効率にかかってくる。たとえば、この位置関係がデータの物理的な位置に直接対応しているとするとき、一群のデータを並行して移動する作業が必要である。そのような機能を実現するには、Random Access Memory と Shift Register の両方を合せたもの（ここでは SRAM とよぶ）を作ればよい。

図3.1 に SRAM のブロック図を示した。このメモリは外から見ると、アドレス線、データ線は普通の RAM と同じであるが、オーダーは READ, WRITE の他に PUSH と POP の計4種である。PUSH は指定された番地から高位番地の方へ、データを上ずつシフトする命令、POP は反対に指定された番地まで上ずつ詰める命令である。メモリ内の各セル M_i は、与えられたアドレスと自分の位置との大小関係が判別でき、また隣のセルとの連絡線 T_i, T_{i-1} をもつ。各命令と、大小判別の結果、各セルの状態の変化を図3.1(b) にまとめてある。(XEP の所は変化なし。) このメモリは通常の RAM と外部仕様を似せて、ビット・スライスに作ることができるのか、論理回路でセルを作るとすると、1 セル当り筆者らの試算では 70 ~ 80 ゲートになる。

次に重複のない m 個のデータ d_0, d_1, \dots, d_{m-1} のソート方法を簡単に示す。

アルゴリズム

1. $M_0 \leftarrow 0$ (最小数), $M_1 \leftarrow X$ (最大数)。

(d_0, \dots, d_{m-1} は 0 と X を含まない)

2. $d_0 \sim d_{m-1}$ を順次とり出し、各 d_j について次の操作を行う。

2-1. M_0 から M_{j+1} までをバイナリ・サーチして、
 $M_{k-1} < d_j < M_k$

となるような k をみつけよ。 (READ)

2-2. 2-1 で求めた k を利用して、
 $\text{push}(k)$ を行う。

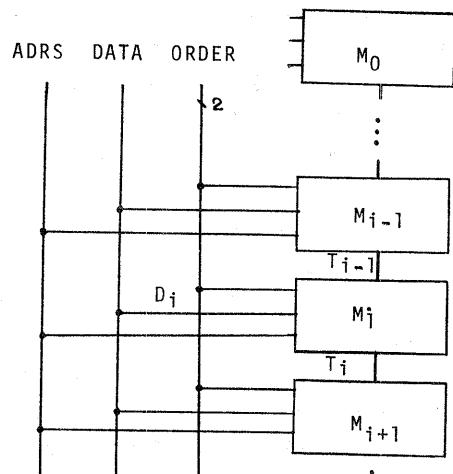
2-3. $M_k \leftarrow d_j$ (WRITE) 以上。

この方法でソートに要する命令の数は、

$$\sum_{j=1}^m (2 + \log_2 j) \text{ のオーダーである。}$$

また m 個のデータがすでに蓄積されていて、その中の 1 データを消去するには

$2 + \log_2 m$
の命令ですむ。



(a) Block Diagram

Compare	ADRS	ADRS	ADRS
Order	< i	= i	> i
READ	X	$D_i \leftarrow M_i$	X
WRITE	X	$M_i \leftarrow D_i$	X
PUSH	X	$M_i \leftarrow T_{i-1}$ $T_i \leftarrow M_i$	X
POP	X	$M_i \leftarrow T_i$ $T_{i-1} \leftarrow M_i$	X

(b) State Diagram

Fig. 3.1 Shifting RAM (SRAM)

このような特殊メモリは、前に述べたように高集積度のものを作るのが困難で、高価なものにつくのが難点であるが、システムの一部に使用するには有用であろう。

3.2. 並列プロセッサによる方法

3.1節の SRAM ではデータ間の大小関係を、そのまま物理的位置に対応させたため、データの移動を各セル上で並列して行う必要が生じた。ここでは、データの移動を減らすため、ポインタを利用することを考える。すなわちデータをいくつかにブロック化し、ブロック間の大小関係はポインタで示すことにする。このようする方法はソフトウェア上ではよく行われており、

B-treeによる方法が効率が高いことが知られています¹⁵⁾。そこでB-treeを原理的には用いることにして、これを2.2節で述べた分散記憶の形態を利用してハードウェア的に高速化することを考えた。図3.2にこのブロック図を示してある。ここで各メモリ・ブロック(MB)に付属しているSRAMは、全体として(横1列で)もSRAMになっている。MBは普通のRAMであり、プロセッサ(PE)によって読み書きされ、またSRAMに転送される。このメモリ上で各MBに対して同じアドレスをえたときにえらべるデータの組(図で破線で示してある)をB-treeのノードに対応させ、この中ではデータはリートされてしまうようとする。また各データは、ノード間の順序を示すポインタ領域(P)をもっている(図3.3参照)。

