

# データベース管理システム向き計算機アーキテクチャについて

ON THE DATA-BASE-MANAGEMENT-SYSTEM-ORIENTED COMPUTER ARCHITECTURE

京都大学工学部 渡 辺 勝 正 萩 原 宏

KYOTO UNIVERSITY

Katsumasa WATANABE

Hiroshi HAGIWARA

現在広く利用されている計算機 (Conventional Computer) は、必ずしも、データベース管理システム (DBMS) あるいは情報検索システムに適しているとはいえない。本論文では、たとえば、記憶装置内におけるデータ構造の実現、可変長データの操作、並列アクセス法の問題をとりあげ、ソフトウェア、ファームウェア、ハードウェアを含めた総合的な面から、大規模 DBMS に適した計算機アーキテクチャを考へ、DBMS の効率向上をはかるとともに、DBMS の問題の本質を追究する。

## 1. データベース管理システムにおける問題

データベース管理システム (DBMS) あるいは情報検索システム (以後とくに区別をせずに DBMS と総称する) の役割は、データの集りのある有機的な関係のもとで定義し、そのデータの集りからある条件に合ったデータを抽出することであり、そのシステムは図1のように示される。

このようなシステムを実現するために、いくつかの問題がある。

### 1) DBMS 向き命令体系

DDL および DML における基本操作は、多くの場合、(数値計算に向いている) 計算機の機械語命令と大まかに似ている。しかも、DDL および DML は、より水準の高い (あるいは、非手続き的な) 言語にたいする傾向にある。これらの言語翻訳の問題と関連して、DBMS 向き計算機の命令体

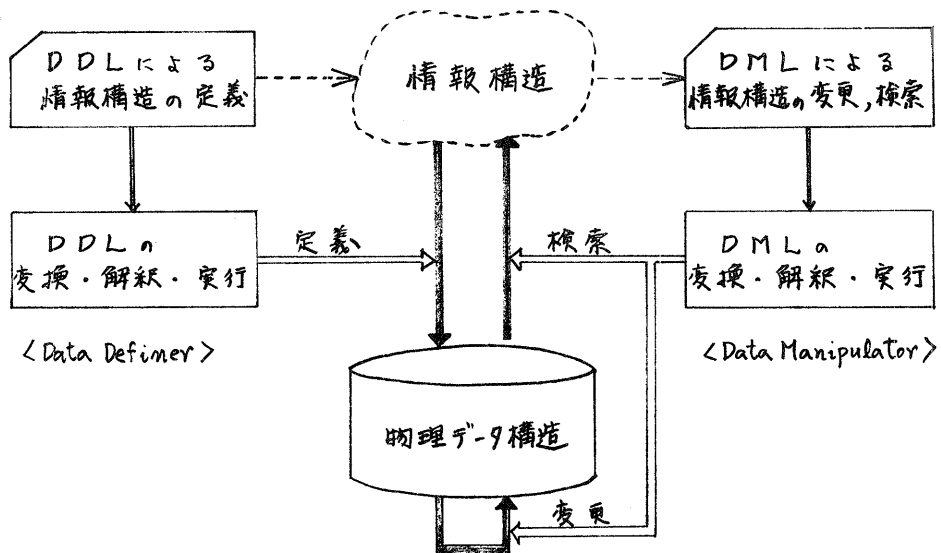


図1. データベース管理システムの構成

系はどのようなものが望ましいかを考える必要がある。

## 2) 記憶の構造

ユーザが頭の中に描いているデータの情報構造と、記憶装置の中で実現されるデータの物理構造とは、必ずしも一致しない。ユーザはその応用に応じて様々な情報構造を定義するが、計算機の中では、1次記憶装置、2次記憶装置の構造に応じて、それを実現せざるを得ない。このとき、アクセスの効率を考慮に入れて、ブロックに分割したり、ある種の問い合わせに対してすばやく応答できるように、データ構造を再編したりすることがある。情報構造と物理構造との間のデータ構造変換の問題と関連して、どのような記憶の仕方がよいか、そのための、どのような記憶装置の機構が必要かを求めねばならない。

