

# 動的マクビットによるリレーショナル演算の並列処理方式とその評価

PARALLEL PROCESSING OF RELATIONAL OPERATIONS WITH  
TEMPORARY MARK BITS AND ITS EVALUATION

青木 康 田中 孝一 上林 憲行 相磯 秀夫

Yasushi KIYOKI, Kouichi TANAKA, Noriyuki KAMIBAYASHI, Hideo AISO

慶應義塾大学 工学部

Faculty of Engineering, Keio University

## 1 はじめに

リレーショナルモデルはその強固な理論的基盤、非構造的なデータ構造、アクセスの対称性、データ独立性、非手続き的な問い合わせ言語の実現の容易性、集合論的な基本演算子の完備性等、従来のデータベースモデルに比べて優れている点があるにもかかわらず、均質なテーブル構造によってデータ間の関係を示す必要から同様の関係を表現するのにデータ量が他のモデルに比べて大きくなる点、さらに集合論的な演算をソフトウェアで処理するのは非常にオーバヘッドを招来する点などが問題とされ、優れたユーザインタフェースを提示する能力があるにもかかわらず、具体的にシステム化するのにこうしたデータ量の増大や処理効率の悪さからもう一つ決め手を欠いていたというのが現実であると思われる。こうした意味でリレーショナルデータベースのハードウェア化による処理効率の向上改善はリレーショナルデータベースの命運を左右する一つの鍵である。そのためリレーショナルデータベースマシンの研究開発はデータベースマシンの大きなハイライトである。リレーショナルモデルを基本モデルとして採用したデータベースマシンはRAPプロダクトのOriginal RAP<sup>[1]</sup>、RAP-2<sup>[2]</sup>、RAP-3<sup>[3]</sup>、CAPS<sup>[4]</sup>、LEECH<sup>[5]</sup>、DIRECT<sup>[6]</sup>、EDC<sup>[7]</sup>等が代表的である。

筆者らは現在、リレーショナルデータベースマシンSPIRIT<sup>[7][8]</sup>の検討を進めているが、本稿ではSPIRITで採用したリレーショナル演算を動的マクビットを効果的に活用した並列処理方式の詳細、それに基づく基本演算アーキテクチャ

とそのハードウェア構成、さらに他の演算方式との比較を通して本稿で提案する方式の有効性があることが確認されたのでここに報告する。

本稿の構成は第2章で動的マクビット方式の基本的な特長を他の方式との比較を通して明らかにする。第3章では具体的な演算アーキテクチャの説明を通して動的マクビットのデータ構造や各基本操作の詳細を明らかにして最後にその実現のための具体的なアーキテクチャを提示している。第4章では動的マクビット方式の評価をRAPおよびDIRECTで採用されている方式を正規化して定量的評価式を与えことにより行い、その結果、本稿で提案する方式の有効性を示している。さらに動的マクビット方式がコード化リレーショナルに対する各種の演算に対して選択的なレコード操作が可能で並列・非同期処理に適した方式であり他の方式に比べて大きな処理能力があることを明らかにしている。

## 2 動的マクビット方式とそれに基づくリレーショナル演算アーキテクチャの特長

本章で説明する動的マクビットを活用した演算アーキテクチャは対象となるリレーショナル演算がSequential Codeで置換されているコード化リレーショナルであることを前提条件として成立しえる方式である。

以下にリレーショナル演算の具体的な各アルゴリズムの方式の比較の基準と動的マクビットを活用した方式の特長を項目別に述べる。

(1) マク情報方式：リレーショナル演算の具体的な方式上の最初の選択肢は各演算ごとにリレーショナル

の再構成を必要とする DIRECT 方式とするが、何らかの中間情報を用いて再構成を回避する方式とするかであるが、一般的に再構成コストは著大であり RAP のような tuple に固定的に付加されているマークビットを使用して再構成を極力抑える方式が有利である。

[2] 動的マークビット方式：RAP の固定マークビット方式に比べて本稿で提案する動的マークビット方式（各 tuple から分離された独自のデータ構造によって属性間およびリレーション間の演算結果を保持する）の有利な点としては ① データベース格納媒体上にマークビット用の格納スペースが不要 ② マークビットは各演算やトランザクション単位で必要に応じて動的にアロケーションされるので固定マークビットで問題であったマークビット数の制限による複雑な演算の遂行能力や tuple の共有によるマルチプログラミング多重度の制限を原則的に解消している。

[3] 属性単位、演算系の確立：主に演算系の基本データ単位としては属性単位 (Attribute 単位) の方式を採用している。属性単位の演算方式が tuple (リレーション) 単位の演算方式に比べて有利な点としては ① 必要最小限のデータで処理が実行可能 ② 格納単位が属性単位である必要は属性が転送されるにデータ転送時間の短縮が可能等が挙げられる。

[4] 並列アルゴリズム：リレーション間演算の projection や join の演算は本質的には逐次的な繰り返しのアルゴリズムを必要とする。このためにこのアルゴリズムを改良しないかぎり並列性を高めても（セル数を増加）セル間通信のためのオーバーヘッドが無視できず効果的な並列処理を期待できないというジレンマを解決するためにコード化値と動的マークビットを体系化することによって、完全非同期に各セルが並列動作可能なアルゴリズムを考案している。

[5] 動的マークビットの体系化：動的マークビットを単なる一次元的な有効 tuple の識別情報 (VTF) として利用するだけでなく、さらにある属性に対する同一種類のマルク化を示す (VTF<sub>g</sub>)、リレーション間の演算結果をリンクテーブルという形で保持する (AIR)、さらにコード化された属性のユニークな種別を定めるためにそれ各セル間での通信用に使用される (LIL) 等に体系化することによってセル間の通信同期コストを 1 tuple あたり 1 ビットで実現可能なアルゴリズムが実現され、繰り返

し演算の演算オーダーを tuple 数 × 種別数 (長) の一般的なアルゴリズムに比べて本方式にすると長 × 長 のオーダーに変換して処理を行なうことが可能となっている。一般に tuple 数に比べて種別数は少ないのでこの効果は大きい。