次に図3.2のメモリ上で、3.1節と同様、流れこんでくるデータをリートしたがる蓄積するアルゴリズムの概略を示そう。ここで、MはMBの個数で奇数である。

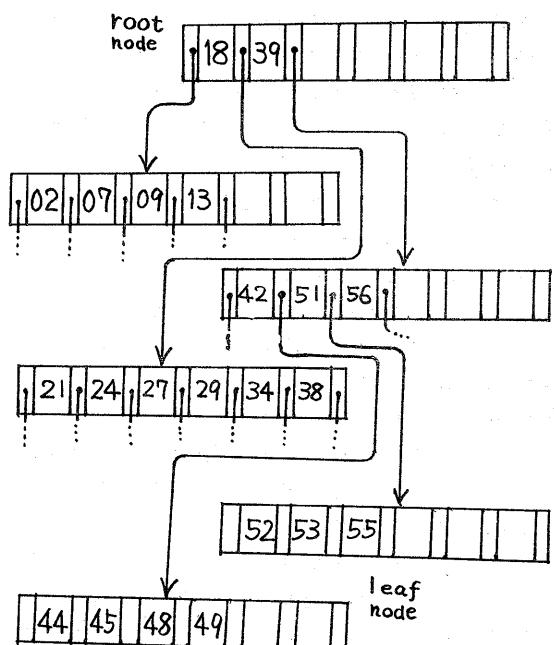


Fig 3.3. B-tree Example

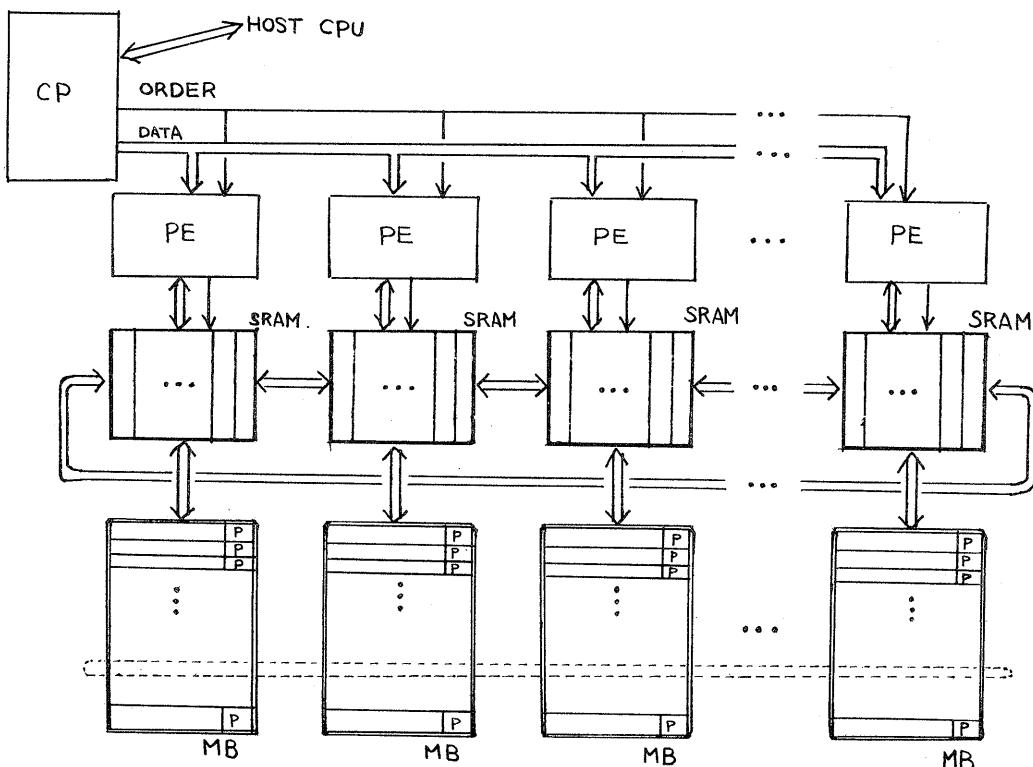


Fig 3.2. Distributed SRAM

アルゴリズム (データ d の挿入)

