

可変構造多重処理データベースマシンにおける ソートモジュール

SORT MODULE
IN
DYNAMICALLY RECONFIGURABLE AND HIGHLY PARALLEL DATABASE MACHINE

喜連川 優
M. KITSUREGAWA

鈴木 重信
S. SUZUKI

田中 英彦
H. TANAKA

元岡 達
T. MOTO-OKA

東京大学

工学部

UNIVERSITY OF TOKYO

FACULTY OF ENGINEERING

1. はじめに

近年、データベースの普及に伴い、ソフトウェアサポートによるDBMSの性能の限界が次第に明らかとなり、データベース管理の諸機能をハードウェアでサポートする事を試みたデータベースマシンの研究が各所でなされる様になった。⁽¹⁾我々も既にリレーショナルデータベースのサポートを中心にセルラロジックタイプのDBMをより発展させたMIMD型マシン可変構造多重処理データベースマシンの開発を行なってきた。⁽²⁾図1,2に示される本マシンのデータ操作部は、多重チャネルリングバスにより多くのプロセッシングモジュール(P)群とメモリモジュール(M)群との間に柔軟な結合関係を生成し、処理負荷に応じた動的プロセッサ割付けを可能としている。又、メモリからプロセッサ群へはデータブロードキャストを行ない、メモリコン

フリプトのない高次の並列処理が可能となし、駆動プロセッサ台数にほぼ比例した処理速度の向上が得られる。メモリモジュールには、マイナ方式の磁気バブルメモリを用い、マークビット管理機構によりデータ実効転送レートの向上を計っている。今回、更に処理の高速化を計る為、ソートモジュールの検討を行なったので報告する。

ソートはデータ処理の基本操作であり、ソートされるとその後の処理は多くの場合容易になる。ここではデータベースマシンへの組込みを意図しているが、ソータ自身としても当然有益と思われる。データベース管理システムの発展が目ざましいが、そこではソートは一層不可欠な機能と考えられる。リレーショナルデータベースでのジョインやプロジェクションはソートされていれば、 $O(N)$ で処理可能であるし、又一般データベースでのインデックス生成などでもソートが

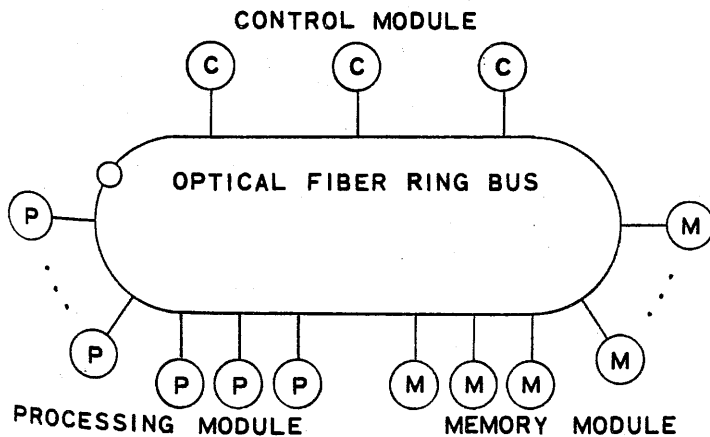


FIG. 1 ARCHITECTURE OF DATA MANIPULATION UNIT

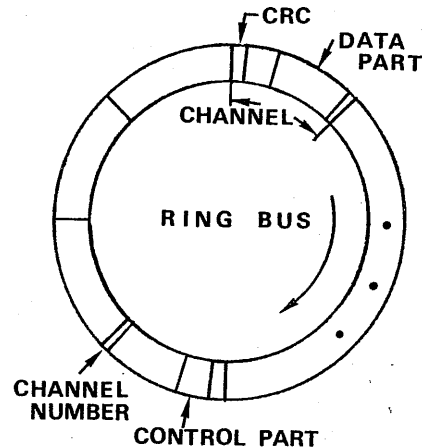


FIG. 2 RING BUS DATA FORMAT

必要とされ、高速ソータの需要は大きいと考えられる。更に二次記憶からのデータレートに合った $O(N)$ ソータは外部記憶に於ける von Neumann Bottle Neck を解消する上でも重要である。

2. ソータ設計指針

ソータ設計に際し、次の様な点に留意した。

- 1) 入カデータストリームに同期したソータ
 一般にソータは N 本のレコードに対し、 $O(N \log N)$ の比較が必要とされ、クイックソート、トナメントソータなど多くの効率の良いアルゴリズムが存在する。しかし、内部ソータでは、ソータすべきデータが処理開始時点で全てコア内に存在する事を仮定している。一方、実際のデータ処理応用分野では多量のデータは二次記憶装置上に存在する場合が多く、主記憶にロードする時点で同時にソータ出来れば理想的である。ここでは、外部記憶装置からの入カデータストリームに遅れることなく同期したソータを行ない、処理と入カとをオーバーラップさせることにより、最後のレコードの入カ終了と共にソータされたレコード出力が得られる様な構成を目標とする。 $(O(N))$ ソータ

- 2) レコード長に対する考慮

キー部のみを扱う場合、或はエンコードされたデータを扱う様な場合には、レコード長は短く、固定的アーキテクチャで問題はない。しかし、一方で一般レコードを扱う場合も考えるべきであり、レコード長の変化に対して柔軟に対処できる構造が好ましい。一般レコードのソータに於いてソータ対象をレコード全体ではなくキー部に限定する事によってレコード長に制限を設ける事も可能であるが、パイプラインソーティングに於いてはキー部、非キー部の分離には意味がない場合が多い。即ち、キー部のみソートし非キー部はソート後結合する事も可能であるが、フィールド毎のファイルではなく、レコード単位の従来のファイル構成の場合には、二次記憶からのレコード取出しレートによってプロセッサの処理能力が決定され、パイプライン処理では非キー部のデータ移動は処理負荷の増大にはつながらない。従って、レコード全体をソータする方が便利である場合も多く、短いレコードのソータ文でなく、一般レコードのソータも考える

事とする。そしてこの場合には、レコード長の変動が激しいと考えられ、柔軟に対処し、なるべくソータ出来る様考慮する事とする。但し、簡単な為、1ストリーム内ではレコード長一定とする。これは可変構造多重処理データベースマシンでは1リレーション内でダブル長に定としている事による。

- 3) 半導体RAMチップの利用

LSI化を考える場合、Kungらの Systolic Array の如くロジックとメモリセルの一体化が自然であるが、より柔軟な制御構造を入れようとするといつ一つのセルが複雑なものになってしまう。セルラロジックではセルの簡素化と構造の柔軟性とは相反する要素と考えられる。ここでは $\log_2 K$ 台 (N : ソータすべきレコード数, K : K way sort) のプロセッサと大容量半導体RAMバンクによる構成とした。 $O(N)$ ソータの最も簡単な手法は N 台のプロセッサを用いる事であるが、ここでは $\log_2 K$ 台にし、プロセッサオーバーヘッドを少なくした。又半導体RAM技術は一層進歩すると思われる。RAMチップを利用する事とした。

