

マルチマイクロプロセッサシステム HARPS

Multimicroprocessor System HARPS

データベースマシン および その他の応用
Emulation of database machine
and Other Applications

田中 譲・田畑 隆司・津田 孝夫
Yuzuru TANAKA, Takashi TABATA, Takao TSUDA

北海道大学・工学部

Hokkaido University

A new approach to relational data base machine is shown as an application of HARPS. HARPS is used to emulate this data base machine called OPREC (Orthogonal Processor for Relational Computation). This system is based on two important ideas, i.e., 1) Parallel access to files, and 2) Parallel processing. This two kinds of parallelism are both applied to domains of a relation in relational data base. Together with the outline of the OPREC's architecture, various problems on data base machines and their design theory are described.

[1] 序論

北海道大学情報工学専攻では、日本ミニコンピュータ(株)と共同開発になるマルチ・マイクロプロセッサ・システム HARPS¹⁾ (Hokkaido University Array Processor System) を現在研究開発中であり、汎用並列計算機の開発を目指した MOSES¹⁾ の他、種々の応用システムの研究開発も開始されつつある。主なものは、DDA (Digital Differential Analyzer)、言語プロセッサの並列処理化、神経回路網のデジタル・シミュレータ、FFT等の特殊計算用高速計算機、データベース・マシン等の研究がある。この報告では、著者等が手掛けているもの内、特に興味を抱いて研究を行っているデータベース・マシン開発への HARPS の応用について述べる。

この研究は、HARPS の柔軟なエミュレーション機能を用いてデータベース・マシンのアーキテクチャを研究するもので、この研究のために、ローカル・メモリーの増設等の HARPS の一部拡張も決定している。現在はアーキテクチャ設計のための基礎研究の段階であり、この諸結果に基づき、具体的なアーキテクチャの提

案をいくつか行い、HARPS によるエミュレーション等を通して問題点の把握を行いつつ、最適システムの設計へと収束させていく予定である。

この研究の基本思想は

1. ファイル・アクセスの並列化
2. 処理の並列化

を押し進めることにより、関係形式データベースを効率良く実現するハードウェア、ファームウェアを開発することにある。この基本思想は

3. 並列化の単位を関係中の各定義域(domain)に選ぶ。

ことにより、関係形式の特長を活かした並列化が可能となる。

同じような方向を目指したシステムとしては RAP²⁾、RARES³⁾、STARAN⁴⁾ 等があるが、並列アクセスの考えが弱く、処理の並列化も充分ではない。同一データの除去、データ管理、マルチ・ユーザーをも考慮したスケジューリングの問題等において、満足な解決がなされていない。この報告では、著者等が目指している OPREC (Orthogonal Processor for Relational Computation) のアーキテクチャの現時点での考え方を示し、OPREC の設計理論やマルチ・ユーザー状況下でのスケジ

ユーザ問題等を概説するとともに、データベース・マシンに対する要請と批判について考察する。

[2] DBMSにおける並列化の可能性

DBMSには大きく分けて3種類のモデルがあるが、これらは、ファイル・アクセス経路決定の3種類の方法と考えられる。これにDBMSの持つべき機能と探索手法を加えいとめると以下のようなものが得られる。

- A. ファイル・アクセス
 - a. 階層モデル
 - b. ネットワーク・モデル
 - C. 関係形式モデル
- B. 機能
 - a. ソート機能
 - b. マージ機能 } 等.
- C. 探索手法
 - a. ハッシュ・テーブル
 - b. Balanced-tree, B-tree etc. }
 - c. 連想記憶 } 等.

データベース・マシンと呼ばれるものは、これらA, B, Cの内の1つないし2つ以上の項目に関し、特殊アーキテクチャにより処理の高速化を図ろうというものである。しかし、この報告では、データベース・マシンはデータベース全体に関連するものを指すとし、従って、Aに関するもののみを考察することにする。

DBMSにおけるデータと数値計算におけるデータの性質を比較すると以下のことが言える。

1. 数値計算では1語データが個別に取り扱われることが多いが、DBMSではデータが集合単位で取り扱われ、1つの集合中の個々のデータは一律な処理を受けることが多い。
2. 数値計算ではデータ間に因果関係を持つことが多いが、DBMSでは処理のスケジュールに影響するような因果関係は個々のデータ間ではほとんどなく、集合間においても統計的に言いつ弱い関係しかない。

以上の2点は上述のA, B, Cにおいて、並列処理の可能性が非常に高いことを意味している。

1は並列処理の対象が充分多いことを意味し、2は並列化の際の競合や、並列化を阻む条件が少いことを保障する。1, 2はまた、並列アクセス

スという数値計算の場合には使われることの少なかつた手法が効果的であることを示している。

[3] データベース・マシン—要請と批判—

DBMSの3種類のモデルの内、階層モデルとネットワーク・モデルは、フォン・ノイマン型計算機の汎用OSのつくる記憶空間の構造に本質的に強く依存しており、純粋にDBMSの高効率化を図った結果の産物とは考えられない。これに対し、関係形式モデルは論理的に整った構造を持ち、推論的検索システムへの発展の可能性を秘め、表現能力が高く、高度な検索要求を容易に与えることができる。データベースの更新も容易で、サブモデルや応用システムも作り易く、データ独立性が優れている。

