

データフロー制御データストリーム処理方式 データベースマシン (基本構想)

田中 譲 北海道大学・工学部

1.はじめに

1970年代には種々のアーキテクチャのデータベースマシンが盛んに研究されたが、80年代に入って、この分野は反省期を迎えているように思われる。主な反省点は、

- 汎用データベースマシンは存在し得るか?
- 大容量データベースをどう扱うか?
- RASISをどう支援するのか?

の3点である。

第1点は、ニーズの分析の結果として生じた疑問というよりも、処理対象データのデータ長の柔軟性と、ジョイント等の複雑な演算の高速処理とを同時に満たすようなマシンを開発することの難しさに根差した疑問である。利用者のニーズとしては、IRやバンキングシステムにDBMSを用いたいとの動きも始めている。マシンの分化を考えることの重要であるが、一方では、

- 可変長データが扱え、
- タガル数が 10^6 程度に及ぶ大規模な2つの関係のジョインを $O(n)$ で処理でき
- ジョイン属性の値として、100文字以上の文字列が現われることも許す

というようなマシン・アーキテクチャの研究が必要である。本論文ではそのようなマシンを提案する。

データベースマシンの大容量化に関する反省は、最下層メモリとして CCD や磁気ハーフルメモリを用いるマシン・アーキテクチャに向けられたものである。本研究では大前提として、

「最下層メモリには移動ヘッドディスクを用いる」

という条件を課す。これに伴い、関係のセグメンテーションが必要となる。

もう1点は、RASISの支援について考慮したマシン・アーキテクチャが皆無であることにについての反省である。これは、従来の研究の重点が、データベース処理全体の内の一部である関係代数演算の処理に置かれ過ぎたためである。本研究では、セグメント化された関係をオブジェクトと見做し、オブジェクト指向アーキテクチャで、他のデータベース処理モジュールを制御することを考えている。

本研究では、関係の各属性値をデータ・タイプ毎にコード化したユード化関係を用いて関係代数演算を行うことを提案し、この演算の高速処理のためツバガシステムと、符号化と復号化のみぞれのためのサブシステムのアーキテクチャの概要を示す。

2.アーキテクチャの概要

2.1. 方式に関する基本的な考察

本研究の指針となる、基本的な考察を以下に列举する。

a. 普通の移動ヘッドディスクを最下層メモリとして用いるという仮定のもとでは、セレクションの実行時間を、インデックス、ファイルを持ったソフトウェア後置型データベースマシンと、これを持たないハードウェア後置型データベースマシンとで比較しても大差ないと予想される。

b. 更新処理に於て過大な負担となるのが限りは、データの冗長度を上げて

ぐり、処理の局部性大、処理系の単純化を図ることによつて、総合的なコスト・パフォーマンスを上げるべきである。

c. 関係モデルでは、どの属性が一様に扱われるか、現実には、検索条件式中に現わることから、出力の際にのみ必要となる属性がある。このよろづ *read-only* の属性を、関係代数演算の処理系に持ち込むことは、効率を下げるだけや、何等利点がない。

d. 性別のように、男と女の2値しか現れないよろづな属性を持つ関係 R(X, sex) は、R-male(X) と R-female(X) のよろづに2つの関係に分けて管理処理すればよい。

e. データベースマシンとホストコンピュータの間、ホストコンピューターと入出力機器、およびネットワークとの間、入出力機器と利用者との間のいずれぞれにチャネル容量の上限があり、これらを無視して、データベースマシンの高速化を図つては無駄である。データベースマシンの入出力のチャネル容量は、セカンド 1~10 MB/sec と見做してよいであろう。これは、検索結果の出力が 5,000~10,000 タブル/sec で得られるスループットに対応する。

f. 検索結果の統計解析結果のみを出力する場合には、そのスループットの値はあつてはまらない。こゝよろづな統計解析用の検索処理のスループットは大きければ大きい程良い。

g. 属性値のデータ型を、数値と文字列、順序の定義の有無、出力の降順リストの指定を免けるか否か、不等式によるセレクション演算を免れるか否か、不等式によるジョインリストリクションの対象となるか否か、四

