

データ分配ネットワークを用いた
データベース計算機の構成法

Data Base Machine Architecture using
Data Partitioning Network.

小田 泰充

Yasumitsu Oda

日本電信電話公社 武藏野電気通信研究所

Musashino Electrical Communication Laboratory, N.T.T.

1. まえがき

近年の情報化社会の発展にともない、データベース（以後単にDBと記す）に対する高機能化、高速化の要求が高まっている。しかし、従来のノイマン型計算機では、① 高機能化を果さうとすると、ソフトウェアの負担が増大する。② 高速化を果さうとすると、補助記憶と主記憶および主記憶とCPUとの間のデータ（および命令）の転送が問題となる。という欠点があつた。

これに対し、VLSI技術の進歩により急速に価格が低下しつつあるハードウェアを利用して、DB専用のデータベースをもつた計算機（Data Base Machine 以後単にDBMと記す）を作成しようとする試みが活発に行なわれている。その主なものに

(1) ベル論理装置型のもの^{(1),(2)}

RAP, CASSM, RARES, DBBC
DIRECT, EDC⁽⁵⁾

(2) ハッシュビットマップを用いたもの
CAFS⁽³⁾, LEECH

(3) 動的ビットマップを用いたもの
SPIRIT⁽⁴⁾

(4) 分散データベース指向のもの
DIALOG⁽⁶⁾

(5) 国際演算専用プロセッサを用いたもの
BSDBM (ソートエンジン, サーチエンジン)⁽⁹⁾

Systolic Array による国際代数演算器⁽⁷⁾

データ圧縮リードによる国際代数演算器⁽⁸⁾
などがある。

これらのDBMは、いずれもCoddが提唱した関係モデル(RDB)を効率よく実現する目的で(少なくとも目的の1つとして)開発されてゐる。しかし、(1)のタイプのものは、検索、更新、挿入、削除(以後これを基本演算と呼ぶ)は効率よく実行するが、RDBの特長である、JOIN, PROJECTIONなどの集合と集合の差的演算(これを基本演算と区別して、集合演算と呼ぶことにする。)には向いていない。(2)のタイプのものでは、ハッシュ技法を集合演算の前処理としてのみ利用しており、演算そのものは、ホスト計算機で実行しなければならない。(3)のものでは、ビットマップを作成する手間が余分にかかることと、場合により作成したビットマップが巨大になるという欠点がある。

(4), (5)のものは、いずれも集合演算を高速に実行するための専用ハードウェアをもなしてゐる。特に、(5)のものは、 $O(n)$ (n : タップル数) で集合演算ができる。しかし、ここれらのハードウェアを m 台用意し、 m 並列入力を行つても、これらを單に、直列、並列、あるいは網状に結合するだけでは、 $O(n)/m$ とすることはできない。

1つの解決策は、 m 並列入力を直列入力に変換して集合演算器に入力することであるが、そのためには、読み出し速度の伝導の転送速度と、それに応じて高速動作する集合演算器が要求されるので、この方法には限界がある。

そこで、ここでは、基本演算ではセル論理装置を用い、集合演算では、データの入力速度に追隨して集合演算を行ふ集合演算器を用い、こ

れとセル論理装置を、データ分配ネットワーク (Data Partitioning Network: 以後OPNETと呼ぶ) で結合することにより、セル論理装置の並列読み出しに応じた並列度で、基本演算も集合演算も行えるDBMを提案する。本稿のOPNETは、セル論理装置から集合演算器へのデータ転送時間を利用して、一種のダイナミッククラスタリングを行ふものである。

2章でシステム概要を、3章でシステム構成要素を述べた後、4章で評価と考察を行う。

2. システム概要

2.1 構成方針

本稿が提案するDBMは、次の方針のもとに構成した。

(1) DB用DBMとする。

(2) DBの類機としては、戈がバイトオーダーとするが、個々のベースリレーションは数十メガバイト以下の中のものが大部分であるとする。

(3) DBの格納媒体は、価格の点から磁気ディスクを前提とする。

(4) 同一ベースリレーションは、同一シリオ群に集める。

(5) 言語インタフェイスは、関係代数レベルとする。(表2-1)

(6) シリング群の並列読み出し速度に応じた処理速度で、基本演算も集合演算も行えるDBMとする。

表2-1 提供する関係代数演算

演算タイプ		関係代数演算
基本演算	OPφ	RESTRICTION, INSERTION UPDATE, DELETION
集合演算	OP1	PROJECTION AGGREGATION OPERATION (max, min, average sum, cardinality, count)
	OP2	JOIN (implicit JOIN, explicit JOIN) UNION, DIFFERENCE INTERSECTION
	OP3	DIVISION

(7) DBの完全性チェック、Queryの最適化はソフトで行う。

(8) 同時アクセス制御は、関係表レベルのロックで実現する。(本稿では、これについては省略する。)