しかし反面において、(A)これを従来の汎用OSの作る記憶空間内で実現する際の効率の悪さと、(B)従来作られたデータベースが、階層モデルやネットワーク・モデルにおいて一部インバーティッド・ファイルを用意することや効率良く実現できなかつたという事実ゆえに、本格的な大規模データベースにおいて関係形式が採用されることは少なかつた。(B)はしかしながら、今後関係形式の研究が進むにつれ状況が変化すると考えられる。問題は(A)であり、ここにデータベース・マシン研究の大きな動機がある。

データベース・マシンという考え方に対する批判も多い。それは次のような理由からである。

- A. コスト・パフォーマンスと技術の面から、特殊計算機を開発する意義があるか？
- B. 既存のデータベースを有効に利用可能か？
- C. ホストとデータベース・マシン、あるいはデータベース・マシンと端末を結ぶチャンネル・スピードがボトルネックにならないか？
- D. マルチ・ユーザに対して対応し得るか？
- E. 新しいメモリー・デバイスが開発されれば、データベース・マシンの構想も変化するのではないか？

これらの批判に対し、著者等は次のように考える。英文字は上述の批判のそれに対応する。

A. DBMSを始めいくつかの分野においては、計算機システムそのものが、計算をするものから、データベース的なものに変化しつつあ

り、将来主従関係が逆転することも考えられる。技術的には、マイクロコンピュータの発展が実現性を増している。

- B. 関係モデルそのもの問題である。むしろ今後この分野の研究が進むものと思われる。
 - C. 検索コマンドの高度なマクロ化と中間結果の独自処理により解決できる。
 - D. 2章で述べた理由により、数値計算に対するスケジューリングよりも容易である。効率向上も大きい。
 - E. OPREC は、例えば CAM (Contents Addressable Memory) の使用により、一層の高速性を得ることが出来る。
- 以上述べたように、データベース・マシン開発の意義は大きい。

[4] 関係形式データベース⁵⁾

関係形式データベースでは、情報ないしは情報の構造情報が表の形で表現される。

4.1. 定義

集合列、 D_1, D_2, \dots, D_n に対し、直積集合 $\Pi_{i=1}^n D_i$ の部分集合 R を関係と呼び、同じ $\Pi_{i=1}^n D_i$ の部分集合の族 \mathcal{R} を関係の型という。 n を関係の型の次数といい $\text{deg}(\mathcal{R})$ で表す。 R の要素を組¹⁾といい、 R の要素を組数²⁾といい、 $\text{tuple}(R)$ で表す。各 D_i は定義域と呼ばれ、 D_1, D_2, \dots, D_n はすべし異なる必要はなく、これらを区別するために A_1, A_2, \dots, A_n なる名前をつけ属性と呼ぶ。 A_i は才³⁾ i 属性である。 R が A_1, A_2, \dots, A_n なる属性を持つことを $R(A_1, A_2, \dots, A_n)$ で表現する。 R の属性全体を $\Omega(R)$ で表す。 $A \subset \Omega(R)$ が $A = \{A_{i_1}, A_{i_2}, \dots, A_{i_k}\}$ のとき、 $\lambda(A) = (i_1, i_2, \dots, i_k)$ とし、これを A のインテックス集合、 i_h をインテックスと名付ける。関係 R が

$$R = \{t : t(i) = 1, 2, \dots, \text{tuple}(R)\}$$

$$t(i) = (r_{i1}, r_{i2}, \dots, r_{in})$$

なるとき、 R は才¹⁾ i 行才²⁾ j 列に r_{ij} を書いた $\text{tuple}(R)$ 行 n 列の表によって書き表すことができる。このように表において、 $r_i \neq r_j$ ($i, j, i \neq j$) が言えるとき、同一データの重複がない⁴⁾ といふ、このときさらに、どの属性もモデルと

してそれ以上分割して考えられることかないとき、この表現を才¹⁾ 正規形という。関係形式データベース \mathcal{R} とは関係の型の集合であり、ある時点でのデータベース \mathcal{R} は関係の集合である。 \mathcal{R} に現れるすべし属性の集合を $\Omega(\mathcal{R})$ で表す。

4.2. 集合演算と関係計算

関係は集合であるから、 \cup (集合和)、 \cap (集合積)、 $-$ (集合差)、 \times (直積) 等の2項集合演算が可能である。

関係計算には Codd の提唱による関係代数があるが、この報告ではそれとは異なるより一般的な表現法を示す。

2つの関係の型 $R(A_1, A_2, \dots, A_n), S(B_1, B_2, \dots, B_m)$ に対し、

$$B_i = f_i(A_{i_1}, A_{i_2}, \dots, A_{i_{l_i}})$$

により、 R を S に移す関数 F を考えることができる。 F を

$$F \equiv (f_1(i_1, \dots, i_{l_1}), \dots, f_m(m_1, \dots, m_{l_m}))$$

で表現する。 $F = (1, 2)$ は才¹⁾ 1、才²⁾ 2 属性への射影を意味する。

関係 R に対し、 R から R の部分集合への写像 C を考えることができる、制約条件と呼ぶ。 C は一般に、

$$C \equiv (V(\bigwedge_{i=1}^n g_i(i_1^i, \dots, i_{l_i}^i) \theta v_i^i) \wedge (\bigwedge_{j=1}^m h_j(j_1^j, \dots, j_{l_j}^j) \theta k_j(j_1^j, \dots, j_{l_j}^j)))$$

