

コネクションマシン CM-2 による大規模関係データベース処理とその評価

岡田 英明 * 松本 和彦 喜連川 優

東京大学 生産技術研究所

{okada,matsu,kitsure}@tkl.iis.u-tokyo.ac.jp

概要

無共有型並列計算機による関係データベース処理性能向上への試みがこれまでに数多くなされている。一方、コネクションマシンをはじめとするデータパラレルマシンと呼ばれる SIMD 型の並列計算機を関係データベース処理に利用するという研究は、今のところほとんど報告されていない。これまでに、我々は主記憶上での関係データベース処理アルゴリズムを提案、実装してきた。本論文では外部記憶装置を含めた大規模関係データベース処理をデータパラレル記法により記述可能であることを示すとともに CM-2 に実装した性能評価に関し報告する。

Large relational database processing and its evaluation on the Connection Machine CM-2

Hideaki OKADA Kazuhiko MATSUMOTO Masaru KITSUREGAWA
Institute of Industrial Science, University of Tokyo

Abstract

Many attempts have been made to accelerate the performance of the relational database processing by utilizing shared nothing parallel machines. On the other hand, few researches that utilize SIMD type machines called data parallel machines such as Connection Machine for the relational database processing. So far we proposed and implemented relational database processing algorithms on the CM-2. In this paper we show data parallel model can describe relational database processing. Then we implement large relational database processing on the Connection Machine CM-2 with secondary storage and evaluate its performance.

*現在、三菱電機情報システム研究所

1 はじめに

これまでに、関係データベースシステムの性能向上のために数多くの研究がなされてきた。中でも GAMMA[1, 2], SDC[4]などのような shared nothing の並列データベースシステムは最も有効なアプローチの一つと考えられており開発、評価等が進められている。また、NCUBE, Sequentなどの汎用並列計算機上に関係データベースシステムを構築して性能向上を達成しているものもある。

一方でコネクションマシンをはじめとするデータバラレルマシンとよばれる、制御レベルではなくデータレベルの並列性を利用する並列計算機が注目されている。ここでいうデータバラレルとは並列計算機アーキテクチャについて述べているのではなく、計算モデルの意味で述べている。データバラレルモデルはプロセッサ数の増加などのシステムの拡張に容易に対応できるという大きな利点を持っている。

現在のところ、このようなデータバラレルマシンの応用分野はデータ並列性を有する大規模数値計算や画像処理が主であるが、データ並列性を有し、大容量データを扱う関係データベース処理も有効な応用分野となりうるのではないかと考えられる。

しかし、我々の知る限りではデータバラレルマシンによる関係データベース処理の研究はほとんど報告されていない。そこで我々はデータバラレルマシンの関係データベース処理への有効性を確認するため、まず主記憶上でのジョイン処理アルゴリズムを提案、実装した[5]。本論文では、外部記憶装置を含めたシステムの性能を概観するために、そのアルゴリズムを拡張し、主記憶上では処理しきれない大規模なジョイン処理について述べる。

第2節においてコネクションマシンシステムについて、第3節では主記憶内のジョイン処理アルゴリズムについて概観する。第4節において大規模データベース処理アルゴリズムについて説明し、その性能測定結果を第5節で述べる。最後に第6節で結論を述べる。

2 コネクションマシンシステム

コネクションマシン CM-2 コネクションマシン CM-2[6] は SIMD 型の超並列計算機である。CM-2 は 1bit の PE を 16PE 単位にハイパーキューブ網で結合しており、ローカル、グローバル通信ともにこれを用いている。各 PE はそれぞれローカルメモリを持ち、最大構成時では 65,536PE で構成される。各 PE の命令はフロントエンドからシーケンサと呼ばれるハードウェアを通してブロードキャストされる。各 PE はその命令を内部状態より選択して実行し、SIMD 動作を行う。

仮想マシンアーキテクチャ CM-2 のプログラミングでは PARIS とよばれる並列命令セットが使用可能である。PARIS は仮想プロセッサ (VP) の概念を持つ。仮想プロセッサ数と物理プロセッサ数の比を仮想プロセッサ比 (VP 比) と呼ぶ。各物理プロセッサは VP 比数分の仮想プロセッサの命令を実行する。したがって、プログラム実行時にはシステムの物理プロセッサ数とプログラムの仮想プロセッサ数に応じて動的に仮想プロセッサ比が決定される。PARIS の仮想プロセッサの機能は、1 つのデータ要素を 1 つの仮想プロセッサに割り当てる capability を可能にする。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。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。複数の VP セットを使用することもできるが、制御対象は常に 1 つの VP セットであり、VP の数は物理プロセッサ数の 2 の累乗倍の数でなければならない。メモリ割り当てなどもこの VP セットに対する命令によって行われる。

