

ジョイン演算のための並列実行システム

A PARALLEL EXECUTION SYSTEM FOR JOIN OPERATION

中山 敬 三工 博司 平川 正人 田中 稔 市川 忠男

Takashi Nakayama Hiroshi Mikami Masahito Hirakawa Minoru Tanaka Tadao Ichikawa

広島大学 工学部

Faculty of Engineering, Hiroshima University

1.はじめに

関係モデル^[1]は、データの表現が理解しやすい、データ独立性が高い、演算が数学的に定義されているなどの優れた点を持つており、一般のデータベースへの適用の他に、より高度な機能を備えたデータベース（知識ベース）システムを構築するための基本モデルとしても期待される。

筆者らは、すでに関係モデルに基づく高機能データベースシステムARESを開発している^[2]。ARESでは、データ間の相互関連にしたがって、入力データに対してセマンティックな関連のあるデータを引きだすことができる。これにより、問合せ条件に関する従来の断定的記述に対して、必要に応じてある程度の幅をもたせた柔軟な解釈が可能となっている。

データ間のセマンティックな関連はリレーションの形でデータベースに蓄えられる。これをセマンティック・リレーションと呼び、各定義域ごとに定義する。すなはち、従来のリレーションについては、コンベンショナル・リレーションと呼んでこれを区別する。セマンティック・リレーションをコンベンショナル・リレーションから独立して管理することによって、セマンティックスの定義をユーザに解放することができ、かつ、これに対するシステムの追従性を保証される。

意味的な関連をもつデータの検索要求が行われた時には、まず、検索条件がかかる属性の定義域で規定されたセマンティック・リレーションに対してセレクション操作を行い、次にそ

の結果とコンベンショナル・リレーションとのジョイン操作が行われる。このように、検索条件にあいまい性が付与された問合せもすなはち、従来の関係代数演算の範囲で実行可能となる。

ところでも関係モデルでは、上記のようないくつかある反面、実際の処理系で通常アーキテクチャ上に構築する場合に関係代数演算の実行効率の低さが大きな問題として現われる。データベースの高機能化に伴ってこの傾向はますます顕著となり、したがって関係代数演算の効率的且実行を支援するための研究が今後重要となる。

本論文では、関係代数演算の中でも特に処理負荷の高いジョイン演算に注目し、これを効率よく処理するためのシステムについて述べる。以後このシステムをジョインプロセッサと呼ぶ。

ジョインプロセッサはマルチプロセッサシステムであり、ジョイン操作を並列に実行する個々のプロセッサ（スレーブユニット）とこれらとのプロセッサを管理するための1台のプロセッサ（マスターユニット）から構成される。

各台のスレーブユニットごとにジョイン操作を実行するにあたっては、ハッシュ技法を用いることによって比較演算の負荷を低減させている。すなはち、マスターユニットは、ジョイン演算を行った属性の値についてハッシングを行い、リレーションの各タップルをそのハッシュ値に対応するスレーブユニットに転送する。

これによって、ジョイン条件が "=" の場合には、それぞれのスレーブユニット内ですべての操作が完了する。それ以外の場合には、さらにリレーションの一方についてスレーブユニット

間でタップルの転送を行い、そのつど各スレーブユニットでジョイン操作を行う。但し、ジョイン条件が" $>$ ", " \geq ", " $<$ ", " \leq "の場合には、ハッシングによってクラスタリングされたタップルの集合の間に順序関係が得られるようはハッシュ関数を用いることにより、転送すべきタップルを制限づけることができる。また、データ転送の方向を一方向に限定することができ、したがってハードウェアとの制御が簡略化される。

以下では、ジョインプロセッサのシステムアーキテクチャ、ならびにアルゴリズムについて詳述する。^[5] さらに、マルチプロセッサ上で実際にシステムの開発を行っており、その実験結果を示す。なお、ジョイン演算の対象としては2つのリレーションを、それをソース・リレーションとターゲット・リレーションと呼ぶ。また、各スレーブユニットに分配されたリレーションについては、上記の区別にしたがって、それをソース・サブリレーションとターゲット・サブリレーションと呼ぶことにする。

2. システムアーキテクチャ

関係代数演算においては、一般にリレーションのサイズが小さくないにしたがってその実行時間は短くなる。したがってリレーションをいくつかのサブリレーションに分割し、それらを複数個のプロセッサで並列に処理することによって全体としての実行時間が短縮できる。特にジョイン演算はリレーションの各タップルごとに独立に処理可能であり、したがって並列処理の導入によて期待される実行速度の改善度は大きい。