の形をしている。 \wedge, \vee は and, or を表し、

" v " は具体値を表す。 $C(R)$ に対し、 $g(i) \theta v$ は R のセレクションと呼び、 $h(i) \theta k(i')$ は R のリストラクションと呼ぶ。ここで θ は関係演算子で、 $=, >, <$ 等々がある。例えば $\{1, 2, 3\} \in \lambda(\Omega(R))$ に対し、

$$C(R) = (\{1 = \alpha\} \wedge \{2 = 3\}) R$$

は、 R において、才¹⁾ 属性が値 " α " や、才²⁾ 属性と才³⁾ 属性の値が等しいような組を求めることを意味する。 F は関係に対しても適用でき、関係 R に対し、 F, C を同時に適用することを $(F; C)R$ で表す。 $(F; C)R$ を関係計算と呼び、 $(F; C)$ を FC 作用素、 R を台集合と呼ぶ。

$$(F'; c') R \times (F; C) S$$

を

$$((F'; c') \times (F; C)) R \times S$$

で表す。(F'; c'), (F; C) は、それぞれ台集合 R, S の作用素である。これを台集合 R x S の作用素に拡張する操作を σ で表わすと、

$$\begin{aligned} & (F'; c') \times (F; C) \\ &= \sigma(F'; c') \cdot \sigma(F; C) \\ &= (\sigma[F'] \cdot \sigma[F] ; c' \wedge \sigma[C]) \end{aligned}$$

たとえば、((1, 2, 3); T) · ((1/2); T) = ((1/2); T) 等とする。T は無制約であることも表す。例えば deg(R) = 2, deg(S) = 2 とし

$$\begin{aligned} \text{a. } (F'; c') &= (1/2; 2 = "5"), (F; C) = \\ &= (2; 1 = 2) \text{ とすれば、} \sigma(F'; c') = \\ &= (1/2, 3, 4; 2 = "5"), \sigma(F; C) = \\ &= (1, 2, 4; 3 = 4) \text{ となり、} \end{aligned}$$

$$\begin{aligned} & \sigma(F'; c') \cdot \sigma(F; C) \\ &= (1/2, 3, 4; 2 = "5") \cdot (1, 2, 4; 3 = 4) \\ &= (1/2, 4; 2 = "5" \wedge 3 = 4) \end{aligned}$$

となる。ついでながら、これは、

$$= (1/"5", 4; 2 = "5" \wedge 3 = 4)$$

と簡略化できる。

関係代数のジョイン演算

$$R[A \theta B] S = \{ (r, s) \mid P_A(r) \theta P_B(s) \\ r \in R, s \in S \}$$

$$A \subset \Omega(R), B \subset \Omega(S)$$

$$P_A: \text{属性集合 } A \text{ の射影: } (\lambda(A); T)$$

は、

$$(\perp ; \lambda(A) \theta \sigma[\lambda(B)]) R \times S$$

で表現される。 \perp は恒等写像を表す。

関係代数における商

$$R \div S = \{ \delta \mid (s, \delta) \in R, \forall s \in S \}$$

は、 $S = P_{A_1}(S, \delta), \delta = P_{A_2}(S, \delta)$ とし、

$$(\lambda(A_2) \div \delta[\lambda(\Omega(S))]; \lambda(A_1) \div \delta[\lambda(\Omega(S))]) R/S$$

と表現する。 δ はインテックスにプライム " ' " を 1 つ附加する操作を表す。これにより、各イ

ンテックスが、台集合の一般形

$$\Pi R / \Pi S / \dots$$

において、どの部分のインテックスであるか知ることができる。

関係計算 (F; C) による記法は関係代数に比し複雑に見えるが、これは、関係代数が概念的であるのに比し、関係計算 (F; C) が計算機処理を念頭に置いているからである。

関係計算 (F; C) に対し、以下の関係が成り立つ。

$$\text{a. } (F'; c') \times (F; C) = \sigma(F'; c') \cdot \sigma(F; C)$$

$$\text{b. } (F'; c')(F; C) = (F'F; c \wedge c'[F])$$

$$\text{c. } F(F_1 \div F_2) = (FF_1, F_1 \div F_2)$$

$$(F_1 \div F_2)F = F_1F \div F_2F$$

$$(F_1 \div F_2)(F_3 \div F_4) = F_1(F_3 \div F_4) \div F_2(F_3 \div F_4) \\ = (F_1F_3 \div F_2F_3, (F_3 \div F_4))$$

ここで、 $c[F]$ は c なる制約を、 F を作用する前のインテックスに関する制約に直すことを意味し、例えば $(3=4)[5, 2, 1, 6]$ は $(1=6)$ である。