データポート データポートはコネクションマシンの大容量二次記憶システムである、RAID-2 レベルのディスクアレイである。その容量は単体で 10G bytes であり、20Gbytes まで拡張可能である。転送速度は 25 MB/sec を維持でき、複数のデータポートを並列に使うと合計 100MB/sec を越える転送速度を実現できる。データポートのファイルにも仮想プロセッサの概念があり、コネクションマシン上の主記憶とのデータのやり取りは同じ VP 数の VP セットに対してのみ可能である。

3 主記憶内ジョイン処理アルゴリズム

これまでにデータバラレルマシンの関係データベース処理応用の第歩として主記憶内にリレーションが収まる時のジョイン処理アルゴリズムについて、ソートマージ法を用いたアルゴリズムとハッシュによる分割を用いたアルゴリズムを提案し、コネクションマシン CM-2 上に実装、評価している。本論文ではソートマージ法を用いたアルゴリズム DPSMJ (Data Parallel Sort Merge Join) を使用した。

リレーションの仮想プロセッサへのマッピングは、1 タブルが 1 仮想プロセッサに対応するものとする。

DPSMJ アルゴリズムは大きく三つのステップに分けられ、それぞれソートフェイズ、マージジョインフェイズ、置換フェイズと呼ぶ。それぞれのフェイズにおける動作を図 1, 2, 3 に示す。

ソートフェイズ

最初に、リレーション R と S が併合されて 1 つのリレーション M になる。リレーション M は最終的なリレーション J を得るために準備として一時的に作られるだけのものであり、それ自体では関係データベースのリレーションとしての意味は持たない。したがって、M はキー属性とアドレス、リレーションの ID (R が 0、S が 1 とする) のみからなるタブルとして作成する。次にリレーション M の中でタブルがキー属性に従って昇順にソートされる。同じキー属性を持つタブルについては、リレーション ID を利用して R からのタブルの後に S からのタブルが並ぶという順序にする。

マージジョインフェイズ

リレーション M の上で必要な処理を行い、ジョイン後のリレーション J を生成する。最初に、リレーション J のための空間をアロケートして必要なデータをそこに転送する。このフェイズはデータ送信とデータ複製からなる。

M に属する全てのタブルは、自分と同じキー属性を持つタブルの数をかぞえる。M のそれぞれのタブルは、自分と同じキー属性を持つタブルが、J の中にいくつ存在することになるかを計算し、スキャン加算して J のタブル数を得る。

同じキー属性をもつタブルの中で最初のタブルが、生成される J の数に応じたアドレスにデータを送信し、J の中に必要な数だけそのデータを複製する。

置換フェイズ

マージジョインフェイズが終了した段階では、まだデータの組合せが適切なものではない。求めるリレーションを得るために、S のタブルから送られてきたデータを再配置する。このようにして、最終的に望むリレーションを、J の中に作ることができる。

4 データボルトを用いた大規模ジョイン処理アルゴリズム

4.1 アルゴリズムの概要

各リレーションを主記憶内で処理できるほどの大きさに分割する。これをスプリットフェイズと呼び、分割されたリレーションのそれぞれをパケットと呼ぶこととする。処理対象であるリレーションをパケットに分割する方法としては、ハッシュによる分割法を採用する。キー属性値を入力としてハッシュ関数を適用し、同じハッシュ関数値をもつタブルの集合が 1 つのパケットを構成する。したがって、同じハッシュ関数値を