## 3) データ転送量

取り扱われるデータベースの大きさは、1次記憶で記憶できるデータ量よりはるかに多いため、データベースを適当に分割して2次記憶に記憶しねばならない。様々なタイプの問い合わせを扱って付けたDBMSにおいては、ある種の問い合わせに対して、ある特定のページだけを参照すれば良いようにページを構成することは困難であるので、(しかも、1ページ内にその問い合わせに対して必要な情報は無い場合が多いので)、多くのページを1次記憶に転送しねばならない。このようなとき、大容量の1次記憶と高速転送能力をもつデータチャネルを有する大型計算機を望むよりも、本当に必要なページのみを指摘したりして、2次記憶に検索能力をもたせたりして、転送量を減らすように工夫することの方が好ましいと考えられる。

## 4) データベースの更新

データベースがいくら大きくても、それが固定されたものであるならば、

あらかじめ、各々の応用あるいは問い合わせの型に応じたファイルを別個に準備しておけば、応答時間を短くすることが出来る。この場合には、一つのデータがいくつものファイルの中に異なったキーを持ったレコードとして記録されることがある。このようなときデータベースが変更されると、変更データに関連したファイルを見つけ、それらのいくつものファイルを改めねばならない。一部のファイルを改めた時点で問い合わせがくると、データベース内に不一致が生じた好ましくない状態となり、問い合わせを一時待ち合わせる必要が生じる。データベースの更新・拡張に耐え、しかも、柔軟なアクセス法をもつ物理データ構造の実現 (inverted file, chained organizationとは異なった自由なデータ構造の実現) が必要となる。

このほか、データの入力方法、データの保護、データベースに関する記録等々の多くの問題がある。これらも含めて、上記の問題は、単にソフトウェア、アルゴリズムの検討といった面だけの問題ではない。以下では、これらの問題を頭に入れて、ソフトウェア、ハードウェアの全体的な面から、DBMS向き計算機のアーキテクチャについて考える。

DBMSの役割りの増大をうけて、与えられた計算機システムの上でDBMSを実現する努力とともに、DBMSにおける基本操作および本質的な概念をとらえて、新しい計算機アーキテクチャを求めたことは重要なことである。

## 2. ソフトウェア、ファームウェア、ハードウェアの役割り

新しい計算機システムを設計・製作するにあたって、ある機能を実現するのにソフトウェア、ファームウェア、

ハードウェアのいづれによるかを決定することは重要な問題である。それは機能の性格によって一意的に定まる場合もあるが、実現する技術的背景とシステムの汎用性・柔軟性に対する考え方や大きく影響される。一般に、そのような影響のもとに、1つの計算機システムにおいて、ある機能をいづれで実現するかという傾向は、歴史的に下記のようになっているとみられる。

- 1) 計算機が用いられた当初は、まさにハードウェアそのものであった。
- 2) ハードウェアの高速・高性能化に伴って、汎用という目標のもとに、システムプログラム、ソフトウェアパッケージとしてのソフトウェアの増大がみられる。
- 3) マイクロプログラミング方式の採用により、応用に依りて命令体系を変換し、可変制御構造をもった計算機として、ハードウェアの機能をファームウェアで実現した。
- 4) 汎用性の利用と共にソフトウェアの比率が増々大きくなったが、ICによる書き換え可能な制御記憶装置の実現が容易になって、処理速度の向上、ソフトウェアの単純化(高水準言語計算機、DBのファームウェア化)をめざして、ファームウェア化する傾向がみられる。
- 5) LSI技術の発展につれて、回路の標準化、並列性の利用と行ったことから、再びハードウェアを中心に考え、ソフトウェアの比率の増大をふせこうとする努力がみられる。

以上の傾向をまよひて考えると、ソフトウェアは我々が問題を考え、それを表現しやすいように支援し、ハードウェアは、回路および設計技術からの制限を受けながら、汎用性・融通性ということにも留意して、ソフトウェアを実現するために最適な基本要素を提供し、ファームウェアは、両者の間において、基本要素の組み合わせによっ

て“表現”を実現し、システムの柔軟性と設計の簡易性を保つ役割りを果たしているものといえる。

したがって、ソフトウェア・ファームウェア・ハードウェア方式では、一般に、基本的なハードウェアに対して、プログラム表現に近い機械語命令をファームウェアによって解釈することで、ソフトウェアに依存した仮想機械を柔軟に構成することが出来る。このように、ソフトウェアによる(外部に於ける)アルゴリズムの表現とファームウェアに対する(内部に於ける)表現とが近接していることは、それらの間の変換が容易であり、命令数が少なくて命令取り出しの回数が減少するという利点だけでなく、我々が計算機構成の頭の中に描きながらアルゴリズムを考へることが出来るという利点をもつ。