[6] 選択型の tuple 処理：動的マークビットの体系化によって次に処理すべき有効 tuple の情報かまえて判断可能となっているので、固定マークビット方式に比べて有効 tuple のみを選択的に処理することが可能で毎回全 tuple を検索するのに比べて効率向上が著しい。

[7] 機能分散方式：本方式では動的にコード化リレーションから動的マークビットの集合 (Association 情報) を生成し、それらの間の演算を通して最終結果を導く方式をとっており、本質的に Association 間の処理とコード化リレーションから Association 情報に変換する処理を機能分散することが可能である。総合的に処理効率の向上に寄与する。

[8] リソース管理の必要のないアルゴリズム：処理プロセスのメモリ容量に限界がある場合ページ単位にリソース管理が必要である演算子が存在するが、本方式では各データは基本的には全部 Association 情報に変換されその意味を保持するのでコード化リレーションの入れ替えが必要となる。これによってページ単位のリソース管理の甚大なコストが軽減される。

### 3 動的マークビット方式の処理アルゴリズム

処理アルゴリズムを論理的レベル、並列処理レベル、ハードウェアレベルでその詳細を述べる。

#### 3.1 論理アーキテクチャ

##### 3.1.1 論理アーキテクチャの基本概念

論理アーキテクチャを図 3.1 に示す。Association 集合は Entity 間の関係を示した情報の集合であり、その要素は複数のマークビット組である。Entity 集合は tuple の集合でありリレーションという部分集合に類別される。複数のマークビット列はル

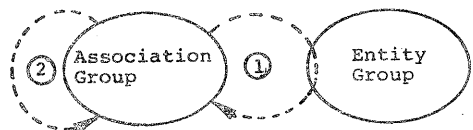


図 3.1 論理アーキテクチャ

ション内の tuple の関係を示し、その1つのマークビット列はリレーション内の有効 tuple を示している。次印はこのデータ構造間の演算を示している。①は演算対象が Association によって選択された Entity であり、結果が Association であることを示す。②は Association を対象とする演算で結果も Association であることを示す。

### 3.1.2 データ構造

Association には次の3種類のデータ構造が用意されている。

(1) VTF (Valid Tuple Flags) : マークビット列によってリレーション内の有効 tuple を示すものでありそのリレーション内で定義される。tuple の関連を示す基本的なデータ構造で他のデータ構造の要素ともなるが、単独でも Selection, Restriction 等の結果として生成される。また次の演算の際に対象とする tuple を選択する働きを持つ。このように、有効 tuple をマークビット列によって示す方法はデータが固定長であり処理が単純であるためハードウェア化が容易である。

(2) VTF<sub>g</sub> (VTF group) : 同一リレーションで定義されたVTFの並びで各VTFは共通の意味を持つ tuple であることを示し、VTF<sub>g</sub>全体でリレーションを分類している。各VTFをインデックスと考えればVTF<sub>g</sub>は配列ファイルと見れば、Group by, Compound Projection 等の演算時には生成される。このようにマークビットを高度に体系化することによって単なるマークビットでは不可能な演算が可能となる。

(3) AIR (Associations for Inter-Relation) : 異なるリレーションで定義された2つのVTF<sub>g</sub>の並びである。2つのVTF<sub>g</sub>に対応するVTF同士は互いに Join される tuple を示している。VTF<sub>g</sub>とVTF<sub>g</sub>を結びつけることによってリレーション内およびリレーション間の演算も新たなリレーションを再構成することも容易に実行可能である。

$VTF = \{ 1001 \dots 01 \}$  : マークビット列

$VTF_g = \{ VTF_1, \dots, VTF_n \}$  : 各VTFは同一リレーションで定義されている。

$AIR = \{ VTF_{g1}, VTF_{g2} \} = \{ (VTF_{11}, \dots, VTF_{1n}), (VTF_{21}, \dots, VTF_{2n}) \}$  :  $VTF_{ij}$ と $VTF_{zi}$ は互いに Join される tuple 群を示す。

### 3.1.3 基本操作

① Entity から Association 情報への変換  
VTFによって選択された Tuple から VTF, VTF<sub>g</sub>, AIR を生成する。

\* (リレーション, VTF) → VTF : Selection, Restriction

\* (リレーション, VTF) → VTF<sub>g</sub> : Group by, Projection

\* (リレーション, VTF, リレーション, VTF) → AIR : Join

② Association 間の演算

\*  $VTF \oplus VTF \rightarrow VTF$

\*  $VTF_{g1} \oplus VTF_{g2} \rightarrow VTF_{g3} : A_{2i} \oplus A_{3i} \rightarrow A_{3i} \quad i = 1 \dots n$   
 $A_{ij}$  は VTF,  $\oplus$  は bit 2 の論理演算

\*  $VTF_g \rightarrow VTF : A_1 \oplus A_2 \oplus \dots \oplus A_n \rightarrow VTF$

\*  $VTF_{g1} \otimes VTF_{g2} \rightarrow VTF_{g3} : VTF_{g1} = (A_1, \dots, A_m), VTF_{g2} = (A'_1, \dots, A'_n)$

• Compound Projection に利用する

$VTF_{g3} = (A''_1, \dots, A''_{mn}) \quad A''_{ij} = A_{ij} \otimes A'_{ij}$

• Division に利用する

$VTF_{g3} = (A''_1, \dots, A''_n) : A_{ij} \otimes (A'_1, \dots, A'_n)$

との間で演算を行ない結果によって  $A_{ij}$  と  $A'_{ij}$  とで出力するかどうかを決定する。

\*  $VTF_{g1} \otimes AIR \rightarrow VTF_{g2} : VTF_{g1}$  を定義しているリレーションを別のリレーションに変換し、リレーション間の Compound Projection を行うようにする。  
 $VTF_{g1} = (A_1, \dots, A_m), AIR = \{ (A'_1, \dots, A'_n), (b_1, \dots, b_n) \}, VTF_{g2} = \{ b_1, \dots, b_n \}, b_i = (A_{i1} \cdot A_1) \otimes \dots \otimes (A_{in} \cdot A_n)$   
 $\otimes : OR \quad \text{内積 } (A_{ij} \cdot A_j) \otimes b_j = \begin{cases} A_{ij} \cdot A_j = 0 \\ \text{なら } b_j \end{cases}$