図1にジョインプロセッサのアーキテクチャを示す。システムは1台のマスタユニットと共に連結されたn台のスレーブユニットから構成されている。すなはち各スレーブユニット SU_i ($0 \leq i \leq n-1$) は隣接するスレーブユニット $SU_{(i+1) \bmod n}$ と $SU_{(i-1) \bmod n}$ にそれを接続しており、隣り合うスレーブユニットにデータの転送が行えるようになっている。

マスタユニットは、ジョイン条件が記述された属性に対してもハッシュを適用し、そのハッシュ値に対応したスレーブユニットへソース・サ

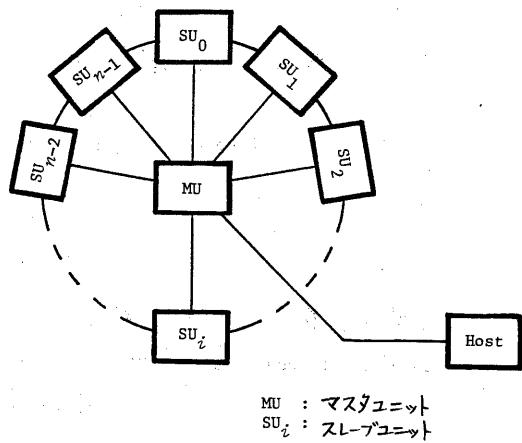


図1 システムアーキテクチャ

ブリレーションならびにターゲット・サブリーションを転送する。さらに、ホストコニピュータとの入出力の管理も行う。スレーブユニットでは、対応づけられた個々のソース・サブリレーションとターゲット・サブリレーションについてジョイン操作が独立に行われる。また、必要に応じてターゲット・サブリレーションはスレーブユニット間で転送され、再びジョイン操作が行われる。

3. アルゴリズム

3.1 リレーションの分配

ソース・リレーションならびにターゲット・リレーションE, それらをいくつかのソース・サブリレーションならびにターゲット・サブリレーションに分割するためにはハッシュを用いる。この場合に適用されるハッシュ関数は次の条件を満たすように決定される。

- (1) 各スレーブユニットにはべく均等な数のタップルが割りあてられる。
- (2) ジョイン条件が指定された属性の値 i, j に対して、それらのハッシュ値をそれぞれ $h(i), h(j)$ とするとき、 $i > j$ ならば必ず $h(i) \geq h(j)$ となる。

但し、(2)はジョイン条件として " $<$ ", " \leq ", " $>$ ", " \geq " が指定された場合にのみ要求される。

図2にハッシュ値とその値をもつタップルが転送されるスレーブユニットとの対応を示す。ジョイン条件が " $=$ ", " \neq ", " $>$ ", " \geq " の場合にはハッシュ値の順序に対応して昇順に " $<$ ", " \leq " の場合には降順に逆ようにスレーブユニットを割りあてる。

ハッシュ関数にこのような条件づけを与えることは次のようす利点をもつ。

(1) データの流れを一方向に限定することができる。したがって制御メカニズムを簡略化することができます。

(2) ジョイン演算の実行にあたって、スレーブユニット間のデータ転送ひらびに比較演算を最小限にすることができる。すばやち。

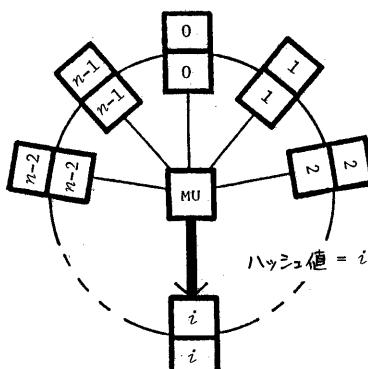
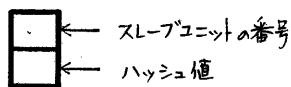
i) ジョイン条件が " $=$ " の場合には、最初にデータが割りあてられたスレーブユニット内の処理だけでジョイン演算が完了するので、スレーブユニット間でデータ転送を行う必要がない。

ii) ジョイン条件が " $>$ ", " \geq ", " $<$ ", " \leq " の場合には、自分より番号の大きいスレーブユニットにだけジョイン条件を満たすデータが存在する。したがって、自分より番号の大きいスレーブユニットにだけデータを転送すればよい。

