

大規模RDB向きデータベースマシン アーキテクチャに関する一考察

佐藤清実、北村正、中村敏夫、井上潮、武田英昭
(日本電信電話公社 横須賀電気通信研究所)

1. まえがき

近年の情報化社会の発展に伴い、データベース(DB)システムは、利用形態の多様化、大規模化への対応に向けた高機能化、高性能化が切実な要求となってきている。

また、LSI技術の進歩に伴うハードウェア・コストの低下を背景に、DB処理専用のアーキテクチャを有する計算機(データベースマシン; DB-M)の研究も活発に進められている。

DB-Mの研究に際して考慮されるテータ・モデルの多くは、DBC⁽¹⁾、GDS⁽²⁾等一部のものを除き(これらは、テータ・モデルを特定しない)、関係モデル(RDB)である。今後予想されるDBの高水準化動向を考えると、今後とも、RDB処理の高速化が、DB-M研究の重要テーマと言える。

このRDB処理を大きく分類すると、RESTRICTION、PROJECTIONのように単一リレーションに対する演算と、JOIN、DIVISIONのように複数リレーションに対する演算に分けることができる。

これまでに提案されているDB-Mを、この観点から整理すると以下のようになる。

(i) 単一リレーション演算の高速化

① 二次記憶装置の専用化

- ・RAP⁽³⁾ ・CASSM⁽⁴⁾
- ・EDC⁽⁵⁾

② 記憶装置コントローラのインテリジェント化

- ・CAFS⁽⁶⁾ ・DBC
- ・サーチスロセッサ⁽⁷⁾

(ii) 複数リレーション演算の高速化

① ハッシュ化ビットアレイによる 事前ふるい落とし機能

- ・CAFS ・LEECH⁽⁸⁾

② ソータの利用

- ・DSDBC⁽⁹⁾ ・GRACE⁽¹⁰⁾

大規模RDBの実現に当って、これらのDB-Mで使用されている手法を適用することを考えた場合、

- ・ DB規模の増大に伴い、システム・コストの占める磁気ディスク装置等ファイル系コストの割合が増加するため、上記の(i)①のような二次記憶装置を専用化する手法は、システム・コストを増加させるため得策でない。

- ・ (ii)②のようなソータを利用する場合は、ソータの容量が、関係演算の単位となるリレーション・サイズに左右されるため、設定すべきソータ容量が、ハードウェア利用率面から重要である。

このように、従来、DB-Mの研究においては、「大規模化」(DB量、トラフィック等)に対する配慮が、必ずしも主要な議論のテーマとされなかった。

本稿では、リレーション・サイズが数100MB(10⁶タプル程度)、DB量が数10GB程度の大規模システムを対象に、実現性が高いと考えられるDB-Mアーキテクチャを提案する。

本稿で提案するアーキテクチャは、①既存の汎用可動ヘッド磁気ディスクを二次記憶媒体とする、②可能な限り、演算対象となり得ないテータは、事前に削除する機構を充

実させる。

③ 機能プロセッサ群を、同報機能を有するスイッチで結合することにより、ビルディング・ブロックによる柔軟なシステム構成を可能とする。

④ 演算機能を、言語/制御系機能から独立させ、演算系の並列/パイプライン処理による高性能化を図る。

等の特徴を有する。

2. 大規模RDB実現上の課題

ここでは、従来のDB-M研究の動向を踏まえ、大規模RDBをコスト/パフォーマンス的に妥当なものとしてシステム化するために考慮すべき主要な項目を、DB量、リレーション・サイズの両面から整理する。

この他、大規模化に当っては、トラフィック面での考慮が必要となるが、これは双方の項目に含めて述べることにする。

2.1 DB量

DB量の増加に伴いシステムとして重要となるのは、DB実体を格納する二次記憶装置や、コントローラ等のファイル類の量(即ちコスト)である。

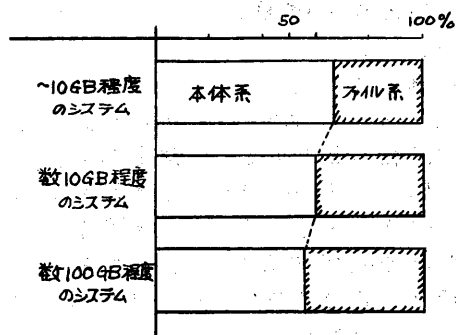


図 2.1 システムコスト構成例

DBシステムのコストを、CPU、メモリ等の本体系コストとファイル系コストに分類し、残ったの大規模システムについて、ラフに試算すると図2.1のようになる。ファイル系コストのシステムコストに占める割合は、DB量の変化とともに、全体の30%~50%程度の範囲となる。