則演算の定義の有無や余類すると表 2-1 のようになる。id# といふのは

	order	sort	selection	join	restriction	小 x - +
Number	v	v	v	v	v	x
id#	v	v	v			-
date	v	v	v	v		*
word	v	v	v			
text						

* : "1982-01-31" + 1 = "1982-02-1"
 "1982-Jan-31" + 1 = "1982-Feb-1"

表 2-1. データ型の分類

識別番号のよろづなデータ型で、date は年月日大、時分秒等のデータ型を指し、text型は順序が定義されない。word は date 型と text 型以外の文字列を指す。date 型以外の文字列型は、不等式によるジョインリストリクションの対象になることはない。date 型は、number 型へ変換することができる。表 2-1 より、number 型と date 型、id# 型と word 型が同じ性質を持ち、これらと text 型の3つの型に分類することができることがわかる。それは

A型 : number, date

B型 : id#, word

C型 : text

と分類する。

h. データベースマシンのアーキテクチャを考える上でも大きな問題となるのが可変長データの取り扱いである。可変長データといふのは、表 2-1 の word 型と text 型に現われる。word 型と text 型を、これらに対して定義されない順序を無視して、固定長のデータへと 1 対 1 に符号化できること

すると、ジョインリストリクエストの処理は、符号化された世界で処理できることが表記よりわかる。プロジェクションの処理における重複タブルの除去が、符号化された世界で処理できる。word型とtext型に対する演算も定義されていないから、このような型の属性について平均や強制的計算を要求されることはない。この型に対する統計演算は数え上げのみであり、これが符号化された世界で処理できる。復号化は、出力に対してのみ行えばよい。符号化された世界での関係代数演算処理は、固定長データの処理になり、このサブシステムのマシン・アーキテクチャを簡略化することができる。

2.2 符号化 / 復号化アーキテクチャ
図2-1に、2.1節のとく述べた符号化関係の処理による要求処理の概念図を示す。図2-2には、処理の具体例を

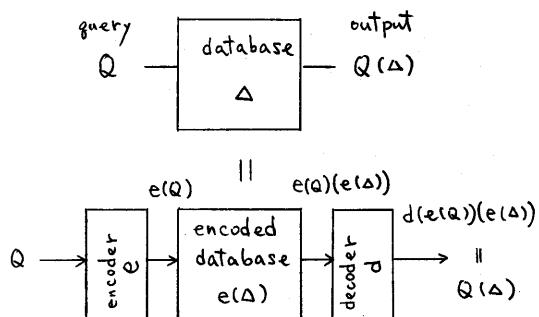


図2-1. 符号化データベースを用いた関係代数処理

ます。データベースの定義時に、各属性には固有のデータタイプが与えられるものとする。ジョインリストリクエストでは、同じデータタイプを持つ属性の間の比較のみが許される。符号化は、各データタイプ毎に行われる。