3. ソータアルゴリズム

3-1 アルゴリズム

ソータのアルゴリズムはマージソータのパイプライン化であり、8)、9)によって既に知られている。

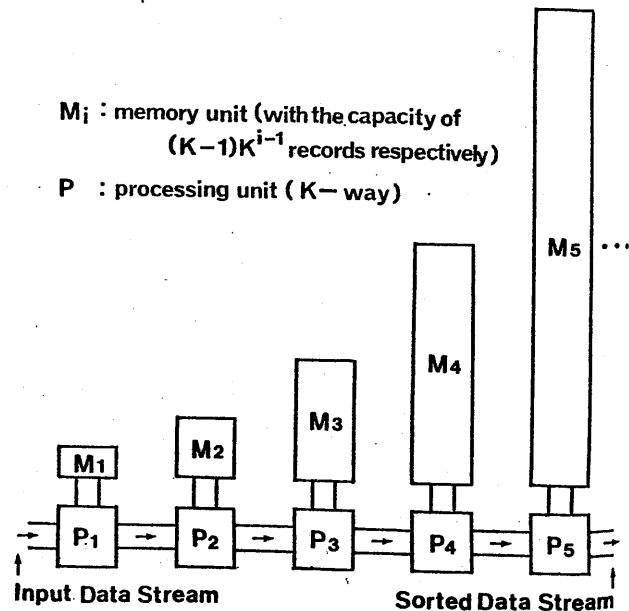


Fig. 3 Global Architecture Of Stream Driven Sorter

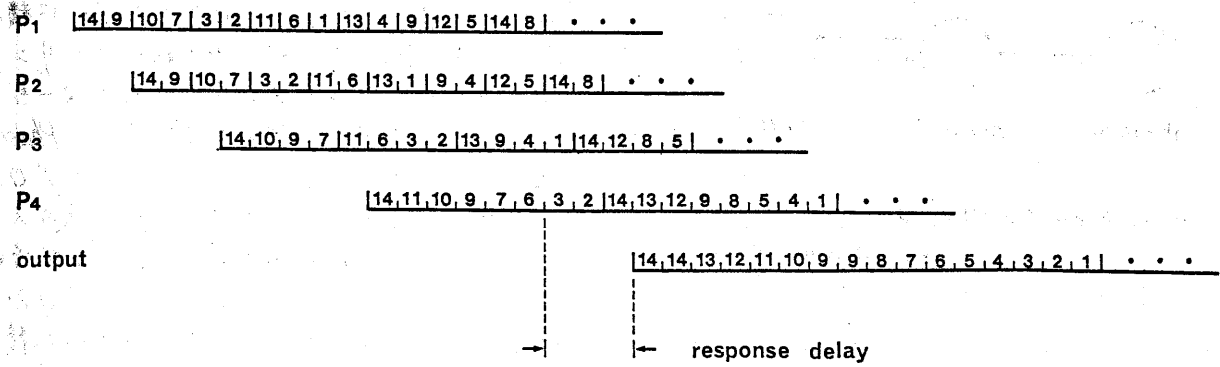


Fig. 4 Sorting Process Overview Of Stream Driven Sorter

necessary processors ----- $\log_k N$
 response delay ----- $\log_k N - 1$
 total sort time ----- $2N + \log_k N - 1$

今、 $N (= k^l)$ のレコードのソートを行なうものとする。 K -way mergeを行なうプロセッサを k 台用意し、一次元状に結合する。この際、 i 番目のプロセッサには $(k-1)k^{i-1}$ レコード分のメモリを付加する。全体の構成は図3に示される。 N 個のレコードはシリアルに i 番目のプロセッサに入力され、 i 番目のプロセッサは $i-1$ 番目のプロセッサから送られてくる k^{i-1} レコードからなるソートされたストリングと k 本マージして、 k^i レコードからなる1本の長いストリングを生成し、 $i+1$ 番目のプロセッサへ出力する。図4の如く、各プロセッサはマージすべき k 本のストリングに関して、 i 番目のストリングの最初のレコードが到着した時点でマージ処理を開始できる。パイプラインのセグメントは簡単の為、上図では、レコードとしているが、これはバイトレベル、ビットレベルにまで落とす事が可能である。(我々はバイトレベルを採用)

3-2 諸特徴量

- 入力レコード"ストリーム"が終了してから、ソートされたレコード"出力"が可能になるまでの遅れ

$$\log_k N - 1$$

これは最終レコードが直列に結合されたプロセッサ間を通り、最終段のプロセッサまで到達するには必要な時間に対応している。この様に遅れは極めて僅かといえる。

- 入力レコード"全体"のソートに必要なとされる時間

$$2N + \log_k N - 1$$

本アルゴリズムでは全レコードをソートに一度入れ、その後とり出す為、最低 $2 \cdot N$ 時間は必要となる。(昇順ソートで、最小レコードが一番最後にやってくるかも知れない為) 時間的オーバーヘッドは上で述べた $\log_k N - 1$ だけである。

• 所要メモリ容量 N

各段でプロセッサはそれぞれ、 k^i レコードからなるストリングを $k-1$ 本貯え、最後の1本の到着とともにマージを開始する。それ以後は1レコードの入力とともに1レコード出力する為、各段では $(k-1) \cdot k^i$ レコード分のメモリ容量が必要となり、全体としては、 $N (= k^l)$ 即ち、全レコード分のメモリ容量があれば良い事がわかる。

3-3 プロセッサ内処理

プロセッサ内での処理は、基本的には既に述べた如く、前段からの k 本のストリングをマージして1本にし、次段へ出力する事であるが、より細く見ると以下の如く3つのPhaseに分かれる。

- Phase 0 最初のストリング"ロード"の時であり、 $(k-1) \cdot k^{i-1}$ 本のレコードを当該プロセッサ内メモリに貯えるPhase。次段のプロセッサへはまだデータを出さない。 $k-1$ 本のストリング"入力"後 Phase 1 へ移る。
- Phase 1 $(k-1)$ 本のメモリ内ストリングと入力されてくる k 番目のストリングをマージして出力するPhase。当該ストリングの入力が終了すると Phase 2 へ移る。
- Phase 2 Phase 1 で当該ストリングの入

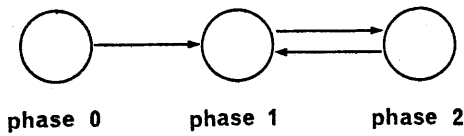


Fig. 5 Phase Transition Of Stream Driven Sorter

力が終わると、次の新しいストリングの入力が始まる。しかし、この段階では前のK本のストリングマージは完結しておらず、K本合わせて $(K-1) \cdot K^{i-1}$ レコードがメモリ内に残っている。Phase 2では、この残っているK本のストリングマージを行ないつつ、次のストリングの入カを行なう。K-1本のストリングを入カするとマージ処理は完了し、Phase 1へ移る。(図5)

