

## データベース処理のハードウェア化

工藤 哲郎

富士通株式会社

汎用コンピュータシステムにおいて、リレーショナル・データベース処理をハードウェア化するための技術について述べる。磁気ディスク制御装置内に、サーチ処理を実行するハードウェア機構を設け、ホストコンピュータの負荷軽減及びデータ転送量の削減を図ったものである。実現にあたっては、従来ソフトウェアにて処理されていたデータベース構造を、直接ハードウェアロジックによって処理することを可能にした。また、複数の処理を並列に実行するとともに、ハードウェア化に適した選択条件判定方式を採用することにより、データベース処理の高速化を実現した。

## *A HARDWARE FOR DATABASE PROCESSING*

*Tetsuro Kudo*

*Fujitsu Limited*

*The method to make the relational database processing hardware in the general purpose computer system is described. We installed the hardware mechanism which executed search operation in the magnetic disk controller to achieved the load reduction of the host computer and the reduction of the amount of the data transfer. In the method, the database structure so far processed by software can be directly processed by the hardware logic. Futher, we adopted the parallel execution of two or more processings and a suitable selection condition judgment method for hardware. As a result, the speed-up of the processing with hardware was achieved.*

## 1. まえがき

近年、ホストコンピュータ上で行われていたデータベース処理を、専用のハードウェアによって実施することにより、ホストコンピュータの負荷軽減及び処理の高速化を図ることが注目されている。

本稿では、ホストコンピュータとデータベースが格納されている磁気ディスク装置との間に、専用のデータベースアシスト機構（以下、DBA）を設け、リレーショナルデータベース処理をハードウェア化するための技術を紹介する。

ハードウェア化を実現する上では、いくつかのデータベース機能が想定出来るが、今回は一般的で使用頻度の高く、効果が期待出来るサーチ関連機能の実現について述べる。

## 2. システム構成

図1にDBAのシステムにおける位置付けを示す。

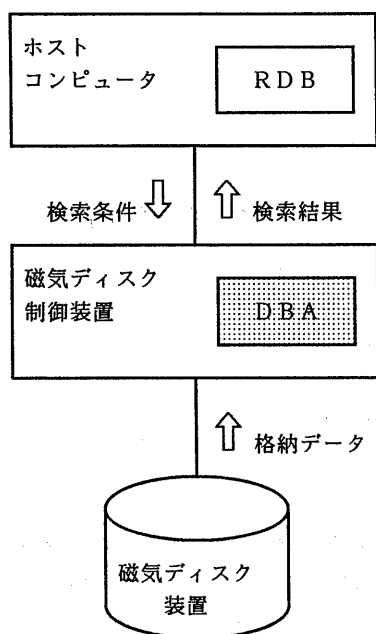


図1 システム構成

DBAは、ホストコンピュータと磁気ディスク装置の間に位置する磁気ディスク制御装置に搭載され、ホスト上に常駐するリレーショナルデータベース・プログラム（以下、RDB）と連携して動作する。

磁気ディスク制御装置内にデータベース処理機構を搭載することは、1つの磁気ディスクサブシステムにおいて、一般のI/O処理、DBAを使用したデータベース処理及びDBAを使用しないデータベース処理を効率良く混在制御することが可能となる。また、独立した装置を設ける場合に比べて、外部インタフェース制御部を初めとする物理資源の共用が可能になるとともに、物理スペース面でのメリットが大きいと言える。

但し、一般のI/O処理との同時制御を行うという観点より、制御の複雑さ及び一般処理への悪影響等を解決する必要がある。

## 3. ハードウェア化機能

今回ハードウェア化の対象としたサーチ関連機能は以下の3つである。

- ・選択：条件に合致するレコードの選択
- ・射影：レコードからの指定カラムの抽出
- ・集約：統計的集計を行い、その結果を通知

上記機能は、磁気ディスク装置から読み出される大量のデータをオンザフライで順次処理することが可能であるため、ハードウェア化には適した機能であると言える。また、ホストコンピュータに送出する際に、膨大なデータ量を大幅に削減することが可能となり、システム効率上のメリットも大きい。

