

多重プロセッサデータベースマシン
におけるデータフローの最適化
OPTIMIZATION OF DATA FLOWS FOR A
MULTIPROCESSOR DATABASE MACHINE

齋藤 直樹 馬場 敬信
Naoki SAITO Takanobu BABA
宇都宮大学工学部情報工学科
Utsunomiya University

1. はじめに

知識ベースや統計データベースなど、データベース（以下、DB）の発展に伴い、DBシステムに対する処理の要求は増大する一方であり、これを支援するために、より高性能のDBマシンが求められている。DBマシンに関しては、従来から、並列処理の活用を目指して数多くの提案がなされており、またその中には試作まで行なわれたものもある。これらの多くは、ディスクにおける並列トラックや、電子ディスクの活用等ハードウェアの特徴を生かして二次記憶からの高速アクセスを目指すと共に、データアクセスの過程でのマルチプロセッサによる並列処理を活用している[1,2]。

このように並列処理を取り入れる際に重要なことは、まずシステムの負荷の状況を調べ、最も効果的な所から並列化を進めることであろう。例えば、DBマシンの設計に当たっては、従来ディスクI/Oのボトルネックに重点が置かれているものが多いが、問い合わせの内容によってはCPU上での処理の負荷が大きくなり、CPUネックとなる場合のあることも実験的に示されている[3,4,5]。

このような背景から、本研究ではまず関係DBシステムを対象に、内部の処理におけるデータフローに着目してモデル化を行なう[6,7]。次に、このモデルのもとで与えられた問い合わせに対して可能なデータフローを生成し、空間的、時間的な並列処理を取り入れて、プロセッサ数ごとの各データフローの応答時間を評価する。この結果を用いて、最も有効なデータフローを選択すると

共にそれをプロセッサに割付けるためのアルゴリズムを求めらる。

2. 機能割付けの基本概念

2.1 DBシステムの基本概念

DBシステムの機能は、多くの基本的な要素から成る。図1はこれを簡単なモデルで示したものである。ターミナルのユーザから発生した問い合わせ（QUERY）は構文解析によって木構造に変換される。この結果生成された問い合わせ木は、コマンドコントローラによって実行されるが、その内容はデータ定義、データ操作、問い合わせ実行の3種類に分けられる。これらは更に、より基本的なルーチン（図では記憶システム）を起動する。

例えば、INGRES[8]における4つのプロセスは、

1. ターミナルモニタ
2. 構文解析
3. コマンドコントローラ、問い合わせ最適化、実行
4. データ定義、データ操作

に対応する。System R[9]におけるRDS (Relational Data System) とRSS (Research Storage System) はそれぞれ、

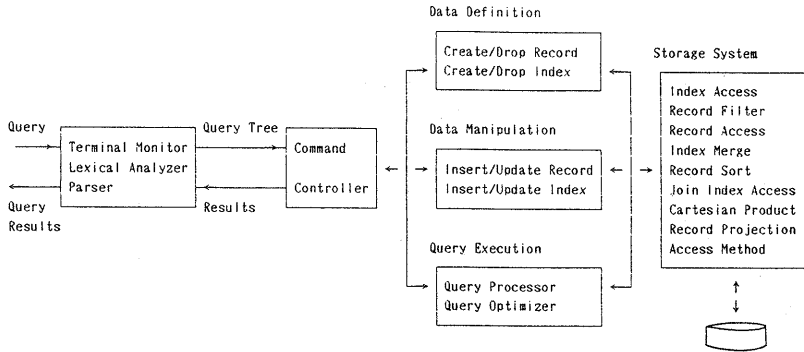


図1 データベースシステム機能の分割
 FIGURE 1: A decomposition of database functions

1. 構文解析、コマンドコントローラ、問い合わせ最適化、実行、選択、直積、射影
2. データ定義、データ操作、インデックス/レコード操作、マージ、ソート、結合

に対応する。

従って、多重プロセッサでDBマシンを実現する場合、これらの基本的な機能をいかにDBマシンに割付けるかが重要な課題となる。図2に、我々の提案する多重プロ

セッサへの割付けの過程を示した。まず問い合わせ木からはデータフローが生成されるが、これは処理のアルゴリズムや、インデックスの有無などにより種々の可能性がある。更に、このデータフローの各々に対して多重プロセッサへの割付けの可能性を検討し、最適な割付けを出力するのが多重プロセッサへの機能割付け (Function Allocator) の役目である。このように、割付けにあたっては、「データフローの生成」と「機能の割付け」の2つが要点となるが、今回は後者を中心に検討する。また、ここでは最適性の規範として、DBにデータフローを与えてから結果が得られるまでの時間 (応答時間) の最小化を目標とする。