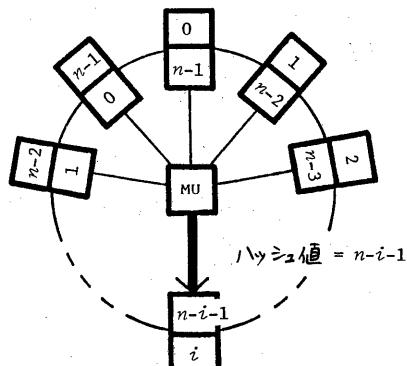
iii) ジョイン条件を必ず満たすデータが存在するスレーブユニットにだけデータを転送する。したがって、スレーブユニット間で転送されたデータに対しては比較演算を行なう必要はない、単に直積をとるだけでよい。

3.2 マスタユニットのアルゴリズム

マスタユニットは、ホストコンピュータからソース・リレーションおよびターゲット・リレーションを入力する。各タップルは、ジョイン条件が記述された属性の値から求められるハッシュ値に対応したスレーブユニットへ逐次転送される。その後、処理の終了したスレーブユニットから結果を受け取り、ホストコンピュータへ出力する。このためのアルゴリズムを以下に示す。



(a) ジョイン条件が $=$, \neq , $>$, \geq の場合



(b) ジョイン条件が $<$, \leq の場合

図2 ハッシュ値とスレーブユニットの対応

```

begin
(*ソース・リレーションの転送*)
while ソース・リレーションの終わりでない
  do
begin
  ポストコンピュータからソース・リレーションを1タップル入力する;
  ジョイン条件が記述された属性の値からハッシュ値を求める;
  ハッシュ値に対応するスレーブユニットへのタップルを転送する;
end;
(*ターゲット・リレーションの転送*)
while ターゲット・リレーションの終わりでない do
begin
  ポストコンピュータからターゲット・リレーションを1タップル入力する;
  ジョイン条件が記述された属性の値からハッシュ値を求める;
  ハッシュ値に対応するスレーブユニットへのタップルを転送する;
end;
(*結果のリレーションの転送*)
while 結果のリレーションの終わりでない
  do
begin
  when スレーブユニットから結果のリレーションが1タップル転送された do
    そのタップルをポストコンピュータへ出力する;
end;
end.

```

3.3 スレーブユニット SU_i ($0 \leq i \leq n-1$) のアルゴリズム

スレーブユニットは、マスタユニットから転送されたソース・サブリレーションとターゲット・サブリレーションとのジョイン操作を行う。ジョイン条件が " $=$ " 以外の場合には、図3、図4に示すように、その後スレーブユニット間でターゲット・サブリレーションの転送を行い、直積操作を行う。二つのためのアルゴリズムを以下に示す。

```

begin
(*ソース・サブリレーションの入力*)
while ソース・サブリレーションの終わりでない
  do
begin
  when マスタユニットからソース・リレーションが1タップル転送された do
    そのタップルを受け取る;
end;
(*ターゲット・サブリレーションの入力あり*)
while ターゲット・リレーションの終わりでない do
begin
  when マスタユニットからターゲット・リレーションが1タップル転送された do
    そのタップルを受け取る;
    そのタップルとソース・サブリレーションとのジョインを行う.
end;
(*ターゲット・サブリレーションの転送あり*)
while ターゲット・リレーションの直積 do
case ジョイン条件 of
  ">", " $\geq$ ", "<", " $\leq$ ":
    for  $j := 1$  to  $i$  do
      begin
        if  $i \neq n-1$  then
          ターゲット・サブリレーションをスレーブユニット  $SU_{i+1}$  へ転送する;
        when ターゲット・サブリレーションがスレーブユニット  $SU_{i+1}$  から転送された do
          ソース・サブリレーションとターゲット・サブリレーションとの直積を行う;
      end;
      if  $i \neq n-1$  then
        ターゲット・サブリレーションをスレーブユニット  $SU_{i+1}$  へ転送する;
    else:
      for  $j := 1$  to  $n-1$  do
        begin
          ターゲット・サブリレーションをスレーブユニット  $SU_{i+1}$  へ転送する;

