

6 L - 8

コネクションマシン CM-2 における関係データベース処理

松本 和彦 喜連川 優 高木 幹雄

東京大学 生産技術研究所

1 CM-2 アーキテクチャ

コネクションマシンは、数万個というオーダーのプロセッサアレイとハイパーキューブ形状の通信ネットワークを持つ、単一命令ストリーム多重データ (SIMD) の並列アーキテクチャマシンである。現在 Thinking Machines Corporation が市販している最新モデルは CM-2 と呼ばれ、フルセットで 65536 個のプロセッサを備えている。

CM-2 は、ワークステーションなどのフロントエンドにバスインタフェースを介して結合している。逐次的な処理はフロントエンド上で通常のプログラムコードとして実行され、並列処理の部分だけがバスインタフェースを通じてコネクションマシンに指令される。

CM-2 の操作は、PARIS というライブラリの関数を呼び出すことによって行なう。次節で、コネクションマシンのいくつかの基本概念について述べる。

2 PARIS インタフェース

PARIS とは、コネクションマシンシステムをプログラミングするための、PARAllel Instruction Set のことである。それは、フロントエンドからコネクションマシンプロセッサを指令するための低レベルインターフェースである。PARIS の実体は、関数、マクロ、変数の集合である。関数やマクロは、呼び出されることにより、必要に応じてコネクションマシンのシーケンスに命令を送出する。

2.1 仮想プロセッサ

PARIS は、抽象化されたコネクションマシンハードウェアをユーザに見せることによって、スケーラビリティを拡張している。その核心となるのは、各々の物理プロセッサが複数の仮想的なプロセッサをシミュレートするという、仮想プロセッサの機能である。

プログラムは好きなだけの仮想プロセッサ数を仮定してプログラムすれば良く、プログラムのできた後で、仮想プロセッサは物理プロセッサに割り当てられる。こうすることによって、ソフトウェアが物理的なハードウェアによって変更を強いられることはなくなる。そこに残るのは、単に、物理プロセッサ数と実行時間やメモリ容量との間のトレードオフである。また、1つの物理プロセッサがいくつかの仮想プロセッサをシミュレートするかという数のことを VP 比と呼ぶ。

2.2 VP セット

あるデータ集合が割り当てられている仮想プロセッサの集合のことを、仮想プロセッサセット (VP セット) と呼ぶ。仕事によっては、複数のデータ集合を扱うことが必要となるので、PARIS は複数の VP セットが同時に存在することを許している。例えば、ドキュメント検索のプログラムでは、ある時点では記事 (数千バイト) をデータ要素として扱い、またある時点では記事中の単語 (数十バイト) をデータ要素として扱う。このプログラムでは、それぞれに対応した 2 つの VP セットを使う。単語をデータ要素とする VP セットは、記事をデータ要素とする VP セットよりはるかに大きくなる。VP セットは PARIS への関数呼び出しによって生成される。その大きさ (VP

の数) は生成された時点で決まり、以後変化しない。VP の数は物理プロセッサの倍数でなければならない。

2.3 フィールド

メモリは、フィールドという単位で扱われる。フィールドとは、メモリ中に割り当てられた連続するビット群のことである。フィールドの長さは任意で、割り当てられる時にユーザが指定する。フィールドは、必ずどれか一つの VP セットに所属する。フィールドを VP セットに割り当てるといのは、VP セット中の全ての VP について 1 つずつ、メモリ中の同じ位置に記憶を割り当てるといことに相当する。

2.4 コンテキストフラグ

PARIS の仮想プロセッサの持ついくつかの 1 ビットフラグのうち、特に重要なのはコンテキストフラグである。コンテキストフラグが立っている仮想プロセッサはアクティブであると言われ、ほとんどの PARIS インストラクションはアクティブなプロセッサでのみ実行される。

3 コネクションマシン上でのジョイン処理

コネクションマシン上のアプリケーションには、有限要素解析等の数値計算、コンピュータビジョン、意味ネットワークなど、様々な分野のものが考えられているが、非常に多くのデータアイテムに似通ったオペレーションを行なうデータベース処理は、特に適したアプリケーションであると思われる。

CM-2 上で実際にジョイン処理を行なうプログラムを作成し、その性能評価を行なった。ジョイン処理 (結合演算) は、関係データベース処理の中心的なオペレーションであり、その手法には様々なものがあるが、今回コネクションマシン上にインプリメントしたものは、マージに基づくアルゴリズムである。