これに対して、ソフトウェア・ハードウェアのみの方式では、それらの間のギャップを、いづれかの歩み寄りによってうめるか、トランスレータによって、ソフトウェアによる表現をハードウェア向きに変換しなければならぬ。とくに、DBMS用のソフトウェアと、汎用計算機のハードウェアについては、両者の間のギャップが(たとえば、記号列操作、データ集合の演算、可変長データの取り扱い等の面で)大きいと感じられる。

### 3. DBMSにおけるソフトウェア・ファームウェア・ハードウェア

DBMSにおける問題に対する対策を、それぞれソフトウェア、ファームウェア、ハードウェアを主体とした面から、いくつかの事例<sup>[1, 2, 3]</sup>を参考にしてみよう。これらの事例の中には、後置プロセサ(back end processor)方式として、データベース管理用あるいは情報検索用の専用プロセサを採用するものがあるが、それは汎用計算機

システムの機能を分離したものであって、専用プロセサにおける問題は、結局は、フジのハブメカに関連づけられる。

3. i) ファイルディレクトリの強化  
 与えられた計算機システムの中で、可動ハードディスクのヘッドの動きを少なくし、しかも、2次記憶から1次記憶に転送されるデータ量を最少におよぼすことができる。ファイルの中からある条件を満たすレコードを探し出すには、  
 i) 条件を満たすレコードを含めようと思われたデータブロックを抽出する、  
 ii) そのブロックの中で条件を満たすレコードを順次調べ、  
 という考えに基づく[4]。

そのため、フジのようく、与えられた“条件”から“必要なブロック”が決定できるように、ファイルディレクトリを強化するとともに、同じ条件を満たすレコードをできるだけ同じブロックに含めるようにレコードを記憶する。

レコードは簡単に(属性、値)対の集合であるとする。レコードの集りをファイルとよび、ファイルがいくつか集ってデータベースを構成するものとする。キーはレコードを特徴づける属性-値の対で、問い合わせは、キーの関係を表わされる。

一才、可動ハードディスクのインデックスを検索単位(retrieval unit)とよび、1つの検索単位の中にいくつかのブロックがあるものとする。r番目の検索単位の中の名前bのブロックを、  
 $(r, b)$

で表わし、物理セルという。

ファイルFとキー $K_i$ に対して、物理セルの集合

$$D(F, K_i) = \{ (r, b) \mid \text{ファイルFのレコードで、キー} K_i \text{をもつものが、少なくとも一つ、物理セル}(r, b) \text{に含まれる} \}$$

を定義し、Fに対するキー $K_i$ のディレクトリ・エントリーとよぶ。すなわち、キー $K_i$ が与えられると、そのディレクトリ・エントリーを見て、キー $K_i$ をもつレコードが少くとも一つ含まれている物理セル $(r, b)$ をすべて知ることが出来る。ファイルFのディレクトリは、F上のすべてのレコードに対するキーのすべてのエントリーから成るものである。

同じキーをもつレコードは同一のセル内にあるいは同じ検索単位の中に記憶されるようにファイルを構成する。

問い合わせQに対して、Qを満たすレコードはフジの手順で計算される。

i) Qの中のあるすべてのキー $K_i$ に対する物理セルの集合の組を抽出する。

$$\{ D(F, K_1), \dots, D(F, K_m) \}$$

ii) Qにおけるキー間の関係元から集合演算によってQを満たす物理セルの集合を決定する。

$$D(F, Q) = R [ D(F, K_1), \dots, D(F, K_m) ]$$

iii) 得られた物理セルを1次記憶にとり出して、各レコードを順次調べ、Qを満たすものを選び出す。

ファイルFにキー $(K_1, \dots, K_j)$ をもつレコードを追加する場合kは、

i)  $D(F, K_1), \dots, D(F, K_j)$ を求め、

ii) これらの集合の中で最も多く共通に含まれている物理セル $(r, b)$ を決定し、

$$(r, b) = \text{Max}_{\substack{r, b \\ (r, b) \in D(F, K_i)}} [ (r, b) \in D(F, K_i) ]$$

iii) その物理セルにこの新たなレコードを記憶する。

この方式により、はじめにかかげた目標は達成されるが、ファイルディレクトリに大きなスペースが必要となる。ただし、当然のことながら、ディレクトリのスペースは、ファイルのスペースに比べて少なくてよい。その大ききの割合は、1レコードの長さ、ブロックの大きさ、ファイル中の異ったキーの数、レコードの関係の複雑さによる。