近年、個体ファイルとして、磁気バフル・メモリの用途が進められているが、数10GB以上にも及ぶDB用のバフル・メモリとして使用されることはコスト面から当面なく、磁気ディスク装置(DKU)が主流となると予想される。

DKUは、基本的にDB量に応じた台数が必要となるが、さらに、トラフィック等による性能面からの台数への影響(例えば、アクセス競合を避けるための収容率の低下)が加わる。

従って、DKU等、二次記憶装置そのものへ、RDBのための専用機構を付加することは、ことさらに、ファイル系コストを高くする。

2.2 リレーション・サイズ

リレーション・サイズは、関係演算のためのCPU、メモリ等システムとして高価な処理系リソース量に関係するた

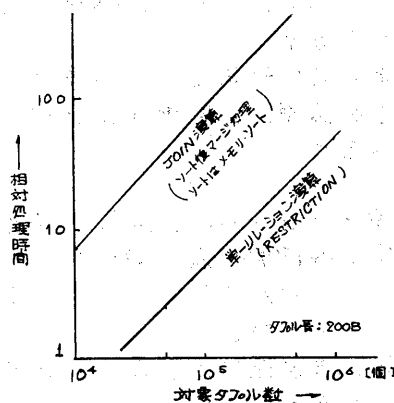


図 2.2 演算の種類と処理時間の傾向 (非インデックスDB時)

め、トラフィック量と合わせて考慮すべき重要な項目である。

特に、JOIN 演算は単一リレーション演算に比べ、1桁以上もの処理時間を要し(図2.2)、それだけシステムでの処理系リソースの使用率が高い。

そのため、JOIN 演算の対象タブルを可能な限り、事前に削除しておくことが、大規模化時に重要となる。

通常、JOIN 演算は、再構成処理時等を除いて、RESTRICTION / PROJECTION のような単純検索と組み合わせられて使用されるので、これらの単純検索を先行して行わせたり、従来から種々提案されているふるい落とし機能⁽⁴⁾⁽⁵⁾を用いることは、有効な手法である。

「リレーションサイズ」に関して、もう1つ重要なことは、一般に、リレーションサイズをシステムとして特定できないため、リレーションサイズの変動に対して柔軟な処理系リソースの確保が必要となる。これは、特に、処理系リソースを多く使用するJOIN 演算時に重要である。

中でも、JOIN 演算の高速化手法として提案されているハードウェア・ソータのハードウェア量は、数100MB(10⁶タブル程度)程度のリレーションを対象とした場合、10³~10⁴クエリ程度にも及ぶと予想される。そのため、大容量ソータを直接利用する方式を、様々なリレーションサイズが存在する現実のシステムに適用することは、リソースの有効利用面から向題となる。

また、リソースの有効利用を図るため、リレーションを分割して、比較的小容量の複数ソータで並列処理させる方式では、ソータへのデータ分配時に、偏りが生ずると、分割した意味がなくなるため、大規模化に当たっては、この分配手法が重要となる。

3. アーキテクチャの概要

3.1. 実現目標

本稿で提案するDB-Mの検討に当たっては、次の目標を設定した。

- (i)リレーションサイズとして数100MB(10⁶程度のタブル数)、システムとして数10GB程度のDBをサポートする。
- (ii)単一リレーションに対する演算、又は、複数リレーションに対する演算のレスポンス・タイムが10秒程度の性能を持つ。
- (iii)トラヒック的には、数100~数1000件/時程度の能力を持つ。