4 アルゴリズム

ジョインの手順を以下に示す。
 1. 2 つの VP セット A、B を作り、それぞれに元となるリレーションをロードする (図 1)。
 2. scan-with-add オペレーションを使い、それぞれのプロセッサが、自分と同じキーを持つプロセッサ (タプル) が自分よりも上方にいくつあるのかを調べる (図 2)。

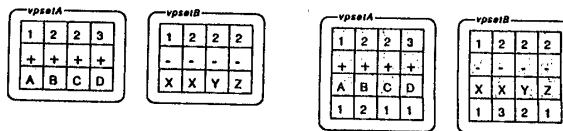


図 1

図 2

3. それぞれのキーの先頭にあたるプロセッサを駆動し、キーの値を宛先アドレスとして自分と同じキーのタプルの数を新しい VP セット L に送信する (図 3)。

⁰Relational Database Processing on the Connection Machine CM-2
 K.Matsumoto, M.Kitsuregawa, M.Takagi
 The Institute of Industrial Science, University of Tokyo

4. Lでは、メッセージを送ってきたタブルがマージによって移動すべきアドレスを計算し、そのアドレスをそれぞれのプロセッサに送り返す(図4)。

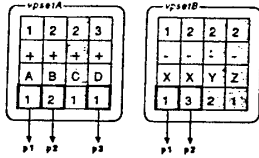


図3

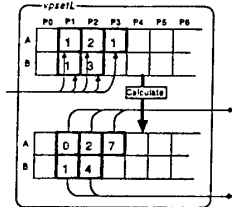


図4

5. AとBでは、返ってきたアドレスをもとに、それぞれのキーの先頭以外のタブルが移動すべきアドレスをセットアップする。

6. AとBの全てのタブルが、新しいVPセットMの中の、移動すべきプロセッサへと送信される。この時点でマージは完了する。

7. それぞれのキーにおいて、Aから送られてきたタブル数とBから送られてきたタブル数を数え、それを掛け合わせる(図5)。

8. 7.で求めた数を使い、scan-with-addオペレーションによって、ジョイン後にそれぞれのキーの先頭タブルが位置するアドレスを計算する(図6)。

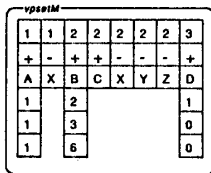


図5

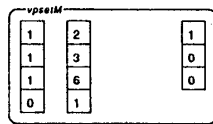


図6

9. 8.の計算結果をもとに、キーと属性を新しいVPセットJの適切なアドレスに送信する(図7)。

10. scan-with-copyで、必要なだけのデータを複製する(図8)。

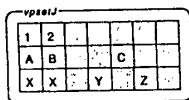


図7

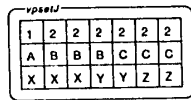


図8

11. 一番下のフィールドで置換を行ない、ジョイン後のリレーションをJ内に完成させる(図9)。

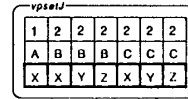


図9

5 性能評価

選択率100%のジョインの場合の性能測定の結果を、表1に示す。タブル長はキー32ビットと属性32ビットの計64ビットである。なお、この時間にはリレーションの入出力の時間は入っておらず、純粋にプロセッサがマージとジョインを行なっている時間のみが含まれる。

タブル数	処理時間	VP比	タブル数	処理時間	VP比
100	56.5	1	10000	267.5	4
200	67.7	1	12000	215.0	4
500	81.8	1	14000	227.7	4
1000	71.5	1	16000	244.0	4
2000	74.0	1	20000	334.5	8
3000	80.3	1	24000	339.2	8
4000	97.0	1	28000	361.2	8
5000	123.7	2	32000	368.7	8
6000	121.2	2			
7000	145.5	2			
8000	145.3	2			

表1: ジョイン処理の時間 (msec)

6 終りに

高並列微粒度マシンCM-2にメインメモリデータベースを実現し、その基本性能を測定した。演算は全て主記憶上でなされており、入出力は含まれていない。高速ディスクDataVaultを使ったより大規模なデータベース処理が、今後の課題である。

参考文献

[1] W.Daniel Hillis, "The Connection Machine", The MIT Press, 1985
 [2] Thinking Machines Corporation, "Connection Machine Technical Summary", 1989