2.2 基本的考え方

ここでは、OPNETを採用した基本的考え方を述べる。

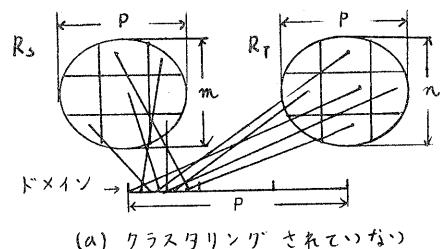
今後、表2-1のOP2タイプの演算をP台のプロセッサを利用して並列処理することを考える。この演算を、OP2(Rs, Rt)と記す。ここで、Rsをソースリレーション、Rtをターゲットリレーションと呼ぶ。Rs, Rtが各々MP個、NP個のセル(回転メモリの1周期分のメモリ)、この場合1セル=1トラック)に格納されており、1プロセッサが1セルのRsと1セルのRtのOP2を実行するのにTc時間かかるものとする。

このとき、Rs, Rtが図2-1のとく、セル単位にP個にカラストラリングされている場合とを並べて場合を考える。

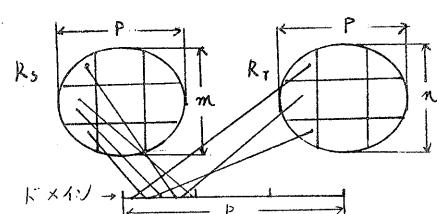
(a) クラストラリングなしの場合

Rs, Rtの全セル同志の比較が必要になりその処理時間Tcは、次式となる。

$$T_c = \frac{m \cdot n \cdot P}{P} T_c = m \cdot n \cdot P \cdot T_c$$



(a) クラストラリングされていない



(b) クラストラリングされていな

図2-1 クラストラリング

(a) クラスタリングありの場合

他のクラスターに属するセル同志の比較は必要でなくなるため処理時間 T_2 は次式となる。

$$T_2 = \frac{\sum_{i=1}^n m_i n}{p} \quad T_C = m \cdot n \cdot T_C$$

このように、クラスタリングされていれば、処理時間を $1/p$ に短縮することができます。図 2-1 (b) のごとくサブをクラスタリングして格納できればよいか、リレーションは 2 つ以上の属性に対して定義されており、どの属性に対しても関係演算が定義できることで、すべての属性に対して本を行っておくことは出来ない。

そこで、処理時間を $1/p$ に短縮するためには演算のたびに、ダイナミックに、退避なく、図 2-1 (b) の状態を作成することが要求される。

本稿で提案する ω PNET は、これを補助記憶から集合演算器へのデータ転送時間を利用して行なうものである。

クラスタリングの手法としては、ソート技術を使用する方法やヒストグラムを知りて分割した結果を平均化するように範囲を決めてやる方法があるが、前者はソーティングに $O(n)$ を要する必要があること、後者はヒストグラムに関する知識が必要となり、しかも演算は部分集合に対して行な

われることが多いため、全体のヒストグラムに対する知識では役立たなくなる恐れがあること、などの問題点がある。そこでここでは、ハッシュ関数を通じてダイナミッククラスタリングを行う方法を採用した。ハッシュ技術を適用する新たな発明する問題点として、

① データの分布を本当にランダムにすることができるか?

② クラスター間の大小比較が不可能になるか関係演算への影響はないか?

である。①に対しては、入力データの性質によるため、数種のハッシュ関数を用意し、ハッシュ関数を選択できる機能をもたらすことにより対応する。②については、関係代数が集合に対して定義されているため、大小比較が必要となるのは θ -JOIN のみであり、他の演算では問題ない。この θ -JOIN が実用上必要となるのは極めてまれであるので、BM としてのハード的サポートは行なわないことにする。

なお、ハッシュ関数を用い、ランダムなデータ分布にしてから、データの分割を行なう集合演算を行なうとの考えは、R.E. Shaw (11) にも見られるが、彼は関係表の読み込み時間は表の大きさと無関係であるとし、1つの集合演算器だけを用いた場合について論じているので、本稿の