このことから、任意の関係計算は、

$$(F_1 \div F_2, (F_3 \div F_4, (\dots)); C) R_1/R_2/\dots$$

$$R_k = \Pi R_{ki}$$

の形の標準形に変形可能であることがわかる。この計算は、プライム " ' " の多くついでいるインテックスから順に実行して行くという制約があるだけで、プライムの個数の等しいインテックスに関する実行の順序は自由である。関係計算を標準形に直すことは、(F; C) 作用素の簡略化やスケジューリングにとって重要である。

関係 R がテ-グベースに表形式で記憶されているとき、R を一次関係と呼び、R が関係計算の結果生まれるものがあるとき、R を二次関係と呼ぶことにする。

4.3. 閉数従属性

A, B $\subset \Omega(R)$ のとき、R の任意の関係 R, R の任意の 2 組 r_1, r_2 に対して、

$$P_A(r_1) = P_A(r_2) \Rightarrow P_B(r_1) = P_B(r_2)$$

が成立するとき、B は A に閉数従属であるといひ、 $A \rightarrow B$ と表す。

$A \rightarrow \Omega(R)$ や、任意の $A' \subset A$ に対して $A' \rightarrow \Omega(R)$ が成立するとき、 A はキャンディデートキーと呼ばれる。 R のキャンディデートキーの集合を $K(R)$ で表す。 $K(R)$ の中に属性 A_i を含むキーがないとき、 A_i はプライムでないと言われる。任意に選んだ $A \in K(R)$ をプライマリキーと呼ぶ。ある $A \in K(R)$ に対し、任意のプライムでない属性 $B_i \in \Omega(R)$ と任意の $A' \subset A$ に関して、 $A' \rightarrow B_i$ ならば R は A に全関係従属であるという、このとき、 R はプライマリキーに対して全関係従属である。

全関係従属になるように関係を分解し得られたモデルを第2正規形と呼ぶ。 $A \rightarrow B$ や $B \rightarrow C$ のとき、 C は A に推移従属であるという。第2正規形や、任意のプライムでない属性 $B_i \in \Omega(R)$ が、ある $A \in K(R)$ に推移従属でない場合、 R は第3正規形であるという。第3正規形では、プライマリキーに属する属性が、プライマリキーに推移従属することがあり、これをキー破壊と呼ぶ。分解された関係は、ジョイン演算により複元できる。この点において、ジョインは重要である。

[5] アーキテクチャ

- 既に RAP, RARES, STARAN 等が提案されているが、これらに関し以下の問題点がある。
- 演算の結果生じた2次関係は1次関係と同様に取り扱い可能であるべきである。
 - マルチユーザー状況での考察が正しい。
 - 並列アクセスの利点を充分利用していない。
 - 効率に関する評価に乏しい。
 - ジョイン演算の高速化に関する工夫がない。
 - 同一データの除去、ソート、関係の除算の実行が遅い。
 - 可変語長を目指しているが充分でない。
- 等の問題がある。

OPREC はこれらの点を考慮し、定義域毎の並列アクセスと並列処理を基本思想として考えられている。

OPREC は検索要求解析システムのもとも働くがこれを図1に示す。検索要求解析システムの内、スキーマプロセッサの一部を除いては本報告では、このシステムには触れない。スキーマプロセッサは関係計算の解析を行い、これを標

準形に変形し、演算順序の解析を行った後、他のユーザーの関係計算に対する結果と比較し、同時処理の可能性も解析した後、この結果に基づいて各関係計算を基本演算に分解し、OPRECに送る。OPRECを図2に示す。関係は、1つ1つの組が $\text{MP} \rightarrow \text{MEM} \rightarrow \text{PMSW} \rightarrow \text{MP} \rightarrow \text{PMSW} \rightarrow \text{MEM} \rightarrow \text{MDSW}$ の多重経路を各定義域毎に分けられて通過していく間に処理される。