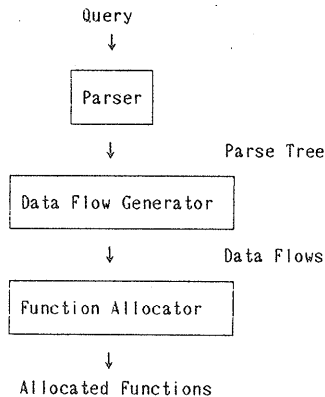


図2 基本機能割付けの基本概念
 FIGURE 2: Allocation of data flow to multiprocessors

2.2 最適化の前提

最適化に当たって考慮すべきことは、基本的な機能をいかにDBマシンに実現するかである。もし、ハードウェアやファームウェアで実現する場合には、割付けられた機能に応じてプロセッサを特殊化することが考えられる。この時には、各機能の実行は速くなるが、1つのプロセッサに実現される機能は制限されると考えるのが妥当である。一方、ソフトウェアで実現する場合には、現在のメモリICの価格と性能から考えて、1台のプロセッサにはDB処理に必要な全機能が含まれると仮定してよいであろう。

前者の立場に対する最適化の方法については既に検討済み[10]のため、ここでは後者の立場に立って検討を行なう。

多重プロセッサDBマシンのアーキテクチャとしては、

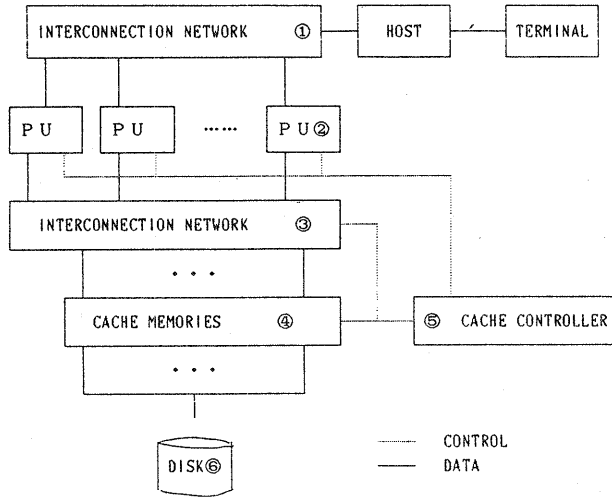


図3 データベース・マシンのアーキテクチャ
 FIGURE 3: Database machine architecture

ごく標準的な図3のような構成を前提とする。

P U-P U間ネットワーク①には、ターミナルからの問い合わせを処理するホストと、実際にDB処理を行なうバックエンド型のP U②が結合されている。ホストは最適化を行ない選択したデータフローをP Uに割付ける。更に、ディスク⑥に対してキャッシュ④を設け、ディスクI/Oのボトルネックを防ぐ。そして、ネットワーク③がP U-キャッシュ間を結合しており、キャッシュコントローラ⑤がそれらの制御を行なう。

3. データフローとそのコスト

3.1 データフローグラフ

DB処理における基本機能を節点に対応させ、さらに基本機能間のデータの流れを、方向を持つ枝に対応させると、処理の過程は有向グラフとして表現できる。典型的な問い合わせに対する基本機能の一覧を表1に示す。

例えば、Restriction Index(RI)は、与えられた条件に従ってそれを満足するレコードへのポインタ集合を出力する。そのコストは $unit(RI) \cdot n + dio(RI)$ で表わされ、 $unit(RI)$ はRIに対するCPU上での単位演算時間、 $dio(RI)$ はディスクアクセス時間を表わす。

表1 基本機能の動作とコスト
 Table 1. Data flow and cost model of elementary operations.

Elementary Function (f)	Contents	Processing Time C(f) (see footnote)
RI (Restriction Index) ---○---	according to input restriction predicate. RI accesses an index file and produces a set of pointers to records satisfying the predicate	$unit(RI) \cdot n + dio(RI)$
JI (Join Index) ---○---	according to input join value. JI accesses an index file and produces a set of pointers to records which contain the same value for join attribute	$unit(JI) \cdot n + dio(JI)$
RA (Record Access) ---○---	inputting a set of pointers. RI accesses a database and produces a set of records	$unit(RA) \cdot n + dio(RA)$
S (Sequential Scan) ---○---	given a relation name. S sequentially scans a relation and produces a set of records	$unit(S) \cdot n + dio(S)$
ST (Sort) ---○---	ST sorts and merges a set of input records so that they are grouped by ascending join values	$unit(ST) \cdot n \cdot \log(n)$
P (Projection) ---○---	P extracts specified attributes of input records	$unit(P) \cdot n$
RF (Restriction Filter) ---○---	according to restriction predicate. RF selects records satisfying the predicate	$unit(RF) \cdot n$
JF (Join Filter) ---○---	JF receives two inputs: records grouped by a join value and join value processed by the other JF. JF checks the match and produces the records with matched value	$unit(JF) \cdot (n1 + n2)$
C (merge and Concatenation) ---○---	C inputs two sets of records grouped by join values and concatenates them with matching join values	$unit(C) \cdot n1 \cdot n2$
I (Intersection) ---○---	I receives two sets of pointers, one of them grouped for a join value, and produces their intersection	$unit(I) \cdot (n1 + n2)$

$unit(f)$: unit CPU operation cost of function f.
 $dio(f)$: disk I/O time for an input data segment.
 n, n1, n2: number of records or pointers to be processed.