4. ソータの実装方式

4-1 メモリ管理

K番目のプロセッサは K-1本のストリングを入カした後、K番目のストリングに対し、1レコード入カしては比較し、1レコード出力する為、常に有効なメモリ容量は一定である。ところが、マージ処理時には、ストリングの先頭を取る順番はランダムであり、当該レコードの出力によって空となった所に順次入力していくのでは、ソートされた順序を維持する事は出来ず、何等かのレコード管理機能が必要となる。これに対し、ここでは以下の如き3つの手法を検討した。(以下簡単な為 *2-way sort* を仮定する。)

4-2 Double Memory Method

2-way の場合、最も簡単な手法は2倍の容量を用意し、マージすべき2つのストリングに対し、各々のエリアを割り当てる事である。(図6)

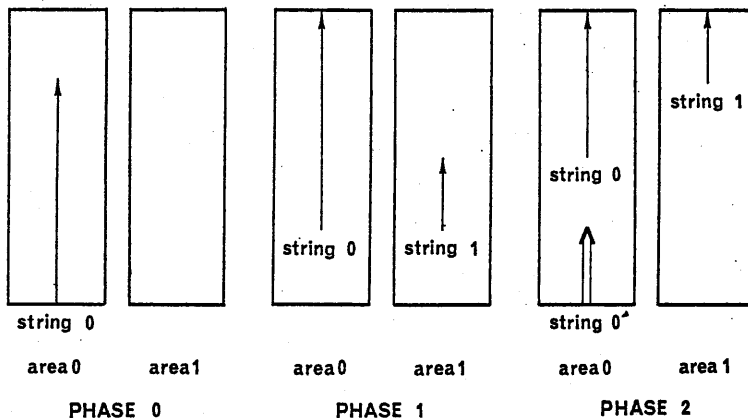


Fig. 6 Sorting Process In Double Memory Method

Phase 0では、ストリング0をエリア0に積み上げてゆく。Phase 1では入力されるストリング1をエリア0に存するストリング0とマージするとともにエリア1に積み上げてゆく。Phase 2ではエリア0, エリア1に残ったストリング0及びストリング1をマージすると共に、次のストリング0をエリア0に積み上げる。いかなる時点に於いても、ストリング0, 1, 0の合計は2ⁱ⁻¹レコードから成り、従つて、エリア0に積み上げる事によってストリング0が残存するストリング0に上書きする事はない。即ち、元来、ある時点でプロセッサ内には2つのマージすべきストリング及び、次ストリングの3つのストリングが現在する事になるが、2本分のエリアで足りるのである。この手法は制御が大変簡明ではあるが、この段階では2レコードの容量を有し、これが大きくなると50%の無駄は無視出来ない。

4-3 Pointer Method

各レコードにポインタ領域を持たせ、ポインタによってレコード間のシーケンスを維持する手法。この手法はレコード部に対してポインタ部のオーバーヘッドが無視出来る場合には、制御が容易で有効である。

4-4 Block Division Method

Double Memory Method と *Pointer Method* の中間策であり、メモリを固定長のブロックに分割する。ブロック内ではデータはソートされており、ブロック間の前後関係を別の構造として持たせるブロック単位のメモリ管理手法である。この手法に関して次の事実が成立する。

「当該プロセッサに必要とされるメモリ領域を l ケの等しいブロックに分割すると、 $l+2$ ケのブロック分のメモリ領域を用意する事によって、1ブロック分のレコード入カ後、必ず1つの空ブロックが生成される。」

この事実を利用し、以下の如く実行する。Phase 0でストリング0を l ケのブロックに積み上げた後、Phase 1では、入力されるストリング1とストリング0とマージしながら適当に選んだ空ブロックに積み上げていく。ストリング1を1ブロック入カし終える毎に

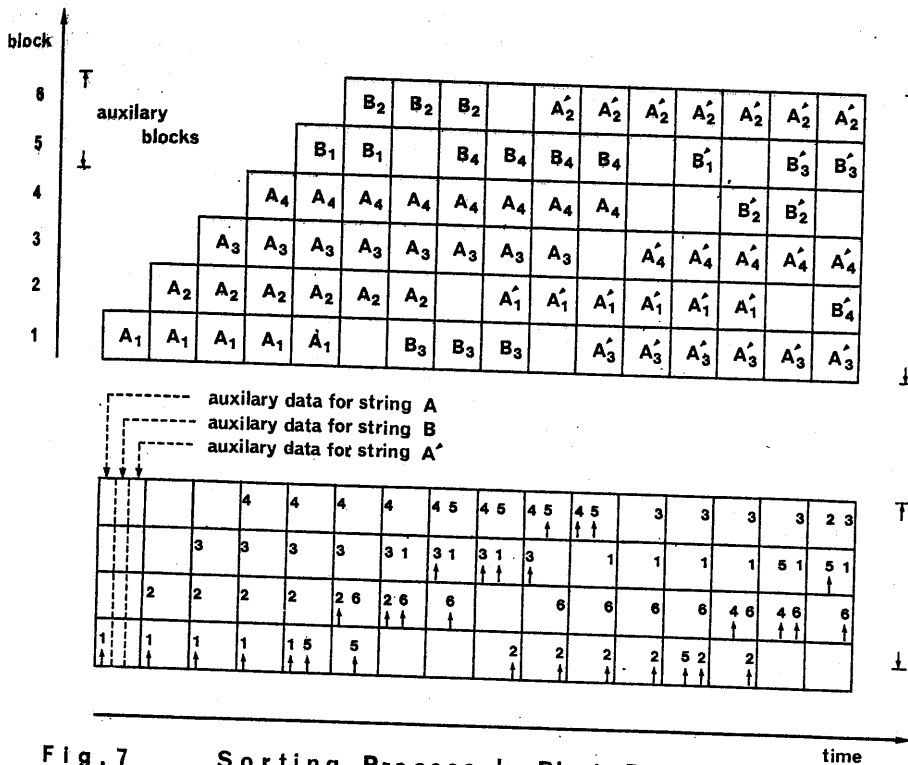


Fig. 7 Sorting Process In Block Division Method (Block Fetch Sequence)