\*  $AIR_1 \otimes AIR_2 \rightarrow AIR_3 : \text{リレーション A-B 間の } AIR_1 \text{ と B-C 間の } AIR_2 \text{ から A-C 間の関係を求める } AIR_3 \text{ を生成する。}$

$AIR_1 = \{ (a_1, \dots, a_m), (b_1, \dots, b_n) \}, AIR_2 = \{ (b'_1, \dots, b'_n), (c_1, \dots, c_n) \}, AIR_3 = \{ (a_1, \dots, a_m), (c_1, \dots, c_n) \}.$   
 $(c_1, \dots, c_n) \leftarrow (b_1, \dots, b_n) \otimes AIR_2$   
 $a_i, b_i, c_i$  は異なるリレーションとする。

### 3.2 動的マークビットによる並列処理化

#### 3.2.1 並列処理方式

集合演算を並列処理するためにリレーションを tuple 単位に均等に分割して各セルに分担させ、各セルではその部分リレーションを対象とする処理を行う。マークビット群を対象とする部分だけ生成・操作する。この結果、処理量が増える代わりに必要とする最大マークビット量も減少する可能性がある。

例えば 1000 tuple の 11レシオンでは VTFg の最大量は  $1000 \times 1000 = 10^6$  bit であるが 10セルに分割すると  $(100 \times 100) \times 10 = 10^5$  bit とする。

### 3.2.2 非同期並列処理

各セルでの処理を可能な限り非同期に実行するために UIL, Compact UIL (CUIL) という補助データ構造を導入する。これらは各セルで非同期に生成される VTFg の要素である VTF のつらかりを保持するためのものである。UIL は 11レシオンのある属性中に存在するアイテム値を全種類並べたものである。各アイテム値はコード化されていることを前提としているから、VTF がマークビットによって有効 tuple 番号を示したように UIL をマークビット列によって存在するアイテム値を示すこととする。この結果各セルの UIL は共通の意味を持つのでこの UIL の各アイテム値に対応させてこの属性に対する VTFg (転置マシ)内の VTF を生成可能な VTF 間のつらかりを保持する。CUIL は、アイテム値と関係なく単に VTF のセル間でのつらかりを保持するものでセルの各 VTF が全体の VTFg の何番目の VTF の一部分であることを、マークビット列で示している。(図 3.2)

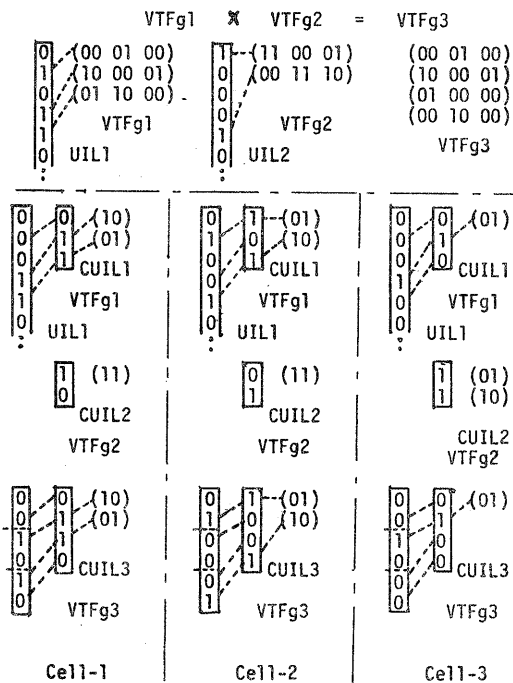


図 3.2 VTFg の並列処理

### 3.2.3 セル間通信

この処理方式でセル間通信とは、各セルで生成されるセル内 UIL (すなわち CUIL) の全ての OR をとることによって、グローバル UIL (すなわち GUIL) を作成することである。これは従来の通信とは異なるものであるといえる。この通信は VTFg, AIR を生成する演算に比する必要がある。通信量は UIL 長 (bit 長) であるが伝わる実質的な情報量はこの数字倍と考えられる。また、この通信とセル内の処理は、ほとんどの場合分離できるので、機能分散処理化が可能である。

### 3.3 物理アーキテクチャ

#### 3.3.1 ハードウェア構成

動的マークビット方式による並列処理演算は次のようなハードウェア構成で効果的に実現できる。

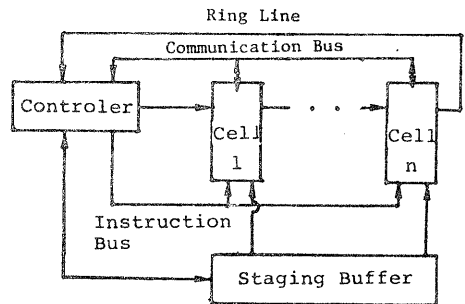


図 3.3 動的マークビット方式のハードウェア構成図

・コントローラ：各セルにブロードキャストで命令を送り、その完了は各セルで終了した時点で認識できる。コード値データをセルとステージングバッファの間で管理する。

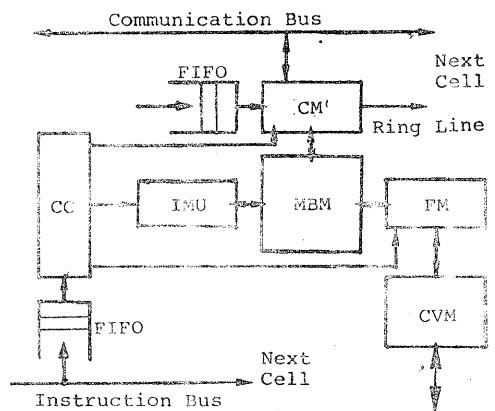


図 3.4 Cell の構成

- CC : 命令のシフトを考慮し、各ユニットが並列動作が可能になるように命令を各ユニットに配る
- FM : コード値からマクビット樹を作成し、コード値の更新等のコード値に関する処理を全て行う
- IMU : マクビット間の単純な論理演算を行う
- CM : ①セル間での通信をワイヤードORを利用して行う。ここでは、全セルが同期して動作する。②命令の完了標号、関数演算の結果の累積等をセルからセルへ方向に伝達し、最終的にコントローラに送る。
- MBM : VTF, VTFg, UIL等の全てのマクビットを格納するメモリでIMU, FM, CMの各ユニットからアクセス可能。
- CVM : リレーションの属性単位のコード値が各セルに均等に分割されて格納されているメモリ。セル内はFM, IMU, CMに機能別のプロセッサが用意されている。それらは各々①コード値集合に対する処理②マクビット間の処理③通信に関する処理に大きく大別し、各機能を専用化して処理と実行する。