データフローとしてここで考慮するのは、多変数に対する代表的な処理アルゴリズムとして知られているネストドループジョインとソートマージジョインの2つである。

<例1>ジョインインデックスを用いたネストドループジョイン (図4)

2変数の一方について、インデックスのアクセス (RI-2) を最初に行う。他方の変数についてインデックスのアクセス (RI-1) を行ない、レコードをアクセス (RA-1) し、ジョインインデックスの共通項が2つの流れからとられる (I)。そして、レコードをアクセスした後 (RA-2)、レコードの結合を行ない (C)、フィールドの取り出し (P) が、行なわれる。

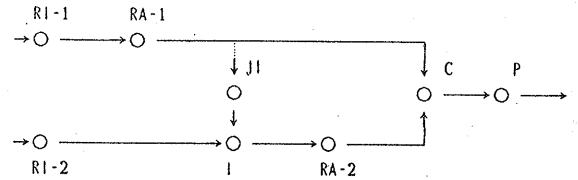


図4 J Iを用いたネストドループジョインのデータフロー
FIGURE 4: Data flow of nested loop join with join index

<例2>ジョインインデックスを用いないネストドループジョイン (図5)

インデックスをアクセス (RI-1) した後、レコードをアクセス (RA-1) する (インデックスがなければシーケンシャルスキャン (S))。他方の変数についてもインデックスのアクセス (RI-2) を行ない、レコードをアクセス (RA-2) し、これを与えられた条件によって選択 (RF) した後、2つの流れのレコードの結合を行ない (C)、必要なフィールドの取り出し (P) を行なう。

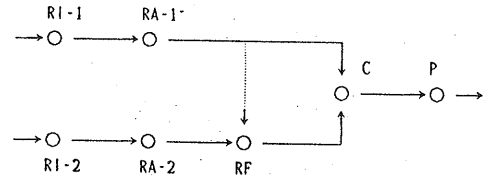


図5 J,Iを用いないネストドループジョインのデータフロー
FIGURE 5: Data flow of nested loop join

<例3>ソートマージジョイン (図6)

2変数の両方について、インデックスのアクセス (RI-1, RI-2) を最初に行い、レコードをアクセス (RA-1, RA-2) する。次にそのレコードをソート (ST-1, ST-2) し、結合を行ない (C)、フィールドの取り出し (P) が行なわれる。

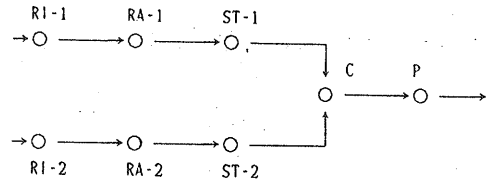


図6 ソートマージジョインのデータフロー
FIGURE 6: Data flow of sort and merge join

3.2 処理時間の計算

各データフローの処理時間の計算は次のように行なう。

まず、基本的な時間として単位データ (セグメントと呼ぶ) に対する処理時間を定義する。標準的なセグメントとしては、1 ページ分のレコードなどが考えられる。

コストをディスク I/O 時間、CPU 時間、転送時間に分割して次のように計算する。f は基本機能を表わす。他の各変数は表 2 に示した。

表2 変数表

TABLE 2: Database system parameters

Parameter	Meaning
C(f)	DIO(f)+CPU(f)+T(f)
DIO(f)	disk I/O time for a unit data segment
CPU(f)	CPU time for a data segment
T(f)	transfer time for a data segment between processors
d	average disk access time for a data segment
c	average cache access time for a data segment
rq	average time to wait for I/O. (seek and latency)
h	cache hit ratio
unit(f)	unit operation time of function f
t	unit transfer time
r(f)	selectivity of function f
wd	cache-processor bus contention time
wt	processor-processor bus contention time
N	number of segments
P	number of processors
Np	N/P (number of data segments to be processed in a processor)
k	number of records in a data segment
pipeC	pipeline processing time for a data segment
repCi	parallel processing time for a data segment, using i processors

$$C(f) = DIO(f) + CPU(f) + T(f) \quad (3.1)$$

ここで、各コストは以下ようになる。

$$DIO(f) = \begin{cases} rq + \{h \cdot c + (1-h) \cdot (d+c)\} + wd \\ 0 \end{cases} \quad \text{f がディスク I/O を行わない時} \quad (3.2)$$

$$CPU(f) = \text{unit}(f) \quad (3.3)$$

$$T(f) = t \cdot r(f) + wt \quad (3.4)$$

ディスク I/O 時間は、シーク時間と回転待ち時間、キャッシュ又はディスクから P U への転送時間、転送ネットワークの競合による待ち時間からなる。次に CPU 時間は、基本機能に応じた処理時間からなる。転送時間は、基本転送時間、転送ネットワークの競合による待ち時間からなる。