```

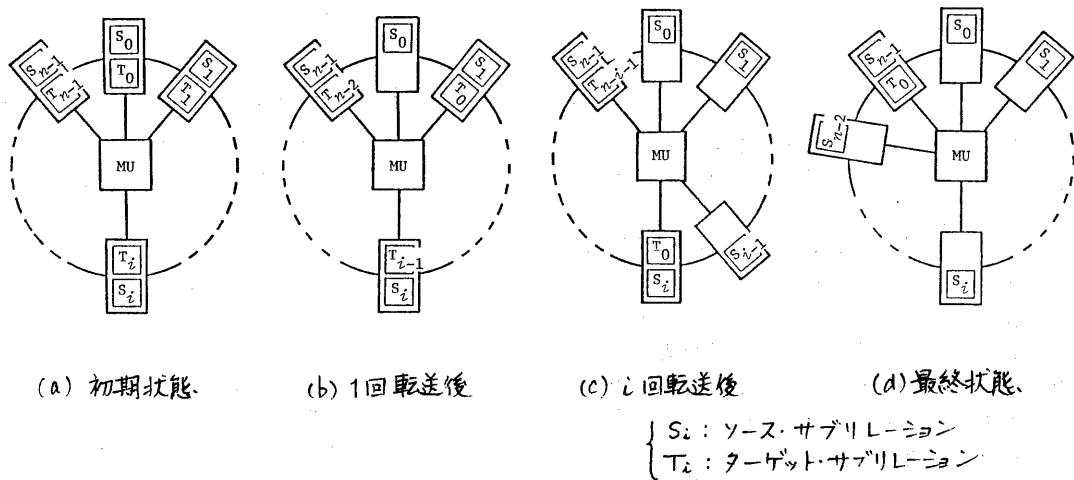


図3 θ -join ($>$, \geq , $<$, \leq) の場合のデータ転送

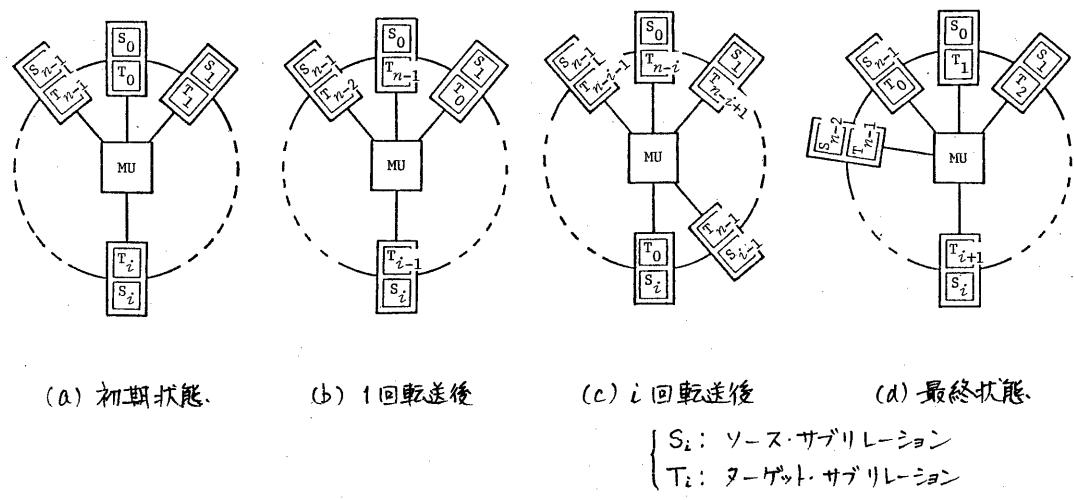


図4 θ -join ($=$) の場合のデータ転送

```

when ターゲット・サブリレーション
  ガスレーブユニット  $SU_{i-1}$  から転送さ
  れた do
    ソース・サブリレーションとターゲッ
    ト・サブリレーションとの直積を行う;
end;
end;

```

(*結果のサブリレーションの転送*)
while 結果のサブリレーションの終わりで
 then do
 マスクユニットへ結果のサブリレーション
 を1タップル転送する.
end.