### 3.3.2 ハードウェア基本命令

基本的な命令を表3-1に示す。このうちから簡微的なハードウェア命令をいくつか説明する。

#### • Generate UIL (ATL, VTF#1; UIL#1)

ATLとVTF#1で選択されるCVM内のアイテム値をフェッチし、そのコード化された値自身によってアドレスされるUIL内のビットをマークする。VTFg, ATRを作成する際の前処理である。

#### • Generate VTF (ATL, VTF#1, UIL#1; VTF#2)

ATLとVTF#1によって選択されるアイテム値自身によってアドレスされるUIL#1内のビットがマークされている場合は、そのアイテム値をそのtupleを有効としてVTF#2にマークする。

#### • Generate VTFg (ATL, VTF#1, UIL#1; VTFg#1)

ATLとVTF#1によって選択されるアイテム値自身によってアドレスされるUIL#1内のビットがUIL#1の先頭から何番目のマクビットにおいて、アイテム値に対応するVTFを認識し、そのVTFにそのアイテム値をそのtupleを有効としてマークする。

#### • Global OR UIL (UIL#1; UIL#2)

UIL#1のビット列をジョイン・イン・ジョインバス上でワイヤードORを行い、結果のグローバルUILをUIL#2としておける。

#### • Global OR UIL (UIL#1; CUIL#2)

UIL#1のビット列をジョイン・イン・ジョインバス上でワイヤードORを行い、結果のグローバルUILをCUIL#2としておける。各セル内のUIL#1と比較することによってC-UILを生成する。

#### • Global OR UIL (UIL#1, UIL#2; CUIL#1, CUIL#2)

UIL#1のグローバルUILとUIL#2のグローバルUILとのANDにおいて、Joinされるアイテム値を示すUIL'を作成する。

$$UIL' \#1 \leftarrow UIL \#1 \text{ and } UIL \#2$$

$$UIL' \#2 \leftarrow UIL \#2 \text{ and } UIL \#1$$

C-UIL#1とC-UIL#2はそれぞれUIL#1とUIL#2に対応するC-UILとなる。

### 3.3.3 基本的なQueryの展開方法

① データ駆動型制御：コントローラはデータフローで表現されたセルへの命令列を待ち、下位メモリからのステータスがマップにおいて、実行に必要なデータが揃った命令からセルへ送る。この結果処理中に最適な形で命令の実行順を決定でき、各セルにおいては、即座に実行が可能となる。

② セル内の並列処理：セル内の各ユニット(FM, IMU, CM)は並列動作が可能であるのでデータを生成されていくが互いに並列処理を併行

表3.1 基本命令

	source	destination	unit
CS	ATL, VTF, OP, Value	VTF	FM
CS	ATL, ATL, VTF, OP	VTF	FM
G. UIL	ATL, VTF	UIL	FM
G. VTF	ATL, VTF, UIL	VTF	FM
G. VTFg	ATL, VTF, UIL	VTFg	FM
G1. OR UIL	UIL	UIL	CM
G1. OR UIL	UIL	CUIL	CM
G1. OR UIL	UIL, UIL	(UIL, CUIL)*2	CM
Multiply	VTFg, VTFg	VTFg	IMU, CM
Multiply	VTFg, VTFg, VTFg	VTFg	IMU, CM
Select VTF	VTFg, VTFg	VTFg	IMU, CM
B. VTF	VTF, VTF	VTF	IMU
B. VTFg	VTFg	VTF	IMU
B. VTFg	VTFg, VTFg	VTFg	IMU
Total	ATL, VTF	ACC	FM
I. Total	ATL, VTFg	ACC	FM
Count Bit	VTF	ACC	IMU
Write	Value, VTF	ATL	FM

CS: Conditional Search, ATL: Attribute Location  
G.: Generate, G1: Global, B.: Boolean Operation  
I.: Image, OP: Operator

```
SELECT UNIQUE S.CITY, J.CITY
FROM S, SPJ, J
WHERE S.S# = SPJ.S#
AND J.J# = SPJ.J#
```

SEQUEL

G.UIL(S.S# ; UIL1)  
 G.UIL(S.PJ.S# ; UIL2)  
 G1.OR UIL(UIL1,UIL2 ; UIL1,CUIL1,UIL2,CUIL2)  
 G.VTFg(S.S#,UIL1 ; VTFg1)  
 G.VTFg(S.PJ.S#,UIL2 ; VTFg2)  
 G.UIL(J.J# ; UIL3)  
 G.UIL(S.PJ.J# ; UIL4)  
 G1.OR UIL(UIL3,UIL4 ; UIL3,CUIL3,UIL4,CUIL4)  
 G.VTFg(J.J#,UIL3 ; VTFg3)  
 G.VTFg(S.PJ.J#,UIL4 ; VTFg4)  
 Multiply(VTFg2,VTFg4,VTFg3 ; VTFg5)  
 G.UIL(S.CITY ; UIL6)  
 G1.OR UIL(UIL6 ; CUIL6)  
 G.VTFg(S.CITY,UIL6 ; VTFg6)  
 Multiply(VTFg6,VTFg1,VTFg5 ; VTFg7)  
 G.UIL(J.CITY ; UIL8)  
 G1.OR UIL(UIL8 ; CUIL8)  
 G.VTFg(J.CITY,UIL8 ; VTFg8)  
 Multiply(VTFg7,VTFg8 ; VTFg9)