4. 最適化

4.1 最適化の操作

複数の基本機能を多重プロセッサ上に実現する場合、分散の仕方によって、以下の 3 手法に分けられる。

- (1) 1 台のプロセッサ上での逐次処理（複数の基本機能を一つにまとめる）
- (2) 異なるプロセッサ間の時間的な並列処理（パイプライン）
- (3) 異なるプロセッサ上の空間的な並列処理（同じ逐次処理を複数のプロセッサで実行する）

以下では単に (2) をパイプライン処理、(3) を並列処理と呼ぶ。

基本機能を A、B、プロセッサ数を 2 としたときのパイプライン処理と、並列処理の例を図 7 に示した。以下で、3 手法によるコストを定義する。

(1) 逐次処理

1 つのプロセッサ内でのデータ転送を単にパラメータの受渡しで行なうものとし、その時間を無視すると逐次処理のコストとして以下の式が定義される。ただし、 f_1, f_2, \dots, f_n は基本機能である。

$$C(g(f_1, f_2, \dots, f_n)) = \sum_{i=1}^n \{CPU(f_i) + DIO(f_i)\} + T(f_n) \quad (4.1)$$

(2) パイプライン処理

パイプライン処理においてセグメントは、P U から P U へ次々と流れていく。ここでは、簡単化のために各セグメントの処理時間を一定と考える。この場合、最も時間のかかる逐次処理を行なう P U がネックとなり、パイプライン処理のコストもその時間で決定される。従って、各 P U の逐次処理のコストを $C(g(f_1, \dots, f_i)), \dots, C(g(f_j, \dots, f_k)), C(g(f_{k+1}, \dots, f_n))$ とすると、パイプライン処理のコストは、次の式で表わされる。

$$\text{pipeC} = \max\{C(g(f_1, \dots, f_i)), \dots, C(g(f_j, \dots, f_k)), C(g(f_{k+1}, \dots, f_n))\} \quad (4.2)$$

(3) 並列処理

基本機能を m 台のプロセッサに複写した時、m セグメントを同時に処理することが可能なので処理時間 (DIO + CPU) は $1/m$ になる。ただし、ソートを並列処理で実行する場合、ソートしたセグメントをマージしなければならなくなる。このように、並列処理によって新たなコストを生じることもあるが、ソートを除きそのコストは無視できるものとする。

基本機能の逐次処理を $g(f_1, \dots, f_n)$ とすると、m 台のプロセッサを用いた並列処理のコストは、以下のように定義される。

$$\text{repC}_m = 1/m \cdot \sum_{i=1}^n \{CPU(f_i) + DIO(f_i)\} + T(f_n) \quad (4.3)$$

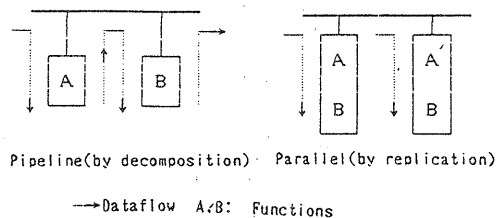


図 7 パイプライン処理と並列処理
FIGURE 7: Parallel and pipeline dataflow

4.2 最適化のアルゴリズム

前節までの定義を用いることにより、次のことが言える。

「同一数のPUを用いて基本機能の割付けを行なう場合、並列処理は常にパイプライン処理より効果が大きい。」

この証明は付録に記した。

従って、常に並列処理を優先させて適用していけばよいのだが、処理のアルゴリズム上、並列処理だけでは済まないものがある。その代表的なものには先に述べたソートとそれに続くマージである。ソートを並列処理で実行すると必然的にマージとの間でパイプラインを形成することになる。

以上の議論をまとめると、ある与えられたデータフローに対する割付けのアルゴリズムは以下ようになる。

ルール1)

基本的には、1台のプロセッサに、全基本機能を載せる。そして、全プロセッサで並列処理を行なう。

ルール2)

ソートマージのように、アルゴリズム上PU間のデータ転送が必要となるものについては、パイプライン処理を行なう。

4.3 ソートマージジョイン

最適化において特別扱いするソートマージについて、更に具体的に考えると次の2つの方法が考えられる。

【方法1】NセグメントをP台のPUが各々N/Pセグメントをソートした後、パイプラインを形成し、ソートしたセグメントを次々とマージする(図8)。

【方法2】始めからソートを行なうPUとマージを行

なうPUに分けてパイプラインを形成し、ソートしたセグメントを次々とマージする(図9)。

ここで、1セグメントのソートコストをCs、2セグメントのマージコストをM2とする。そして、mセグメントとnセグメントのマージコストは、 $M2 \cdot \max\{m, n\}$ と仮定する。

方法1についてコストは以下ようになる。

まず、1セグメントをkレコードとした時のコストを $Cs = k \cdot \log k$ とする。各PUは、 $Np (= N/P)$ セグメントのソートを行なうので、その時のソートコストは次のようになる。

$$k \cdot Np \cdot \log(k \cdot Np) \quad (4.4)$$

ソート終了後、次に、PU間でマージを行なう過程がある。まず、P/2台の内容を他のP/2台に転送(図8(a))し、空になったP/2台のうち、P/2-1台を用いて、パイプラインを形成(図8(b))する。