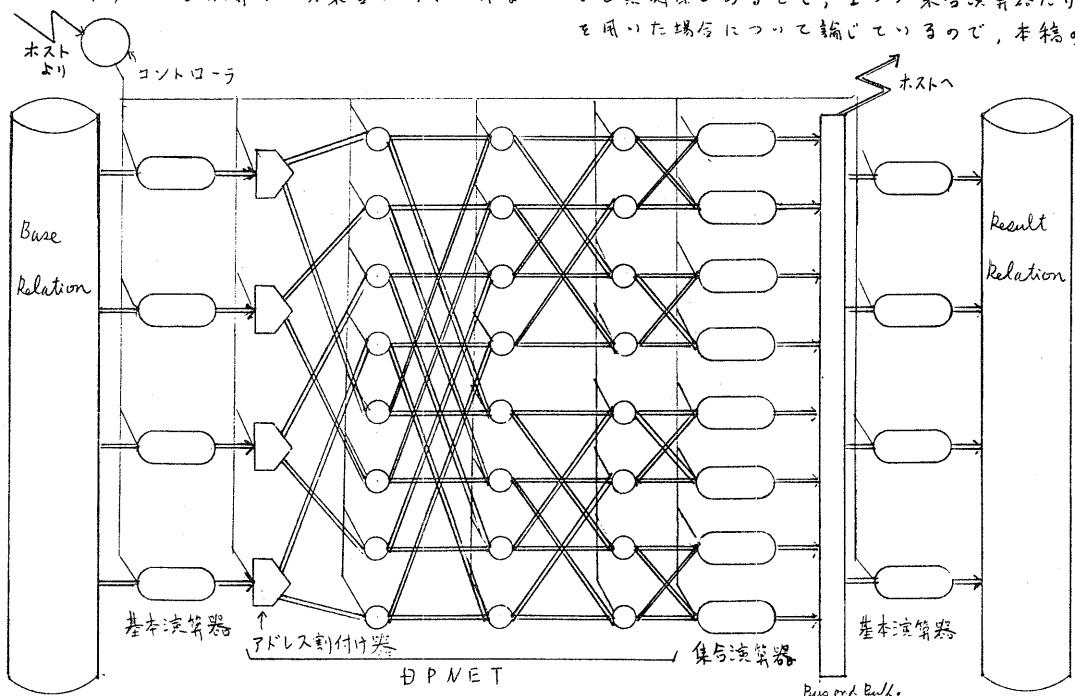


図 2-2 システム構成図

場合（表の大きさと読み込み時間は比例、演算器は複数台並列処理）とは、大きく異なる。

2.3 全体構成

図2-2に、OPNETを用いたBMのシステム構成図を示す。以下、全体の動作概要をOP2(R_S, R_T)演算を例に説明する。

① コントローラは、基本演算器にR_Sの送出を、OPNETにR_Sの転送準備を、集合演算器にOP2演算の準備を、それぞれ指示する。このとき、R_S送出と同時にRESTRICTIVE演算が可能なら、その旨基本演算器に指示を与える。

② 基本演算器は、補助記憶から関係表を読み込み、集合演算器が必要とするタップル（または、タップルの必要部分のみ）を選択し、OPNETに送出する。

③ OPNETでは、演算対象となる属性をキーとしてハッシュ函数を計算し、その結果からタップル選出先の集合演算器を決定して、当該集合演算器にタップルを送り込む。

④ 各集合演算器は、OPNETから送られて来るタップルを、これまでの集合演算器のオーバーフローするまで、蓄えている。

⑤ コントローラは、基本演算器がS-R_S終了通知を受けたか、集合演算器がS-R_Sオーバーフロー通知を受けるかしたS、基本演算器にR_S送出準備を、OPNETと集合演算器にR_T受け取り準備を各自指示する。

⑥ R_Tに対し、基本演算器、OPNETとも、②、③と同様の動作を行なう。

⑦ 集合演算器は、R_Tの1タップルを受けとるたびに、OP2演算を行なう結果を出力する。結果はバッファに蓄えられ、まとめてセル論理装置、またはホスト計算機に送出される。

⑧ コントローラは、集合演算器がS-R_Tの終了通知を受けたと、全R_Sの読み込みが終了していればOP2演算を終了し次の演算を、そうでなければ、R_Sの残りを読み込み②～⑧を行なう。これを、基本演算器、OPNET、集合演算器の各々に指示する。

以上が、システムの全体構成と処理の流れの概要である。OP1, OP3タイプの演算については、次章で述べる。

3 システム構成要素

3.1 基本演算器

図3-1に、基本演算器のブロック図を示す。図3-1に示すごとく、ここでは基本演算器に必要な比較器の他に、送出済タップル判定ユニット(TDU)と、タップル転送許可信号バッファ(TSB)をもっている。

TSBは、次段から送られてくるタップル転送許可信号(TS)を蓄えておくためのバッファであり、次段でのタップル受け付け可能数を表示している。

TDUは、TSBが0にかつたとき（処理速度が読み込み速度に追いつかなくせつたとき）セルの空読みを行なう、次の回転で前回の読みきりタップルを読みだしたり、R_Sの読み込み中に集合演算器が一杯になりR_Sの読み込みを中止した後、再びR_Sを前回の読みきりから読みだしたりする（2.3節の④参照）のに必要なユニットである。

タップル送出ユニット(TSU)は、比較器からの比較結果信号(1), TDUからの未送出タップル信号(2), TSBからの転送許可信号(3)のAND条件が満足されたら、TSU内のタップルを次段に送出する。ここで、比較器での条件比較、TDUの処理、TSUへのタップル格納は並列に行なわれ、しかも、タップル長に比し、演算対象となる属性長ははさりので、上記判定を行なっても、転送遅れが生じることはない。