3. 2 中央処理装置の強化

フォームウェアあるいはハードウェアによって、記号処理向きの命令や理想処理機能を実現して、データ処理能力を向上させようとするものである。

(1) よく使われる手続きを命令として実現する

記号列の比較、転送、挿入、削除、変更、さらには、サーチ、ソート等記号処理に頻繁に使われる手続きを命令として与える。これによって、命令取り出し回数が減少するとともに、機械語命令とユーザの用いた操作との関係が緊密になり、言語の翻訳が簡単になる。

(2) 可変長データの取り扱い

名前(ファイル名、フィールド名)、定数(10進数、記号列)およびフィールドの値は、一般に、その長さが様々であるが、これをそのままバイト単位の可変長データとして取り扱い、データの記憶形式をユーザのみぞ形式(少しでも)近づける。それとともに、命令形式も可変長のオペランドを持つ可変長の形式とする。したがって、たとえば、DDMにあるいはDDMで書かれたプログラム中の名前、定数は命令の中ではそのまま直接オペランドとして使われる。

命令コード { 長さ・名前 記号 定数 }\*

オペランドの個数は、個々の命令に依存する。

(3) 記号処理向きのアキュムレータ

記号処理(あるいは非数値処理)のための編集領域(アキュムレータ)を連想記憶または(8ビット並列の)高速レジスタで実現して、前の2層をハードウェアの層からサポートする。

(4) ディレクトリの階層構成

データアクセスのための辞書を階層的に構成して、データ構造との対応を明確にしたり、データの記憶位置の決定を速くしたりする。これには2つの

応用が考えられる。

その一つは、名前(ファイル名、フィールド名等)をデータ構造の階層構成にあわせて登録し、名前の階層関係を明らかにする。これは、2層と関連して、可変長の名前をそのまま扱う場合に重要となる。

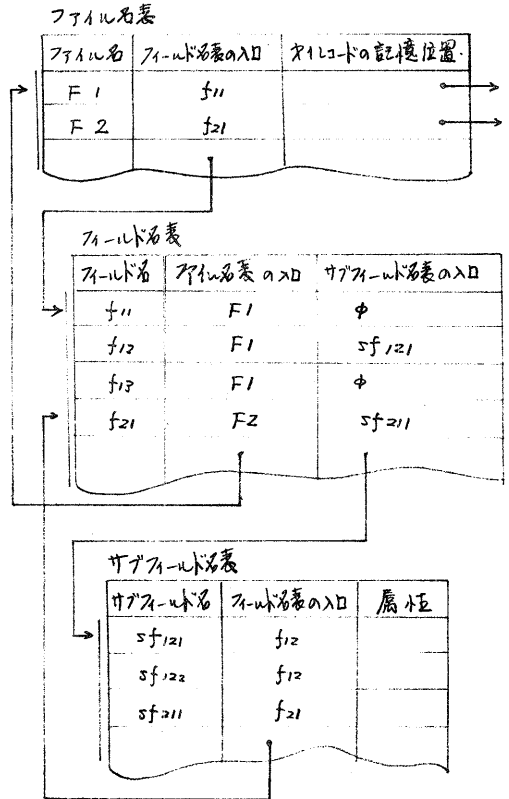


図2. 名前の階層構成

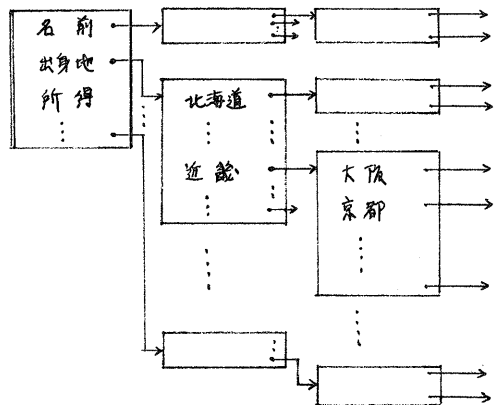


図3. キーの階層構成

もう一つは、ディスクの物理セルを  
つまとめるに当って、キーの階層構造  
に従ってキーディレクトリを構成し、  
所要の記憶位置を決定するものである。

オアおよびオア4項の莫から、ある程  
度の容量をもつ本格的な連想記憶装置  
の実現が得られる。

その他、すでに実現されているもの  
もあり、DBMSにおいてソフトウェア  
エディタではあまりにもオーバーヘッド  
が大きく、ぜひファームウェア、ハード  
ウェアによる支援が望まれるものに  
つぎのものが考えられる、