図 3.5 問い合わせの展開例

#### 4. 構造的マージット方式の評価

前章ではバリエーション的マージット方式 (TMB方式) による集合演算-マージットを RAP, DIRECT との対比を平比して評価を行った。TMB方式の特徴は (1) 選択的理想処理方式 (2) 並列処理方式 (3) 集合演算の結果の保持方法 (4) 処理の並行性の最大化に対する配慮などの点である。これらは、リレーショナル演算と完全、かつ効率良く処理することを可能としている。

##### 4.1 各種演算方式の評価式と評価方法

一般にリレーショナルに対する集合演算は次のようにステップで実行されると仮定できる。① 演算対象 tuple の認識 ② アイテム値の参照、および比較操作 ③ 演算結果の保持及び書き込み、である。①に関連して RAP では tuple に付随した固定マージットを判断することによって行われる。DIRECT では演算対象 tuple と他の tuple を識別する特別なデータ構造は用意されていない。例えば対象リレーショナルはすべて有効 tuple であるように再構成が行われるので①の処理時間を必要としない。②のアイテム値の参照、及び比較操作は①の結果に従って行われる。③の演算結果の保持については TMB方式では Association 情報に変換される。RAP の場合はリレーショナル演算に対してはビットのマージットを、Join 等の演算では新しいリレーショナルを生成する。DIRECT ではすべての演算の結果は新しいリレーシ

ョンによって保持される。代表的なリレーショナル演算については①②③の項目別に各演算方式の評価を示す。

##### 4.1.1 評価式

評価式内の係数  $A_i, F_j, W_k$  は各々①②③に対応する係数であり評価式の一般形式(1)式のように表わされる

$$E = A_i \cdot G_1(t, n, k) + F_j \cdot G_2(t, n, k, d) + W_k \cdot G_3(t, n, k, d) \dots \dots \dots (1)式$$

(1)式の  $G_1, G_2, G_3$  は  $t, n, k, d$  の関数であり各パラメータとマージットの特徴に応じて異なる。各係数は生データ処理に要する時間、コストあるいはある基準値に対する比、等の意味をもつ。係数の定義を次に示す。(コストは、時間その他の単位をとる)

- TMB方式の評価式の係数
- $A_n$  : 1 マージットを参照するのに要するコスト
- $A_{UIL}$  : G.UILにおいて1 マージットを参照するのに要するコスト
- $A_{VTFg}$  : G.VTFg " "
- $F_{12}$  : 1 アイテム値をマージット比較操作するのに要するコスト
- $F_{UIL}$  : G.UILにおいて1 アイテム値をマージットするのに要するコスト
- $F_{VTFg}$  : G.VTFg " "
- $F_{G1ORUIL}$  : G1.OR UIL を実行するのに要するコスト
- $F_{AND UIL}$  : AND UIL " "
- $F_{VTF}$  : 1 VTF をマージットするのに要するコスト
- $W_n$  : 1 マージットを VTF に登録するのに要するコスト
- $W_{UIL}$  : G.UILにおいてUILに1 マージットを登録するのに要するコスト
- $W_{VTFg}$  : G.VTFg " VTFg "
- $W_{G1ORUIL}$  : G1.OR UILにおいてUILを各セルへ書き込むのに要するコスト
- $W_{AND UIL}$  : AND UIL " "
- $W_{SV}$  : 1-CUIL を書き込むのに要するコスト
- RAPの係数  $A_n$  : 1 マージットを参照するのに要するコスト
- $F_{12}$  : 1 アイテム値をマージット比較するのに要するコスト,  $F_{13}$  : 1 アイテム値をマージットし各セルへ登録するのに要するコスト,  $W_k$  : 1 マージットを登録するのに要するコスト,  $W_{k2}$  : 1 タプルを書き込むのに要するコスト
- DIRECTの係数  $F_{02}$  : 1 タプルをマージットするのに要するコスト
- $F_{02}$  : 1 タプルをマージットし各セルへ登録するのに要するコスト
- $W_{02}$  : 1 タプルを書き込むのに要するコスト

各係数は物理的要素で決定される値であり、変数部は各処理方式の方式上の特徴が反映される。各演算については TMB, RAP, DIRECT の各方式の評価式を示す。ただし RAP の場合、演算結果はマ

ビットに併替され、Joinの場合にはレコードを生成する。DIRECTでは各演算後にレコードの再構成を行なう。ゆえに常に全tupleが有効な演算対象となる。(ET: TMB式, ER: RAP, ED: DIRECT)

(1) Selection

$$ET\text{-select} = AR_1 \cdot t/n + FR_2 \cdot d \cdot t/n + WR_1 \cdot d \cdot t/n \dots\dots\dots (1-1) \text{式}$$

$$ER\text{-select} = AR_1 \cdot t/n + FR_2 \cdot t/n + WR_1 \cdot d \cdot t/n \dots\dots\dots (1-2) \text{式}$$

$$ED\text{-select} = F'_{D2} \cdot d \cdot t/n + W_{D2} \cdot d \cdot t / (n \cdot k) \dots\dots\dots (1-3) \text{式}$$

(2) Projection

$$ET\text{-Proj} = (A_{Eu} + A_{Ev}) \cdot t/n + (F_{Eu} + F_{Ev}) \cdot d \cdot t/n + F_{Eu} + (W_{Eu} + W_{Ev}) \cdot d \cdot t/n + W_{Eu} \dots\dots (2-1) \text{式}$$

$$ER\text{-Proj} = 2AR_1 \cdot k \cdot t/n + F'_{R2} \cdot k + FR_2 \cdot (t \cdot k/n) \cdot ((1 + (1/k))/2) + WR_1 \cdot (k + (d \cdot t - k)/n) \dots\dots (2-2) \text{式}$$

$$ED\text{-Proj} = F'_{D2} \cdot k + F_{D2} \cdot (d \cdot t \cdot k/n) \cdot ((1 + (1/k))/2) + W_{D2} \cdot k \cdot (d \cdot t/n) \cdot ((1 + (1/k))/2) \dots\dots (2-3) \text{式}$$