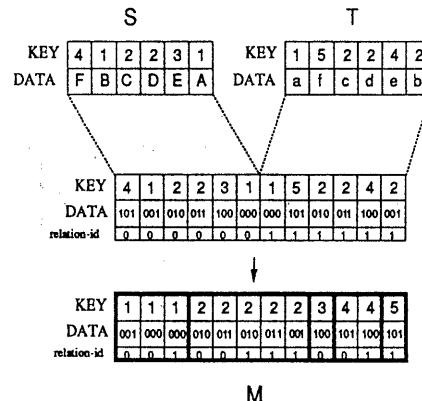


図 1: ソートフェイズ

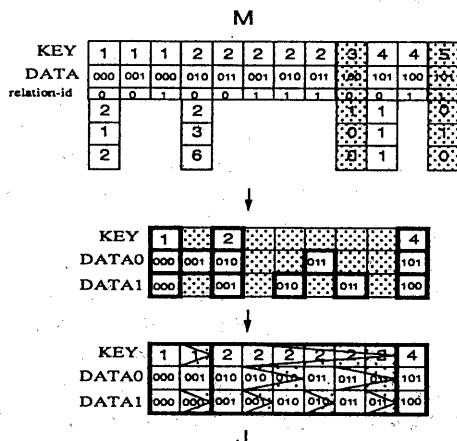


図 2: マージジョインフェイズ

KEY	1	1	2	2	2	2	2	4
DATA0	000	001	010	010	010	011	011	011
DATA1	000	000	001	001	010	010	011	100

KEY	1	1	2	2	2	2	2	4
DATA0	A	B	C	C	C	D	D	F
DATA1	a	a	b	c	d	b	c	e

↓

J

図 3: 置換フェイズ

もつバケット同士の主記憶内でのジョイン処理を全バケットについて行うことにより、全体のジョイン処理を完了する。これをジョインフェイズと呼ぶこととする。

ここでは、ジョイン処理されるリレーションを R 、 S とする。スプリットフェイズにおいて N 個のバケットに分割されるとき、リレーション R 、 S から生成された i 番目のバケットを R_i 、 S_i とする。リレーション R 、 S はそれぞれ、長さ l_R 、 l_S のタブルからなり、 C_R 、 C_S 個のタブルを持つものとする。また、キー属性の長さを l_{key} と表す。

はじめに、リレーション R 、 S はデータボルトにファイルとして書き込まれている。このデータボルト上のファイルにもコネクションマシンの主記憶と同じく VP 数という概念があり、ある VP 数である VP セットにより書き込まれたファイルは同じ VP 数の VP セットによってしか読み込むことができない。したがって、もしデータボルト上のファイルと異なる VP セットを読み書きする場合には、データボルト上の VP セットに等しい VP セットを確保しておき、それを I/O バッファとして用いる。例えば、読み込み時には I/O バッファに読み込んだリレーションを send 命令を用いて、そのリレーションを必要としている VP セットに送る必要がある。このデータボルト上の VP セットの VP 比をリレーション R 、 S についてそれぞれ V_R 、 V_S とする。また、物理プロセッサ数を P とする。

以下に、ジョインアルゴリズムの詳細について述べる。

スプリットフェイズ

スプリットフェイズでは、リレーションはキー属性値のハッシュ関数値に応じたバケットに分割され、データボルトに書き込まれる。リレーション R 、 S についてそれぞれ独立に同じ処理を行なうので、以下の説明ではリレーション R のみについて説明する。

スプリットフェイズでは、スプリットバッファ、バケットバッファという 2 種類のバッファが少なくとも必要となる。スプリットバッファは読み込んだリレーションのキー属性に対してハッシュ関数値を計算して、各タブルがどのバケットに属するかを決定する。バケットバッファには、同じハッシュ関数値を持つタブルがスプリットバッファより送られる。それらをデータボルトに書き込む。

このスプリットバッファとバケットバッファの割り当て方法にはに次の 2 つ方法が考えられる。1 つは図 4 に示すように、巨大なスプリットバッファとバケットバッファを 1 つだけ用意して共用する。もう一方は、図 5 に示すようにファイルの VP セットと同じ VP セットであるスプリットバッファを用意して、バケットバッファは分割数と同じだけ用意する。

ここでは、次のような理由から前者の方法を採用した。

ハッシュ関数の値に応じてタブルをバケットバッファに送る時に、送信先になれるバッファは 1 つしか存在しない。すなわち、後者の方法を探って複数のスプリットバッファを小さくとった時には、分割数に応じたバケットバッファが存在していることになるが、スプリットバッファからすべてのバケットバッファへ同時にタブルを送ることは、コネクションマシンの VP セットの制約により許されていない。したがって、バケットバッファを多く確保することによる大きな利点を活かせないからである。

前者の方法を選択すると、スプリットフェイズではリレーション R を読み込むための I/O バッファを VP 比 V_R で確保し、 V_R の倍数でありかつ、できるだけ大きい VP 比 V_{split}^R でスプリットバッファを確保する。バケットバッファに送られるタブル数は、平均的にはスプリットバッファに入るタブル数の分割数 N 分の 1 となるので $\frac{V_{split}^R}{N}$ 、もしくはキー属性値の分布を考慮して $\frac{V_{split}^R}{N}$ より大きい VP 比 V_{bucket} のバッファを確保する。VP 比と主記憶割り当てる関係や分割数の決定方法については後に触れる。

以上をまとめるとスプリットフェイズでは、以下の処理を $\frac{C_R}{PV_{split}^R}$ 回行なうことになる。

1. 次の処理を最大で $\frac{V_{split}^R}{V_R}$ 回行なう
 - (a) データボルトから入力バッファに VP 比 V_R

でタプルを読み込む

- (b) 入力バッファからスプリットバッファへ send 命令によりタプルを送る
2. スプリットバッファにおいてハッシュ関数値を計算する
 3. 次の処理を i が 0 から $N - 1$ のときについて N 回行なう
 - (a) ハッシュ関数値が i であるバケットを数え、バケットバッファに送る
 - (b) バケットバッファに送られてきたタプルをデータポルトに R_i バケットのファイルとして書き込む