3.2 構成概念

2章で述べた課題の解決に向けて、本稿で提案するDB-Mの構成概念を図3.1に示す。

主たる特徴は、以下の点である。

- (1)DB用のバルク・メモリとして汎用のDKUを使用する。
レスポンスの向上を図るため、リレーションを複数のDKUに分割格納し、並列処理する。
- (2)単一リレーション演算機構(SER)は、複数のDKUで共用することを考慮して、DKUコントローラ相当部分に付加する。
このことにより、DB量に直接関係するDKUの専用化によるコスト増加を避けることが可能となり、かつ、トラヒック見合いの演算機構の設置が可能となる。
この演算機構は、(1)の条件下で並列動作する。

(3) JOIN 演算は、3ステージに分割して行う。即ち

- ① 結合可能性の無いタプルをふるい落としとともに、複数のソータに、リレーションを分割する処理 (FL)
- ② 分割されたリレーションをソートする処理 (SOR)
- ③ 結合処理 (一部の EXC) の3ステージで行う。

タプルのふるい落としでは、両リレーションを並列処理させ、かつ、両リレーション・サイズ差によって生ずるアイドル時間を削減するために、クロス・ハッシュ方式(後述)を採用する。

ソータへのデータ分配法では、従来行われているハッシュ法⁽¹⁰⁾とは異り、統計的処理法を採用する。

(4)各プロセッサは、スイッチで結合する。

従来、結合方式として増設の容易性、コスト面から光ループのような、リング・バスによる結合方式が利用されている。

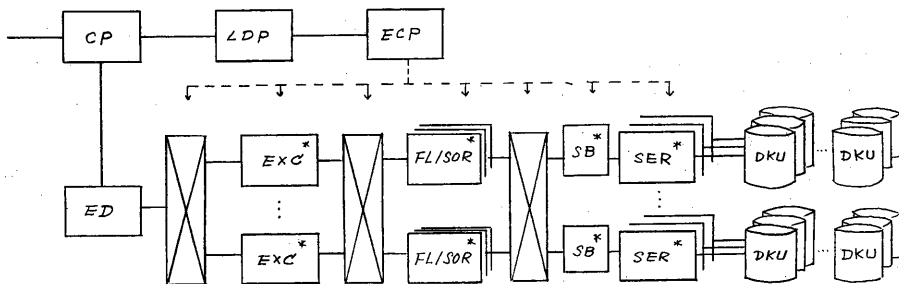
現状の光ループの転送レートは数10 MB/s 程度である。そのため、例え単一リレーション演算機構等により転送データ量が削減されるにしても、リレーション・サイズ数100 MB、数1000件/時のトラフィックを処理するシステムでは、転送路ネックが生ずる。

従って、性能的にスイッチによる結合方式が必要となるが、コストの観点からの考慮が必要である。

(5)言語処理/制御機構と演算機構を分離する。

10⁶タプル程度の大量タプルを処理する場合には、ディレクトリによるマッピング処理、命令の最適化等を行う言語処理時間の全体に占める割合は極めて小さく、数%以下と想定される。

従って、マルチ・トランザクションを処理するシステムでは、言語処理機構は、マルチ・トランザクションで共用し、演算機構は、1トランザクション対応に構成可能とすることが、システム構成上効果的である。



CP : Communication Processor
 LDP : Language & Directory Processor
 ECP : Execution Controller
 ED : Editor

EXC : Executor
 FL/SOR : Filter & Sorter
 SER : Searcher
 SB : Selected data Buffer

図 3.1 構成概念

特に、演算枝構をスイッチ結合とし、各プロセッサ・リソースを演算に先立って、あらかじめスイッチ制御により割当てすることは、テータが、そのパスを通過するだけで良く、タブル単位の切替制御が不要となるため、通信オーバーヘッドの削減とともに、パイプライン処理の容易化が期待できる。

4. 構成要素の概要

ここでは、3章で述べた構成概念の要素のうち、ふるい落としと結合を実現する要素を中心に概要を述べる。

4.1 検索枝構 (SB & SER)

検索枝構は、RESTRICTION / PROJECTIONの単一リレーションに閉じる演算を行うSER (Sercher) と、その演算結果を蓄積するSB (Selected data Buffer) から構成される。