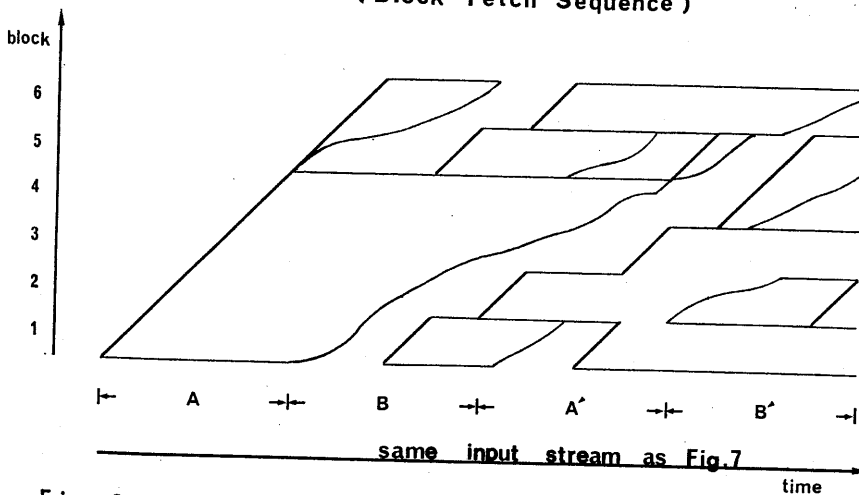
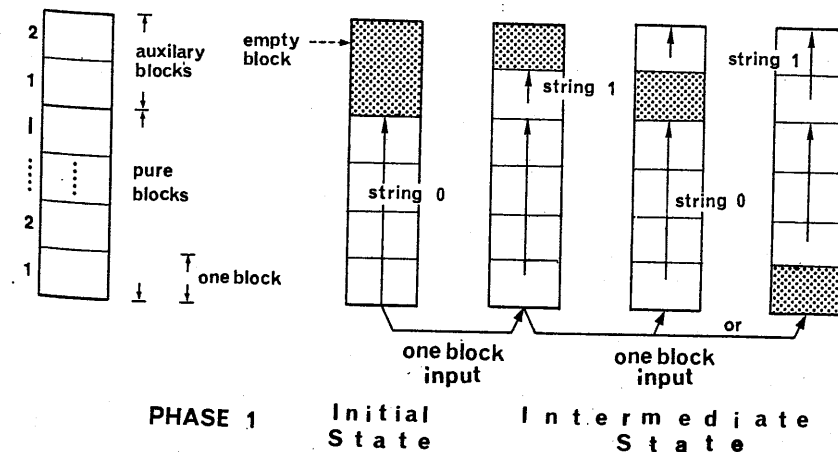


Fig. 8 Sorting Process In Block Division Method (String Top Position Trace)



先の事実から少なくとも1つ空ブロックの存在が保証されており、そこへ入かして行く事が出来る。Phase 2でも同様に、string 0を1ブロック入力する毎に、空ブロックが生成され、string 0と1のマーヅと並行してstring 0の貯蓄がなされる。この際、空ブロックの生成される順番はデータに依存しランダムであるから、ブロック間のシークは補助的データによって維持されているものとする。(図7, 8, 9)

さて、前記の事実は、以下の点から明らかである。即ち、もし、1つの空ブロックもないとすると、そもそもメモリ内には、完全につまみつかないブロックは、高々一つしか許されないから、少なくとも1ブロック分のレコードと、それに加えてレコード分必要になってしまう。これは各段での必要とされる一定のメモリ容量と一致しない。

Fig. 9 Block Division Method. Empty Block Generation After Each Block Input

• ブロック管理用補助データの構造

Block Division Method ではブロック間の順序を維持する為の補助的データが必要となる。このブロックレベルでのシーケンス維持に対しても、今まで述べて来たレコードレベルでの手法を再び適用する事が可能である。即ち補助データの補助データを考える事が出来る。しかし、このデータ容量は元来、真のデータに比べると僅かな量であるから、ここでは容量を節約するよりもむしろ簡単な構造を持たせる方が賢明である。即ち、Double Memory か又は All Pointer Method が好ましい。さて All Pointer Method はそもそもレコード長に比べてポインタ長が短い場合に有効な手法と考えられるが今回は、実レコードがブロック No. でありポインタ長はそれに等しくなる。依つて必要なメモリは両者大差なく、ここではむしろ Double Memory Method を採用する事が最も容易と考えられる。

4-5 K-way への拡張

これまでは 2-way Sorter の場合について説明して来たが、これは容易に K-way に拡張することが出来る。K を大きくする事により、1つのプロセッサはより複雑となるが、一定容量を実現する為のプロセッサ数、メモリバンク数は少なく済み、又遅延時間をわずかに減少出来る。逆にメモリ容量は 2-way の場合に比べて一台増す毎に大きく変化するが、これに関しては最終段のプロセッサが $2 \leq K' \leq K$ とし、そのパワーを落として調整することによって適切な容量を実現できる。

K-way でのメモリ管理方式は次の様に変更される。

Double Memory Method

K本のストリングをマージする K-way ソータでは、 $2 \times (K-1)$ 本のストリング容量のメモリエリアを用いる事によりメモリ管理可能である。

Block Division Method

K-way ソータでは、当該プロセッサの1本のストリングを l ヶの等しいブロックに分割すると、 $(K-1) \cdot l + K$ ブロック分のメモリ領域を用意する事により、1ブロック分のレコード入力後、必ず1つの空ブロックを生成できる。

4-6 Block Division Method に於ける補助メモリのコスト評価

Block Division Method に於いては真に必要なメモリ容量の他に Kヶのブロック及びブロック間順序を維持する為の補助メモリが必要となる。既に述べた如く、補助メモリとしては、Double Memory Method を利用する事にする。ここでは、各段の補助メモリコストの評価式を求め、最適ブロック分割を行なうとともに、真のデータに対するコスト比を計算する。尚ブロック分割数、way 数 K は簡単な為 2 のべき乗とする。i 番目のプロセッサについて考えると補助メモリのコストは以下の如く表わされる。

- K-way $K = 2^a$
- レコード長 M (bit)
- 真に必要なメモリ容量 (レコード数) $(K-1) \cdot K^{i-1} = (2^a - 1) 2^{a(i-1)}$
- ブロック分割数 $l = 2^b$
- 1ブロックの大きさ $K^{i-1} / l = 2^{a(i-1)-b}$
- 補助ブロック数 K
- 補助ブロックの大きさ $K^{i-1} / l \cdot K = 2^{a(i-1)-b}$ 以上より
- ▶ 全実データブロックの大きさ $(K-1)l + K$ (ブロック) $= (2^a - 1) 2^{a(i-1)} + 2^{a(i-1)-b}$ (レコード) $= M \cdot [(2^a - 1) \cdot 2^{a(i-1)} + 2^{a(i-1)-b}]$ (bit)
- 実データブロック数 $(K-1)l + K = 2^{a(i-1)-b} + 2^a - 2^b$
- 補助データ用ブロックポインタ長 $\lceil \log_2 \{ (K-1)l + K \} \rceil$ (bit) $\begin{cases} = a + b & (a \leq b) \\ = a + b + 1 & (a > b) \end{cases}$
- 補助データ用メモリエリア数 $2 \cdot (K-1)$ $= 2 \cdot (2^a - 1)$ (Double Memory Method)
- 1エリア内レコード数 $l = 2^b$ 以上より
- ▶ 補助データ容量 $\lceil \log_2 \{ (K-1)l + K \} \rceil \cdot 2^{a(i-1)-b}$ $\begin{cases} = (a+b) \cdot 2 \cdot (2^a - 1) \cdot 2^b & (a \leq b) \\ = (a+b+1) \cdot 2 \cdot (2^a - 1) \cdot 2^b & (a > b) \end{cases}$
- ▶ 全メモリ容量 $C(M, K, l, i)$ (bit) $= C(M, a, b, i)$ $= M \cdot \{ (2^a - 1) \cdot 2^{a(i-1)} + 2^{a(i-1)-b} \}$ $\begin{cases} + (a+b) \cdot 2 \cdot (2^a - 1) \cdot 2^b & (a \leq b) \\ + (a+b+1) \cdot 2 \cdot (2^a - 1) \cdot 2^b & (a > b) \end{cases}$

