

3L-3

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 =

喜連川優 鈴木重信 田中英彦 元岡達
東京大学工学部

§0 はじめに

可変構造多重処理データベースマシンは Cellular Logic Type の DBM をより発展させた MIMD 型マシンである(図1)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(PM)群とメモリモジュール(MM)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インフリクトのない高次の並列処理を行なう。今回、更に処理の高速化を計る為、ソートモジュール(SM)の検討を行なったので報告する。

§1 ソートモジュールの設計指針

1) 入力データストリームに同期したソート --- MM は割当てられた当該チャネルに対しデータを連続的に出力する。SM はメモリからのデータ入力に遅れる事なくソートを行ない処理と入力をオーバーラップさせ、最後のデータ入力終了と共に、すぐに、ソート出力が得られる構成が好ましい。

2) タプル長に対する柔軟性 --- タプルデータをコード化し、対象データを一定長に制限する事は設計を容易にするが、ここではより一般的にタプル長の変更に対しても、柔軟に対処できる様にする。但し、簡単の為1リレーションのストリーム内ではタプル長は固定とする。(ドメイン選択バブルメモリを仮定)

3) 半導体RAM 技術の進歩への期待 --- LSI 化を考える際 Systolic Array²⁾の如くロジックとメモセルの一体化が自然であるが、より一般的な制御構造(Programmability)を入れ易くする為、メモリチップとプロセッサチップに分離することとし、 $\log_2 N$ (N : ソートタプル数, K : K -wayソート) 台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。

§2 ソートモジュールの構成

2.1 ソートアルゴリズム --- アルゴリズムは 3) の如くによって既に知られているパイプラインパラレルマージソートであり、図2に示される如くプロセッサは各ステージで等しく、メモリは 2-way の場合 i ステージでは 2^i 個のタプル分の容量をもつ。 i ステージ目のプロセッサは $i-1$ ステージ目のプロセッサから与えられる 2^i 個のタプルからなる 2 本のソートされたストリームをマージして 2^{i+1} 個のタプルからなる 1 本のストリームとして $i+1$ 番目のプロセッサへ出力する。前のストリームに対するマージと次のストリームの入力をオーバーラップさせパイプライン処理する事により、 $\log_2 N$ 台のプロセッサにより $O(N)$ の時間でソート出来る。

2.2 メモリ管理方式 --- i ステージ目のプロセッサは 2^i 個のタプルからなるストリームを入力した後、次のストリームに対し、1タプル入力しては比較後1タプル出力する為常に有効なメモリ容量は一定であり、 2^i タプル分あればよい事になる。ところが、ストリームトップをとる順序はランダムであり、空いた所に入力していくのではそのソートされた順序を維持出来ない。これに対し、ここでは次の様な方式を検討した。

a) Double Memory Method --- 最も単純な管理手法は 2 倍の容量を用意し、各々独立した FIFO として、2つのストリームに割当てる事である。しかし、 i 段目では、 2^i 個のタプルを入力する為、 i が大きくなると、この無駄は無視出来ない。

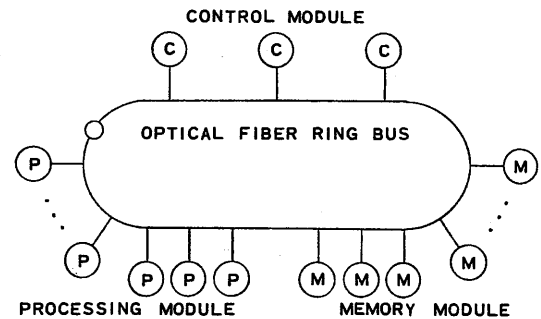


FIG.1 ARCHITECTURE OF DATA MANIPULATION UNIT

b) *Pointer Method* ---- ポインタによって、タプル間のリンクをとって行く手法

c) *Block Division Method* ---- a) と b) の中間策であり、メモリを固定長のブロックに分割し、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もたせる。今当該ステージでその必要なメモリ域をいくつかのブロックに分割すると $l+2$ 々のブロック分のメモリを用意する事によって 1 ブロックのデータ入力後必ず 1 々の空ブロックが生成されるという点に着目し、ブロック単位のメモリ管理を行なう。