リレーションは、そのテータ量に応じた適当な台数のDKU上に

- ①各DKUのテータ量が極力平均化し、
- ②DKU内では、極力同一シリンダ内に、
- ③同一シリンダ内では、極力連続トラックに、

という管理アルゴリズムで格納される。SERは、DKUの回転速度と同期して各タブルについて演算を施し、その結果得られたタブルをSBに書き込む。

SBの必要量は、あらかじめ予想することができないため、SERの演算条件によって容量不足が生ずることがあるが、その場合は、SBをセグメント化しておき、ECP経由で新しいSBの割付けを受け、テータ転送先を変更することとしている。

4.2 事前ふるい落とし / 分配ソート枝構 (FL / SOR)

FL / SOR (Filter & Sorter) では、JOIN 演算対象タブルに対する事前ふるい落とし枝能、ソート枝能及び、複数台のソータを動作させるためのテータ分配枝能とその分配に必要なテータを収集する統計解析枝能を実現する。(図4.1参照)

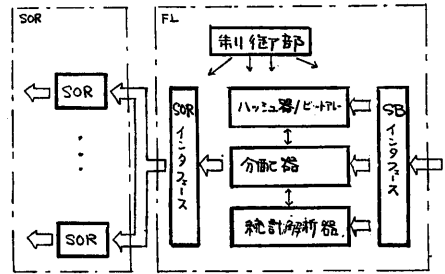


図4.1 FL/SORの枝能ブロック概念

(I) 事前ふるい落とし枝能

JOIN 演算対象リレーションに前処理を施し、演算量を削減する方式としては、従来、ハッシュ化ビットアレイを用いた以下の2方式が提案されている。

(i) 片ハッシュ方式

ハッシュ器、ビットアレイを各々1台用意し、一方のリレーションのみをふるい落とす

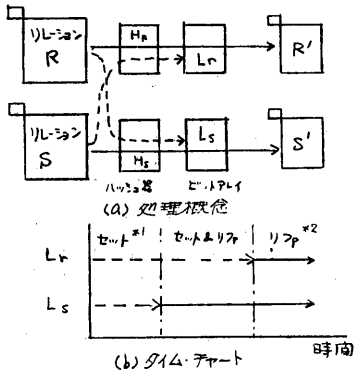
(ii) シリアル・ハッシュ方式

ハッシュ器及び2台のビットアレイを用い、2つのリレーションのふるい落としを逐次に行う。

片ハッシュ方式は、制御が単純であるが、片方のリレーションのみのふるい落としであるため、後段のソート、結合(マージ)の負荷削減の効果が少ないこと、どちらのリレーションをふるい落としの対象とすべきかの迷いが困難であるという問題がある。

シリアル・ハッシュ方式は、処理が逐次のため、ふるい落とし時間が長いという欠点がある。

本稿での方式は、SBとFLを同報機能を持つスイッチで結合し、2つのリレーションのふるい落としを並列に実行するもので、クロス・ハッシュ方式⁽¹⁾と呼んでいる。方式概念を図4.2に示す。



*1: ビット・アレイの該当位置に「1」を立てること
*2: 「セットリフ」を参照し、キーが結合対象かをチェックする。

図4.2 クロスハッシュ方式概念

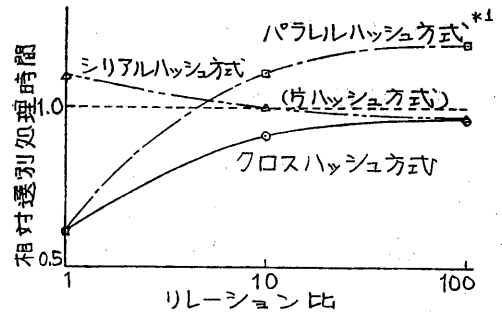
アルゴリズムの概要は、以下のとおりである。(図4.2参照)

- ① リレーション R, S ($|R| < |S|$ とする) のタブルをリレーション対応に並列に読み出し、結合属性をハッシュングして、Rをビット・アレイ L_S に、Sをビット・アレイ L_R にセットする。
- ② Rのセットが終了した時、Sのセット中のタブルを記憶するとともに L_S のモードを切替え、Sの結合属性値を同報機能により2つのハッシュ器に送り、 L_R はセットを続け、 L_S はリファを行う。
- ③ Sのセットが終了した時、 L_R のモードを切替え、Rの全タブルを順次リファし、 L_S はSの先頭からモードを切替えた時、Rのタブルまでをリファする。