データ分布の不均一性などにより 3(a) 段階において、バケットバッファの溢れが生じる場合がある。そこで、溢れが生じる場合にはバケットバッファに入る切るタプルは出力し、入り切らないものはスプリットバッファに保存しておく。リレーション読み込み時には I/O バッファからスプリットバッファの空いている領域へタプルを転送する。

ジョインフェイズ

ジョインフェイズでは、図 6 に示すようにスプリットフェイズで生成されたバケット R_i, S_i をジョインバッファと呼ばれる 2 つのバッファにそれぞれ読み込んで主記憶内でのジョイン処理を行ない、結果をデータポルトに出力することを分割数の数だけ繰り返す。主記憶内でのジョイン処理アルゴリズムには、DPSMJ アルゴリズムを使用する。

スプリットフェイズにおいて V_{bucket}^R の VP 比でデータポルト上に書き込まれているバケットを読み込むには、 V_{bucket}^R の VP 比を持つバッファを用意しなければならない。そのバケットに対する書き込みが 1 度であれば、そのバッファがジョインバッファとなることができるが、複数回の書き込みを行なう必要があれば、そのバッファは単なる入力バッファとして使用され、バケットのサイズに応じた VP 比 V_{join}^R でジョインバッファを新たに確保する必要がある。

スプリットバッファよりもリレーションのサイズが大きい場合には、バケットバッファの書き込みが複数回行なわれる。したがって多くの場合、このような入力バッファが必要となる。

ジョインフェイズでは以下のようないくつかの処理を N 回行なう。

1. バケット R_i をすべて読み込むまで、次の処理を繰り返す
 - (a) バケットを入力バッファに読み込む

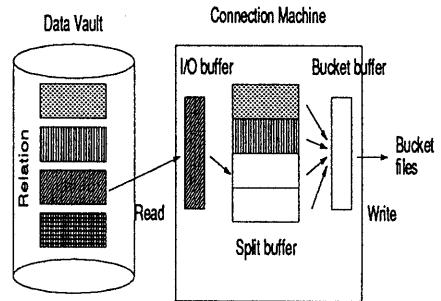


図 4: スプリットフェイズバッファ割り当て法 1

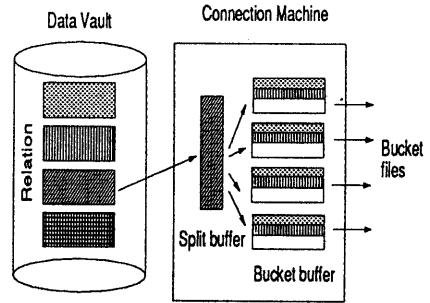


図 5: スプリットフェイズバッファ割り当て法 2

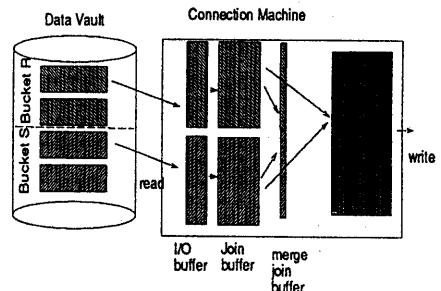


図 6: ジョインフェイズバッファ割り当て法

- (b) 入力バッファからジョインバッファへ転送する
2. パケット S_i についても同様の処理を行なう
 3. リレーション R,S から生成されたパケットファイルの VP 比を V_{join}^R, V_{join}^S とすると、キー属性とアドレス、リレーション ID からなる、
- $$V_{join}^R + V_{join}^S \leq V_{merge}$$
- をみたす VP 比 V_{merge} のリレーション M を用意し、 R_i, S_i のキー属性、アドレス、リレーション ID を M に送る。
4. リレーション M において DPSMJ を部分的に利用してソートを行なう
 5. DPSMJ により生成される結果リレーションのタプル数を求め、その数に応じた VP 比 V_{result} で長さが $l_S + l_R - l_{key}$ のバッファ J を確保する。確保できない時は、利用可能な主記憶量だけ確保して、以下の処理を分割して行なう。
 6. リレーション M の内容に応じて、 R_i, S_i からデータを J に送り、J 内で完全な結果リレーションの生成を行なう
 7. 結果リレーションをデータボルトに出力する

4.2 主記憶の割り当て

このデータボルトを用いたジョイン処理アルゴリズムにおける各バッファの主記憶上の割り当てについて考察する。主記憶量を M 、リレーション R, S のサイズを M_R, M_S と表す。