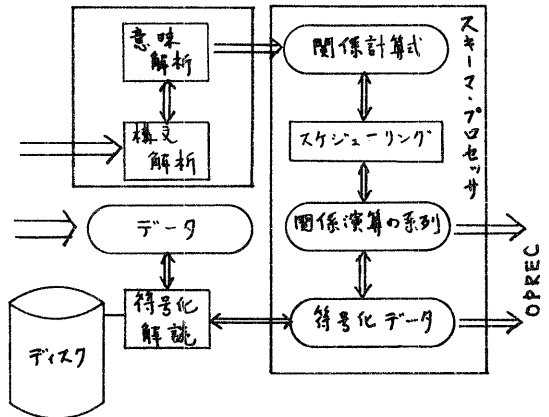


図1: 検索要求解析システム

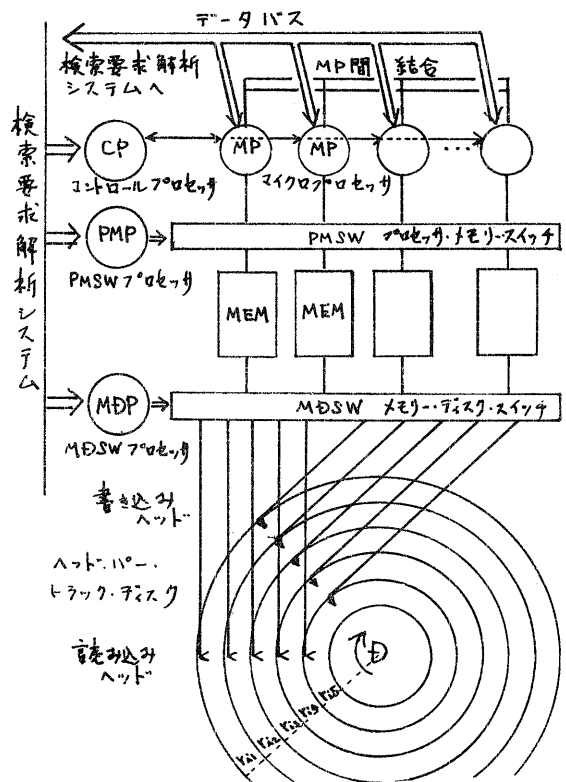


図2: OPRECの構成図

この処理は回転記憶装置での回転に遅れることなく繰り返される。処理の単位は1組ついでなく、バッファ・メモリー-MEMのサイズをBとしてB語がつまめて処理される。関係Rは基本的には回転記憶装置のトラックに/定義域が割り当てられ、回3のように記憶される。回3におけるMEMのチェーンはMDSWとPM SWで作られ、 N_B 台のMPに割り振られる。読み込みヘッドと書き込みヘッドを回のように配置することにより、データの更新が可能である。読み込みトラックと書き込みトラックが異なることもある。 N_B は通常1である。

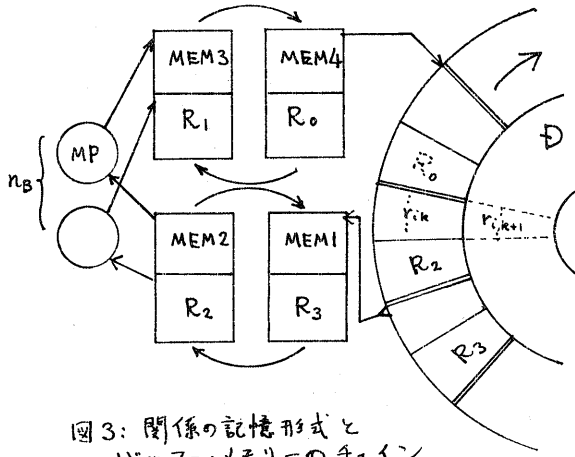


図3: 関係の記憶形式とバッファ・メモリーのチェーン

ディスクは複数台に分かれてよい。 $R=R_1 \cup R_2$, $R_1 \cap R_2 = \emptyset$ のとき, R_1 と R_2 は2台のディスクに分配可能であるが, $A \cup B = \Omega(R)$, $A \cap B = \emptyset$ なる A, B を分配することは許されない。これはディスクが異なれば回転の同期がとれないからである。定義域内のデータ長が, OPRECの基本語長の W_D 倍あるとき, D を記憶するためのトラック数を T_{RD} , 処理用バッファ・メモリー数を M_D , 処理用プロセッサ数を P_D とすると, $P_D = M_D = T_{RD} = W_D$ とすることにより可変語長のデータが取り扱える。 $P_D = M_D = N_B T_{RD} = n W_D$ ($n = k \cdot N_B$) による処理形態を (n, n) way方式, (n, n) 方式を n way方式と呼ぶ。これらは処理スピードを向上する。OPRECは関係Rの処理結果に従い, Rの各組にマークをつける機能(例えば $R(A_1, A_2)$ に対し ("mark", 1, 2; 1 > 2))を持つが, マーク用には関係Rのための P_R, M_R, T_{RR} とは別に P_M, M_M, T_{RM} が取られる。 ($n > n_R$) way方式では $P_M = M_M =$

$N_B T_{RM} = n W_M$ であり, W_M は自由に選べる。
 R, S を1次関係ないしはディスクに既に表形式で生成されている2次関係とする。OPRECに対するコマンドは, $(F; C)R$, $(F; C)R \times S$ ($F \div F'; C$) R/S , $DUP(R)$ 等がある。 $DUP(R)$ はRにおける同一データの除去を行う。これを一般に $(F; C)R$ と表わすと, $d(F; C)$ や F, C のいずれかに関係するインデックスの集合を表すとき, MDSWは $P_d(F; C)$ 足を行い, 選ばれた定義域に対するトラックを空いている $(MP; MEM1, MEM2, MEM3, MEM4)$ の組に結びつけ, $MEM1$ とこのトラックを結合する。 $(MP; MEM1)$ の組は PMSW によって作られる。処理結果が, 別のトラックと異なるトラックに書き込まれることも許し, $MEM4$ と書き込みトラックの結合が MDSW で行われる。 (F, C) 足の実行中のある時点では, $MEM1$ に対し読み込みトラックからのデータ転送が起っており, 同時に $MEM4$ から書き込みトラックへのデータ転送がおこなわれている。一方 $(MP; MEM2, MEM3)$ を1つのプロセッサ・システムとする密結合を持ったアレイプロセッサ・システムでは, $MEM2$ に書かれたB語のデータに対する (F, C) 足の実行がプロセッサ間の密結合を使っている。処理結果は $MEM3$ に書き込まれていく。これらの過程は同時に終了し, $MEM1$ と $MEM2, MEM3$ と $MEM4$ が交換され次のB語の処理に進む。検索要求に対する最終結果は, $MEM3$ に書き込まれない。検索要求解析システムに並列に送られる。
 $(F, C)R$ は1回転, D をトラック長として $(F, C)R \times S$ は D/B 回転で処理される。 $(F \div F'; C)$ には, 最初の D/B 回転で $P_d(F \div F'; C)(R)$, $P_d(F \div F'; C)(S)$ の同一データの除去が行われ, 次の D/B 回転でRに対する中間結果 R' の内, 中間結果 S' に対し条件を満足するものに対し, R' に付加したマーク・エリアに1を加えていく。次の D/B 回転で直前のRの中間結果 R'' に対し, $P_d(F)$ に対し同一データの除去を行うと共に, それらのマーク・エリアは加え合わせる。最後の1回転でマーク・エリアが $stop(S')$ となっているものを標す。結局 $(D/B) + 1$ 回転必要である。
 OPRECはその他, カウント命令, ソート命令 Min, Max 命令等を持つものとする。 Min, Max