▶ 補助メモリの全メモ
容量に対するコスト比 $R(M, a, b, i) = C(M, a, b, i) / M(2^a - 1) \cdot 2^{a(i-1)} - 1$

今、レコード長 M 及び K が与えられると各段で $R(M, a, b, i)$ を最小にするブロック分割数 l を求めることが出来る。図10, 11 は 2-way の場合について各段でブロック容量を変化させたときのコスト比の値と、その際求められた最適ブロック分割数及びブロック容量を示している。ここでレコード長は 8 byte とした。キーのエンコードを行なう応用分野では、この程度の長さで充分と考えられるが、Pointer Method で実現しようとするとき、この様に短いレコードの場合には、そのオーバーヘッドは無視出来ない。一方 Block Division Method では図12に示す如く、16段のソータで 1~2% 程度のオーバーヘッドに押さえられる事がわかる。

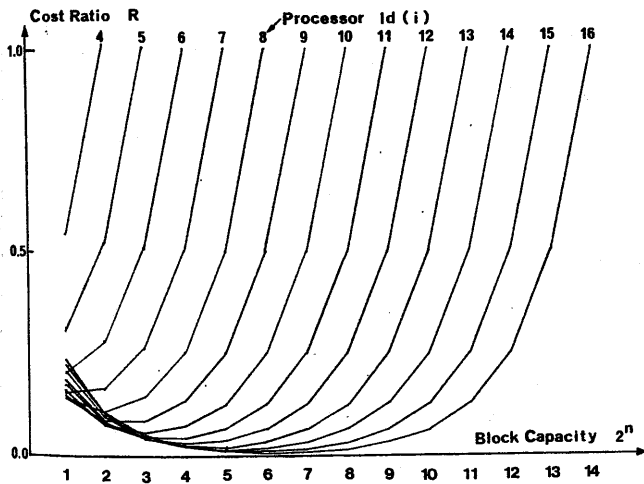


Fig. 10 Auxiliary Memory Cost Ratio For Each Block Capacity

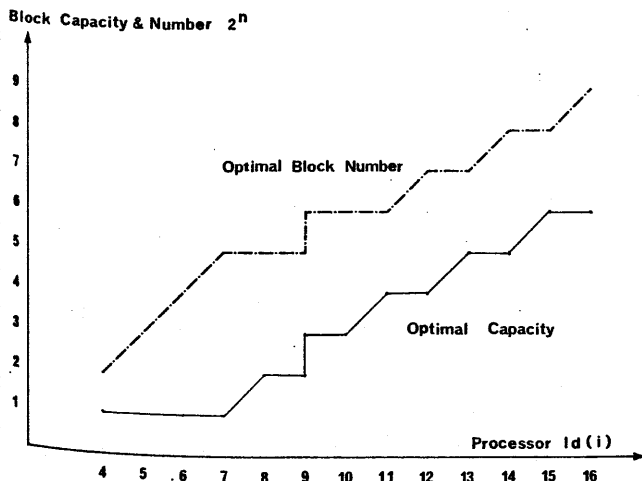


Fig. 11 Optimal Block Capacity & Number

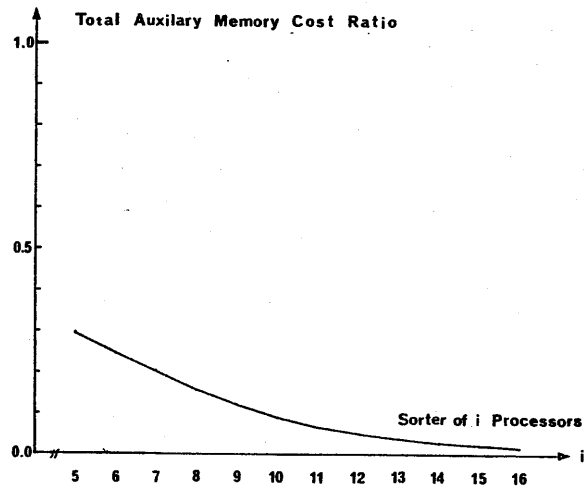


Fig. 12 Total Auxiliary Memory Cost Ratio For Optimal Stream Driven Sorter of i Processors

5. ソータの柔軟性

本節では、ソータの使用される環境及びその下で有すべき柔軟性について検討する。

5-1 レコード長とメモリ管理方式

a) レコード長が短い場合

キー部のみを高速にソートする場合、又それがエンコードされている様な場合には、17フルはかなり短いと考えられる。(数バイト) 従って Pointer Method を用いるとポインタ部だけでも 2~3バイト必要となり無視出来ない。この場合には Block Division Method を採用する事により効率化が図れる。この方式はブロックを大きくした極限で Double Memory Method に対応し、逆にブロックサイズを小さくした極限で 1 ブロックが 1 レコードに対応し、Pointer Method となる。

b) レコード長が長い場合

一般に、データの格納はレコード単位に行われ、フィールド分割は行なっていない。この様な場合、既に述べた如くパイプラインソータでは時間的ロスがない為、キー部の分離には意味がなく、フルレコードを直接ソートした方が便利である。レコード長は、かなり長くなると考えられ、ポインタ部の占める割合は僅かである為、制御の容易さから Pointer Method で良い。尚この場合には、レコード長の変化に対して効率よく対処する必要があり、次節以下に述べる。

5-2 レコード長変動に対する効率化 (String Length Tuning)