4.2.1 スプリットフェイズ

スプリットフェイズでは、入力バッファを VP 比 V_R で、スプリットバッファを VP 比 V_{split}^R で割り当てる。パケットバッファは、キー属性値の分布によるパケットの溢れを考慮した変数 α を用いて表すと、VP 比

$$V_{bucket}^R = \alpha \frac{V_{split}^R}{N} \quad (1)$$

で割り当てる。 α はパケットの溢れを大きく考えるほど 1 より大きくなっていく。

したがって、入力バッファ、スプリットバッファ、パケットバッファを確保するのに必要な主記憶量 $M_{split-inbuf}^R, M_{split}^R, M_{bucket}^R$ はそれぞれ、

$$M_{split-inbuf}^R = l_R P V_R \quad (2)$$

$$M_{split}^R = l_R P V_{split}^R \quad (3)$$

$$M_{bucket}^R = l_R P V_{bucket}^R = \alpha \frac{V_{split}^R}{N} l_R P \quad (4)$$

となる。これらが、主記憶内に収まるには

$$M_{split-inbuf}^R + M_{split}^R + M_{bucket}^R < M \quad (5)$$

を満たせばよい。

$M_{split-inbuf}^R$ は処理するリレーションによって決まっており、 V_{bucket}^R と V_{split}^R には、式(1)の関係が成立するため、

$$M_{split}^R (1 + \frac{\alpha}{N}) < M - M_{split-inbuf} \quad (6)$$

を満たす M_{split}^R が割り当てる可能である。分割数 N はジョインフェイズに対する考察により決定されるが、 $\alpha \leq N$ であるから、

$$M_{split}^R < \frac{M - M_{split-inbuf}}{2} \quad (7)$$

である M_{split}^R は割り当てる可能である。したがって、 M_{split}^R は、 $\frac{M}{2}$ 程度には割り当てる可能であるといえる。

4.2.2 ジョインフェイズ

ジョインフェイズでは、 R_i, S_i 用の入力バッファ、ジョインバッファ、マージジョインバッファ、結果出力バッファが必要となる。入力バッファはそれぞれ VP 比 $V_{bucket}^R, V_{bucket}^S$ で割り当てる。

ジョインバッファは、VP 比 V_{join}^R で割り当てる。

マージジョインバッファは、 $V_{join}^R + V_{join}^S \leq V_{merge}$ を満たす VP 比 V_{merge} で割り当てる必要があるが、VP 比は 2 の幂乗であるので

$$V_{merge} = 2 \max(V_{join}^R, V_{join}^S) \quad (8)$$

で割り当てる。

結果出力バッファの VP 比 V_{result} については、選択率を考慮する必要がある。しかし、ジョイン処理時の選択率は動的に決定され、結果リレーションを生成するバッファも動的に確保するが、このバッファをどの程度の大きさにするかを決める必要がある。そこで、仮の選択率とも言える β を

$$\beta = \frac{V_{result}}{V_{merge}} \quad (9)$$

と定義する。100% ジョインは $\beta = 1$ のときに対応する。実際に生成されるタプル数が PV_{result} を越えた場合には、結果リレーションの生成、出力過程を分割して行なう。

したがって、入力バッファ、ジョインバッファ、マージジョインバッファ、結果出力バッファに必要となる主記憶量 $M_{join-inbuf}^R, M_{join-inbuf}^S, M_{join}^R, M_{join}^S$ 、

M_{merge} , M_{result} は、アドレスを表すのに必要な領域の長さを l_{adr_s} と表すと、

$$M_{join-inbuf}^R = l_R P V_{bucket}^R \quad (10)$$

$$M_{join}^R = l_R P V_{join}^R \quad (11)$$

$$M_{merge} = 2(l_{key} + l_{adr_s}) P \max(V_{join}^R, V_{join}^S) \quad (12)$$

$$M_{result} = 2(l_R + l_S - l_{key}) P \max(V_{join}^R, V_{join}^S) \beta \quad (13)$$

となり、

$$\begin{aligned} M_{join-inbuf}^R + M_{join-inbuf}^S + M_{join}^R \\ + M_{join}^S + M_{merge} + M_{result} < M \end{aligned} \quad (14)$$

を満たすようにすればよい。

$l_{key}, l_{adr_s} \ll l_R, l_S$ と仮定すると、

$$\begin{aligned} M_{join-inbuf}^R + M_{join-inbuf}^S + M_{join}^R \\ + M_{join}^S + M_{result} < M \end{aligned} \quad (15)$$

とできる。

ここで、

$$N_{split}^R = \frac{V^R}{V_{split}^R} \geq 1 \quad (16)$$

である β を定義する。 β は、キー属性値の分布が等しい場合にはスプリットフェイズにおいてパケットに対して書き込みがあった回数、すなわち V_{join} が V_{bucket} の何倍になるかを示している。このとき、

$$M_{join}^R = N_{split}^R M_{join-inbuf} \quad (17)$$

となる。

したがって、式(15)は

$$M_{join}^R (1 + \beta + \frac{1}{N_{split}^R}) + M_{join}^S (1 + \beta + \frac{1}{N_{split}^S}) < M \quad (18)$$