4. 実験と評価

4.1 実験システム

広島大学工学部計算機工学研究室で開発されたマルチマイクロプロセッサシステムUNIPを用いて実験を行った。

UNIPの構成を図5に示す。UNIPは1台のマスタユニットとn台のスレーブユニットから構成されており、マスタユニットと各スレーブユニットは並列バスによって接続されている。また、隣接するスレーブユニット間は互いに直列ポートによって接続されている。各ユニットのハードウェア仕様を表1に示す。

プロセッサ間の通信はすべて同期的に行われる。すなはち、送信側、受信側いずれかが先に通信を開始した方が待ち状態に入り、相手が通信を開始するまで待ち続ける。

4.2 実験結果と評価

今回の実験では12台のスレーブユニットを使用した。ただし、スレーブユニット SU_0 と SU_1 は直接的には接続されていないため、この間のデータ転送はマスタユニットを介して行う。また、各スレーブユニットに均等な大きさのリレーションが分配されるようデータを用い、それがループの長さは20バイトとした。

ホストコンピュータ(PC-8001)上でジョイン演算を行った場合とジョインプロセッサによってジョイン演算を行った場合について、それぞれの処理時間と実測した。ジョイン条件として“=”、“>”、“キ”を設定したときのタップル数と処理時間との関係を、それぞれ図6、図7、図8に示す。ホストコンピュータのみで処理を行う場合と比べて、ジョインプロセッサでは明らかに実行速度の改善がなされている。

次にジョインプロセッサの場合について、タップル数を600としたときの処理時間の内わけを表2に示す。それぞれデータの転送時間と待ち時間の合計が57.4%，71.9%，84.7%を占めており、転送時間の短縮と待ち時間の削減が重要である。

以下、付録にジョインプロセッサの性能評価式を示す。

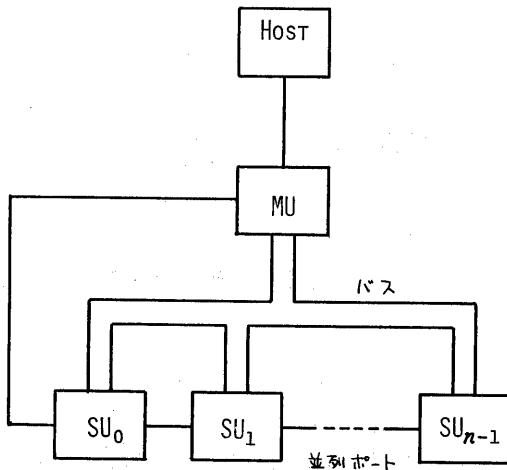
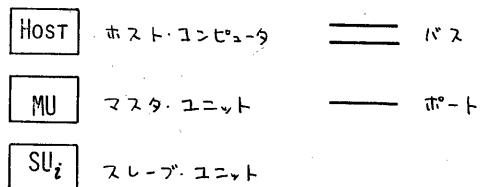


図5 UNIPの基本構成

マスタユニット

CPU	Z80A CPU (4MHz動作)
メモリ	ROM 2kバイト, RAM 14kバイト実装
並列バス	Z80A PIO
直列ポート	Z80A SIO
カウンタ/タイマー	Z80A CTC

スレーブユニット

CPU	Z80A CPU (4MHz動作)
メモリ	RAM 16kバイト実装
並列バス	Z80A PIO

表1 マスタユニットおよびスレーブユニットのハードウェア仕様

処理時間(秒)

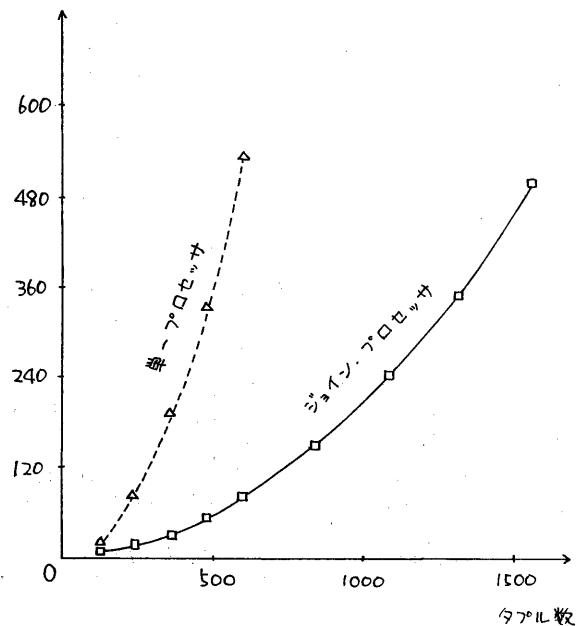


図6 単一プロセッサとの比較(ジョイン条件: =)

処理時間(秒)