そこでマージ過程のコストは、

(P/2台の内容転送)

$$\begin{aligned} &+ (TOP \text{ のPUにセグメントが到着するまでの時間}) \\ &+ (TOP \text{ のPUが全セグメントを処理する時間}) \\ &= N \cdot t/P + (M2 \cdot t) \cdot \log(P/2) + (N-1) \cdot M2 \end{aligned} \quad (4.5)$$

となる。以下で式の意味を説明する。

マージ過程では、パイプラインを形成するために、まずP/2台の内容を他のP/2台に転送してP/2台を空にする(第1項目)。

パイプラインが形成されると、パイプラインの葉に相当するPU(例えば、図8(b)ではPU1)から根に相当するPU(図8(b)ではTOP)へセグメントが流れて

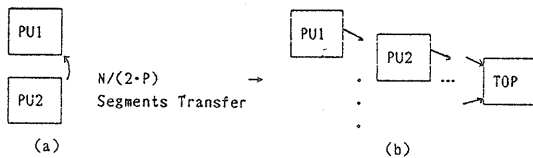


図8 ソートマージ(方法1)
FIGURE 8: Sort merge method 1

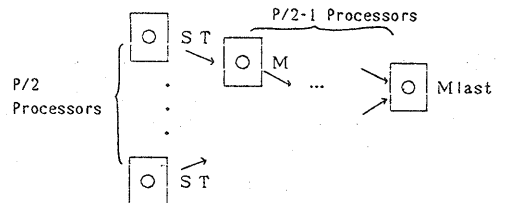


図9 ソートマージ(方法2)
FIGURE 9: Sort merge method 2

行く。TOP の P U に最初のセグメントが到着するまでの間、マージと転送がパイプラインの木の高さから TOP の P U を除いた $\log(P/2)$ 回行なわれる (第 2 項目)。

そして、TOP の P U に最初のセグメントが到着した後には TOP の P U でマージが行なわれる。ここでのマージはマージ済みのセグメントと流れてくるセグメントの大小関係が保証されているので、すでにマージ済みのセグメントとマージする必要がなく、 $N-1$ 回のマージ (第 3 項目) で全セグメントのソートマージを完了する。

従って、(4.4) と (4.5) から全処理コストを $SMJ1$ とすると、以下の式で表わされる。

$$SMJ1 = k \cdot Np \cdot \log(k \cdot Np) + N \cdot t/P + (M2 + t) \cdot \log(P/2) + (N-1) \cdot M2 \quad (4.6)$$

方法 2 についてのコストは以下のようになる。

パイプラインの木の高さは、 $\log(P/2) + 1$ であり、ソートとマージに用いる P U 数はそれぞれ $P/2$ 台、 $P/2-1$ 台である。

パイプラインではソートに用いた $P/2$ 台が葉に相当するので、 $P/2$ セグメントが 1 回のデータフローの単位となる。そして、流れた単位の数をステージ数と呼ぶ。

方法 1 と異なりセグメント内のレコードはソートされているが、セグメント間ではソートされていないため、すでにマージ済みのセグメントとマージする必要が生じる。従って、最後の段 (マージ) のコスト $Mlast$ は、ステージ数の増加に比例して、以下のようになる。

$$(P/2 \text{ セグメントのマージ}) + (\text{前回までのマージ済みセグメントと } P/2 \text{ セグメントのマージ})$$

ここで、第 1 項目は次のようになる。

$$(P/2 \text{ セグメントのマージ}) = (P/2-1) \cdot M2$$

次に第 2 項目の (前回までのマージ済みセグメントと $P/2$ セグメントのマージ) では、ステージ数とコストに以下の関係が成り立つ。

ステージ数	1	2	i
コスト	0	$(P/2-1) \cdot M2$	$(i \cdot P/2-1) \cdot M2$

従って、 $Mlast$ は常に他のマージよりコストが高いので、(4.2) から Cs と $Mlast$ の大小関係でパイプライン処理のコストが決まる。すなわち、このパイプライン処理のコストは、 $\max\{Cs, Mlast\}$ となる。そこで、パイプライン処理のコストを決めるために Cs と $Mlast$ の大小関係を比較しなければならない。j 回目のステージの $Mlast$ は、以下のようになる。

$$(P/2-1) \cdot M2 + \sum_{l=1}^j (l \cdot P/2-1) \cdot M2 = M2 \cdot \{P/2-1 + j \cdot (P \cdot (j+1)/4 - 1)\} \quad (1 \leq j \leq 2 \cdot N/P)$$

j 回目のステージで、 $Cs < Mlast$ になったとすると、逆に $j-1$ 回目までは Cs の方がコストが高いと言える。(j は上式と Cs を比較して決定する)