とできる。

M_{join}^R, M_{join}^S より一度に主記憶内でジョインできるタプル数が確定する。そして、分割数が、

$$N = \max(\frac{M_R}{M_{join}^R}, \frac{M_R}{M_{join}^S}) \quad (19)$$

ときまる。

5 性能評価

以上のようなデータボルトを用いたジョイン処理アルゴリズムをコネクションマシン CM-2 に実装した。ここでは、その性能評価、解析を行なう。

使用したのは 8K プロセッサ、ローカルメモリ 8KB、全体では 64MB の主記憶を持つ CM-2 である。この測

定では、タブルのキーは、1 からタブル数までの値を重複なく設定している。すなわち、100% ジョインになる。測定時間は、コネクションマシンの処理時間のみであり、フロントエンドの処理時間は含まない。測定に当たり、キー長は 4 バイト固定とし、タブル長は 64,128,256, 512 バイトの 4 種類の値をとるものとする。両方のリレーションは同じタブル長、同じタブル数であるとし、リレーションのサイズを 8,16,32,64MB と変化させた。タブル数は、タブル長とリレーションのサイズに応じて決定され、16K から 1024K の値をとる。

主記憶の割り当てについては、 $\alpha = 1, \beta = 1, N_{split} = 1 or 2$ として行なった。スプリットフェイズでは 32MB をスプリットバッファに、パケットバッファ、入力用バッファは 4MB とした。ジョインフェイズにおいては、入力用バッファ、ジョインバッファはリレーション毎に 4MB とし、出力バッファは 8MB とする。これらは第 4 節において導いた式から得られる値と比べると小さくなっているが、これはユーザには見えてこないメモリ領域、例えばソートルーチンでの使用、などにより M が主記憶量より小さくなっているためである。

図 7,8,9 にそれぞれ、全体の処理時間、スプリットフェイズの処理時間、ジョインフェイズの処理時間を示す。また、リレーションサイズが 32MB のときのスプリットフェイズ、ジョインフェイズの処理時間の内分けを図 10,11 に示す。