(3) Join

$$ET\text{-join} = (A_{Eu} + A_{Ev}) \cdot (ts + tp)/n + (F_{Eu} + F_{Ev}) \cdot (ds \cdot ts + dp \cdot tp)/n + F_{Eu} + F_{Ev} + (W_{Eu} + W_{Ev}) \cdot (ds \cdot ts + dp \cdot tp)/n + (W_{Eu} + W_{Ev} + W_{U}) \dots\dots (3-1) \text{式}$$

$$ER\text{-join} = 4AR_1 \cdot \frac{ts}{n} + 3AR_1 \cdot ts \cdot tp/n + F'_{R2} \cdot ts + FR_2 \cdot ts \cdot tp/n + 2WR_1 \cdot ds \cdot ts + WR_1 \cdot ds \cdot ts \cdot \alpha \cdot dp \cdot tp / (n \cdot kp) + (WR_{2s} + WR_{2p}) \cdot \alpha \cdot ds \cdot ts \cdot dp \cdot tp / kp \dots\dots (3-2) \text{式}$$

$$ED\text{-join} = F'_{D2} \cdot ds \cdot ts + F_{D2} \cdot ds \cdot ts \cdot dp \cdot tp/n + (W_{D2s} + W_{D2p}) \cdot \alpha \cdot ds \cdot ts \cdot dp \cdot tp / kp \dots\dots (3-3) \text{式}$$

(4) Implicit Join

$$ET\text{-I-Join} = A_{Eu} \cdot ts/n + A_{Ev} \cdot tp/n + F_{Eu} \cdot ds \cdot ts/n + F_{Ev} \cdot dp \cdot tp/n + W_{Eu} \cdot ds \cdot ts/n + W_{Ev} \cdot \alpha \cdot dp \cdot tp/n \dots\dots (4-1) \text{式}$$

$$ER\text{-I-Join} = AR_1 \cdot (ts + tp)/n + F'_{R2} \cdot ts + FR_2 \cdot ts \cdot tp/n + WR_1 \cdot ds \cdot ts + WR_1 \cdot ds \cdot ts \cdot \alpha \cdot dp \cdot tp / (n \cdot kp) \dots\dots (4-2) \text{式}$$

$$ED\text{-I-Join} = F'_{D2} \cdot ds \cdot ts + F_{D2} \cdot ds \cdot ts \cdot dp \cdot tp/n + W_{D2} \cdot ds \cdot ts \cdot dp \cdot tp / kp \dots\dots (4-3) \text{式}$$

(5) Division

$$ET\text{-Division} = A_{Eu} \cdot (ts + tp)/n + A_{Ev} \cdot 2tp/n + F_{Eu} \cdot (ds \cdot ts + dp \cdot tp)/n + F_{Ev} \cdot 2dp \cdot tp/n + 2F_{Gou} + F_{VTF} \cdot B_g(n) \cdot k_s \cdot B_{pcn} \cdot k_p + W_{Gu} \cdot (ds \cdot ts + dp \cdot tp)/n + W_{Ev} \cdot 2dp \cdot tp/n + 2W_{Gou} + W_{sv} \cdot k_p \dots\dots (5-1) \text{式}$$

$$ER\text{-Division} = 2AR_1 \cdot kp \cdot tp/n + 2AR_1 \cdot ts \cdot kp + AR_1 \cdot ts \cdot tp \cdot kp/n + F'_{R2} \cdot kp + FR_2 \cdot (tp \cdot kp/n) \cdot ((1 + (1/kp))/2) + F'_{R2} \cdot ts \cdot kp + FR_2 \cdot ts \cdot tp/n + WR_1 \cdot (kp + (dp \cdot tp - k)/n) + WR_1 \cdot (2ds \cdot ts \cdot kp + dp \cdot tp/n) \dots\dots (5-2) \text{式}$$

$$ED\text{-Division} = F'_{D2} \cdot kp + F_{D2} \cdot (dp \cdot tp \cdot kp/n) \cdot ((1 + (1/kp))/2) + F'_{D2} \cdot ds \cdot ts \cdot kp + F_{D2} \cdot ds \cdot ts \cdot dp \cdot tp/n + W_{D2} \cdot kp \cdot (dp \cdot tp/n) \cdot ((1 + (1/kp))/2) + W_{D2} \cdot dp \cdot d'p \cdot tp/n \dots\dots (5-3) \text{式}$$

4.2 TMB方式の性能評価

(1) 選択的理想処理方式：TMB方式ではマ  
ークビットがtupleと物理的に分離されており  
あらかじめ特定の有効tupleの列を選択的に操  
作することが可能である。この点、各式のFの  
項におけるdの値(d ≤ 1)によって表現されている。  
図4.1はSelectionにおける処理時間をRAP  
を基準とした比によって示したものである。有効  
tuple数が1/(種別数)(=d)であると仮定  
するとSelectionに要する処理時間の比はdの  
値が小さくなる程TMB方式及びDIRECT方式  
はこの効果が著しい。他の演算においてもdの  
寄与する項によって選択的にレコード操作を行なえる  
ので効率向上に大きく寄与する方式であるといえる。

(2) 並列処理の効果：一般にSelection  
(条件検索)は並列処理により1/(セル数)の処  
理時間となるが、Projection, Join, Implicit  
Join, Division等の複雑なセル間交換およ  
び多回の繰り返し比較演算が大きな負担とな  
る。つまり単純にセル数を増加して通信コストを軽  
減する対策が必要である。TMB方式では前章  
で述べたように通信コストを1アイテムに対してビット  
でかつ複雑な制御を伴わない方式を考案して