データベースの形式は、DBAの有無に係わらずデータベース構造の共通化を図るために、特にハードウェア化を意識した構造とはしなかった。

したがって、ページ形式はもっとも一般的なスロット方式を採用し、レコード形式は、IDや長さを表す固定部を先頭に、固定長部、可変長部、さらに、各カラムにIDを有するカラムワイズ部により構成される構造とした。

(1) 述語比較

述語判定は，“＝”，“<”，“>”及びこれらの否定である“≠”，“≥”，“≤”を基本とし，文字タイプでは部分文字列比較（％），任意文字指定（＿）及びE S C A P E文字をサポートした。

(2) データ形式

データ形式については，以下のものをサポート対象とした。

- ・文字形式（1，2バイト）
- ・I N T E G E R形式
- ・D E C I M A L形式
- ・N U M E R I C形式
- ・F L O A T形式

(3) 選択条件

選択条件判定は，述語のA N D / O Rによる任意の論理式にもとづいて判定する。

(4) ナル処理

目的のカラムが存在しない場合やN U L Lタグ値がゼロで無い場合には，カラムをナルと見なした上，述語比較及びカラム抽出を実施する。

(5) カラム抽出

カラム抽出では，指定された任意のカラムを抽出し，新たなレコードを作成する機能をサポートした。

(6) 集計処理

集計処理については，取り合えずもっとも効果の期待出来る以下の機能を実現した。

- ・C O U N T（＊）
- ・C O U N T（項目）
- ・S U M（項目）

但し，F L O A T形式については，精度の問題があり，S U M（項目）の実現は見送った。

#### 4. ソフトウェア・インタフェース

ハードウェアにてデータベース処理を実施するためのソフトウェアとのインタフェースを下記に述べる。

(1) カラム情報

サーチ処理を実施する場合には，ハードウェアにおいてデータベースをカラム単位まで解読する必要がある。レコード単位の処理はページ内に格納されたレコードポインタを参照することにより可能であるが，カラム単位の処理を実施するためには，カラム属性，データ形式及びデータ長等を表す情報が必要となる。

したがって，予めR D Bから受取った複数のカラム情報を装置内に格納しておき，データベース読出し時に，対応するカラム情報と照合することによって，カラム単位の識別を可能とした。

(2) セレクション情報

述語判定を行うためには，比較対象カラム毎に，比較演算子，比較データ，比較データタイプ，位取り等の情報が必要である。さらに，選択条件判定のためには，述語のA N D / O Rによる論理式を何らかの形でR D Bより受け取る必要がある。今回は，D B A処理の依頼時に，セレクション関連の情報を受け取ることとした。

論理式の条件判定は，ハードウェアに適した高速処理方式を検討した結果，事前にC N F / D N F形式に展開されたものを受け取りそれをハードウェアにて処理する方式を採用した。

(3) プロジェクション情報

カラムの抽出処理を実施するためには，セレクション情報同様，R D Bより事前に抽出条件を受け取る必要がある。

固定長部及び可変長部については，格納順に並べられたビットマップにより抽出条件を指定し，カラムワイズ部については，抽出対象となるカラムI Dを指定することとした。

尚、集計処理の対象カラム指定は、インタフェース情報量を削減する目的で、プロジェクト情報を併用することとした。

(4) 処理結果通知

DBAの有無に係わらず、RDB側での処理の共通化を図るために、サーチ結果の通知形式は、格納時と同一な構造とした。したがって、DBA側は処理結果を、スロット方式により構成されるページ形式にまとめた後、ページ単位にRDBへ送出する。

5. ハードウェア構成

図2にディスク制御装置の内部構成図を示す。チャンネルアダプタ(CA)はホストとのコミュニケーションを分担するモジュールであり、デバイスアダプタ(DA)は磁気ディスク装置を制御するモジュールである。