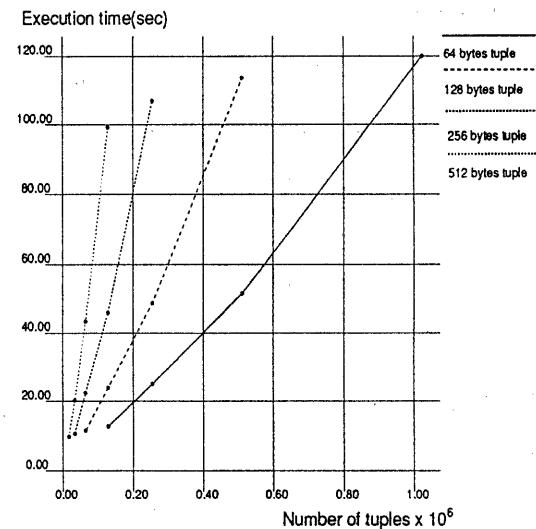


図 7: データボルトを用いたジョイン処理プログラム全体の処理時間

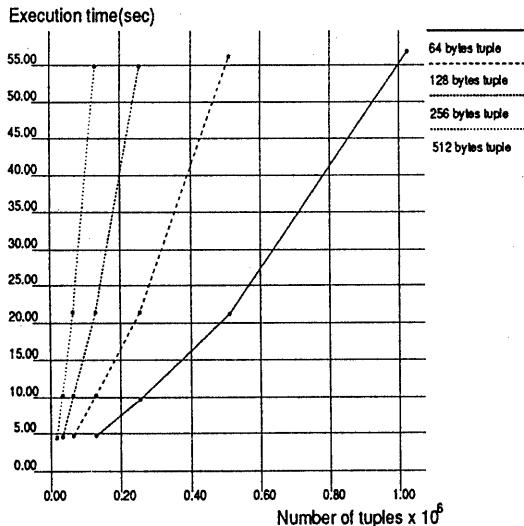


図 8: スプリットフェイズの処理時間

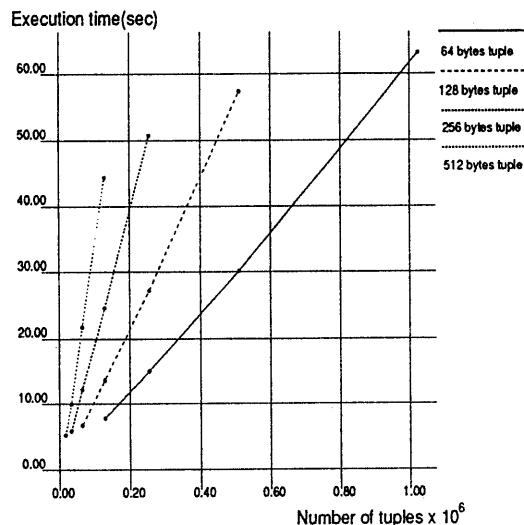


図 9: ジョインフェイズの処理時間

“DISK I/O”はデータboltとの入出力時間を表す。スプリットフェイズにおける、“send to bucket”はスプリットバッファからバケットバッファへのデータ転送時間、“send to split buffer”は入力バッファからスプリットバッファへの転送時間、“hash”はハッシュ関数値を求める時間を表す。また、ジョインフェイズにおける“sort”はキー属性のソート時間、“make result”はキー属性のソート結果に基づいて結果のリレーションを生成する時間、“send to join buffer”は入力バッファからジョインバッファへのデータ転送時間、“merge join”はソート処理と結果タプルの生成以外のDPSMJの時間をそれぞれ表している。

スプリットフェイズでは、処理時間の7割以上がデータboltとの入出力によって占められていることがわかる。ジョインフェイズではタプル数に依存してソートの処理時間によって割合は変わるもの処理時間の約半分がデータbolt入出力に費やされている。データboltとの入出力性能はタプル数ではなく容量に依存している。

このジョイン処理でのデータboltとのI/O性能を転送速度で表すとスプリットフェイズで6MB/s程度、ジョインフェイズで8-10MB/sと、最大転送速度の1/2から1/4しか引き出せていないことが分かる。これは、入出力のアクセス単位が4-8MBとあまり大きくないことが原因として考えられる。ジョインフェイズの方が転送速度が大きいのはそのためであり、結果リレーションの出力時のアクセス単位が最も大きいからである。

また、リレーションのサイズが32MBを越えると性能が落ちていくことがわかる。これらは、次のような理由による。

リレーションのサイズがスプリットバッファのサイズより大きい場合に、我々はパケットバッファをディスクに書き込む際にVP比を小さくして書き込んでいる。これは第4.1節でも説明したように、データboltを用いたコネクションマシンのファイルシステムはVPセットそのものを書き込むために、パケットバッファに空きが多いと入出力性能に影響を及ぼすからである。また、中身のつまつたパケットを読み込んでジョインバッファに送る時には、send命令ではなく、次に述べるような分散した、負荷の小さい処理で置き換えられ、データ移動の時間が短縮されるからである。したがって、64MBのリレーションを処理する時にはスプリットバッファが32MBで分割数が16であるから、平均的なパケットのサイズは2MBとなるため4MBではなく、半分の2MBで処理している。図8において、32MBリレーションに対する処理のデータを境にして直線の傾きが変化しているのは、2MBというデータboltにとおな書き込み単位で出力している

ため、転送速度が小さくなっているからである。

スプリットフェイズのバッファの割り当て方法について、我々は図4に示した方法を選択をした。その結果、小さなVP比のI/Oバッファから大きなVP比のスプリットバッファへのデータの移動が必要になる。このようなデータの移動はスプリットフェーズのバッファの割り当て方法の選択理由で述べたように、データの移動先の物理プロセッサの偏りが発生することになる。このデータの移動は本来は、物理プロセッサのメモリ上でのデータの移動にすぎないことから、我々は、移動先の物理プロセッサが分散するよう`send`命令の移動先のアドレスを変更している。

次に、処理時間のディスクI/O以外の時間に注目してみると、VPセット間でのデータ転送の時間が大きいことがわかる。スプリットフェイズでは24%が、ジョインフェイズでは15%がそれらの時間に割かれている。スプリットフェイズでのデータ転送とは、スプリットバッファからバケットバッファへの転送が大部分を占めている。残りは、入力用バッファからスプリットバッファへの転送時間である。ここにも、ファイルにおけるVPセットの弊害が現れている。計算が主であるハッシングの部分は非常に小さい。ジョインでフェイズも、結果のリレーションを得るための`send`処理と、入出力用バッファのデータ転送時間が大きな時間を占めている。

全体の処理の中で、ディスクやプロセッサ間でのデータのやりとりを除いた処理は、スプリットフェイズのハッシングとジョインフェイズのソートである。ハッシングの占める割合は小さい。これはその処理の内容から当然といえる。ここで、ソートの時間が大きいといえる。しかも、このソート処理では、ソートキーとアドレスのみを扱っており、タブルのデータ部分の転送時間は含まれていない。このことからも、ソート処理の負荷の大きさが良くわかる。これに関しては、我々はソート自身に関する研究を進めている[10]。

6 おわりに

これまでにデータパラレルマシンの関係データベース処理への有効性に関する研究はほとんどなされていなかった。そこで我々は外部記憶装置とのI/Oをも含めたデータパラレルジョインアルゴリズムをデータポートと呼ばれるディスクアレイを備えたコネクションマシンCM-2に実装し、その性能について評価した。