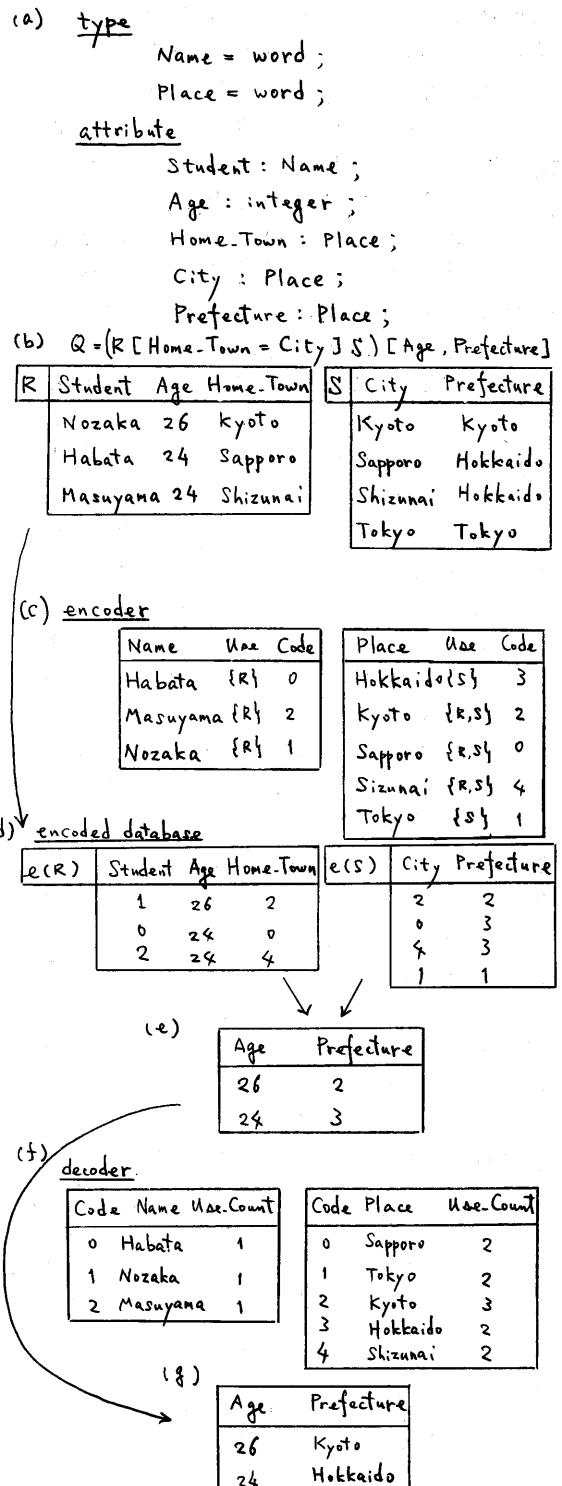


図2-2. 符号化データベースの処理例

符号化が行われるのは、桁数の大きい id#型, date型, word型, text型である。date型は、適当な関数を用いて 1対1に number型へと変換できる。桁数の大きい id#型は、word型として扱う。word型とtext型は変換テーブルを用いるが、word型では、変換テーブル中の word型データは辞書式順序を保つて管理される必要がある。(たとえば word型は B-tree で、text型は hash テーブルを管理する)。図 2.2 (c) に見るように、word型であっても符号化データが、符号化前のデータ間の順序を保つ必要はない。0-ジョインや 0-リストリクションを行うことはないからである。R[Name > 'Ikeda'] のような 0-セレクションは、(c) の変換表より、Masuyama, Nozaka が条件に該当することを knj(e(R)[Name = 2]) ∨ (e(R)[Name = 1]) を計算すればよい。

符号化モジュールと復号化モジュールでの処理は、基本的には、変換テーブルを関係と見做してセレクションの処理を行うことである。2.1 節の A の考察から、これら 2 つのモジュールは、ソフトウェア複雑型データベースエンジンを多重化することで実現するのかよると考える。多重化は符号化テーブル(図 2.2 (c))と復号化テーブル(図 2.2 (f))のセグメント化を意味する。符号器はデータタイプに定義されれる順序やハッシュ値の順序でセグメント化されるのが望ましい。これに対し、復号器は、符号の値の順序でセグメント化されることが望ましい。図 2.1 で更新処理を考慮すると、図 2.3 のよう

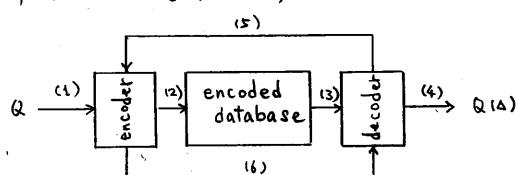
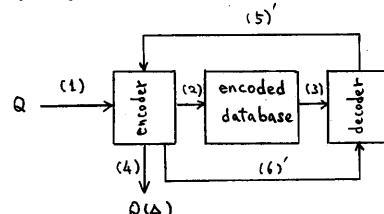


図 2.3 変換テーブルの更新

にデータバスの追加が必要である。(5) は、検索条件を満足するタブルを削除するような場合に必要である。そのようなタブルが符号化されて (3) より得られると、各属性の符号値が復号器で調べられ、その符号値をデータベース中の何箇所で用いてるかを記憶している use-count (図 2.2 (f)) が 1 だけ減られる。その結果、どこで用いられていないとかかった符号値は、そのデータタイプで未使用の符号値を管理していく場所へと戻され、符号器から消去される。(6) のデータバスは、新しく関係を登録するような場合に必要である。

