

# 結合演算を高速に実行する関係データベースマシン

7X-4

大月知之

濱田 喬

東京大学大学院工学系研究科

学術情報センター

## 1 はじめに

関係データベースは従来事務処理に用いられることが多かったが、昨今人工知能、CAD/CAM等応用分野も広がり、その高速処理への要求は高まる一方である。その演算としては選択、結合、射影等があるが、これらのうちの結合は処理負荷が重く、また、最近の応用分野において繰り返し使われることが多いため、その処理コストの削減が特に重要な課題となっている。現在、このような結合演算を高速に行うデータベースマシンについて研究を進めているが、本稿ではこのマシンに関し、その基本思想、機能、構成、演算実行の流れ、マシンの持つ特徴について述べる。

## 2 基本思想

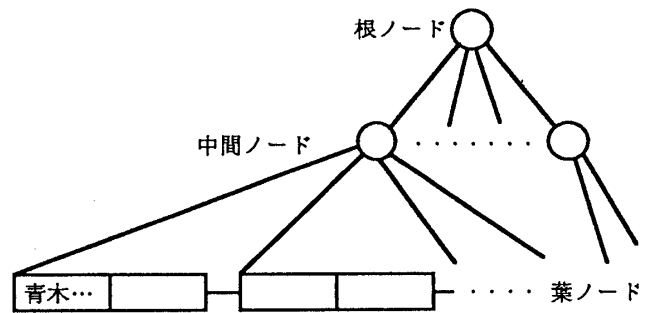
関係データベースマシンにおいて演算の実行にかかるコストには大きく分けて、2次記憶のアクセスコスト、データの転送コスト、処理コストの3つがあり、それぞれのコストを抑えるために、インデックスの利用、高性能ネットワークの採用、並列処理の利用等の方法がとられている。これらのうちインデックスは、データベースマシンに採用されているものの殆どが主に選択演算の実行に威力を発揮するものであり、結合演算等の実行時には2次記憶へのアクセスが無駄に行われていることが多い。しかし、現在提案されているインデックスの中には結合の実行にも有効なものが存在する([1]、[2])。本稿で提案するマシンは、このようなインデックスの1つ(Bc-tree[2])と並列処理(ハッシュとソートを用いる手法[3])の利用により結合の高速化を実現するものである。以下、このインデックスと並列処理について簡単に説明を行う。

まず採用するインデックスのBc-treeについて。これは、インデックスを関係ごとにではなくドメインごとに作るものであり、その内容は図1-aの表のようなものである。但し、タプル位置の情報とは具体的には属性値と対応する複数のタプル識別子(TID)を記したページ(TIDページ)へのポインタである。実際のBc-treeはこの表に属性値に基づくナビゲーション用の中間ノードを付した木構造をとる(図1-b)。インデックスとデータの構成は全体として図2のようなものである。

次に、このインデックスを用いた演算の実行例を紹介する。

属性値	関係R		関係S	
	タプル数	タプル位置の情報	タプル数	タプル位置の情報
青木	3	Page A	2	Page B
青山	0	null	5	Page C

(a) 内容



(b) 構造

図1 Bc-tree

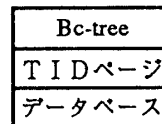


図2 インデックスとデータの構成

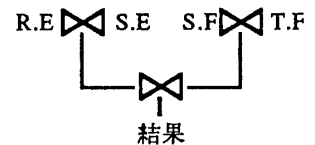


図3 複雑な処理の実行例

<例1> 選択演算  $R.A = \text{const}$

この場合、まず属性Aの属するドメインのBc-tree、次にTIDページを参照することにより必要なTIDを獲得し、TIDに記された格納ページをもとに最終的なタプルを読み出す。

<例2> 結合演算  $R.E = S.E$

この場合、まず属性Eの属するドメインのBc-treeを参照し、結合可能な属性値を求める。次にこの属性値に対応するそれぞれの関係のTIDページを参照することにより結合すべきタプルを決定し、最後に必要タプルを読み出し、結合を実行する。

上記の2例から分かるように、このインデックスを用いると、選択、結合ともに必要なタプルのみを取り出すので2次記憶への無駄なアクセスが避けられる。

次に、より複雑な問い合わせの実行における並列処理の利用を例により説明する。

### <例3> R.E=S.E $\wedge$ S.F=T.F

この場合、まずRとS、SとEの結合に必要なタプルをBc-tree、TIDページを用いて調べる。次に、両方の結合に関し得られたSのタプルをそのTIDに関し結合する。これにより、最終的に必要なSのタプルを決定し、SのタプルをもとにR、Tについても必要なタプルの決定を行い、これらタプルを読み出して結合の実行を完了する(図3)。