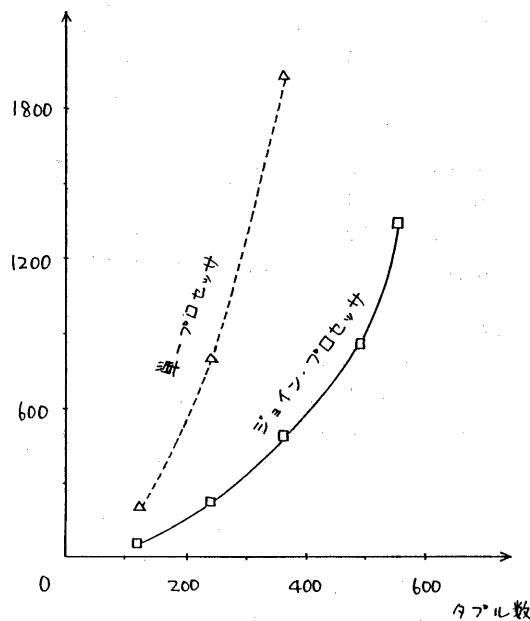


図7 単一プロセッサとの比較(ジョイン条件: >)

処理時間(秒)

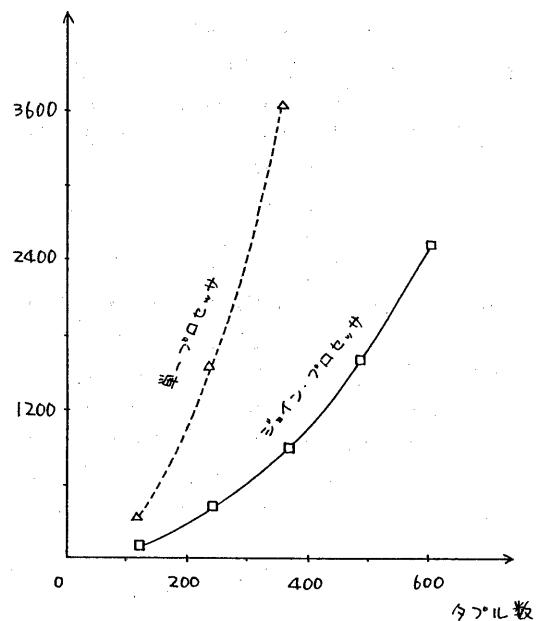


図8 単一プロセッサとの比較(ジョイン条件: ≠)

比較演算子 内わけ	=		>		≠	
	処理時間(秒)	割合(%)	処理時間(秒)	割合(%)	処理時間(秒)	割合(%)
ソース・リレーションおよびターゲット・リレーションの観察時間	7.00	8.6	7.00	0.5	7.00	0.3
ハッシュингに要する時間	10.34	12.7	10.34	0.9	10.34	0.4
比較演算に要する時間	24.26	29.9	40.48	3.0	49.98	2.0
結果ソリューションの転送時間と待ち時間	39.69	49.8	950.82	70.7	2114.11	83.8
スレーブ・ユニット間の転送時間	—	—	9.09	0.7	14.30	0.6
直結に要する時間	—	—	326.26	24.3	326.26	12.9
合計時間	81.30	100.0	1344.00	100.0	2522.00	100.0

表2 処理時間の内わけ(タップル数:600)

5. おわりに

本論文では、ジョイン演算を並列に処理するためのシステムアーキテクチャ、すらびにアルゴリズムについて述べた。システムは、1台のマスター・ユニットと、それと連絡されているN台のスレーブ・ユニットから構成されている。さらに、各スレーブ・ユニットは隣接するスレーブ・ユニットと接続されており、データの転送が行える。マスター・ユニットではまずリレーションのそれぞれのタップルについてハッシュングを行い、次にそのハッシュ値に対応したスレーブ・ユニットにタップルを転送する。各スレーブ・ユニットは独立に動作し、各自がジョイン操作を行う。

ところで、ジョイン条件が "=" の場合、その条件を満足するタップルは必ず同一スレーブ・ユニット上にある。したがって、それまでのスレーブ・ユニット内で処理が完了し、あらたにスレーブ・ユニット間でデータ転送を行ふ必要がない。それ以外のジョイン条件をもつ場合にはスレーブ・ユニット間でデータの転送が起こる。しかし、"≠", ">", " \geq ", "<", " \leq " の条件をとるジョイン演算については、ハッシュバケット間に順序性を持つようなハッシュ関数を適用することによって、転送すらびにスレーブ・ユニット