図 2.3 の構成では、(5) と (6) には生データが転送され、しかも符号器と復号器のセグメンテーションの仕方が異なっているために (5), (6) のネットワークのスループットを上げることが難しい。そこでこれを図 2.4 のように改善する。復号器には、生データを持



encoder P.A. : physical address

P.A.	Name	use	Code
A1	Habata	{R}	0
A2	Masuyama	{R}	2
A3	Nozaka	{R}	1

P.A.	Name	use	Code
A1	Hokkaido	{S}	3
A2	Kyoto	{R,S}	2
A3	Sapporo	{R,S}	0
A4	Suznai	{R,S}	4
A5	Tokyo	{S}	1

decoder		
Code	P.A.	u.c
0	A1	1
1	A3	1
2	A2	1

decoder		
Code	P.A.	u.c
0	A3	2
1	A5	2
2	A2	3
3	A1	2
4	A4	2

図 2.4 符号化 / 復号化方式

へのではなく、生データが符号器で管理されている物理アドレス（物理アロケーション番号）を持つ。これにより、(5)' (6)' のネットワークを流れろデータが固定長にでき、変換テーブルの記憶に安やす記憶領域を小さくできる。この場合も、符号器と復号器とのセグメントーション方式は異なるので、(5)' (6)' では、後述する動的クラスタリングを行った必要がある。

2.3 動的クラスタリングとリート

2.2 節の符号器/復号器を例にと、説明する。両者との物理アドレス（P.A.）と符号（C）の関係を情報として持つ。しかし、符号器では、P.A の順に並べられており、復号器中では C の順に並べられている。符号器では P.A の順に処理することが望ましく、符号器より復号器へ送られるデータは P.A の順になる。一方、復号器では C の順に処理することが望ましい。したがって、符号器から復号器への転送では、C についてリートする必要がある。符号器と復号器がそれを車一モジュールであればこれでよいが、各々が複数モジュールより成る場合には、符号器の各モジュールからの出力を 1 つにまとめてリートし、これを C について各復号器に分配する方法では多数のモジュールに対して転送路が 1 本となり、隘路を生じる（図 2.5(a)）。逆に、C の値で先に分配し、各復号モジュールに転送されながらリートする場合には、分配に際してデータ転送の衝突を生じる（図 2.5(b)）。

これに対し、リートと分配をネットワークの機能として統合した SD ネットワークを提案する。まずはソータで、著者等の開発したリートエンジン（SOE）を基本モジュールとする。分配は M-bus を用いて各符号モジュール毎に行われ、M-bus によって、各復号モジュール毎