パイプライン処理のステージ数 $j-1$ 回目までのコストは、ソートが決定し、

$$(j-1) \cdot (Cs+t) \quad (4.7)$$

となる。残りステージ数 j 回から $2 \cdot N/P$ 回までのコストは、マージ ($Mlast$) が決定するので、

$$(P/2-1) \cdot M2 \cdot (2 \cdot N/P - j) + \sum_{l=j}^{2N/P} (l \cdot P/2-1) \cdot M2 = M2 \cdot (2 \cdot N/P - j) \cdot (N/2 - P \cdot j/4 - 1) \quad (4.8)$$

となる。

従って、(4.7) と (4.8) から全コストを $SMJ2$ とすると、以下の式となる。

$$SMJ2 = (j-1) \cdot (Cs+t) + M2 \cdot (2 \cdot N/P - j) \cdot (N/2 - P \cdot j/4 - 1) \quad (4.9)$$

5. 適用例

関係を以下のように定義する。

SUPPLIERS (SNO, SNAME, ITEM, CITY)

各10,20,20,10レコード, 計60レコード

ORDERS(NO, NAME, ITEM, QUANTITY)

各10,20,20,30レコード, 計80レコード

関係SUPPLIERSは、部品供給者の番号、名前、品名、都市名を示し、関係ORDERSは発注の番号、発注者名、品名、量を示している。この二つの関係に対して、「Brooks氏の発注した品物を保有し、New Yorkにいる供給者」を求める問い合わせは以下ようになる。

```
SELECT SNAME, SUPPLIERS.ITEM,
FROM SUPPLIERS, ORDERS
WHERE NAME='Brooks, B.' AND
SUPPLIERS.ITEM=ORDERS.ITEM, AND
CITY='NEW YORK'
```

ここで、NAMEとCITYに対するインデックス(RI)の有無、及びITEMに対するジョインインデックス(JI)の有無に応じて種々のタイプのデータフローの可能性がある。以下で、NLJはネステッドループジョイン、SMJはソートマージジョインを意味する。

Aタイプ	JIを用いたNLJ (図4に対応)	}	2変数共にRI可 1変数にRI可
Bタイプ	JIを用いないNLJ (図5に対応)	}	2変数共にRI可 1変数にRI可 RI不可
Cタイプ	SMJ (図6に対応)	}	2変数共にRI可 1変数にRI可 RI不可

以下では、2変数共にRI可でJI可とし、PU数を7台とした時の割付けをタイプごとに示す。

<Aタイプ> JIを用いたNLJによる割付けを示す(図10)。

ルール1)により並列処理を適用してデータフローの全開数を全PUに割付ける。ここで、割付けた基本機能の記述順に意味はない。また、括弧はデータフローに表われるものの、他の基本機能と共有できる基本機能であることを示す(例えば、RI-1とRI-2は同じRIである)。

<Bタイプ> JIを用いないNLJによる割付けを示す(図11)。

ルール1)により並列処理を適用してデータフローの全開数を全PUに割付ける。

<Cタイプ> SMJによる割付けを示す。

ルール2)により4.3節で述べた2つの方法による割付けを順に示す。

まず方法1であるか、方法1を用いた割付けでは3つの段階がある(図12)。

1) PU内ソートを完了するまでの段階(図12の(1))。

ここで、各PUはN/Pセグメント(今は、N/7セグメント)のソート処理を実行する。

2) PU間でマージする段階(図12の(2))。

この段階では、各PUがソートしたN/7セグメントが次々にPU間を流れてマージされていく。

3) マージ結果を各PUに分配して、ジョインする段階(図12の(3))。

この段階は、マージの完了したセグメントが1つのPUに集中しているため、各PUに分配して結合しなければならない。ここでは、ソートマージが済んでいるのでルール1)に従い、全PUに並列処理が適用される。

次に、方法2による割付けを行なう。

方法2は始めからソートとマージ間をパイプラインで実現するので、方法1の1、2)の段階が図13となる。ただし、段階3)は方法1と同一である。

PU0~6

RI-1, (RI-2)
RA-1, (RA-2), JI, I, C, P

図10 JIを用いたネストドループジョインの割付け
FIGURE 10: Allocation of functions for nested loop join with join index

PU0~6

RI-1, RA-1, (RI-2, RA-2), RF, C, P

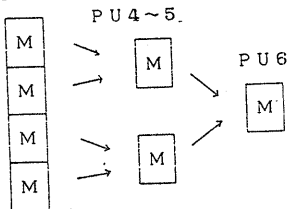
図11 JIを用いないネストドループジョインの割付け
FIGURE 11: Allocation of functions for nested loop join without join index

PU0~6

RI-1, RA-1, ST-1
(RI-2, RA-2, ST-2)

(1) First step

PU0~3



(2) Second step

PU0~3

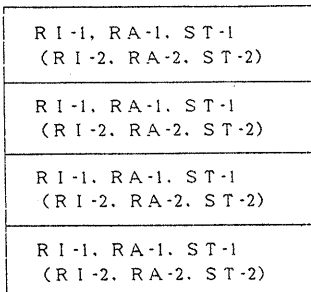
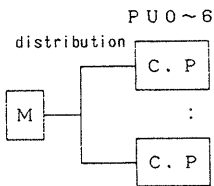