は、対応する定義域を k とし

$$\text{Min}(k; R) = (\lambda(\Omega) \div k', k \leq k')$$

$$\text{Max}(k; R) = (\lambda(\Omega) \div k', k \geq k')$$

に等価である。

次に、B語のプロセッサ・アレイによる処理について述べる。簡単のために、 $(F, C)R$ に対して、リストラクション

$$(1; 2 > 3)R$$

を考える。バッファは、 R の第1, 第2, 第3 属性を図4のように持つものとする。

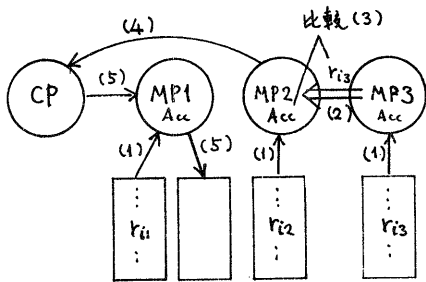


図4: $(1; 2 > 3)R$ の実行

CPから、MP1, MP2, MP3に対し、データのロードが指定される。3台とも同じ番地を指定する。CPはMP3に対しMP2へのデータ転送指令を送り、次にMP2に対しAccと転送データの比較を指示する。比較結果はCPに送られ、これによりCPはMP1にデータ書き込みを指示する。以上が1サイクルで、読み込み番地が加えられた後、次のサイクルに移る。

$(F; C)R \times S$ に対して、簡単のために、

$$(1, 4; 2 = 3)R(A_1, A_2) \times S(B_1, B_2)$$

の実行を考える。これを図5に示す。転送用ポインタ P_x は最初データの先頭にある。比較用ポインタ P_c をMP1, MP2では P_x に、MP3, MP4では $P_x + 1$ にセットする。MP2, MP3で P_x のデータを交換する。このデータを X とする。MP2, MP3において X と P_c を比較し、比較条件を満足するとき、各MPのステータスを1にする。MPアレイのステータスはCPにより常に監視されており、MP2, MP3のいずれかが1になったとき、それぞれに対応して P_x と P_s の2種

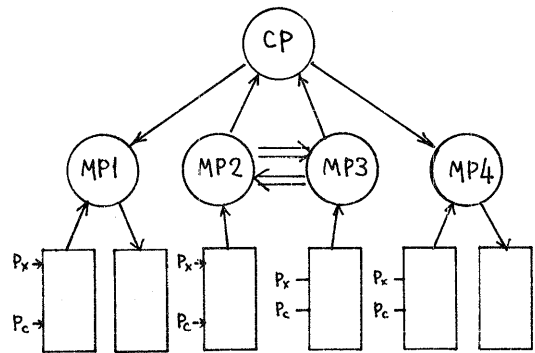


図5: $(F; C)R \times S$ の処理

$$((1, 4); 2 = 3)R(A_1, A_2) \times S(B_1, B_2)$$

のバリスをMP1, MP4に送る。 P_x, P_s の両方が送られることもある。MP1, MP4では、MP2, MP3に対応して、 P_x, P_c が変化しており、 P_x が送られるとMP1は P_c のデータを、MP4は P_x のデータをMEM3に書き込み、 P_s の場合は、 P_x, P_c が逆になる。MEM3が1杯になったときは、PMSWとMDSWによって新しいMEM3, MEM4と新しいトラックがこのMPに割り当てられる。これにより、トラック数はどんどん増すが以後の回転の際に中間結果が使用するトラック数の最少化が行われる。次に P_c が $P_c + 1$ となり、 P_c がバッファサイズをこえたとき、 $P_x = P_x + 1$, $P_c = P_x$ となる。

(n, n_B) way方式では、1度に他の n 個のMPから送られたデータと P_c のさすデータの比較を行う。一回の比較は n 倍の時間がかかるが、同時に n^2 倍の比較がいき、時間は $1/n$ となる。

[6] OPREC の設計理論

OPRECの設計に必要な式を求める。先述のものも含め、以下のような記号を導入する。

関係 R に対して、

$\text{deg}(R)$: 次数, $\text{tuple}(R)$: 組数

$\text{width}(\Theta) = W_\Theta$, $\text{width}(R) = \sum_{\Theta \in \Omega(R)} W_\Theta$

回転記憶装置に対して、

D_l : トラック長 (word)

D_w : トラック数

τ : 1回転に要する時間 (sec)

バッファに対して,

B : バッファ・サイズ (word)

N_{MEM} : バッファ・メモリーの個数

プロセッシング・エレメントに対して,

C : データの比較, 演算に要する時間のオーダー (sec)

Δ : プロセッサ間の 1 word のデータ転送時間 (sec)

C' : 2つの定義域にあるデータ間の比較に要する時間のオーダー (sec) $C' \approx C + \Delta$

N_p : プロセッサの数.