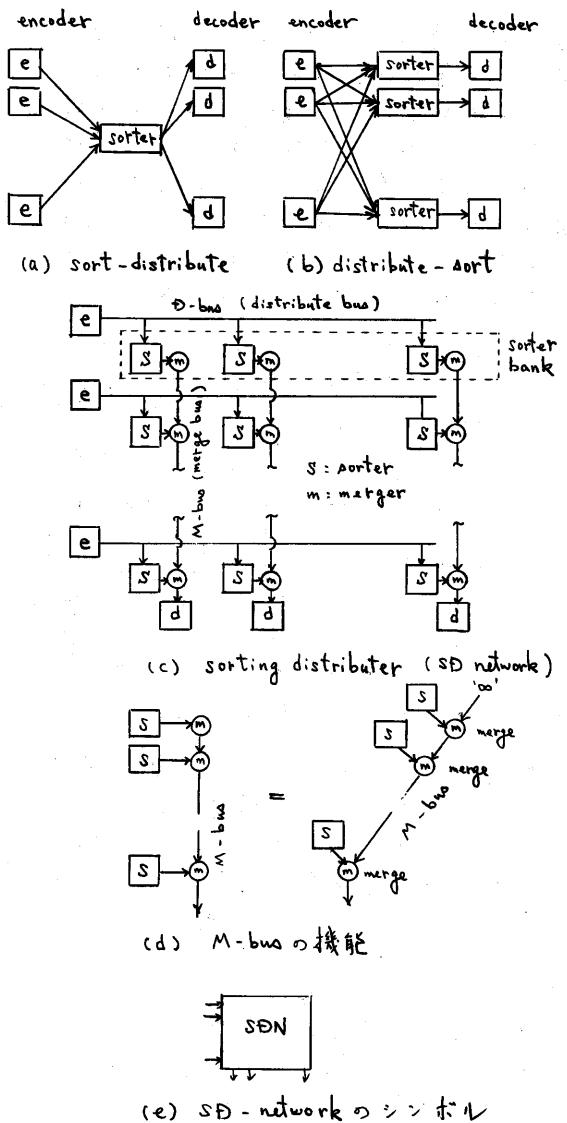


図 2.5 リートと分配の方式

に、各列のリータの出力をマージして復号モジュールへと転送する。SD ネットワークはバッファであり、大容量を必要とするが、後述のように大容量化は通常の RAM を用いて行う。各符号モジュールより一度に送出されるデータ量は、セグメントの大きさや上限が決まる。よって、図の(c)に示すリータバンク毎に要する RAM の大きさは一定である。よって、リータの容量を増

すために使用するRAMバンクは、リータバンク毎に分割することができます。

データ長が固定の場合、2つの処理システム間のデータ転送において、間にリータを介在せることを希望しない場合、各システムが多重化された場合にはSDネットワークを用いることができる(図2.6(a))。ただし、各データ

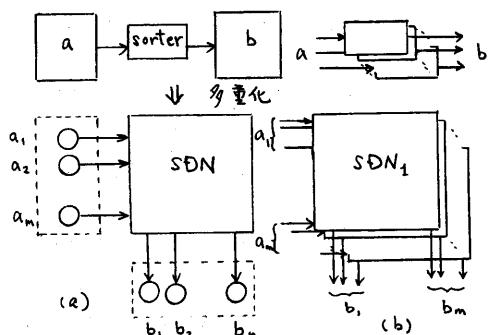


図2.6 リータを介す転送の多重化

の転送先は、モジュラスを計算する等によって、データより計算可能とする。さらには、転送シップの多重化は、同図(b)のように多重度だけSDネットワークを用いることにより実現できる。

2.4 大容量リータ

本節では、読者からリートエンジン(SOE)について既に知っておられるものと仮定する。

SOEはヒート構造を用い、図2.7(a)のような状態でキーを入力すると、同図の次経路上のデータ列にこの経路上で順序正しく新しいキーが挿入される。SOEが図(c)のような時、最下位レベルの(i+1)番目のノードと報ノードを結ぶ経路に次経路を固定し、ヒートの順序と同順序(図では昇順)にリートされたキー列を入力した後、あふれたキーをあふれた順に上から下へ(i+1)番目のノードの下に並べてみると、このアンバランスな本においても根から各葉へ至るキー列はソート

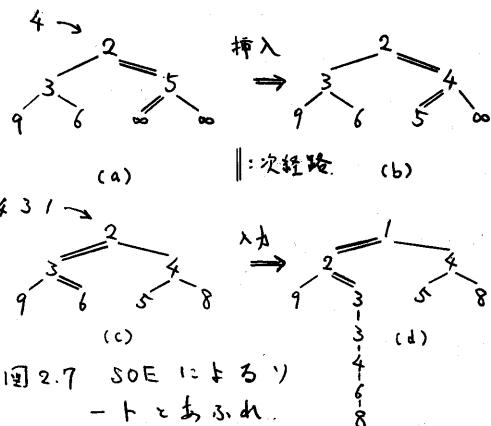
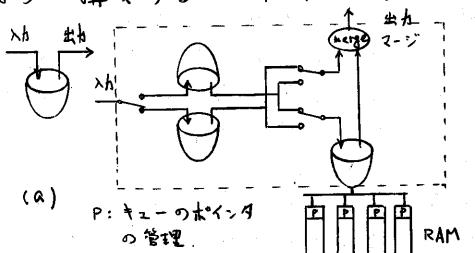


図2.7 SOEによるソート

ーとあふれ

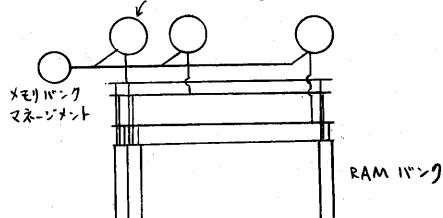
されていく。このあふれの部分をRAMに記憶し、SOEの出力時に(i+1)番目のノードに空孔が生じたとき、あふれ部の上から順に、SOEへと送りこむこと、あふれ部と一緒にリート出力することができる。昇順にリートされたプロック毎にiを変化させれば、ヒートの最下端のすべてのノードの下にあふれ部を設けることができる。今SOEを図2.8(a)のように表わすと、以上のことを用いて、大容量リータを(b)のように構成することができる。RAMは