1. 与えられたデータ d に対して, d にもうとも近い前後の値を含むノード (d を挿入すべきノード) を次の手順でさがす。

1-1. ルート (root) ノードを SRAM にロードする。(CP から各 PE に, ルートのアドレスを与え, SRAM 上へデータを移動。)

1-2. d と, ロードされたデータの大小判定を各 PE 上で行い, 次にロードすべきノードのアドレスを決定する。

もしも SRAM 上のノードが最末端 (leaf) であれば 1. のおり。2 へ。

1-3. 1-2 で決った新しいノード・アドレスのノードを SRAM 上にロードする。1-2 へ。

2. 1 で得られたノードに d を挿入する。

2-1. SRAM 上のデータ列中に d を挿入する。(3.1節のアルゴリズム)

挿入後のデータの総数が m より小さければ, そのまま元の場所にノードを戻して終り。

挿入後のデータ総数が m になったときは 2-2 へ。

2-2. ノードの後半分 ($m/2$ 個) のデータを新しいノード・アドレスを作って, そこへ移す。(同じ MB 内の, もととは別の番地への格納による) 2-3 へ。

2-3. 元のノードの中央値 ($m/2 + 1$ 個目) のデータを, 1 段階上のノードへ挿入するため, この値を新たに d とし, 元ノードをポイントしてリ 1 段階上のノードを新ノードとして SRAM にロードする。2-1 へ。

以上

上のアルゴリズムは B-tree の場合とほぼ同じであるが, ノードの開始点になる MB が 2ヶ所でさるため, SRAM の全体はサイクリックにしておく必要がある。また 2-2 のステップで作られる新しいノードの生成や, ノードとアドレスの管理は CP に行わせなければならぬ。

図 3.2 のアーキテクチャは, データの消去に関するても, やはり B-tree の手法で行われる。

上のアーキテクチャでは, 長さ m のノードの内部での大小判定や移動は, M 個の PE によって常に M 重に並列処理されてリ。いま N 個のデータが蓄積されてリるときに, あるデータの検索, あるいは新ノードの生成を伴わない追加を行なうときの命令の数は

$$(2 + \log_2 m) \cdot \log_m N$$

のオーダーである。また確保された記憶領域中の有効な領域量は平均 75 %, データの追加時に新ノードの生成処理が必要になると確率は

$$2/m - 1$$

なあ上のアーキテクチャでは PE にクレジット処理を負わせることにより, ノード中の各データ長を可変でさるし, leaf ノード中のボイラー・エリアの削除など, ソフト的な B-tree より高い機能をもたせ, また記憶領域を有効に使う工夫も可能である。

最後に本稿では MB として RAM を想定しておいたが, 図 3.2 のアーキテクチャは, 同時アクセスという性質上, 各 MB を CCD, 磁気バブルなどのシフト形式のメモリにあきかえても, なお有効であることを付記しておく。

4. ハードウェア・ソートによるデータベース演算子の実現

本章では前章までに述べたハードウェア・ソートの機能を利用して, データベースの基本的な演算がどのように実現されるかについて述べる。本章で述べるデータベースは, 2 章のドメイン単位蓄積方式を想定しており, またデータベース・マシンのアーキテクチャとして, 3 章のハードウェア・ソートにもとづくチャネル型の並列処理マシンを考えている。(Codd モデルは, ドメイン単位やその他, 任意の形でデータを蓄積できるし, そのようにしても記憶領域のむだは少ない。)

ここでは Codd の提案した集合演算⁽¹³⁾を考えよう。

• Projection

関係 R は domain D_1, \dots, D_m から構成されてリとする。projection

$$R [D_{r_1}, \dots, D_{r_k}]$$

をとるには, tuple $(d_{r_1}, \dots, d_{r_k})$ を順次のように蓄積すればよい。

1. $(d_{r_1}, \dots, d_{r_k})$ がすでに蓄積されてリるか検索する。もしすでにあれば, 何もしないで次の tuple をよむ。

2. まだ蓄積されていないければ, ソートしな

が"う3. 2節のアルゴリズムで蓄積する。

Projection の結果の集合はデータベース・マシン上に残ったものである。

•Join

関係 R, S はそれを domain $D_1, \dots, D_m, E_1, \dots, E_n$ から構成されていとする。

R, S の θ -join

$$R [D_r \ \theta \ E_s] S$$

をとるには次のようにすればよい。

1. R, S のうち tuple 数の少ない方 (仮に R とする) について, join の対象になる domain (D_r) の具体値と, 重複を許して Y-join してから蓄積する。

このとき, 具体値と共に TID (tuple identifier, すなはちここでは R 中の tuple へのポインタ) をデータとリッショナリ化する。

2. もう一方の関係 (S) の tuple を順次データベース・マシンに送りこみ, 各 tuple の join のフィールド (E_s) と, 1. で蓄積した (D_r の) 具体値との間で, θ で結びつくものを検索し, (S の) tuple と (D_r の) TID の組を返す。