いる。図-4.2, 図-4.3, 図-4.4, 図-4.5はそれぞれ Selection, Projection, Join, Divisionの処理時間とセル数の関係を各標式から算出の値を代入して導いたものである。Selectionでは並列セル数の増加が処理効率を正比例的に向上させる。DIRECTではリレーションの再構成コストを示す 4.3 式の第2項の効果で約100倍処理時間がかかる。ProjectionではRAP, DIRECTの場合(2.2)式の  $F_{22}$  と  $k$ , (2.3)式の  $F_{22}$  の値の項に漸近するので  $k$  の値が大きい場合はセル数を増しても並列性の効果が著かたなくなる。TMB方式は  $k=512$  の場合 RAP, DIRECT方式と比べて各々約10倍~約100倍, 数1000倍の効率向上がある。Joinではその結果をTMBで保持するので新しいリレーションを生成する必要がない。RAP, DIRECTでは 9.2 式第7項, 9.3 式第3項が示すリレーション再構成コストおよび 9.2 式第4項, 9.3 式第2項が示す繰り返し並列の比較操作の効果は図-4.4に表われている。DivisionにおいてRAPでは(5.2)式第4項, 第6項, DIRECTでは(5.3)式, 第1項, 第5項においてグラフは漸近する。TMB方式では第6項がグラフを漸近させる。これはTMB方式のVTF<sub>g</sub>間の複算が種別数の積のみに依存することからである。

(3) 集合演算結果の保持方式

集合演算の結果の保持方式は各処理方式の特徴を良く表わしている。図-4.4に示すJoinの処理時間とセル数の関係はRAP, DIRECTの場合リレーションの再構成コストのためにセル数がある程度になると並列性の効果が低くなる。TMB方式ではリレーションの再構成コストの必要がないのでその点有利となる。