次のような OPREC のモデルを考える.

$C\phi$: 1つのトラックには 1つの定義域のみを記憶する.

$C1$: $Dl > \max_R \text{tuple}(R) \gg B$

$C2$: $Dw, N_{MEM}, N_p \gg \max_R \text{deg}(R)$

R, S を一次関係ないし中間結果として表形式で記憶されている二次関係とし, T_1 を $(F, C)R$ の実行に要する時間 (sec), T_2 を $(F, C)R \times S$ の実行に要する時間とすると,

$$T_1 = \tau \quad (1)$$

$$(B/dl)\tau > BC \rightarrow \tau/dl > C \quad (2)$$

$$T_2 = \tau (dl/B) \quad (3)$$

$$\tau \cdot (B/dl) \approx (C/2)B^2 \rightarrow (\tau/dl) \approx (CB/2) \quad (4)$$

が得られ.

$$T_2 = dl^2 \cdot C/2$$

となり, T_2 は τ, B に依存しない. C, dl はたいたいの値が決まってしまうので, τ, B に関する条件を求めると.

$$\tau > C \cdot dl$$

$$(\tau/B) \approx C \cdot dl/2$$

$$B > 2$$

となる. $C = 10^{-7}$, $dl = 10^4$ とすると,

$$\tau > 10^{-3} \text{ sec}, T_2 = 5 \text{ sec}$$

となる. ディスク使用の際は, $\tau = 10^{-2}$ とすると,

$$B \approx 20$$

が得られる.

$C\phi$ を変えて $(n:1)$ way とすると.

$$T_{2,(n,1)} = T_{2,1}$$

$$(\tau_{(n,1)}/B_{(n,1)}) = n(\tau_1/B_1)$$

$$N_{p(n,1)} = n \cdot N_{p1}$$

$$N_{MEM(n,1)} = n \cdot N_{MEM1}$$

となり, ハードウェアは良くない.

n way では.

$$T_{2,n} = n^{-1} \cdot T_{2,1}$$

$$\tau_n/B_n = \tau_1/B_1$$

$$N_{pn} = n \cdot N_{p1}$$

$$N_{MEMn} = n \cdot N_{MEM1}$$

となる.

以上のことから, B は, τ が実現可能程度に小さくすればよいとわかる. τ の小さい CCD や磁気バブルを使用する試みもあるが, これらのデータ長が $\text{tuple}(R)$ に比し小さい場合は, B を高速メモリより安価なメモリに置き換えるという効果に終る. このことは, しかしながら重要であり, 現在検討中である.

以上のモデルでは, 処理時間は $\text{tuple}(R)$ に関係ないが, 回転記憶装置のバッファ・サイズ毎に初期位置があれば,

$$T_1 \propto \text{tuple}(R)$$

$$T_2 \propto \text{tuple}(R)$$

となる.

さらに.

$C\phi$: 回転記憶装置上に R を格納する際, 1トラックの dl/B 個のブロックに均等に $\text{tuple}(R)/B$ 個のブロックを配分する. 残りのブロックは別の関係に使用する. $(F, C)R$ の処理では回転記憶装置が k_1 ブロック進む間に 1トラックの処理が終了するとし, $(F, C)R \times S$ では k_2 ブロック進む時間を要するとする. ($k_2 \gg k_1$) 処理は n way で行う.

とすると,

$$T_1 \approx [\text{tuple}(R) \cdot k_1 / (n \cdot dl)] \cdot \tau$$

$$k_1 \cdot \tau / dl > C$$

$$T_2 \doteq \frac{C \cdot k_2}{2n} \text{tuple}^2(R)$$

$$c/B = n \cdot (C \cdot dl) / 2$$

が得られ、

$$T_2 \propto \text{tuple}^2(R)$$

とできる。

通常の n way にあつて CAM を 1 つ ファイルとして使用するならば、

$$T_2 = dl^2 \cdot c / (n \cdot B)$$

となり、非常に効率が上がる。ただし、

$$c = C \cdot dl$$

をなければならぬ。

[7] スケジューリング

関係計算の標準形

$$(F_1 \div F_2 (F_3 \div F_4 (\dots))) ; C) R_1 / R_2 / \dots$$

$$R_k = \Pi R_k^i$$

が得られたならば、 F の各レベルに関して簡略化を行う、 C について σ "リ" に関する各レベルで簡略化を行う。その結果 $C = F$ (F は偽を表す) となれば、この検索の実行は必要ない。レベル解析から、いくつかの処理順序に関する制限が得られる。これを商条件と呼ぶ。単一ユーザの場合は、 $R_1 / R_2 / \dots$ の中に同じ関係が現れないか探し、商条件を満たす範囲でそれらに対する処理を同時に行うようにする。マルチユーザの場合は、他の関係計算の中に同時実行可能なものがないかを探す。このアルゴリズムは、商を含む場合は簡単集合を関係データベースに対して決められた順番に並べ替えた直積形に変形し、他ユーザの関係計算中に、同じような並べかたがないかを調べればよい。除算を含む場合は、並べ替えの際に商条件が制約として入ってくるだけで、同じようにスケジュール可能である。集合の順番を決めたならば、 (F, C) 作用素を各2項ないし1項演算に分解する。詳細は別の機会に述べる。