内のジョイン操作のコストを低減することができる。

現在、本学工学部阿江忠教授担当の計算機工学研究室で開発中のマルチマイクロプロセッサシステムUNIPを用いて、実験すらびに評価を行っており、今後データベースマシンの構成についても検討を進める予定である。

謝辞

実験設備の提供と、日頃頂いている有益なご討論に対して阿江忠教授に感謝する。また、日頃の討論と実験協力に対し、計算機工学研究室と情報システム研究室の諸氏に感謝する。

参考文献

- [1] M. Brodie and J. Schmidt (Eds.), "Final Report of the ANSI/X3/SPARC DBS-SG Relational Database Task Group," ACM SIGMOD RECORD, Vol. 12, No. 4, July 1982.
- [2] 平川, 清水, 市川, "ARESにおけるデータベース操作機能の拡張" 情報処理学会論文誌, Vol. 24, No. 4, pp. 513-520, July 1983.
- [3] Special Issue on Data Base Machines, IEEE Computer, Vol. 12, No. 3, pp. 7-79, March 1979.
- [4] 田中, "データベースマシン," 情報処理学会誌, Vol. 23, No. 10, pp. 939-947, October 1982.

- [5] T. Nakayama, M. Hirakawa and T. Ichikawa,
"Architecture and Algorithm for Parallel Execution of a Join Operation," Proc., Conf. on Computer Data Engineering, 1984 (in print).
- [6] T. Ae and R. Aibara, "Experimentation and Analysis of Multiprocessor Systems," Proc., Real-Time Systems Symp., pp.69-80, 1982.

付録 ジョインプロセッサの性能評価式

次のようないづれかのパラメータを用いる。

S : ソース・リレーションのカーディナリティ

イ

T : ターゲット・リレーションのカーディナリティ

R : 結果のリレーションのカーディナリティ

S_i : スレーブユニット SU_i に割り当てられたソース・サブリレーションのカーディナリティ

T_i : スレーブユニット SU_i に割り当てられたターゲット・サブリレーションのカーディナリティ

R_i : スレーブユニット SU_i で生成された結果のサブリレーションのカーディナリティ

α : $S \times T$ と R の比 ($R = S \times T \times \alpha$)

α_i : $S_i \times T_i$ と $S_i \times T_i$ の間でジョイン条件を満たすタップル数の比

t_{comp} : ソース・リレーションとターゲット・リレーションのジョイン条件が記述された属性の値を比較するのにかかる時間

t_{create} : 結果のサブリレーションを 1 タップル生成するのにかかる時間

t_1 : マスター(スレーブ)ユニットからスレーブ(マスター)ユニットへ 1 タップル転送するのにかかる時間

t_2 : スレーブユニット間で 1 タップル転送するのにかかる時間

n : スレーブユニット数

(1) ジョイン条件が "=" の場合

$$(S+T+R) \times t_1$$

$$+ \max_{0 \leq i \leq n-1} (S_i \times T_i \times t_{comp} + S_i \times T_i \times \alpha_i \times t_{create})$$

(2) ジョイン条件が ">", ">=", "<", "<=" の場合

$$(S+T+R) \times t_1$$

$$+ \max_{0 \leq i \leq n-1} (S_i \times T_i \times t_{comp} + S_i \times T_i \times \alpha_i \times t_{create})$$

$$+ \sum_{j=1}^{n-1} (\max_{\substack{j \leq i \leq n-1}} (S_i \times T_{i-j} \times t_{create}))$$

$$+ \sum_{j=1}^{n-1} (\max_{\substack{0 \leq i \leq n-1 \\ i \neq j}} (T_i \times t_2))$$

(3) ジョイン条件が "≠" の場合

$$(S+T+R) \times t_1 + \max_{0 \leq i \leq n-1} (T_i \times t_2) \times (n-1)$$

$$+ \max_{0 \leq i \leq n-1} (S_i \times T_i \times t_{comp} + S_i \times T_i \times \alpha_i \times t_{create})$$

$$+ \sum_{j=1}^{n-1} (\max_{\substack{0 \leq i \leq n-1 \\ i \neq j}} (S_i \times T_{i-j} \times t_{create}))$$

$$k = \begin{cases} i-j & (i \geq j) \\ i-n+j & (i < j) \end{cases}$$

但し、上の式では処理時間と転送時間の重ねあわせないと仮定している。