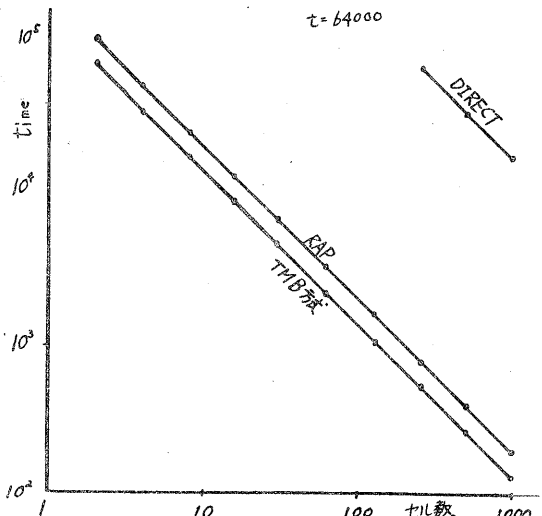


図4.2 Selectionにおける並列処理の効果

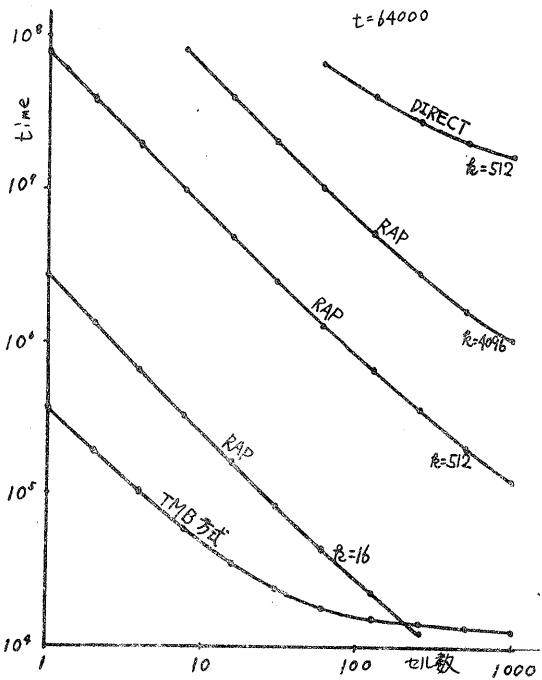


図4.3 Projectionにおける並列処理の効果

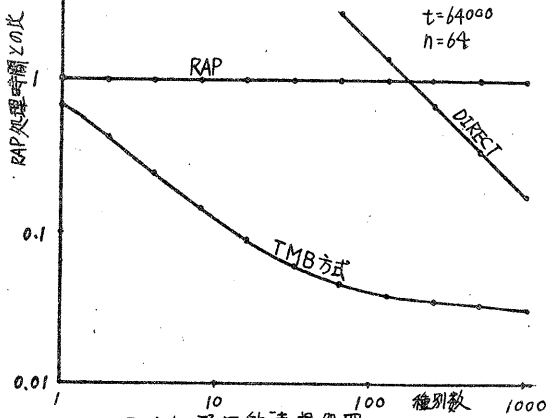


図-4.4 選択的連想処理



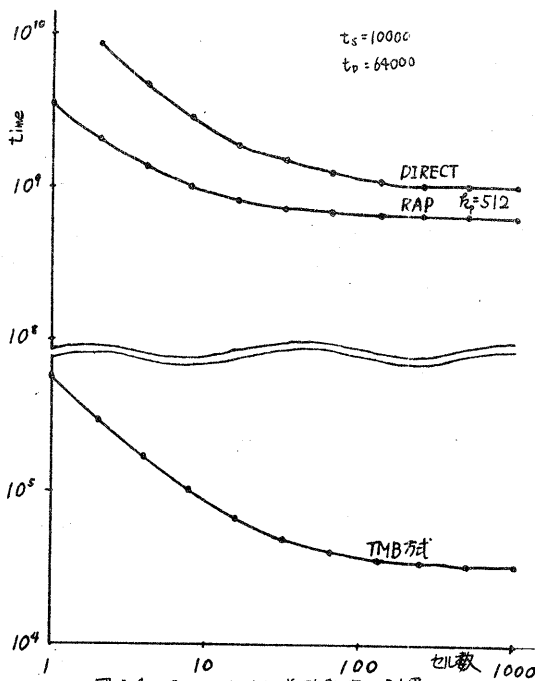


図4.4 Joinにおける並列処理の効果

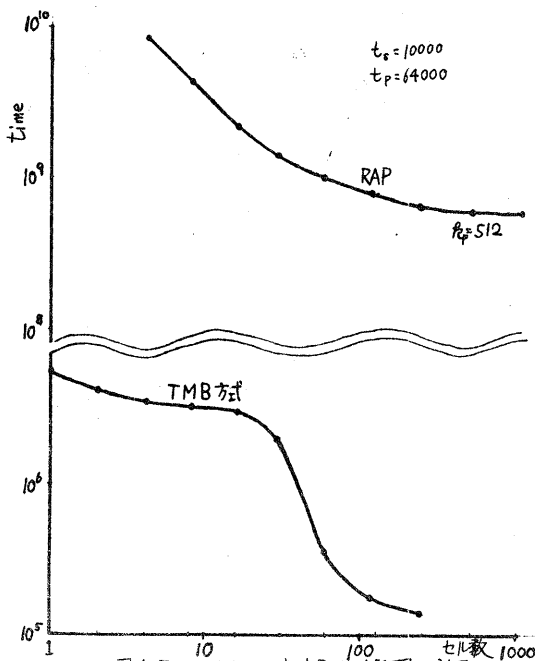


図4.5 Divisionにおける並列処理の効果

(4) 演算に必要なる活性データ量

2回の Selection と Implicit Join, Projection から成る問い合わせを例にして各方式の活性データ量の推移について考察する。

SELECT UNIQUE SN FROM SPJ WHERE J#='J1' AND P# IN (SELECT P# FROM P WHERE COLOR='RED')  
問い合わせ例

$t_p$ : P-リレーション タプル数  
 $t_s$ : SPJ-リレーション タプル数  
 $n$ : セル数,  $k_p$ : オペランド属性 種別数  
 $m$ : RAP, DIRECT における タプル長 (bit)  
 $d_p$ :  $t_p$  中の有効タプルの割合  
 $d_s$ :  $t_s$  " "  
 $\beta_{0j}$ :  $k_p$  中各セル内の 3-値集合がゼロ 種別数の割合の平均値 ( $\leq 1$ )

表 4.1 Xメモリ使用量の推移

	TMB方式	RAP	DIRECT
Selection (obj='RED')	$17 \times t_s$	$(m+4) \times t_s$	$m \times t_s$
Selection (J#='J1')	$17 \times t_p$	$(m+4) \times t_p$	$m \times t_p$
Implicit Join	(1) $64000n + 17 \times t_s$ (2) $64000n + 17 \times t_p$	$(m+4) \times (t_s + t_p)$	$m \times (d_s t_s + d_p t_p)$
Projection	(1) $64000n + 17 \times t_p$ (2) $64000n + 17 \times t_p + \beta_{0j} \times k_p \times t_p$	$(m+4) \times t_p$	$m \times d_p \times t_p$

表 4.1 より TMB方式の最大Xメモリ使用量は

$$64000 \times n + 17 \times t_p + \beta_{0j} \times k_p \times t_p$$

RAPでは

$$(m+4)(t_s + t_p)$$

DIRECTでは

$$m \times (d_s t_s + d_p t_p) \text{ である。}$$

従って RAPでは

$$M_{RAP} < \frac{64000n + 17 \times t_p + \beta_{0j} \times k_p \times t_p}{t_s + t_p} - 4$$

DIRECTでは

$$M_{DIR} < \frac{64000n + 17 \times t_p + \beta_{0j} \times k_p \times t_p}{d_s t_s + d_p t_p}$$

の場合、各々TMB方式より、最大Xメモリ使用量は少なくなる。しかしながら Implicit Join, Projection 等では RAP, DIRECT 方式は TMB方式に比べて処理時間が(2)で述べたように大きくなるので、Xメモリ・レポート獲得TMB方式の方が一般的にかなり小さくなる。また Join が含まれる問い合わせでは RAP, DIRECT では大きなXメモリ領域を必要とする可能性がある。セル内の記述領域に限界がある場合、TMB方式は他の方式に比べて本質的違いがある。それはすべての演算が3-値集合から Association 情報に変換する方式

を基本的に採用している。そのためにコード化レシジョンに対しては繰り返し逐ってページをロードすることが原則的に必要なくなり、処理プロセスのメモリ容量に制限がある場合（処理データがすべてメモリにロード出来ない場合）でも、多回のリフレッシュメントによる処理効率の劣化を未然に防いでいる。

## 5. おわりに

本稿では筆者らが現在検討を進めている競合リレーショナルデータベースマシン(SPIRIT)で採用されているリレーショナル演算アーキテクチャについてその特長、処理方式の詳細、具体的な実現可能なレベルのアーキテクチャ、そして他の演算方式との比較を定量的に行ない、本稿で提案している動的マクヒット方式の有効性を明らかにした。

問題点としては、コード化レシジョンを前提していることであるが、この点に関してはあるデータ圧縮効果による演算対象データ量とそれに転送時間を軽減することか可能であり、たとえ実行時への変換コストを考えても十分に効率の良い方式である。

## 謝辞

本稿をまとめる上で貴重な御教示と御討論をいただいた本学大学院 小沢裕之君、瀬尾和男君、加藤洋君に敬意を覚りて感謝致します。

## 参考文献

- [1] E.A. Ozakaraham., et al.: "RAP - An associative processor for data base management", NCC, (1975)
- [2] S.A. Schuster., et al.: "RAP. 2 - An Associative Processor for Databases and Its Applications", IEEE Trans. on computers, Vol. C-28, No. 6, (June 1979)
- [3] E. Babb.: "Implementing a Relational Database by Means of Specialized Hardware", ACM Trans. on Database systems, Vol. 4, No. 1, (Mar. 1979)
- [4] D.R. McGregor., et al.: "High performance for database systems," in Systems for Large Databases, North-Holland, Amsterdam, 1976, pp. 103-116
- [5] D. J. De Witt.: "DIRECT - A Multiprocessor Organization for Supporting Relational Database Management Systems", IEEE Trans. on computers, Vol. C-28, No. 6, (June 1979)
- [6] 園分, 他.: "磁気バブルデータベース計算機EDCのアーキテクチャ", 信学技報, EC78-46, (1978)
- [7] N. Kamabayashi., et al.: "SPIRIT: A New Relational Database Computer Employing Functional-Distributed Multi-Microprocessor Configuration", The 1st International Conference on Distributed Computing Systems, (Oct. 1979)
- [8] 上林, 他.: "リレーショナルデータベース演算のハードウェアアルゴリズムと、これに基づくデータベースマシンアーキテクチャ", 情報処理学会計算機アーキテクチャ委員会 35-4 (1979.9.12)