K-way ソータをレコード長 (L) に対して設

計した時、実際のレコード長 $X (\geq L)$ の大小によらず、高いメモリ使用効率を保ち、出来る限り多くのデータをソートする為には如何なる構成にすればよいかがこの問題である。一番目のプロセッサについて、高々 $K \cdot L$ 長のストリングを出力するとすれば、実際のレコード長 X に対し、メモリ効率 η_k は

i) $X \leq L$ (短い時)

$$\eta_k = \frac{X}{L}$$

ii) $K \cdot L \geq X > L$ (長い時)

$$\eta_k = \frac{\left[\frac{K \cdot L}{X}\right] \cdot X}{K \cdot L}$$

と表わされる。

今、実際のレコード長が設計値よりも長い場合を考えると $y = K \cdot L / X$ とおくことにより $\eta = [y] / y$ と表わされ、最低50%のメモリ使用効率となることがわかる。平均使用効率 $\bar{\eta}_k$ は

$$\begin{aligned} \bar{\eta}_k &= \frac{1}{(K-1)L} \int_L^{K \cdot L} \frac{\left[\frac{K \cdot L}{X}\right] \cdot X}{K \cdot L} dx \\ &= \frac{K}{K-1} \sum_{i=1}^{K-1} \int_{\frac{K \cdot L}{i+1}}^{\frac{K \cdot L}{i}} \frac{1}{y^2} dy \quad y = \frac{K \cdot L}{X} \\ &= \frac{K}{2(K-1)} \sum_{i=1}^{K-1} \left(\frac{1}{i} - \frac{1}{(i+1)^2} \right) \end{aligned}$$

例えば、 $\bar{\eta}_2 = 0.75$, $\bar{\eta}_3 = 0.78$, $\bar{\eta}_4 = 0.81$ を得る。

さて、もしストリングの長さ (1ストリングを構成するレコード数) を変える事が出来れば、原理的には次段のプロセッサは $K \left[\frac{K \cdot L}{X} \right] X$ 丈の長さを有するストリングの代わりに $\left[\frac{K^2 \cdot L}{X} \right] X$ 丈の長さのストリングを出力出来る筈であり、メモリ効率は

$$\eta'_k = \frac{\left[\frac{K^2 \cdot L}{X}\right] \cdot X}{K^2 \cdot L} \quad \text{を得る。}$$

先と同様に平均効率を求めると

$$\begin{aligned} \bar{\eta}'_k &= \frac{1}{(K-1)L} \int_L^{K \cdot L} \frac{\left[\frac{K^2 \cdot L}{X}\right] \cdot X}{K^2 \cdot L} dx \\ &= \frac{K^2}{2(K-1)} \sum_{i=K}^{K^2-1} \left(\frac{1}{i} - \frac{1}{(i+1)^2} \right) \quad \text{となる。} \end{aligned}$$

例えば、 $\bar{\eta}'_2 = 0.84$, $\bar{\eta}'_3 = 0.92$, $\bar{\eta}'_4 = 0.98$ を得る。実際、

$$K \left[\frac{K \cdot L}{X} \right] \leq \left[\frac{K^2 \cdot L}{X} \right] < K \left\{ 1 + \left[\frac{K \cdot L}{X} \right] \right\}$$

$$\left[\frac{K^2 \cdot L}{X} \right] = l_1 \left[\frac{K \cdot L}{X} \right] + l_2 \left\{ \left[\frac{K \cdot L}{X} \right] + 1 \right\}$$

$$K = l_1 + l_2$$

と表わす事が出来、プロセッサは K 本のストリングを次プロセッサに出力するが、この際、 l_1 回は $\left[\frac{K \cdot L}{X} \right]$ 本の長さのストリングを、 l_2 回はもう1レコード長くしたストリングを送出すればよい事がわかる。この様にストリング長の増加は l_2 回にわたってなされる為、1回のマージで必要とされる余分の領域は1レコード分だけである。即ち、必要とされる容量は $\lceil \lceil \frac{K^2 \cdot L}{X} \rceil \cdot X \rceil$ を用いて

$$\max_x \left[\left\lceil \left\lceil \frac{K^2 \cdot L}{X} \right\rceil \cdot X \right\rceil \right] = \frac{2 \cdot K^2}{K+1} \cdot L$$

と表わされ、従って、実際

$$\frac{K^2}{K+1} L - (K-1)L = \frac{L}{K+1}$$

丈のわずかな増加で済む事がわかる。

この様に一番目のプロセッサだけが、ストリング長の調整をする場合でも (1次 Tuning) メモリ効率は可成り改善されたが更に2番目のプロセッサもストリング調整を行なうとすると

$$\left[\frac{K^3 \cdot L}{X} \right] = m_1 \left[\frac{K^2 \cdot L}{X} \right] + m_2 \left\{ \left[\frac{K^2 \cdot L}{X} \right] + 1 \right\}$$

$$= m_1 \left\{ l_1 \left[\frac{K \cdot L}{X} \right] + l_2 \left(\left[\frac{K \cdot L}{X} \right] + 1 \right) \right\}$$

$$+ m_2 \left\{ (l_1-1) \left[\frac{K \cdot L}{X} \right] + (l_2+1) \left(\left[\frac{K \cdot L}{X} \right] + 1 \right) \right\}$$

と表わされ (l_1, l_2) を m_1 回実行後 (l_1-1, l_2+1) を m_2 回実行することで、 $\eta = \left[\frac{K^3 \cdot L}{X} \right] \cdot X / K^3 \cdot L$ を得ることが出来る。2段目以後のプロセッサは、 k -way merge を行ない初段のプロセッサが K 以下でマージ数を変える事によって全体の調整をすることになる。

以上の方式を一般化する事により更に、高次の調整が可能となる。即ち、

$$\left[\frac{K^i \cdot L}{X} \right] = l_1^{i-1} \left[\frac{K^{i-1} \cdot L}{X} \right] + l_2^{i-1} \left\{ 1 + \left[\frac{K^{i-1} \cdot L}{X} \right] \right\}$$

$$\vdots \quad \quad \quad \vdots$$

$$\left[\frac{K^2 \cdot L}{X} \right] = l_1^2 \left[\frac{K^2 \cdot L}{X} \right] + l_2^2 \left\{ 1 + \left[\frac{K^2 \cdot L}{X} \right] \right\}$$

$$\left[\frac{K^2 \cdot L}{X} \right] = l_1' \left[\frac{K \cdot L}{X} \right] + l_2' \left\{ 1 + \left[\frac{K \cdot L}{X} \right] \right\}$$

と表わされ、 $i-1$ 番目までのプロセッサが、ストリング長を調整することにより、

$$\eta_k^{i-1} = \frac{\left[\frac{K^i \cdot L}{X} \right] \cdot X}{K^i \cdot L} \quad \text{を得る。}$$

又平均値は

$$\bar{r}_k = \frac{K^i}{2(K-1)} \sum_{j=K^{i-1}}^{K^i-1} \left(\frac{1}{j} - \frac{1}{(j+1)^2} \right)$$