シェアドストレージ(SS)はキャッシュ等を使用される大容量メモリであり、カラム属性等を定義するカラム情報はSS内に保持されている。

一般のI/O処理及びDBAを使用しない処理では、DA⇔SS⇔CAを経由してホストコンピュータへデータを送出するが、DBA使用時にはSSに格納されたデータがDBAにおいて一度加工された後、ホストコンピュータへ送られる。

一つのディスク制御装置において、同時に複数のデータベース処理が実行可能なように、本ディスク制御装置では複数のDBAが搭載可能な構造をとっている。

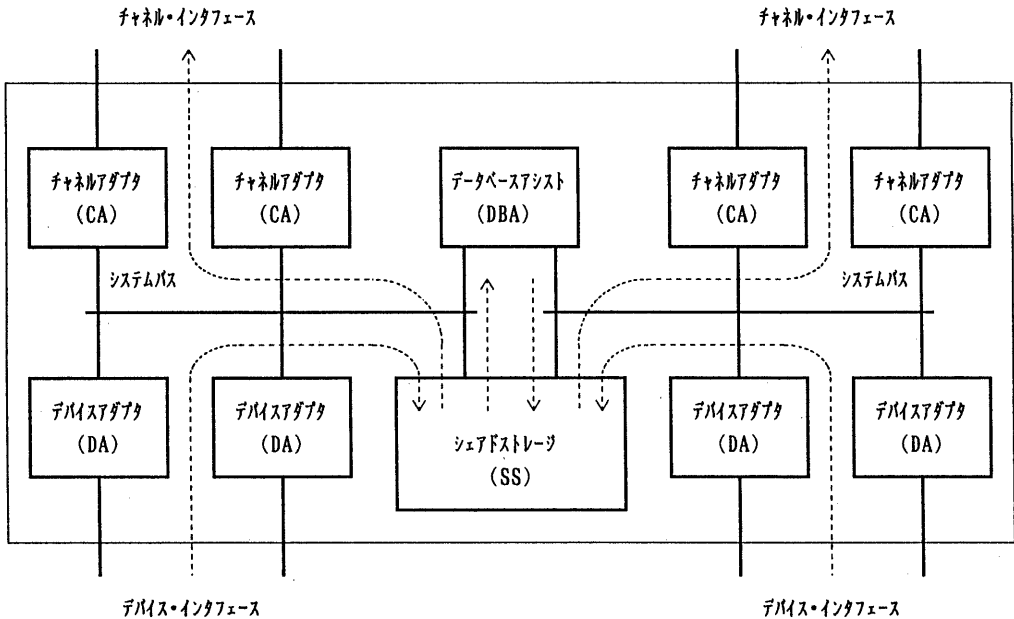


図2 ディスク制御装置内部構成

DBAの内部は、図3に示すように以下の論理ブロックにより構成されている。

- ・入力ページバッファ
- ・出力ページバッファ
- ・解析ロジック
- ・抽出ロジック
- ・格納ロジック
- ・選択ロジック
- ・集計ロジック

(1) 入力ページバッファ

サーチ処理を実施すべき入力データが、システムバスを經由してページ単位に格納されるバッファであり、最大64Kバイト分の入力ページが格納可能である。

入力ページ以外にも入力データに対応したカラム情報が最大64Kバイト格納される。

(2) 出力ページバッファ

サーチ処理結果である出力データがページ単位に格納されるバッファであり、最大64Kバイト分の出力ページが格納可能である。

出力ページ以外にも対象となるサーチ処理に対応したセレクション情報、プロジェクション情報及び集計結果が格納される。

(3) 解析ロジック

解析ロジックは、バッファ内のカラム情報をもとに、バッファに格納された入力ページについてレコード及びカラムの解析処理を行う。解析された入力ページはレコード単位、カラム単位に分解され抽出ロジックへ送られる。