(a) P:キーのポイント
の管理

(b) 大容量リータ

(b) 図の破線で囲まれた部分



(c) RAMバンクの共用

図2.8 大容量リータヒータバンク

(1) Y.Tanaka et al. IFIP Congress 80 pp.429-432.

多パンク構成でないか、同時に複数パンクをアクセスすることはないので1パンクやさよ。多パンク構成の場合には、各パンクの割り当てを管理することにより、多パンクのRAMモジュールを複数のリータイ共用し、リータパンクを構成することができる(図2.8(c))。

2.5. 处理方式

符号化データベースの処理は、

1. タブルセレクタ (TS)
2. タブルメーラ (TM)
3. アグリゲートオペレータ (AO)

の3つのモジュールを行なう。TS, TM, AOは図2.9のように結合されている。

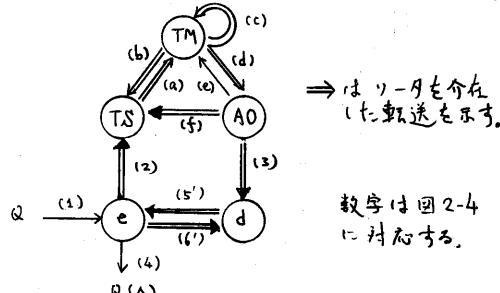
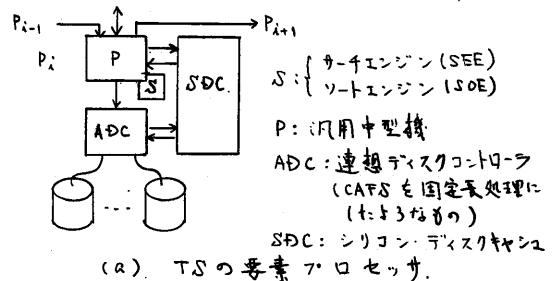


図2.9. 抽象アーキテクチャ。

関係 $R(A_1, A_2, \dots, A_n)$ の各タブルにはタブル番号を定義することができるのべ、これを@R で表し、Rがあたかも @R という属性をもつてかのように見做し、 $R(@R, A_1, A_2, \dots, A_n)$ と表すことにする。Rは符号化されているとする。

TSは、各関係毎に、その関係のread-only である属性毎に、 $R(@R, A_1, \dots, A_n)$ の A_i と @Rへのプロジェクション $R[A_i, @R]$ を A_i についてソートしてファイル $R[\uparrow A_i, @R]$ を持つ。TSが多重プロセッサ化される場合には、各 $R[\uparrow A_i, @R]$ 型のファイルは、 A_i が number型 or date型の場合には、 A_i の値をレンジ分けし、レンジ値のモジュラスの値でプロセッサに分配する。それ以外は、 A_i の値の

モジュラスでプロセッサに分配する。プロセッサ内やさらに $R[\uparrow A_i, @R]$ を分割する場合は、 A_i の値の順序で分割する。各プロセッサは、図2.10に示すよう構成で、プロセッサ間は同図(b)のように結合される。TSでは、セレ



(b). TSの要素プロセッサ間の結合。

図2.10. タブルセレクタの要素プロセッサ

クションを行って、該当するタブルのタブル番号を TMへ送り、たり、等号によるジョインを各プロセッサ毎に独立に処理し、ジョインされたタブル対のタブル番号対の集合を TMに送る。セグメンテーションと符号化の仕方により、等号ジョインはプロセッサ毎に独立に処理することができます。@-ジョイントは、図2.10(b)の結合セットワークを用いて処理する。TSから TMへの転送はリータを介在して行われる。多重プロセッサの場合には SDCネットワークを介する。