例3においてSのTIDに基づく結合にはインデックスは使用できず、他の何らかの手法で結合を行う必要があるが、このようなTIDに基づく結合は処理の複雑さに比例して多くなるものである。そこで、本研究ではこれらの結合をハッシュとソートを用いた並列処理により実現することとした。以下、この結合手法につき簡単に説明する。

この手法では、まず結合を行う両処理対象をその結合属性値にハッシュ関数を施した結果に従ってクラスタリングする。この際、等しい属性値を持つタプルは同じクラスタに割り振られるため、各クラスタは独立なものとなる。そこで、これらのクラスタを複数のプロセッサを用いてソート、マージし、結合を実行する。このような手法を用いることにより結合の処理コストを削減することができる。

以上に記したように、Bc-treeを用いて2次記憶アクセスコストを削減し、ハッシュとソートを用いた並列処理により処理コストを削減して高速な処理を実現することが本研究における基本思想である。

### 3 ホストマシンとデータベースマシンの役割分担

現在研究中のマシンはホストマシンに対するバックエンドマシンという位置づけを持つが、それぞれが以下のような機能を持つことを前提としている。

#### <ホストマシン>

- ・ユーザの権利チェック
- ・問い合わせ解析
- ・概念スキーマから外部スキーマへの変換
- ・障害対策

#### <データベースマシン>

- ・各演算の実行
- ・内部スキーマ、概念スキーマの取り扱い
- ・一貫性制約の実現
- ・障害対策

### 4 マシンの構成

マシンの構成は図4に示すようなものである。但し各略号の意味は以下のとおりである。

- ・CM: 制御モジュール。他の各モジュールを制御し、ホストマシンからの命令を実行する。
- ・IM: インデックス・モジュール。Bc-treeを格納する2次記憶を持ち、インデックス処理の一部分を実行する。
- ・DM: ディスク・モジュール。2次記憶を持ち、データベース、TIDページ等を格納する。
- ・MM: メモリ・モジュール。データのステージングを行う。
- ・PM: プロセッサ・モジュール。汎用プロセッサとハードウェア・ソータを持ち、各処理を

実行する。

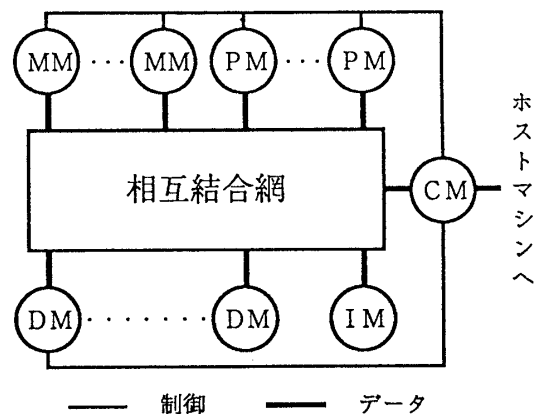


図4 マシンの構成

### 5 演算実行の流れ

問い合わせは以下のようにして実行される。まずIMがBc-treeの必要葉ノードをMMを介してPMに転送する。PMはこれをもとに必要なTIDページを決定し、これをDMに要求する。DMがこれに応じTIDページをMMを介してPMに転送する。これをもとにPMが各素演算(インデックスを用いて完了できる演算)において必要なタプルを決定する。次に、MMを介しながらPMとPMでデータのやりとりをし、タプル識別子の結合を行って、最終的に必要なタプルを決定する。これらタプルをPMがDMに要求し、DMからMMを介してPMにデータを転送することで最終的な処理を完了する。

### 6 マシンの持つ特徴

本研究で構築するマシンは以下のような特徴を持っている。

- ・大規模のデータベースに対応できる。
- ・問い合わせに関わる関係のタプルの総和に比べ、最終結果に必要なタプル数が少ない分野で威力を発揮する。
- ・選択的結合やマルチウェイ結合を高速に処理する。
- ・タプルにおける属性値のサイズ(バイト数)の大きい分野で威力を発揮する。

### 7 おわりに

現在研究中のデータベースマシンに関し、その基本思想、構成、機能、演算の流れ、マシンの持つ特徴について述べた。現在、マシンの持つ具体的な性能の評価、本マシンの適する応用分野の検討を実行中である。

#### <参考文献>

- [1] P.Valduriez: "Join Indices", ACM TODS, Vol.12, No.2, pp.218-246 (1987).
- [2] B.C.Desai: "Performance of a Composite Attribute and Join Index", IEEE T-SE, Vol.15, No.2, pp.142-152 (1989).
- [3] 喜連川, 伏見: "データベースマシン", 情報処理, Vol.28, No.1, pp.56-67 (1987).
- [4] 上林: "データベース", 昭晃堂 (1986).