リレーションの読み出し開始から、

リファ終了までの時間(差別時間)をシミュレーションにより評価した結果を図4.3に示す。

この結果から、クロス・ハッシュ方式は、従来の方式に比べ、最大40%程度の性能向上を図ることができるとがわかる。



*1: パラレルハッシュ方式とは、セット、リファの切替えを、2つのリレーションで同時に行う。

◎ 処理時間は、プロセッサの処理時間と入出力時間の合計値。

◎ リレーションは、セット時はDKU、リファ時はSBから読み込む。

図4.3 クロスハッシュ方式の効果

(2) 分配ソート機能

複数のソータを用いて並列演算を行うためには、結合属性値により、2つのリレーションのタブルを、それぞれソータに分配する必要がある。このために、従来、高速性を生かしたハッシュ関数を用いた手法が提案されている。この方法は処理が簡単なため、振分け時間が短いという利点がある反面、ハッシュ関数とテータの特性の組合わせによっては、分配結果に偏りが生ずる。

本稿での方式は、ソートの前処理として行う事前ふるい落とし処理が、セットとリファの2段階で行われることを利用するもので、セット時に、結合属性値の統計データを収集しておくことにより、SERでの演算後

のタブルの分布状況が推定できる。そのため、リファ・モードへの移行に先立って、収集した統計データを処理（単純な範囲指定程度の情報を得る処理を施す）した後、分配を行うことにより、後続のソートや結合処理の遅れの影響を極めて小さくして、かつ、偏りの問題を避けることが可能となる。

4.3 結合演算機構 (EXC)

図3.1の構成概念中のEXCは、関係演算の他、論理演算等を専用に行う演算プロセッサの総称であり、ここで述べる結合処理も、一部のEXCで行われる。

タブルは、あらかじめソートされているので、単に、マージするだけで良いように思われるが、同一キー値を持つタブルが複数存在する場合には、結合処理のバック・トラックが必要となり、タブル数に比例した処理時間で終わらなくなる。

特に、JOIN演算の最初のステージでふるい落とし処理を行っているため、同一キー値を持つタブルが集中する可能性が高くなる。そのため、バック・トラックを避けるアルゴリズムが必要となる。

そこで、本稿では、タブルを格納するためのスタックと、スタックの深さのレベル毎にタブル連結器を備えたスタック・マージ方式⁽¹²⁾によりマージを行う。(図4.4、図4.5)

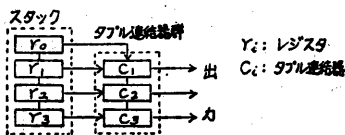


図4.4 スタックとタブル連結器の概念図

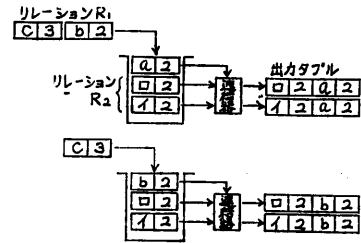


図4.5 スタック・マージ方式の概念

この方式のアルゴリズムは、以下のとおりである。

- ① 2つのソートされたリレーションのタブルを1つつ取出して、キー値の大きさを比較し、キー値が等しくなければ、小さい方のタブルを破棄して次のタブルを取出して比較する。
- ② キー値が等しい場合は、一方のリレーションからタブルを連続して取り出し、異なるキー値が現われる迄スタックに蓄積する。
- ③ スタックに蓄積された複数のタブルと、他方のリレーションの1タブルを、タブル連結器を用いて一斉に結合して出力する。

用意するスタック段数や、リレーションの大きさとスタック対象リレーションの関係をパラメータとして、スタック・マージ方式の性能比較を行った結果を図4.6、図4.7に示す。

この結果から、スタック・マージ方式は、ある程度のスタック段数を準備することにより、重複したキーを持つタブルが多くなっても、比較的安定した処理量となること、リレーションの大きさが異なる場合は、タブル数の小さいリレーションのタブルをスタックに蓄積した方が効率が良いことが判る。