また、カラム内の実データが格納されているアドレス等、カラム比較処理に必要なカラム関連情報を選択ロジックへ送付する。

(4) 抽出ロジック

抽出ロジックは、出力ページバッファ内のプロジェクション情報をもとに、解析ロジックより受け取ったカラムの中から抽出対象となるカラムを選択し、格納ロジックへ送付する。

(5) 格納ロジック

格納ロジックは、抽出ロジックより受け取ったカラムを順次出力ページに格納し、新たなレコードを作成する。

さらに、選択ロジックより通知されたレコードの選択結果をもとに、選択対象となるレコードで新たなページを出力ページバッファ内に作成する。

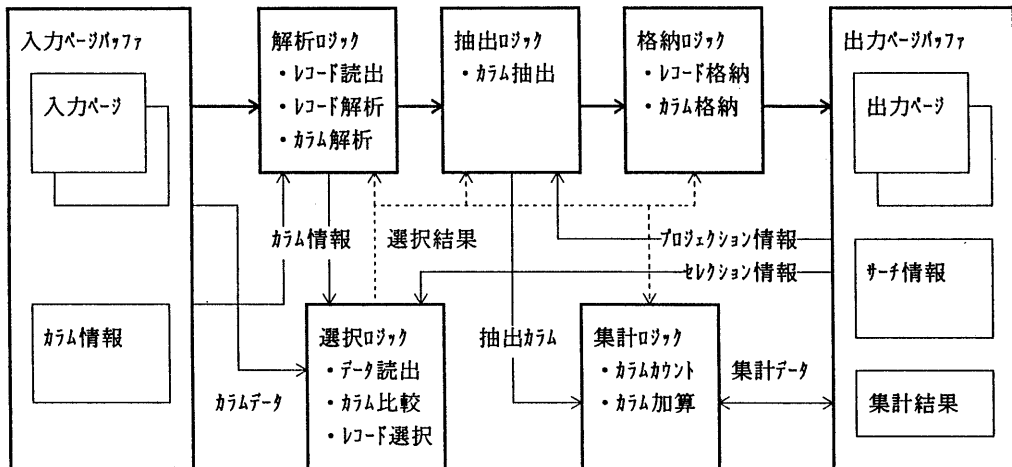


図3 DBA内部構成

(6) 選択ロジック

選択ロジックは、出力ページバッファ内のセクション情報と解析ロジックより受け取ったカラム関連情報をもとに、入力ページ中のカラム内データの比較処理を行い、レコードの選択結果を他のロジックへ通知する。

(7) 集計ロジック

選択対象となったレコードについて、カラム数のカウント処理及びカラム内容の加算処理を行い、集計結果を出力ページバッファ内に格納する。

図4にDBA内部の動作シーケンス例を示す。解析、抽出、格納、選択、集計といった各々の処理動作は、レコード単位に並列に実行される。

本例では、レコードN及びN+2は選択条件を満足したレコードであり、格納処理及び選択処理が完了するのを待って、次のレコード処理へ移行する。レコードN+1は選択条件を満足しないレコードであり、処理途中に選択対象外であることが判明した場合には、そのレコードに対する全ての処理を速やかに中断して次のレコード処理へ移行する。また、集計処理は、抽出処理及び選択処理の完了後、動作が開始される。

各ロジック部は、物理的に3つのLSIにより構成され、システムバスとのインタフェース制御ロジックや全体動作を管理するマイクロプロセッサを含めて、1枚のPCB (Printed Circuit Board) に搭載されている。

表1に各ロジック部のLSIテクノロジーを示す。

表1 LSIテクノロジー

プロセス		富士通製CMOS-LSI (0.8μ, メタル2層)
ゲート数	解析ロジック部	23,894ゲート
	抽出ロジック部 + 格納ロジック部	22,101ゲート
	選択ロジック部 + 集計ロジック部	48,827ゲート + 2KビットRAM
パッケージ		256ピン, PGA