3.2 OPNET

図2-2に示すごとく、OPNETは、アドレス割り付け器と、各ノードにタップル待ち行列用バッファ(TSB)をもったルーライナゲットワークからなる。

図3-2にアドレス割り付け器を、図3-3

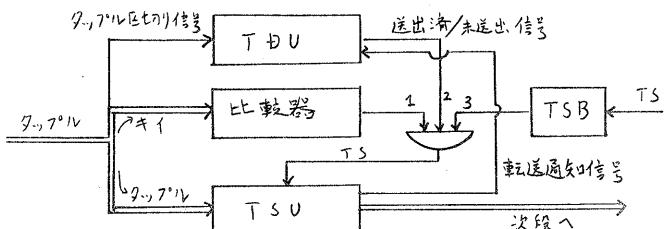


図3-1 基本演算器

に DPNET の各ノードを示す。ここで、TSB, TSU の役割は、基本演算器の場合と同じである。

図中、SW選択は、アドレス割り付け

付け（ハッシュ関数）の結果を受けとり、タップルをどろぐノードに出力するがのスイッチ（SW）の切り替え信号の送出のみを行っている。TQB選択は、タップル格納数の多いちのTQBかSからタップルを取り出すためのものである。

また、TQB, TSU では、RS用とRT用の2種のバッファをもつてている。これは、RSの読み込み中に集合演算器が一杯にならなければ、途中状態保存用である。

図2-2で、ルーティングネットワークの入口から、基本演算器の2倍（転送速度を2倍にすれば、物理的に2倍にする必要はない。）になり、でいることに注意されたい。これにより、DPNET の各ノードの、TQBの大きさを数ヶアル分とするだけで、途中でTQBがオーバーフローする（このとき前段のTSBが0になる。）確率を 10^{-3} 以下とすることができる。（4章参照）

3.3 集合演算器

ここで、集合演算器としては、次々に到着するタップルに遅れることなく、表2-1の集合演算を行なうことを望ましい。その意味で、ソーランゲを行なうもの⁽⁴⁾⁽⁹⁾は、シート用にデータを蓄える必要があるので、望ましくない。また、Systolic Array⁽⁷⁾などのようにRS, RTの同時入力を行なうものは、本稿のDBM

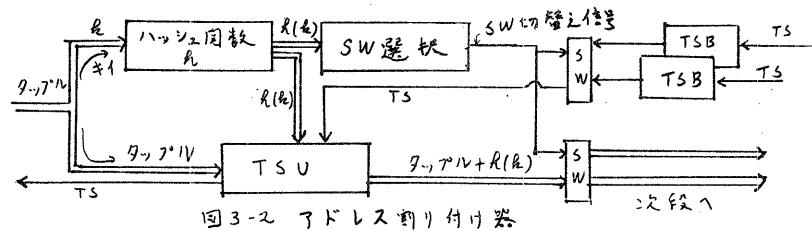


図3-2 アドレス割り付け器

では使えない。

そこで、図3-4に示す集合演算器を提案した。本集合演算器は、

- ①キーサーカロセッサ (KSP)
- ②演算カロセッサ (OP)

から構成されている。

KSPは、連想メモリとプロセッサで構成され、連想メモリには、連想メモリを節約する目的で、OPエタイプの演算を実現する目的で、未登録キーやそのタップルの格納アドレス (RSAD)とともに登録している。

OPは、タップル格納用メモリ (RSM) とバッファ (RTB) をもっており、この2つに格納されたタップルとKSPからの検索結果を使用して、集合演算を行なっている。

RSMは、KSPに格納先アドレス (RSAD) を知らせるためのアドレスカウンタ (AC) を持つている。同時にACは、RSMを節約するため、OPで不要なタップルは、KSPからのタップル不必要信号により、格納しないようにする（次に入力されたタップルを重複書きする。）目的にも使用されている。

RTBも、RSMと同様の目的で、KPSからのタップル必要/不必要信号を受け取っている。

(1) KSP

KSPには、表3-1に示す5つの処理モード

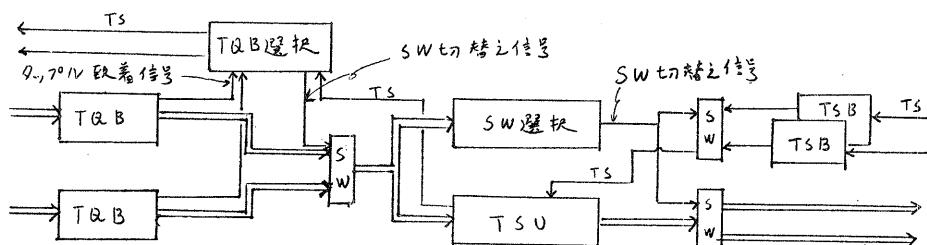


図3-3 DPNET のノード

があり、キー（演算対象となる属性値）と R_{SAD} を受けとり、キーの KSP 内の連想メモリへの登録及び同一キー検索を行っている。

KSP 内では、キーは図3-5の形式で格納されている。ここで TP , TL は、 RSM に格納された当該キーをもつタップルチェイン（チェインは OP の後である。）の、 TOP と $TAIL$ アドレスである。

KSP 内の処理で、検索の他に必要とされる処理は、最大 2 step の格納処理と、キーの読み込み、検索結果の出し、キーの登録であるが、後の 3つは並列に行うことができる。