$R[A_i = B_j] \sqcup [B_k = C_\ell] T$ のような演算では TS で $((R[\uparrow A_i, @R]) [A_i = B_j] (S[\uparrow B_k, @S]))$ $[@S, @R]$ と、同様に $(@S, @T)$ の対とが求められ、両対とが @Sにつきリートするように SDCネットワークに転送することにより、二つの関係のジョインが残る $(@S, @R, @T)$ の3項組を、TMの各プロセッサ中に求めることができる。

TMは、各関係 R をタブル番号のモジュラスで要素プロセッサに分割して

ファイルも、タガル番号の順序でセグメント化して持っている。各要素プロセッサは、図2.10の(a)と同様である。加同図(b)の結合ネットワークは持たない。

$R[A_i]$, $S[B_k]$ に対して、

((R[A_i = 'v']) [A_j = B_k] S) [A_l, B_m]
のように処理では、符号器で v や B_m をコード化され、((R[↑A_i, @R]) [A_i = 'v']) [↑@R]
が TM へ送られ、TM では、この集合を X として、各プロセッサで、Y = ((R[↑@R, A_j]) [OR = @R] X) [A_j] をサーチエンジンを用いて求め、これを A_j バリートして TS に送る。TS の各プロセッサでは、Z = ((R[↑A_j, @R] [A_j = A_l] Y) [A_j = B_k] S[↑B_k, @S]) [OR, @S] を求め、@R バリートして TM に送る。TM では、各プロセッサ毎に独立に、((R[↑OR, A_l]) [OR = @R] Z) [OR, A_l] を求め、これを W として、転送経路(c) を用いて、@S バリートして TM に送る。各プロセッサは独立に ((S[↑OR, B_m]) [OR = @S] W)
[A_l, B_m] をサーチエンジンを用いて処理する。結果は、出力指定ルートへ、A_l か B_m のいずれか、指定がない場合は任意に決め、この属性についてソートして AO に送る。

AO はグルーピ化して統計演算や、重複タガルの除去などを行う。これらプロセッサ毎に独立の多重プロセッサングで処理できる。AO の要素プロセッサは、サーチエンジンとソートエンジンを持った汎用コンピュータである。

(e), (f) のデータバスは、結果として得られた関係の登録や、更新処理に用いる。

結果として得られた関係の各タガルには、AO の各プロセッサ毎に、プロセッサ番号とタガルの順番からなるタガル id がつけられ、各属性毎にバラバラにして、上の例では、A_l とタガル id, B_m エタガル id の対を A_l, B_m の値バリート

して、復号器に送る。復号器では、各符号をこれに対応する生データの入っている物理アドレスに変換し、符号器に転送し、符号器では、各属性毎に、生データとタガル id の対が得られる。符号器から、これらの対をタガル id の順序で 1 個づつ取り出し、タガルを構成することにより結果が得られる。

符号器と復号器は図2.10(a)のような要素プロセッサで構成され、この2つと、TS, TM の各プロセッサでは、検索要求にしたがって、セグメント・ファイルの先読みがなされ、シリコンメモリで構成される大容量(100~500MB)のディスク・キャッシュにステージングされるものとする。先読みと、ステージングがレディーになることをうりの処理のアクティベートは、データフロー制御で行う。処理の単位は、1セグメントであり、これはオブジェクトと見做して管理と処理がなされる。

図2.9 で 2 重矢印の部分には SD ネットワークが用いられるが、SD ネットワークを少くするために、図2.9 のかわりに図2.11 がよいと考える。

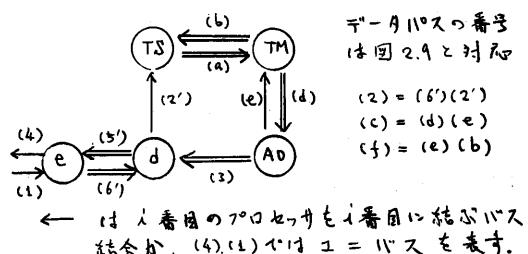


図2.11 改良した抽象アーキテクチャ

3. おわりに。

本マシンは、対象データベースの規模に応じて、各モジュール毎のプロセッサの多重度、SD ネットワークの多重度を自由に設定できる。制御方式の詳細は別稿に譲る。

謝辞 通産省第5世代コンピュータ DDG 委員の方々の議論に感謝します