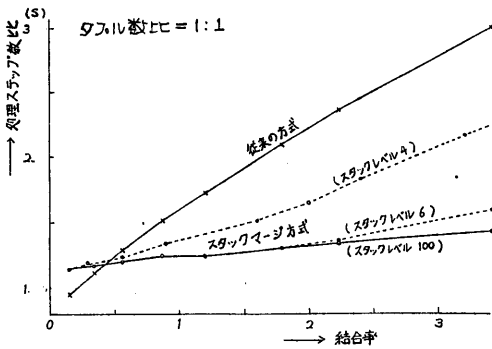
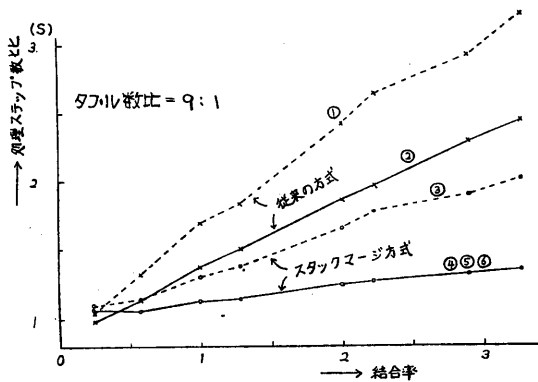


図4.6 スタックマージ方式の性能(リレーションのタプル数が同一の場合)



- ① タプル数の小さいリレーションをマーク
- ② " 大きい " "
- ③ タプル数の大きいリレーションをスタック (スタックレベル6)
- ④ " 小さい " "
- ⑤ " 大きい " " (スタックレベル100)
- ⑥ " 小さい " "

◎性能の尺度として、処理スタック比、タプルのレジスタへのロード、キー値の比較、タプルの連続性等をそれぞれ「コスト」要素としてカウント。

◎タプルのキー値の分布は一様分布とし、同一キー値の重複度合いは、データ結合率 = 結合後のタプル数 / 結合前のタプル数で評価。

図4.7 リレーション入れ替への効果

5. あとがき

本稿では、スイッチ結合を用い、大規模化に向けて柔軟なシステム構成を可能とするアーキテクチャについて、その基本的概念を述べた。

さらに、スイッチにより、データの流制御が容易である3ステージによるJOIN演算方式について論じ、みる

い落ち処理、結合処理のシミュレーションによる評価を中心に、各部の概要を述べた。

本稿で述べたアーキテクチャは、基本構想の段階であり、今後、制御系/更新系処理の詳細化を図り、システムトータルな評価を実施する予定である。

最後に、日頃御指導を頂いている関係各位に感謝致します。

参考文献

- (1) Banerjee J. et al "DBC - A Database Computer for Very Large Databases," IEEE Trans. on Computer, vol.C-28, no.6, pp.414-429, June, 1979
- (2) 箱崎他, "データベースマシン実験システム", 信学会計算機研究会資料, EC-79-68, 1980
- (3) Schuster S.A. et al. "RAP.2-An Associative Processor for Databases and Its Applications," IEEE Trans. on Computers, vol.C-28, no.6, pp.446-458, June 1979
- (4) Su S.Y.W, Lipovski G.J. "CASSM: A Cellular system for very large databases" Proc. VLDB, pp.456-472, 1975
- (5) 植村他, "磁気バブルデータベースマシン", "大型プロジェクトパターン情報処理システム研究成果発表会論文集", pp.263-279, 1980.
- (6) Babb E. "Implementing a relational database by means of specialized hardware", ACM TODS 4.1 pp.1-29, 1979
- (7) 佐藤他, "データベース処理のための専用サーチプロセッサ", 情報論文誌 vol.23, no.4, July, 1982
- (8) McGregor D.R. et al. "High performance hardware for database systems," Syst. for Large Data Bases, North Holland Publishing Co., 1976
- (9) 田中 "データストリーム方式のデータベースコンピュータ", 情報記号処理研究会資料, 12-14, pp.97-103, 1980
- (10) 喜連川他 "HashとSortによる関係代数マシン", 信学技報, EC 81-35, 1981
- (11) 武田他 "並列的い落ちアルゴリズムに関する一考察", 情報第26回全国大会, 4F-13
- (12) 井上他 "大規模データベース向き結合演算方式", 昭和58年度信学会総合全国大会 S1-3