図13 ソートマージジョインの割付け(方法2)
FIGURE 13: Allocation of functions for sort merge join method 2

図12 ソートマージジョインの割付け(方法1)

FIGURE 12: Allocation of functions for sort merge join method 1



(3) Third step

前記の割付けを基本として、下記の値を用いてシミュレーションを行なった。値の設定に当っては、文献[11, 12]を参考とした。図14は、対象とする2つの関係SUPPLIERS (1レコードは60バイト)とORDERS (1レコードは80バイト)のセグメント数(1セグメントは10レコード)を共に100とした時の結果を示す。この時、総レコード数はSUPPLIERS、ORDERS共に1000である。図15は、2つの関係のセグメント数を共に10000とした時の結果を示す。この時、総レコード数はSUPPLIERS、ORDERS共に100000である。

rq : 100 μ s
h : 0.6
d : 1 μ s/byte
c : 0.5 μ s/byte
k : 10 レコード
unit(f) : 0.1 μ s/byte
Cs : unit(f) \cdot k log k μ s/byte
M2 : 2 \cdot unit(f) \cdot k μ s/byte
t : 0.1 μ s/byte
r(RI-1) = r(RI-2) = 0.1, r(P) = 0.3,
r(JI) = 2, r(I) = 0.2, r(RF) = 0.4,
wd = wt = 0

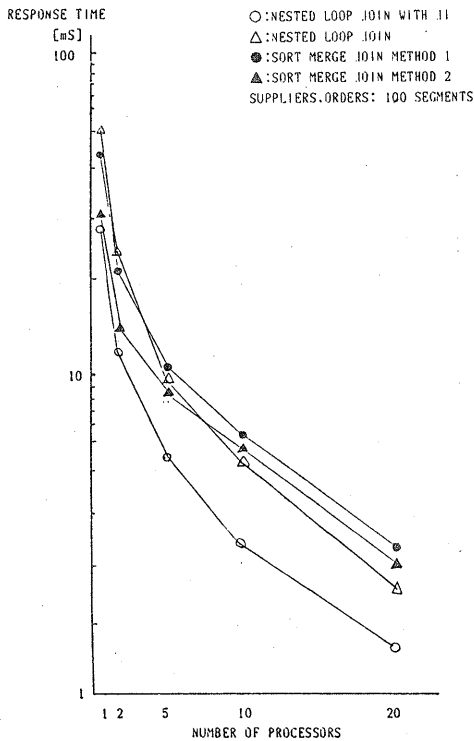


図 14 データフローとプロセッサに対する
応答時間の比較 (セグメント数 100)

FIGURE 14: Response times of the algorithms
versus the number of processors

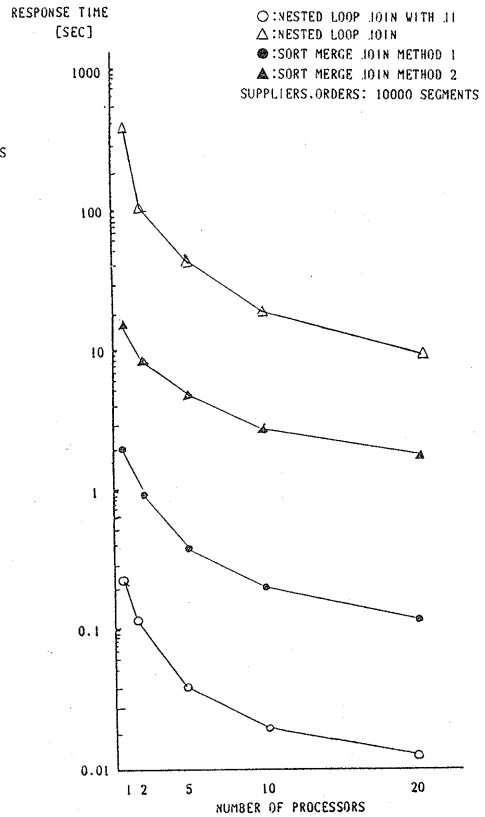


図 15 データフローとプロセッサ数に対する
応答時間の比較 (セグメント数 10000)

FIGURE 15: Response times of the algorithms
versus the number of processors

結果を以下に要約する。

1) 4つの手法の比較

J I を用いることが可能ならば、J I を用いた N L J が最適なデータフローであると言える。しかし、J I を用いることができない場合、他の3つの手法の優劣については、P U 数、データ量に依存して変わるため一概に言えない。

2) J I の効果

J I を用いた N L J は、J I を用いない N L J より結合するデータ量が少なくなり、応答時間も短くなる。そして、処理するセグメント数が大きくなると J I の効果は顕著になる。

3) S M J での2つの方法の比較

セグメント数が多いと、方法2では最終段のマージの負荷が格段に増大して、そのためパイプラインの流れが滞り、応答時間は方法2 > 方法1となる。従って、方法2では最終段のマージの高速化を行わないと方法1に対抗できない。しかし、セグメントが少ない場合、方法2では最終段のマージの負荷が低く、パイプラインが有効に働くために方法1より有利になる。