となる。これを図13に示す。

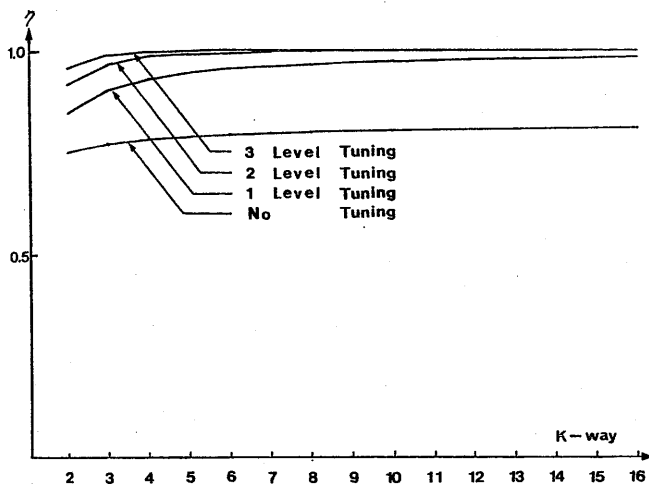


Fig.13 String Length Tuning & Memory Utilization Efficiency

以上の如く、レコード長しに対して設計したソータにより長さの異なるレコードを取り扱おうとすると、そのままの形ではメモリ使用効率が低くなってしまふ。初段での効率が次段へと伝搬してゆき、全体としてのメモリ効率は初段のそれによつて決定されるからである。これに対し、ここでは始めの数台のプロセッサに余分のメモリを少々加え、ストリングの長さを調整することにより後段でのメモリ効率を高めた。メモリ容量は段を増すに従つて急増し、メモリ効率は最終段でほぼ決定される為、最初のプロセッサのメモリ増量は問題にならない。又、実レコード長の範囲を、 $L < X \leq K \cdot L$ としたが、 $X > K \cdot L$ の場合には、2番目のプロセッサが初段として振舞うことにより、同様の制御が可能である。逆に、実レコード長の方が短い場合には元来、 $\log_k \frac{X}{L} N$ 台のプロセッサが必要であるにもかかわらず、現実には $\log_k N$ 台しかなく、そもそもプロセッシングパワーが足りない為効率が落ちる。従つて、短いレコードも効率よく処理する為には、余分のプロセッサを先頭に付加すればよい。この様な方式により、任意のレコード長に対してそのゆらぎを吸収し効率よくソート出来る事がわかる。

5-3 ストリーム長変動に関して

m 台のプロセッサからなる K -way Sorter では、 K^m 個のレコードのソートが可能となる

が、これより短いストリームの場合、最も簡単な手法は、ダミーレコードをストリームの最後に付加し、 K^m 個のとしてソートする事である。しかし、これではソータの容量が増大するにつれて、短いストリームに対する出力遅れが甚だしくなる。そこでこれに対しては、次の手法を採つた。

最初のプロセッサには、最終レコードの先頭にストリング「エンド」のフラグ(eos)を付加して入力し、以後のプロセッサは、ストリーム内の最終ストリングの先頭レコードに eos を付加して出力する事とする。又、プロセッサは eos の検出により、マージを開始する。この様になると eos フラグはプロセッサを介する毎にストリーム内での相対位置が前方にずれて来る事になる。そして最後には、eos がストリームの先頭 SOS (ストリームスタートフラグ) に一致し、それ以後のプロセッサはメモリを介さずに次段のプロセッサへバイパスさせればよい。こうする事により遅れを最小化できる事がわかる。尚 eos, SOS の他に、最終ストリングの最終レコードを識別する為の eor フラグも必要である。これらのフラグはプロセッサ間での通信情報であり、例えば、ポインタ方式によりプロセッサを構成したとすると、レコード毎にポインタを操作する為のオーバーヘッドタイムが生ずるが、この時間内に埋め込む事が可能である。

6. ソートモジュールとデータベースマシンに於けるその位置付け

本節ではこれまで述べて来たソータのデータベースマシンの中での位置付けについて述べる。ソータを用いたソートモジュール内での処理及びソートモジュールの駆動環境の2つの立場から以下説明する。特に後者では Hash とソートの利用により効率のよい $O(N)$ 処理方式の概容が示される。(図14)

6-1 ソートモジュール内での関係代数処理

a) Projection

重複除去は、処理負荷の重いオペレーションであり、SQL など陽に指定しないと実行しないシステムも多い。ソータを用いる事によつて $O(N)$ で処理出来る事は明らかである。

SORT MODULE

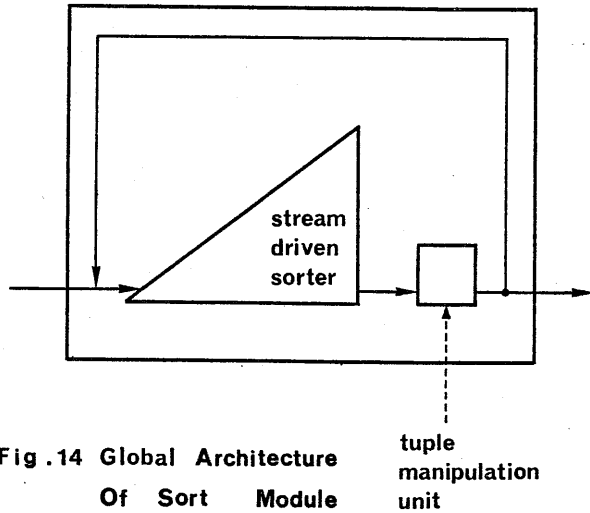


Fig. 14 Global Architecture Of Sort Module

b) Join

Joinは直接実行しようとするると両リレーションのカーディナリティの積に比例した処理時間が必要とされるが($O(N_1 \times N_2)$)、両リレーションがソートされていれば、略 $O(N_1 + N_2)$ で実現出来る。Joinは基本的にはソータを2つ用意し、各々に2つのリレーションを割当て、それをマージすればよいが、実際には、ソータを2つ用意する事は、ソータ内で一段プロセッサを増加する事に等しく、一台のソータにより、Joinを実現する事が望ましい。これはレコードのキーと非キー部の間にリレーションid (0, 1) を検索してソートする事によって解決する。又Join処理はマージと変わらないが、直積をとる際バッファオーバーフローが生ずる可能性がある。即ち、ソータ出力の2つのリレーションをマージし、直積をとる部位のバッファ内にいずれか一方のリレーションの同一ジョインキーを有するタプル集合が入り切らぬ場合には問題はないが、入り切