キー検索を、LSI 連想メモリ用いれば 1 回のメモリアクセスで、また並列ハッシュ演算器 ⁽¹²⁾ を用いても、ほぼ 1 ～ 2 回のメモリアクセスで、検索を終えることができる。このことと、通常キーはタップルに比し短かることを考慮すれば、 KSP の処理時間はタップルの R_{TB} , RSM への格納時間と同程度とし、両者を並列に処理することは、充分可能である。

(2) OP

OP は、 KSP が 5 受け取った検索結果を RS

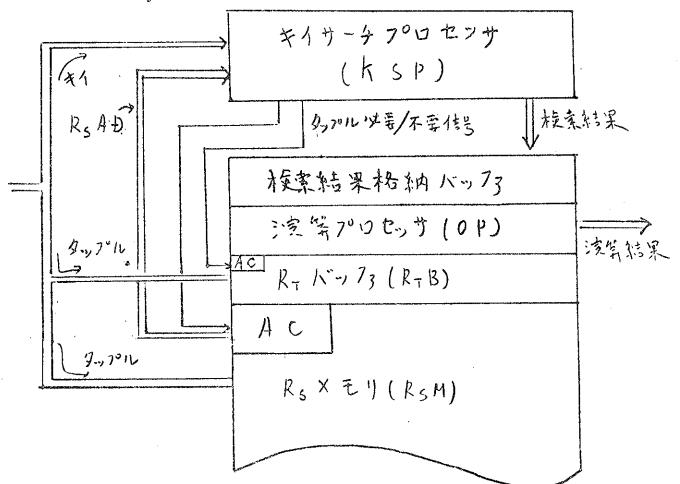


図3-4 集合演算器

M , R_{TB} に格納されているタップルから、次の並列からの処理を行う。

① R_SAD タップルのチェイン (CHAIN)

KSP から受け取った R_SAD と TL を用い、 R_SAD と TL で示されるタップルのチェインをオブジェクト部に格納する。

② R_T と R_S の結合と出力 (OUT($R_T R_S$))

R_{TB} のタップルと RSM のタップルを結合し、演算 (explicit JOIN) の結果として出力する。

表3-1 KSP の処理内容

処理モード	入力	処理			タップル必要信号	検索結果の出力形式
		検索	キー操作	キー登録		
mode 1 (無条件登録)	R_{SAD}	Y	$R_{SAD} \rightarrow TL$	△	不要 必要	$\begin{array}{ c c c c }\hline \text{キー} & TP & TL & R_{SAD} \\ \hline \end{array}$
		N	$R_{SAD} \rightarrow TP, TL$	O	不要 必要	$\begin{array}{ c c c c c }\hline \text{キー} & R_{SAD} & R_{SAD} & NULL \\ \hline \end{array}$
mode 2 (新キー登録)	R_{SAD}	Y		X	必要 不要	$\begin{array}{ c c c c }\hline \text{キー} & TP & TL & Y \\ \hline \end{array}$
		N	$R_{SAD} \rightarrow TP, TL$	O	必要 必要	$\begin{array}{ c c c c c }\hline \text{キー} & R_{SAD} & R_{SAD} & NULL \\ \hline \end{array}$
mode 3 (同一キー検索)	キー	Y		X	必要 不要	$\begin{array}{ c c c c }\hline \text{キー} & TP & TL & - \\ \hline \end{array}$
		N		X	不要 不要	$\begin{array}{ c c c c }\hline \text{キー} & TP & TL & - \\ \hline \end{array}$
mode 4 (新キー検索)	キー	Y		X	不要 不要	$\begin{array}{ c c c c }\hline \text{キー} & TP & TL & - \\ \hline \end{array}$
		N		X	必要 不要	$\begin{array}{ c c c c }\hline \text{キー} & TP & TL & - \\ \hline \end{array}$
mode 5 (キー有無検索)	キー	Y		X	必要 不要	$\begin{array}{ c c c c }\hline \text{キー} & TP & TL & Y \\ \hline \end{array}$
		N		X	必要 不要	$\begin{array}{ c c c c }\hline \text{キー} & TP & TL & N \\ \hline \end{array}$

/ ; 処理せず Y ; 同一キー有り N ; 同一キー無し O ; 登録する X ; らせず TL を変えた

R_{SM} のタップルがチェックされていれば、チャインされているタップルのすべてに対して、 R タップルとの結合を行ふ。

出力時に、 JOIN 結果に対する RESTRICT ON があれば、それを同時に行ふ。

(3) R_T 出力 ($OUT(R_T)$)

R_{TB} が取り出したタップルをそのまま出力する。

出力したタップルは、 OP_2 タイプの演算では、それをそのまま演算結果となるが、 OP_1 タイプの演算では、中間結果となり、全国係表入力終了後、中間結果を R_S として再入力する必要がある。

(4) R_S 出力 ($OUT(R_S)$)

R_{SM} に格納されている全タップルを演算結果として出力する。

(5) AGGREGATION 演算 (AG 演算)

AG 演算では、表3-2のごとく、 R_S タップルに演算結果格納部を付加し (ACを格納部だけ余分にカウントする) そこに演算結果を格納する。演算は、表3-2のごとく行ふが、いづれも数stepsで終了する。