まず、データパラレル記法により外部記憶装置を含めた関係データベース処理が可能であることを明らかにした。

現在のシステムでは充分な性能向上は必ずしも得られず、仮想プロセッサ機能によるメモリ、通信、ファ

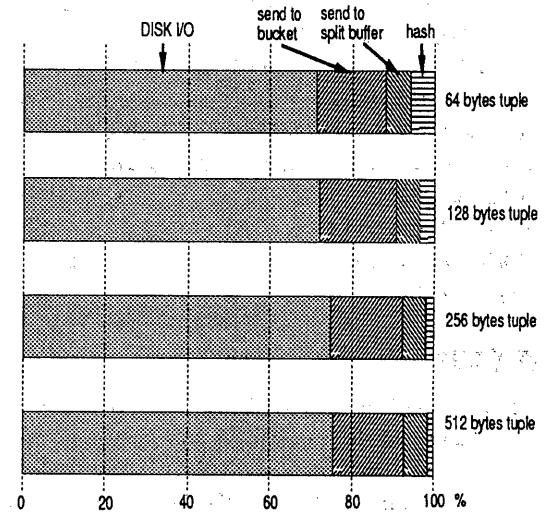


図 10: スプリットフェイズの処理時間の内分け

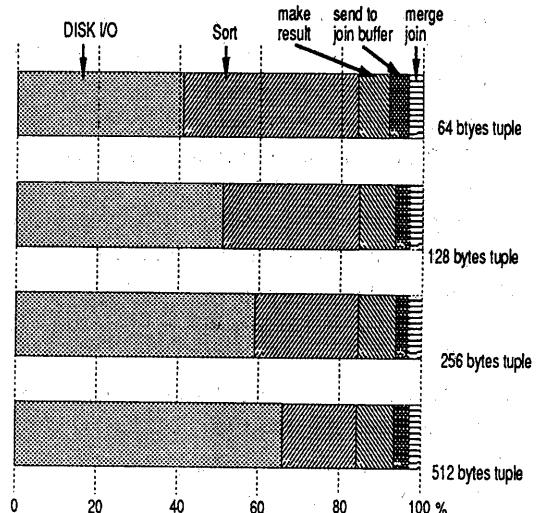


図 11: ジョイン各フェイズの処理時間の内分け

イルに対する制限が大きいことが分かった。

本論文ではコネクションマシン CM-2 への実装を扱ったが、アルゴリズム自身はデータパラレルマシンに汎用的に利用でき、仮想プロセッサに対する主記憶割り当てに関する形式化も有効である。

本論文のプログラムはまだ第一段階にすぎない。データポートにおいては同期 I/O の他に、非同期 I/O の機能を備えている。この機能を用いた性能向上方法やジョイン処理における主記憶割り当て方法についてさらに詳細に検討し、可変長レコードなどの複雑な条件への対応などについて検討したい。また、VP の制限が緩和されていると言われている CM-5 への実装なども行っていきたい。

参考文献

- [1] D. J. DeWitt, R. H. Katz, F. Olken, L. D. Shapiro, M. R. Stonebraker, and D. Wood: implementation Techniques for Main Memory Database Systems, in ACM SIGMOD '84, pp.1-8, 1984.
- [2] D. J. DeWitt, *et al*: A Performance Evaluation of Four Parallel join Algorithms in a Shared-Nothing Multiprocessor Environment, Proc. of SIGMOD, 1989.
- [3] Teradata Corp.: DBC/1012 Data Base Computer Concepts and Facilities, CO2-0001-05.
- [4] M. Kitsuregawa, *et al*: The Super Database Computer (SDC) System Architecture, Algorithm, and Preliminary Evaluation Proc. of the 25th Hawaii International Conference on System Sciences, pp.308-319, 1992.
- [5] M. Kitsuregawa, K. Matsumoto: Massively Parallel Relational Database Processing on the Connection Machine CM-2, The 2nd International Symposium on Database Systems for Advanced Applications (DASFAA'91), pp.226-235, 1991.
- [6] Thinking Machines Corp.: Connection Machine Model CM-2 Technical Summary, Version 5.1, 1989.
- [7] E. Ozkarahan: Database Machines and Database Management, Prentice-Hall, N.J., U.S.A., 1986.
- [8] S. Y. W. Su: Database Computers, McGraw-Hill, N.Y., U.S.A., 1988.
- [9] M. Kitsuregawa, H. Tanaka, and T. Moto-oka: Application of Hash to Data Base Machine and Its Architecture, New Generation Computing, Vol.1, No.1, pp.66-74, 1983.
- [10] 岡田, 喜連川, 高木: 超並列計算機における各種ソートアルゴリズムの実装と評価, 情報処理学会研究報告 92-ARC-97, 情報処理学会, 1992.