[8] データベース・マシンの課題

CAM を使用しないデータベース・マシンは次のような課題を残している。

- A. 同一データの除去
- B. ソートの高速化
- C. 符号化と符号の解読法

ここでは OPREC に関し、A, B を着る。

B に関しては回転記憶装置の外部記憶に対する外部ソート問題の理論に依ればよい。その結果は満足し得るものではなく、特殊な工夫が今後必要になる。この問題は CAM を使用した際に問題となる。

同一データの除去は、ほぼ6章の T_2 だけの時間を要し、大きな損失である。除去はその都度行う必要はないが、中間結果をそのままにしておくと、冗長度はジョインの度に指数関数的に増加する。この問題に対し、いつ冗長なデータが発生し得るか分かれば都合がよい。

関係 R に対し、 $\text{dup}(R)$ を R に同一データの重複があれば ϕ , なければ 1 とする述語とする。 $\text{cand}(A)$ を A が R のキャンディデート・キーがあるという述語とする。これは $\text{cand}(A; R)$ とも書く。 $\text{univ}(A)$ を $P_A(R)$ の要素が1個であるという述語とする。これらを使って、非冗長性の指標規則を以下に示す。

- a. $\text{dup}((1, C)R) = \text{dup}(R)$
- b. $\text{dup}(R[A=B]S) = (\text{dup}(R) \wedge \text{dup}(S))$
- c. $\text{dup}(P_A(R)) = (\text{cand}(A) \wedge \text{dup}(R))$
- d. $\text{cand}(A) \wedge \text{cand}(B) = 1$ のとき、
 $\text{tuple}(A) = \text{tuple}(B)$
- e. $\text{dup}((F; T)R) \geq \text{cand}(\text{dom}(F)) \wedge \text{univ}(F) \wedge \text{dup}(R)$
- f. $\text{dup}((F; C)R) \geq \text{dup}((F; T)R) \wedge \text{dup}((1; C)R)$
- g. $\text{cand}(A; R \times S) \geq \text{cand}(A; R) \wedge \text{dup}(R) \wedge \{R \rightarrow S\}$
 $\text{cand}(A; R \times S) \geq \text{cand}(A; R) \wedge \text{dup}(R \times S) \wedge \{R \rightarrow S\}$

ただし、 $\text{inv}(F)$ は F に逆関数が存在するという述語である。恒等規則はこの他にも多くあるものと考えている。これらを見つけることも、本研究の一つのテーマである。

これらの論理不等式を $(F;C)$ に適用することにより、どこまで非冗長性が保存されているかわかり、保障が切れた時点や除去も行をけよい。むやみに除去を行うより効果は良くなる。 $\text{cand}(A)$ を決定するアルゴリズム、つまり、4章における $K(R)$ を求めるアルゴリズムは存在する。

[9] 結論

この研究は、現在、基礎理論の研究が中心になっているが、近い将来、HARPSによるエミュレーション実験に移行する予定がある。

関係形式データベースは、ファイル・アクセスのみならず、OPRECのスキーマプロセッサに相当する部分についても種々の工夫が考えられるモデルであり、自然言語との関連もつきやすく、今後、さらに重要性を増すものと考える。

OPRECの基本思想である並列アクセス、並列処理は、CAMの使用やプロセッサの仕様の工夫、回転記憶装置の階層化等、今後さらに深い検討を要する。

[謝 辞]

本研究は、その一部を、昭和51年度、文部省特定研究(2)「情報システムの形成過程と学術情報の組織化」における研究活動の一部として、行った。貴重な御意見を頂いた諸先生方に感謝の意を表します。

Bibliography

- 1) Tanaka, Y., and Miyashita, K., et al. HARPS (Hokkaido University Array Processor System): A new hierarchical array processor system. Proc. 2nd EUROMICRO Symp. 1976, North-Holland.
- 2) Ozakarahan, E.A., Schuster, S.A., and Smith, K.C. RAP-An associative processor for data base management. Proc. AFIPS 1975 NCC, Vol. 44, AFIPS Press, Montvale, N.J., pp. 379-387
- 3) Lin, C.S., Smith, C.P., and Smith, J.M. The design of a rotating associative memory for relational database applications. ACM Trans. on Database Systems, Vol. 1, No. 1, 1976, pp. 53-65.
- 4) Davis, E.W. STARAN parallel processor system software. Proc. AFIPS 1974 NCC, Vol. 43, AFIPS Press, Montvale, N.J., pp. 16-22.
- 5) Codd, E.F. A relational model of data for large shared data banks. Comm. ACM 13, 6 (June 1970), 377-387