(6) R_T タップル消去 (CLEAR(R_T))

R_{TB} より、タップルを取り出すだけで何もしない。

以上が、 OP の処理内容であるが、これらは KSP の処理とパイプライン的に結合され、並列処理される。

キー	TP	TL
----	----	----

図3-5 KSPでの
キー格納形式

R_S タップル	AG演算用 附加フィールド
------------	------------------

図3-6 AG演算時の R_S タップル
格納形式

*：新キーフィールドの判定は、表3-1の mode 2, mode 4 の検索結果出力形式より判定する。

(3) 集合演算の実現法と状態遷移

集合演算の実現法を状態遷移図を使用し、図3-7 (OP_1 演算), 図3-8 (OP_2 演算) に示す。また、具体的な演算に対し、 KSP, OP が、各状態でどのような処理を行うかを、表3-3, 表3-4に示す。

なお、 OP_3 の INTERVIEW , PROJECTION , AG 演算の Cardinality Count を組み合せることにより実現できること⁽¹⁴⁾、ここには示してない。

4. 評価と考察

本稿で提案したDBMは、まだ設計段階である。従つてここでは、ハッシュ函数により入力データの分布をランダムにできると仮定した場合の設置するハードウェアと期待される効果についての確率的評価を行つた後、集合演算器の入力追随性と R_{TB} の大きさについて考察する。

4.1 DPNETの各ノードの待ち行列数

DPNETの各ノードの待ち行列数分布を次の前提のもとに求めめる。

① 各ノードへは、下時間ごとにタップルが到着する。

② 各ノードは、下時間で1つタップルを次の段に送ることができる。

すると、第一段目のノードでは、同時に2つのタップルが到着することはないので待ちが発生することはない。

表3-2 AG演算処理

項目番号	AG演算	追加フィールド名	AG演算処理	
			新キーフィールド名	二度目までのキーフィールド名
1	max	MAX	$A_x^{**} \rightarrow MAX$	$if A_x > MAX then A_x \rightarrow MAX$ $else ;$
2	min	MIN	$A_x^{**} \rightarrow MIN$	$if A_x < MIN then A_x \rightarrow MIN$ $else ;$
3	sum	SUM	$A_x^{**} \rightarrow SUM$	$SUM + A_x \rightarrow SUM$
4	Cardinality count	COUNT	$1 \rightarrow COUNT$	$COUNT + 1 \rightarrow COUNT$
5	average	COUNT, SUM AVERAGE	項目番号3, 4を行つた後 $SUM / COUNT \rightarrow AVERAGE$	

** A_x : R_{TB} に格納されたタップルの演算対象属性フィールド名

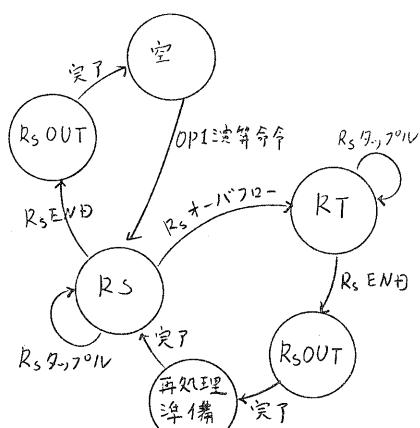


図3-7 OP1 演算の状態遷移図

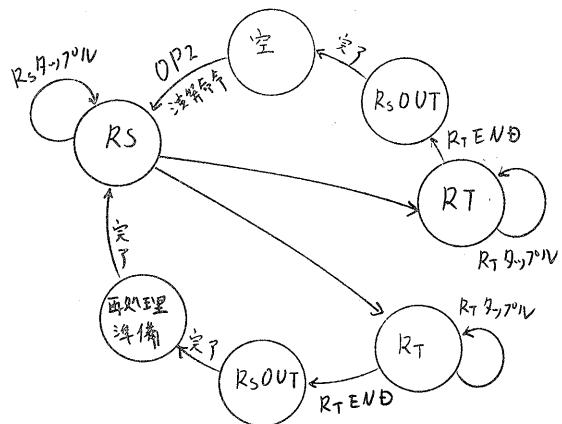


図3-8 OP2 演算の状態遷移図

表3-3 OP1 演算の状態と処理

状態 演算	RS		RT		RS OUT		再処理 準備
	KSP	OP	KSP	OP	KSP	OP	
PROJECTION	mode 2	DEL(R _T)	mode 4	OUT(R _T)		OUT(R _S)	OUT(R _S) として出力された。
AGGREGATION	mode 2	AG 演算	mode 5	Y: 脳演算 N: OUT(R _T)		OUT(R _S)	Y, T, R _T を、 R _S として再入力する。