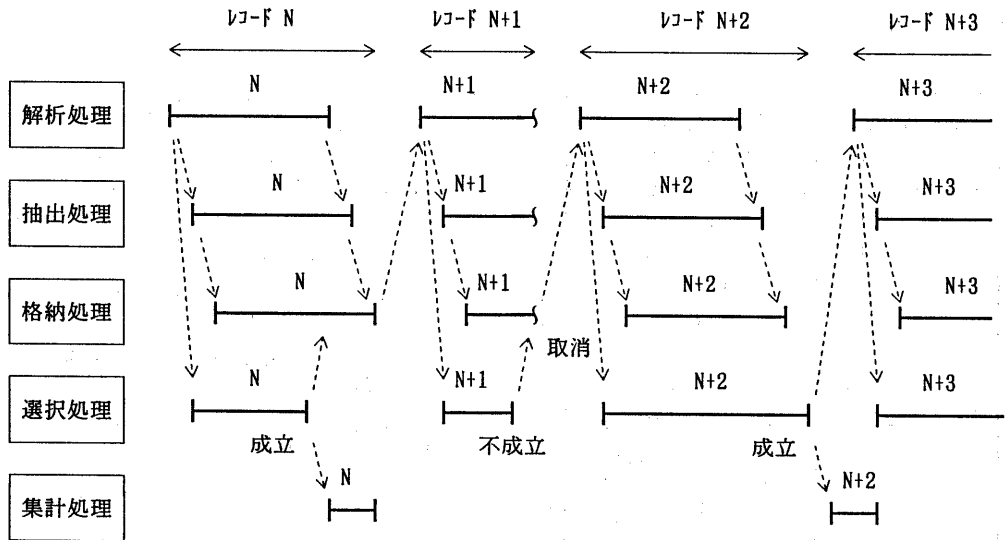


図4 DBA内部動作シーケンス例

## 6. 選択条件判定方式

レコードの選択処理では、ハードウェア回路による高速処理を実現するために、論理式を事前にCNF/DNF形式に展開したものを処理する方式を採用した。但し、一般的にソフトウェアで使用されている方式は、ハードウェアには不向きである。本章では、ハードウェア化に適したレコードの選択条件判定方式について述べる。

(1)は簡単な条件式の例であり、(1)をCNF形式に展開したものが(2)である。

(a1=1989 & a2>10 or a1>1989) &  
(a1=1991 & a2<7 or a1<1992) &  
(a3=XXX )

(1)

(a1=1989 or a1>1989) & ..... ①  
(a2>10 or a1>1989) & ..... ②  
(a1=1991 or a1<1992) & ..... ③  
(a2<7 or a1<1992) & ..... ④  
(a3=XXX ) ..... ⑤

(2)

表2は、(2)の式を実行形式に変換した表であるが、ソフトウェアによる処理を前提とした一般的な構造である。CNF形式では、少なくともある1つのジャンクションについての条件が不成立の

表2 条件判定実行表 (ソフトウェア用)

CNF/DNF種別	ジャンクション数 → 5	
ジャンクション・ポインタ-1	-----	
ジャンクション・ポインタ-2	-----	
ジャンクション・ポインタ-3	-----	
ジャンクション・ポインタ-4	-----	
ジャンクション・ポインタ-5	-----	
a1項目	=条件	1989
a1項目	>条件	1989
a2項目	>条件	10
a1項目	>条件	1989
a1項目	=条件	1991
a1項目	<条件	1992
a2項目	<条件	7
a1項目	<条件	1992
a3項目	=条件	XXX

場合に、全体が条件不成立となる。したがって、ソフトウェアでは、1つ1つのジャンクションを順次処理していく方式が一般的である。

但し、この場合、項目の出現が順不同となり、前のカラムに戻って処理をする必要が出てくる。したがって、これをハードウェアで処理するためには、各カラムの格納位置を保存しておき、再度同一カラムを読み出す等の処理が必要となる。これはハードウェアの物量及び高速化の観点より、ハードウェア化には適していないと言える。