特に1)の結果は、アルゴリズムの優劣が一意的に決まらず、J I の有無、P U 数やデータ量に依存して決まることを示しており、本論文で述べたような最適化の方法が有効であることを裏付けている。

6. まとめ

データベースシステム内部のデータフローのモデル化を行ない、これに基づいて、多重プロセッサ計算機を対象とした、処理機能の割付けを考えた。研究はまだ初期の段階にあり、本論文で述べたことが一般のデータベースシステムに問題なく適合するか否かはこれからの検討課題である。特に、プロセッサの多重化に伴うディスク I/O や結合バスの競合の問題は重要な課題である。今後、実際のデータベースシステム内部の検討を行ない、より実際のモデルを定義するとともに、多重プロセッサシステムへの応用を考えていきたい。

謝辞

本研究を進めるに当たり貴重な御意見を頂いたメリランド大学 S. Bing Yao 助教授、ならびに宇都宮大学 奥田健三教授、山崎勝弘助手に感謝する。

付録

任意の数の基本機能に任意の数のプロセッサを与えた時、常にパイプライン処理より並列処理が有利であることを示す。

基本機能を以下のようにする。(プロセッサ数 < 基本機能数の時は g 個の基本機能群に統合済みとする)

$$(F_1, F_2, \dots, F_g)$$

ここで、プロセッサ P 台 ($P \geq g$) を各基本機能に割付けた場合を考える。まず、パイプライン処理で 1 セグメントを処理するのに必要なコストは、以下ようになる。

$$\begin{aligned} \text{pipeC} &= \max\{C(F_1)/N_1, \dots, C(F_i)/N_i, \dots, C(F_g)/N_g\} \\ &= C(F_m)/N_m \quad (1 \leq m \leq g), p = \sum_{j=1}^g N_j \end{aligned}$$

次に、全基本機能を 1 つにまとめて並列処理を適用した場合、並列処理では P セグメントを P 台のプロセッサで並列に処理できるので、1 セグメントを処理するのに必要なコストは、以下ようになる。

$$\text{repCp} = C(F_1 + F_2 + \dots + F_g)/P = \sum_{j=1}^g F_j/P \quad (1 \leq j \leq g)$$

ここで、 repCp の値を確定すると、パイプライン処理のコスト $C(F_i)/N_i$ が最小値をとるのは、パイプラインの各要素 ($C(F_1)/N_1, \dots, C(F_i)/N_i, \dots, C(F_g)/N_g$) が全て同じ値をとる時である。その時、パイプライン処理のコストの最小値は、

$$\sum_{j=1}^g F_j/P$$

となる。従って、

$$\text{pipeC} = C(F_m)/N_m \geq \sum_{j=1}^g F_j/P = \text{repCp}$$

となり、常にパイプライン処理より並列処理が有利である。

証明終わり。

参考文献

- [1] "Special Issue on Database Machines," IEEE Computer, Vol.12, No. 3, March 1979.
- [2] "Special Issue on Database Machines," IEEE Transactions on Computers, Vol.C-28, No. 6, June 1979.
- [3] C.J.Date: An Introduction to Database Systems, Vol.2, Addison-Wesley, 1983.
- [4] R.Epstein and P.Hawthorn: "Design Decisions for the Intelligent Database Machine," Proceeding of National Computer Conference, pp. 237-241, 1980.
- [5] M.Stonebraker: "Retrospection on a Database System," ACM Transactions on Database Systems, Vol.5, No.2, pp. 225-240, June 1980.
- [6] T.Baba, S.B.Yao, A.R.Hevner, and Z.Shi: "Multiprocessor Database Machine Architectures Based on Data Flow Analysis," Proceeding of the 16th Annual Hawaii International Conference on System Sciences, January 1983.
- [7] 清木、長谷川、雨宮: "先行・遅延評価機構を用いた関係演算処理方式" 情報処理学会論文誌, Vol. 26, pp 685-695, 1985.

[8] H. Stonebraker, E. Wong, P. Kreps, and G. Held: "The Design and Implementation of INGRES," ACM Transactions on Database Systems, Vol.1, No.3, pp. 189-222, September 1976.

[9] M. Astrahan, et al.: "System R: Relational Approach to Database Management," ACM Transactions on Database Systems, Vol.1, No.2, pp. 97-137, June 1976.

[10] 馬場, S.B. Yao, A.R. Hevner: "データフロー解析に基づく多重プロセッサデータベースマシンの設計" 情報処理学会第29回全国大会論文集, pp. 757-758, 1984.

[11] R.k. Shultz and R.J. Zingg: "Response Time Analysis of Multiprocessor Computers for Database Support," ACM Transactions on Database Systems, Vol.9, No.1, pp. 100-132, March 1984.

[12] 菅, 上田, 田中: "ディスク・キャッシュ装置のシミュレーションによる効果測定" 情報処理学会論文誌, pp. 22-28, Jan. 1981.