表3-4 OP2 演算の状態と処理

状態 演算	RS		RT		RS OUT		再処理 準備
	KSP	OP	KSP	OP	KSP	OP	
explicit JOIN	mode 1	CHAIN	mode 3	OUT(R _T , R _S)			R _S の残りを再入力する。
implicit JOIN INTERSECTION	mode 2	DEL(R _T)	mode 3	OUT(R _T)			
UNION	mode 2	DEL(R _T)	mode 4	OUT(R _T)		OUT(R _S)	
DIFFERENCE	mode 2	DEL(R _T)	mode 4	OUT(R _T)			

第2段目のノードでは、

- ① 2つ同時に到着する確率 $\cdots (1/4)^2$
- ② 1つだけ到着する確率 $\cdots 2 \cdot (1/4) \cdot (3/4)$
- ③ 1つも到着しない確率 $\cdots (3/4)^2$

であるので、時刻 t, T における待ち行列数が n である確率を $P_t(n)$ とおくと

$$\left. \begin{aligned} P_t(0) &= (15/16) \cdot P_{t-1}(0) + (9/16) P_{t-1}(1) \\ P_t(n) &= (6/16) \cdot P_{t-1}(n) + (9/16) P_{t-1}(n+1) \\ &\quad + (1/16) P_{t-1}(n-1) \end{aligned} \right\} n \geq 1 \quad 4-1$$

これより、 $t \rightarrow \infty$ の平衡状態における待ち行列の分布函数 $P_{\infty}(n)$ 、待ち行列数の期待値 E は、

$$P_n(x) = (8/9)(1/9)^n, \quad E = 9/64$$

となる。

同様に、3段目以下の各ノードの到着確率も2段目と同じになることが示せる。

4.2 TQBオーバーフロー確率

図3-3のTQBの空きを、各1/2タップル分(TQBは2つあるので、1ノード当たり2タップル分)とし、連続m-T時間オーバーフローが発生しない確率を求めよ。またT時間後にオーバーフローが発生する確率は、4-1式で、 $P_x(q+1)$ となる。これより、 $0 \leq n \leq q-1$ では、4-1式をそのまま使用し、 $\circ x = q$ では。

$$\left. \begin{aligned} P_x(q) &= (6/16)P_x(q) + (1/16)P_x(q-1) \\ \circ x = q+1 &\Rightarrow \end{aligned} \right\} 4-2$$

そして、 $P_x(q+1)$ を求めるとき $P_x(q+1)$ は、またT時間内に1回以上オーバーフローが発生する確率となる。これを計算した結果を図4-1に示す。

図より、RSMを 10^3 タップル分、RTの大きさを 10^6 タップル分としても、 x を各々6, 8個と用意すれば、TQBのオーバーフロー確率も、 10^{-3} (1000回の集合演算で1回の発生)程度とすることができる。

4.3 集合演算器の利用率

RS状態では、それが1つの集合演算器が一杯になると、RSタップルの入力を中止し、RT状態へ切り替えて行っている。このとき、他の集合演算器にどの位のタップルが格納されているのか(これを、ここでは利用率と呼ぶこと

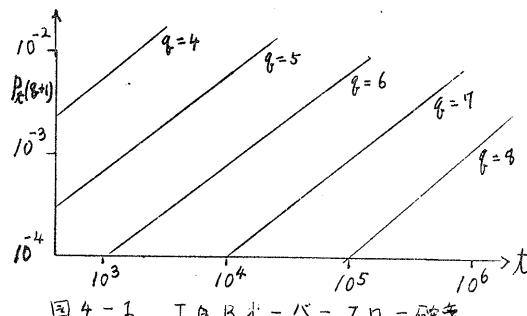


図4-1 TQBオーバーフロー確率

にする。)を、次の仮定のもとに計算する。

- ① タップルは1つずつ到着する。
- ② 到着したタップルは、m個の演算器のどれか1つが、ランダムに選択されてそこに格納される。
- ③ m個の演算機のどれか1つが、m個のタップルを格納した後、入力を中止する。

各個の入力があつたとき、特定の1つの演算器にm個の入力がある確率 $P(m, n, i)$ は2段分布となるので、最高m個の入力があつた時、他の演算器に合計n個のタップルが入力されいる確率 $P(m, n, i)$ は、次式のごとくな。

$$\left. \begin{aligned} P(m, n, i) &= f(m, n+i-1, n-1) \cdot n \cdot (1/n) \\ &= \binom{n+i-1}{n-1} \left(\frac{1}{m}\right)^{n-1} \left(\frac{n-1}{m}\right)^{i-1} \\ \circ n \leq i \leq (n-1) \cdot (m-1) \end{aligned} \right\} 4-3$$

$$P(m, n, i) = f(m, n+i-1, n-1) \times \left\{ 1 - \sum_{j=0}^{i-n} P(m-1, n, j) \right\}$$

これより、その期待値 E_U と集合演算器の利用率 E_U は、次式となる。

$$\left. \begin{aligned} E_U &= \sum_{i=0}^{(n-1)(m-1)} i \cdot P(m, n, i) \\ E_U &= \{E_U / [(m-1) \cdot n]\} \cdot 100 \end{aligned} \right\} 4-4$$