記憶の割付け

ガルバーシ・コレクション

データの保護機構。

### 3.3 ファイル装置の強化

ファイル装置(ディスクあるいはドラ  
ム記憶装置)を、単にデータの記憶  
抽出を行うための補助記憶装置にとど  
めおかず、データの検索・照合機能  
をもたせて、二次記憶から一次記憶に  
送られるデータの量を減らすことを考  
える。検索・照合にあたっては、1語  
ずつを直列に行うのではなく、ディス  
クやドラムの構造からびにその読み書き  
の速度を考慮して、トラック毎に並列  
に行わせる(Logic per Track Device)。

ディスクやドラムのヘッドの各々に  
マイクロプロセサを配置して、与えら  
れたキーに対して、同時にサーチすると、  
ドラムや固定ヘッドディスクの場合には、  
1回転の時間で全領域をサーチする  
ことができる。サーチの結果、必要ならば、  
所要のレコードを含めと判明した物理セル  
のみを一次記憶に送ればよい。一次記憶  
では、そのデータブロックに必要な処理  
をほどこして、場合によっては元の  
ファイルに書き込む。

このように、各ヘッドに特殊なマイ  
クロプロセサをもつ装置を実現するに  
あたっては、各プロセサをどのように

設計するか、どのような機能をもたせ、  
命令セットをどのように定めるかとい  
う問題がある。

また、単純には、ディスクやドラム  
が1回転する間、すべてのプロセサは  
同じ命令に従うと考えられるが、

- 同一の命令を実行するプロセサの一  
部を選択することができ、
- 各々独立に別個の命令に従わせる  
こと、

- 1つの命令の実行結果を(次の命令  
に送ることで)他の任意のプロセサに  
送ることができ、

といったセルラ型プロセサ構成上の  
問題もある。

これらの問題が(ある程度)解決可  
ければ、もっと速度の高い回転型メモ  
リ(磁気バブルやCCD)あるいは半導  
体集積回路によるシフトレジスタに  
拡張して適用することができ、そのよ  
うな場合には、記憶内容にデータタグ  
をつけて、純粋のデータも、命令(プ  
ログラム)も、スタックやテーブルの  
ような制御用データもすべて同じメモ  
リ内に入れて、1命令実行中に次の命  
令を取り出す制御方式が可能となる。  
よって、前節オア項で述べたような可  
変長の命令形式であっても同様に扱  
えるとともに、記号レベルでメモリ  
やファイルが管理できる。

例. このように、1回転で全トラック  
の内容を検索できるように強化された  
ファイル装置をもつ場合の演算時  
間を、各プロセサの様々な特性条件  
のもとで考える。

演算: 2つのデータ集合XとYの共通  
集合  $Z = X \cap Y$  を求める。

前提: 1レコードは1つの値から成り、  
Xは  $M_x$  個のレコードを持ち、 $k_x$  トラ  
ックにまたがって記憶されている。  
Yは  $M_y$  個のレコードを持ち、 $k_y$  トラ  
ックにまたがって記憶されている。  
結果のZは、 $M_z$  個のレコードの集

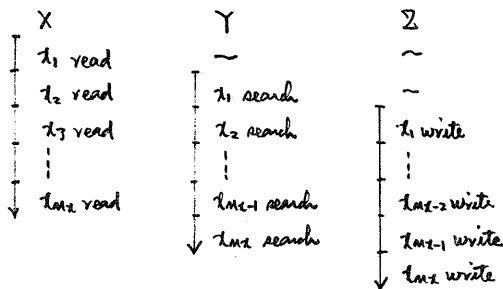
合とわり、 $k_2$ トラックに記憶できるものとする。ここに、

$$0 \leq m_2 \leq \min(M_x, M_y).$$

a. 各ヘッドにつけられたプロセサが互いに独立に異なった命令を実行し、かつ、他の任意のプロセサと通信できる場合:

- X のレコード  $x_i$  を読み、
- $x_i$  を Y の  $n$  番目のレコードと比較する、
- $x_i = y_j$  とした  $y_j$  が存在すれば、 $x_i$  を  $Z$  に加える、