2.3 レコード長による管理方式の選択

A) 短いデータのソート ---- 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、1 タプルはかなり短くなる (数バイト)。従って b) の手法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。この場合は c) の手法を採用する事により、効率化が図れる。最適ブロック分割数を各段で決定する事が出来、必要なメモリ容量は 1~2% 増に過ぎない。

B) 長いデータのソート ---- 一般的にデータの格納はレコード単位に行なわれ、フィールド分割は行なっていない。従ってデータ出力レートは一定であり、パイプラインソータでは時間的ロスがない為、キー部と非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な事が多い。この様な場合には、ポインタの占める割合はわずかであり、制御の容易さから b) の手法が好ましい。

2.4 ポインタ方式でのメモリ効率 ---- 2.3 の B) ではリレーション間でタプル長が変化しても、メモリ効率を高くし、成べく多くのタプルをソート出来る事が望ましい。K-way ソータに於いて (K が大きいと、ステージ数が少なくなり、メモリバンクの細分化が僅かです。) タプル長を L として、各ステージのメモリ容量を決定したとする。この時、実際のレコード長 X に対し、メモリ効率は $\eta = \frac{K \cdot L}{X} \times K \cdot L$ で表わされ、 $L < X < K \cdot L$ で最低 0.5 となる事がわかる。これに関しては最初の数段のメモリを $K \cdot L$ だけ増加し、最初のプロセッサがストリーム長を調整する事にすれば、 $K = 8 \sim 16$ の時平均 97~98% の効率が得られる。

§ 3. 関係演算の処理

当該リレーションを占める複数台の MM からのデータ出力はハッシュを施されて、互いに独立した幾つかのバケットに分割される。SM は各バケットに割付けられ、データの入力に従ってソートを行なう。データの分割が終了すると直ちにソート出力が可能であり、各モジュールは独立に処理を開始する。SM はソータと関係演算処理部からなり、Projection では $O(n^2)$ で重複除去を行なえ、又 Division では Divisor の Cardinality が大きい場合一部で割り、篩い落としをすると共に、商の可能性のあるタプルを再びソータの入力とする事で処理出来る。以上の如く、データがソートされている場合には比較的負荷の重い Join Projection Division も高速に処理出来る事は明らかである。

§ 4. むすび

現在ソータの LSI 化の設計を行なっている。今後 SM の設計、システムレベルの制御手法の研究を進める予定である。

- 参考文献 1. 喜連川他 信学技報 EC 80-51, 1980 2. H.T.Kung Advances in Computers vol 19 1980 P65
3. Shimon Even CACM April Vol 17 No 4 1974 p202 4. S.Todd IBM J.RES.DEVELOP Vol. 22 No 5 1978 p509

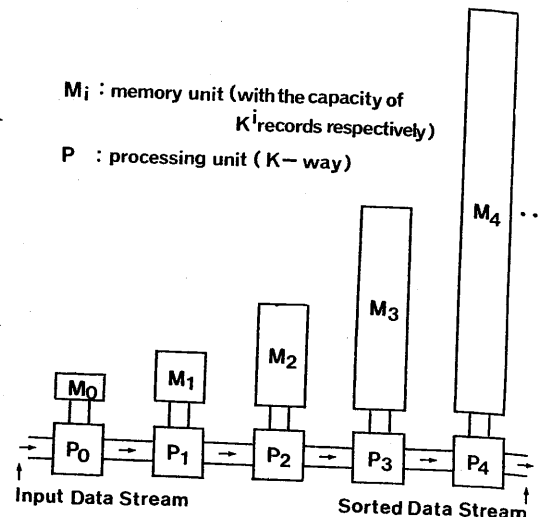


Fig. 2 Global Architecture Of Stream Driven Sorter