図4-2に、 E_U の計算結果を示す。こぞより、mを数十、nを数百とすれば、80%程度の利用率が得られることが判る。

4.4 集合演算器の入力追随性とバッファサイズ

集合演算器の入力追随性が問題となるのは、NPSとOPの処理がハイブリッド化されているRS、RT状態に於いてであるが、OPがUT(R_TRS)を行なう時以外は、OPも(データ転送+数steps)の処理しか実行していないので

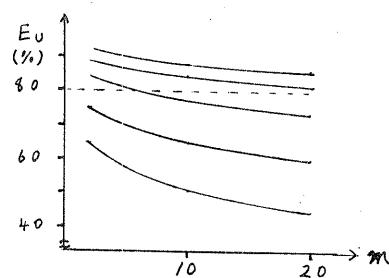


図4-2 集合演算器利用率

データの流れが止ることはない。

$OUT(R_T R_S)$ により、OPの処理が遅れるなど、 $R_T B$ オーバーフローし、D/PNETがデータを集合演算器に送れなくなる。このとき、3.1で述べたごとく、セル論理装置は1回転分の空読みを行ふ。この間OPは、 $R_T B$, $R_S B$ のタップオールに対して $OUT(R_T R_S)$ を行つてゐる。従つて、 $R_T B$ の大きさを、空読み時間とOPの処理時間が一致する程度にしておけば、 $DUT(R_T R_S)$ でも、無駄な遅延がOPに生じることはない。

$R_T B$ オーバーフローが恒常に発生するのは、タップオール割り当てる率とサービス割り当てる率が等しいが、前者の方が大きくなつた時である。 $OUT(R_S R_T)$ の処理時間は、出力量に比例するので、これは、入力量の2倍以上の出力があるとき発生する。従つて、 $R_T B$ の大きさは、補助記憶の1セル分の記憶容量の1/2より大きくする必要はない。

5. あとがき

本稿では、セル論理装置の並列読み出し速度に応じてダイナミックリラストリングを行つて PNETと、タップオールの転送速度に遅れることがある関係代数演算を行つて集合演算器を用ひることにより、 $O(n)/P$ (n :タップオール数, P :セル論理装置の並列読み出し数) で集合演算が行えるDBMを提案した。

補助記憶としてディスクを使用した場合、D/PNETの入出力数は、32~64個であり、また、各ノードに必要なバッファは10タップオール分程度である。集合演算器を、RSIMに64KByte, $R_T B$ に数KByte, KPSのキヤ松納部に16KByte程度用意すれば充分である。また、ディスクからの転送速度を1M~1.5Mbps/secとし、1タップオール100バイトとすると、100μsec~66μsecで1タップオール分の処理ができるよう。

残された問題点として、① explicit JOINの性格と $R_T B$ の大きさに関する詳細な解析、②ハッシュ因数の選択法、③集合演算器の数とデータの実現値の種類数不同程度が、後者の方が多い場合の処理、④障害復旧法の検討などである。

これについては、本稿で述べたDBMの技術の詳細化とともに、今後検討を進めて行く予定である。

〈謝辞〉 日頃御指導頂く第一研究室山下室長、春藤調査役をはじめ、種々御討論頂いた第一研究室の諸兄に深謝する。

〈参考文献〉

1. IE³ TRANSACTION on COMPUTERS VOL.28 NO.6 1979
2. IE³ COMPUTER VOL.12 NO.19 1979.3
以上2つは、IBMの特集号である。中の個々の論文は省略する。
3. E. BABB Implementing a Relational Database by means of specialized hardware. A.C.M. Trans. on Data Base Systems VOL.4. 1. 1979.3 PT.1~24
4. 青木他 多目的エーカベットによるリレーショナル演算の並列処理とその評価 信学会計算機研究会 EC79-66
5. S. Uemura et al. The Designing and Implementation of a Magnetic-Bubble Data Base Machine Proc. IFIP 80
6. B.W. Web et al. DIALOG - A distributed processor organization for data base machine Proc. AFIPS NCC 1980 P243
7. H.T. Kung et al. Systolic (VLSI) Arrays for Relational Data Base Operation Proc. ACM SIGMOD 1980 P105
8. 上野 ソーティングハードウエアを用いたデータベースシステム 京大大型計算機センター開拓セミナー報告 1980.3
9. Y. Tanaka Pipeline Searching and Sorting Modules as Components of Data Flow Database Computers Proc. IFIP 80
10. M. Maekawa Quick Parallel JOIN and Sorting ALGORITHMS 13th. IBM Conference Symp. 1979.10 P II-11
11. D.E. Shaw A Relational Data Base Machine Architecture 5th workshop on Computer Architecture for nonnumerical processing ACM SIGARCH, SIGMOD, SIGIR, P84
12. 小倉他 16ビット超高速メモリLSI 信学会 EC80-44
13. E. Gallo et al. Parallel Hashing Algorithms Technical Rep. of Tokyo Univ. TR-78-16
14. 吉川 データベースの知的アクセスに関する研究 豊橋研報告 第604号
15. 植村, 前川著 データベーススキン 情報処理叢書1