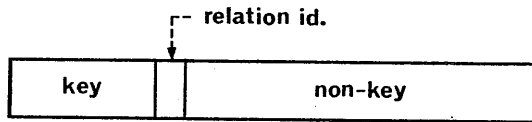


Fig. 15 Modified Tuple Format

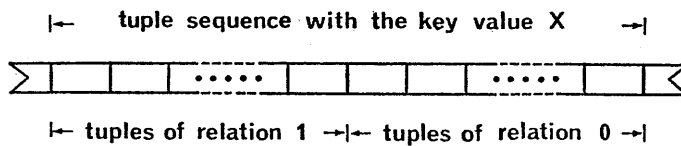


Fig. 16 Output Tuple Sequence Of Stream Driven Sorter

らない場合の方策を考える必要がある。本ソータではソートストリーム出力時にも内部処理が施されており、一端出力したデータは、再利用する事は出来ない。ここではバッファオーバーフローに対し、処理対象として残すべきタプルを再びソータに入力し、2回目のマージ時に処理する事とした。尚、処理済のタプルは捨てられる為2回目のスキャンでは不要なタプルはかなり篩い落とされるが、オーバーフローし、再入力したにもかかわらずJoin出来ないタプルに対しては、リレーションidを反転する事によって捨て去る事が出来る。(図15,16) 以上の如く、Projection, Joinはソータにより、高速処理可能であることがわかる。その他、集合演算操作 Intersection や Difference 更に、Divisionなど、処理負荷の重いオペレーションも同様に1つのソートモジュール内の処理に関しては高速に実行可能である。

6-2 Hashとソートによる関係代数処理

リレーション全体をソータのサイズのブロックに分割し、ソータ内では既に述べた如く $O(N)$ で、ブロック間では $O(n^2)$ で処理する事も可能であるが、ここではHashを利用する方式を考察したので、その概略を述べる。Hashの利用は10), 11) らによって既に指通されているが、これらはJoinやProjectionの候補を篩い出す前処理として利用するもので、実際の演算自体を高速化するものではない。ここでの方式はHashによるクラスタリング機能を利用するものである。Joinの場合、両リレーションにHashを施し、各々全体を幾つかのバケツに分割すると、結合のとれるタプルは互いに同じ番号を持つバケツのタプル同志に限られ、異なる番号のバケツ同志のマッチングは必要なくなる。この様：リレーションを幾つかの小さいサブリレーションに分ける際、単に分割するだけでは、そのレベルで $O(n^2)$ の結合処理が必要となるが、Hashを用いるとそのクラスタリング特性から $O(n)$ ですむ。そしてバケツ内でのマッチング処理は本ソータを用いれば、 $O(N)$ で処理が可能となり、従って全体として $O(N)$ のJoinを実現出来る事になる。処理手順は以下の如くなる。リレーションはディスク上に貯えられているものとし、Restriction Filter Processor を通して、必要なタプルのみが磁気バブル

メモリからなるメモリモジュールにステージングされる。この際、メモリモジュールはジョインタリビュートに関し、Hash を施し、動的にバケツId を生成すると共にマークビットによりこれを管理する。ステージングが終了するとメモリモジュールはバケツ毎にタプルを送出する。改良されたM/m 方式磁気バブルメモリを用いる事により特定のマークが施されたタプルのみを高速に送出する事が可能となる。⁽⁵⁾ この出力をソートモジュールが入力するが、当該バケツの出力が完了した時点で、ソータはソート出力可能となり処理が開始される。一方メモリは次のバケツの出力を始めるが、これは別のソートモジュールに割り付けられる。この様にしてバケツシリアルに処理を進めてゆく事が出来る。バケツの大きさはスイッチング"オーバーヘッド"によって決定され、又大きさのゆらぎに関してはオーバフローの生じない様 *Bucket Size Tuning* を行ない、プロセッサ使用効率を高いている。

以上、1つのオペレーションとして、 $O(N)$ の高速処理が期待出来る事が判るが、更に *Query Tree* 上オペレータ間でもパイプライン可能である。即ち、当該演算処理中の出力タプルはメモリモジュールに送出される毎に Hashing が施され処理フェーズと次オペレーションの準備の Hashing フェーズとを重畳する事が出来る。

この様にソータレベルの低次のパイプラインから関係オペレータレベルの高次のパイプラインまで種々のレベルの並列性を効果的に抽出する事により本マシンは高い処理能力を發揮すると考えられる。ここではバケツシリアルな処理方式について述べたが、バケツパラレルな方式⁽⁴⁾及び、パイプラインの擾乱効果、性能評価等、詳細は、別の機会に譲る。

7. おわりに

マージソートのパイプライン化を行なう *Stream Driven Sorter* について、メモリ管理方式及びその構成法に関する検討を行なった。現在ポインタ方式によるソータのLSI化を考えている。プロセッサ間は8bit中では結合され、ポインタ部は8バイト長である。レコードの1バイト処理は2回のメモリアクセスにより実現され1M byte / Sec 程度の入力レートに追従出来ると考

えられる。プロセッサ内レジスタ容量に関しては、レコード全体を入れようとする、レコード長の最大値を制限する事になり、好ましくない為、キー部のみをレジスタ内に置き非キー部はメモリに入れてポインタにより管理する事とした。キー部もメモリに入れ、キー長の制限を取ら事も可能であるが、メモリアクセスが増加し処理速度が低下する。非キー部は比較が不要な為余分のメモリアクセスを伴わずメモリに置いて速度が低下する事はない。

又Hashとソートによるデータベースマシンは、Hashによるクラスタリング機能とバケツ内ソート処理を結合し、処理負荷の重い関係代数演算子と $O(N)$ で処理出来る事が示された。更に本マシンは種々のレベルの並列性を抽出しており、従来のマシンに比べ高い性能が期待出来る。

参考文献

- 1) 植村, 前川著 "データベースマシン" 情報処理叢書1 情報処理学会 1980
- 2) 喜連川他 "可変構造多重処理データベースマシンの構成" 信学技報, EC80-51, 1980
- 3) 喜連川他 "多重チャネルリングバスに於けるブロードキャスト伝送制御手順" 情報処理学会分散処理研究会5-2, 1980
- 4) 喜連川他 "可変構造多重処理データベースマシンのシステム構成, ソートモジュール" 情報処理第22回全国大会, 3L-3, 1981
- 5) 鈴木他 "可変構造多重処理データベースマシンのメモリモジュール" 全上 3L-4
- 6) H.T.Kung et al "Systolic (VLSI) Arrays for Relational Data Base Operation" Proc ACM SIGMOD 1980
- 7) H.T.Kung "The Structure of Parallel Algorithms" Advances in Computers Vol 19 Academic Press 1980
- 8) Shimon Even, "Parallelism in Tape Sorting" CACM vol 11 No 4 April 1974
- 9) S. Todd, "Algorithm and Hardware for a Merge Sort Using Multiple Processors" IBM J. RES. DEVELOP vol 22 no 5 1978
- 10) E. Babb, "Implementing a Relational Database by Means of Specialized Hardware" ACM TODS vol 4, no 1 1979
- 11) D.R. McGregor et al "High Performance Hardware for Database Systems" in Systems for Large Data Bases. North-Holland, 1976