上の方法で, TID で指定された R の tuple と S の tuple を結合した集合が本当に必要ならば, マスター・スペース上で, この情報をもとに併合する (merge) しかない。しかし, 通常は 1 つのファイル (S) からもう一方 (R) が直接参照できればよいので, 上記演算で十分である。

•Division

Joinと同じ前提で division

$$R [D_r \div E_s] S$$

をとるには次のようにすればよい。

1. R の tuple の, D_r 以外の domain の具体値の組と, TID を, 具体値の組でソートしながら蓄積する。

2. 1 の集合と具体値の組が共通なものごとに 1 つの部分集合を作り, この部分集合ごとにマスター側 (GPC) へ帰き出す。

3. S の tuple の, E_s domain の具体値をソートしてから蓄積する。

4. 2 でえられた部分集合ごとにデータベース

・マシンに D_r の具体値を送りこみ, S の具体値を全部尽していいかどうか調べる。尽して「れば」, その部分集合の共通な (D_r 以外の) 具体値の組が, 求める division の要素である。

この他, restriction は, チャネル型のデータベース・マシンである以上, 問題がないし, Union, Intersection も, ハードウェア・ソートができるれば, そのまま処理できる。

5. まとめ

本稿では, ハードウェア・ソートを基礎におくデータベース演算子の実現およびデータベース・マシンの形成について, 筆者らの DB マシンの構想紹介も兼ねて, できるだけ一般性のある議論をしてきた。データベース・マシンの計画は, 現在電子技術総合研究所や北海道大学でもすすめられており,¹⁷⁾ 今後活発な進展が期待できる。

本報告を終えるにあたり, 日頃有益を示唆をいただりしてある電子技術総合研究所の面野博二部長, 植村俊亮主任研究官はじめ, 大型工業技術研究開発連絡会議データベース・マシン WG の諸氏に感謝の意を表します。

参考文献

- 1) 植村: "CODASYL 方式のデータベース・システム", 情報処理 Vol. 17, No. 10, 1976.
- 2) 魏鶴, 他: "データベースの関係形式", 同上。
- 3) 関野, 他: "データベース・マシン", 同上。
- 4) E. A. Ozkarahan, et. al.: "RAP-An Associative Processor for Data Base Management.", AFIPS (NCC), Vol. 44, 1975
- 5) M. Edelberg, et. al.: "Intelligent Memory", AFIPS (NCC), Vol. 45, 1976.
- 6) 有澤: "データベース・処理の専用マシン化について" 情報処理学会 17 回全国大会, 1976.

- 7) 有澤・土肥："データベース・マシン向きのデータ処理方式." (投稿中)
- 8) S. Schuster, et.al: "A Virtual Memory System for a Relational Associative Processor." AFIPS(NCC) Vol. 45, 1976.
- 9) 飯塚, 他 "C. mmp マルチ・ミニ・プロセッサについて." 情報処理, Vol. 15, No. 7, 1974.
- 10) E. F. Codd: "A Relational Model of Data for Large Shared Data Banks", Comm. of the ACM, Vol. 13, No. 6, 1970.
- 11) 有澤: "リレーショナル・データベース上の領域間の決定関係についての考察" 情報処理, Vol. 17, No. 11, 1976.
- 12) 有澤・土肥: "リレーショナル・データベース・マシンについての考察" 昭52年電子通信学会総合全国大会。
- 13) E. F. Codd: "Relational Completeness of Data Base Sublanguages" Courant Computer Science Symposia 6, 1971.
- 14) P. L. Gardner: "Functional memory and its microprogramming implications" IEEE Trans. C-20, No. 7, 1971.
- 15) D. Knuth: "The Art of Computer Programming" Vol. 3, pp 473- Addison-Wesley
- 16) P. A. Alberg: "Space and Time Saving through Large Data Base Compression and Dynamic Restructuring." Proc. IEEE, Vol. 63, No. 8, 1975.
- 17) 本研究報告と同時に発表される。