この操作を並列にパイプライン的に実行すれば  $(M_x + 2)$  回転で完了する。



$$t_a = (M_x + 2) \times ROT$$

ROT は 1 回転の時間。

b. 他の任意のプロセサとの通信はできないが、すべて (あるいは一部) のプロセサが同じ命令しか実行できない場合:

- $x_i$  の取り出し、
  - Y の要素との比較、
  - $Z$  への書き込み、
- を順に行わねばならないので、 $m_2$  の値によって、 $2M_x$  回転から  $M_x$  回転かかる。

$$2M_x \times ROT \leq t_b = (2M_x + m_2) \times ROT \leq M_x \times ROT.$$

c. 他のプロセサとも通信できない場合:

- X を 1 次記憶に転送する  $\sim k_1$  回転
  - $x_i$  を順次ファイルに送って、比較し、結果によって  $Z$  に加える操作をパイプライン的に行う  $\sim (M_x + 2)$  回転
  - $Z$  を 2 次記憶に転送する  $\sim k_2$  回転
- $$t_c \leq (k_1 + k_2 + M_x + 2) \times ROT.$$

d. 比較を 1 次記憶で行う場合:

- X, Y を 1 次記憶に転送する  $\sim k_1 + k_2$  回転
  - $x_i$  に対して Y のレコード  $y_j$  を 1 つずつ比較する  $\sim (M_x \times \text{FETCH}) + (M_x \times M_y \times \text{COMPARE})$
  - $Z$  を 2 次記憶に転送する  $\sim k_2$  回転
- $$t_d = (k_1 + k_2 + k_2) \times ROT + (M_x \times \text{FETCH}) + (M_x \times M_y \times \text{COMPARE}).$$
- FETCH は 1 次記憶から 1 レコードを取り出す時間 (せいぜい  $1 \mu\text{sec}$ ).  
COMPARE は 1 レコード同士の比較時間 (せいぜい  $10 \mu\text{sec}$ ).

レコード長、トラック容量にもよりますが、多くの場合、 $t_a < t_c < t_b$ .

ROT が  $10 \text{ msec}$  のオーダーであれば、レコード数が  $500 \sim 1000$  位まであれば、 $t_a < t_c$ 、それ以上にすると、 $t_a < t_d$  になる。

1 次記憶での比較照合に連想記憶を用いると COMPARE の回数は  $M_x$  になる。よって、 $k_1, k_2$  が  $M_x$  くらいに充分小さければ、 $t_d < t_c$  となる。

回転型記憶装置上での記号列操作、あるいは、そのための命令体系は、読み出し・書き込みの機構によって大きく左右される。固定長データのみを取り扱う限り、読み出し (データ検索) のみを問題にする場合には、あまり困難はないが、とくに、可変長データを挿入・変更などの書き込みも含めて取り扱うことは、現在のディスク、ドラムの読み出し・書き込み機構だけでは無理がある。そのような操作に対しては、ランダムアクセスメモリの、インテリジェント型バッファを補助的に用いて書き込み操作をその上で行うか、あるいは、ディスクやドラムの読み出し・書き込み機構の改良を加える必要があると考へる。

#### 4. おわりに

本論では、DBMSの問題点を整理し、これらの問題を解決するための対策をソフトウェア・フレームワーク・ハードウェアの機能分類という面から、事例を交えて検討した。これらの対策は、局所的なもので、計算機アーキテ

クチャある1台計算機システムとして  
まかすところなものでない。今後、  
以上の考察をもとに、より具体的な  
検討をすすめると共に、システム全体  
としてまかすものを導き出したい。  
その際、とくに、可変長データの直接  
操作、検索・照合などの並列処理がホ  
イントになると考えられている。

#### — 参考文献 —

- [1] 関野陽, 植村俊亮: データベース・マシン; 情報処理, Vol. 17, NO.10 (10月, 1976), pp 940~946.
- [2] 植村俊亮: データベースマシンのアーキテクチャ; 電気四学会連合大会, 昭和51年.
- [3] 関野陽: 非数値処理アーキテクチャ会議に出席して—DBMSのハードウェアサポートの研究; 情報処理学会データベース研究会資料 DB29-1, 昭和51年7月.
- [4] R. I. Baum & D. K. Hsiao: A Data Secure Computer Architecture; COMPCON 76, pp 113~117.