表3は、ハードウェアによる処理を前提とした実行形式表である。本表は、同一項目に関するものをまとめた構造をとっている。この構造のもとでは、項目順に比較処理を行い、その結果を対応したジャンクションに反映していくことになる。そして、CNFの場合、全ジャンクションのAND条件を常に監視し選択条件を判定する。

本方式では、CNF形式の場合、全てのジャンクションの条件が成立した場合にのみ、全体の条件が成立となる。したがって、本来のCNF/DNF方式の目的とは異なってくる。しかしながらジャンクション内は少なくとも1つの条件が成立すれば十分であるため、条件式によっては、本方式の方が速いケースもある。また、もっとも重要なポイントは、本方式では、格納順にカラムを処理していくことが可能になるとともに、1つのカラムについて同時に複数の比較処理を実行することも可能になる点である。

表3 条件判定実行表 (ハードウェア用)

CNF/DNF種別	ジャンクション数 → 5	
a1項目	比較データ数	
=条件	1989	ジャンクション番号 → 1
>条件	1989	ジャンクション番号 → 1
>条件	1989	ジャンクション番号 → 2
=条件	1991	ジャンクション番号 → 3
<条件	1992	ジャンクション番号 → 3
<条件	1992	ジャンクション番号 → 4
a2項目	比較データ数	
>条件	10	ジャンクション番号 → 2
<条件	7	ジャンクション番号 → 4
a3項目	比較データ数	
=条件	XXX	ジャンクション番号 → 5

## 7. ハードウェア性能

DBAのデータ処理能力は、最低限磁気ディスク装置のデータ転送速度を満足する必要がある。現在、一般的な大型磁気ディスク装置の最大データ転送速度は4.5MB/Sであるが、磁気ディスク装置のテクノロジーの進歩、ディスクアレイを初めとする新しいアーキテクチャの確立により、データ転送速度は飛躍的に向上する可能性がある。したがって、処理速度は平均10MB/S以上を目標値とした。また、十分な処理能力があれば、1つのDBAで複数の処理を時分割的に実施することも可能になる。

以下に、DBA単体としての性能評価結果を示す。表4は使用したデータベースの仕様であり、表5は、表4のデータベースをもとにした処理速度である。

測定結果からは、プロジェクション条件は処理速度には影響しないことが判る。これは、サーチ条件に係わらず、全カラムについて解析処理を実施するためである。処理速度に影響する条件は、

表4 評価用データベース仕様

ページ長	4 Kバイト
ページ数	5
レコード数	10 (ページ当たり)
カラム数	100 (レコード当たり)
データ形式	文字, 数値混在

## 8. まとめ

DBAの単体性能は、当初の目標値をクリアしているが、システムとして十分な性能値を得られるかが、今後のポイントである。また、同一装置内における一般I/O動作との混在動作時や複数のDBA動作時の性能評価も必要である。

ハードウェア化については、今後一層大規模な

実際に比較されるカラムの数である。今回の評価では意図的に全ての比較対象カラムに対する処理が完了しないと全体の条件が判定できない設定としたが、実際には全ての比較対象カラムを処理する必要が有るケースは少ないため、実際の処理速度は表5に示した結果よりも若干向上すると言える。

DBA内部での処理速度は、セレクションの条件によって大きく左右されるが、一般的なサーチ条件下では、9~12MB/Sの範囲内に収まると言える。

表5 処理速度測定結果

サーチ条件		測定結果	
プロジェクション対象カラム	セレクション対象カラム	処理時間 (msec)	処理速度 (MB/s)
全カラム	条件無し	0.97	20.6
3カラム	条件無し	0.99	20.2
全カラム	5カラム	1.68	11.9
3カラム	5カラム	1.65	12.1
3カラム	10カラム	2.15	9.3
3カラム	全カラム	3.64	5.5

LSIが使用可能になるにつれて、機能拡張も可能となってくるであろう。また、ハードウェアの弱点である論理変更の困難性については、今後、プログラマブル・ゲートアレー等を活用することにより、データベース構造や機能の変更に対しても、柔軟に対応していくことが可能である。