

# 放送通信機能を用いた分散型データベースシステムの通信処理の実現について

## 要 訳 誌

(財)日本情報処理開発協会)

### 1. 序

分散型データベースシステム(以降DDBSと記す)とは、意味的に関連するデータを有する複数のデータベースシステム(DBS)が、通信網によって相互結合され、利用者に1つの仮想DBS(i.e. 1つのスキーマと1つのアピセス言語)としてのDBSサービスを提供出来るシステムである[TAKIM 78, 79]。DDBSとしては、70年代後半に、SDD-1[ROTHJ84], POLYPHEME[LADISM80]等が開発されている。これは、(i)関係DBSから成る同種システム, (ii)広域パケット交換網(e.g. ARPANET, CYCLADES)の利用, (iii)問合せの分散処理の実現等を特徴としている。これを、第I期DDBSと呼ぶ。ここでの主要な問題の1つは、広域網の伝達性から生じる性能問題である。これらに対して、80年頃開始されている新DDBSは、次の様な特徴を有している。

- (i) 異種DBS(e.g. CODASYL型DBS)の統合(JDDBS[TACH80], MULTIBASE[SHITHJ81])
- (ii) ローカル網(LAN)(e.g. Ethernet)の利用(SIRIUS-DELTA[LEBIJ81])による通信
- (iii) 分散更新処理の実現 } 問題の解決,
- (iv) 従来の大型DBSに加えて、個人用小型DBS(パソコン + DBS + WP)の組み込み
- (v) オフィス情報処理への適用

当協会が1977年より開発を進めているJDDDBS(Jipdec DBS)は、上記5つの目標達成を目指している。これらシステムを第II期DDBSと呼ぶ。第II期DDBSの主要な特徴の1つは、高速度(~10Mbps)ローカル網による、第I期DDBSにおける通信問題の解決を試みていることである。

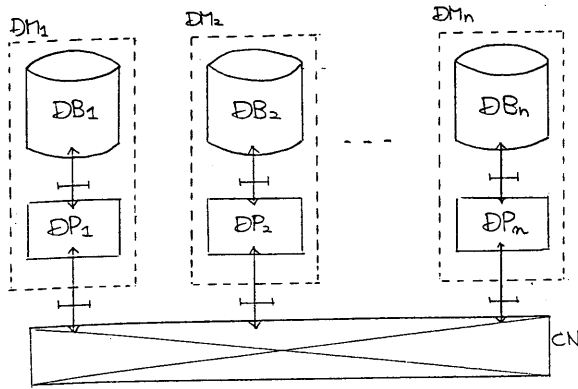
今後の情報システムは、各組織体(例えば会社の部、課)毎に、組織内の個人用DBS、既述大型DBSがローカルに結合されたローカル情報システムを基本単位として、これ等が広域網と結合された形態を取ると考えられる。ローカル情報システム(LIS)内では、各DBS間での大量、頻繁なデータ利用が、大容量LANを通して統合的になされる。一方、LIS間では、大量かつ頻繁なデータ通信処理に対する要求はあまりないと考えられる。又、異な、た組織体間でのデータ統合は、データの意味、伝達手段の差等が問題となり、実現はかなり困難とみられる。従って、DDBSは、まずLIS毎にLANを用いて実現される。

Ethernetの様なLAN、又無線網の様に通信媒体(e.g. 同軸ケーブル、無線)を共有する形態のものは、従来の広域パケット交換網に対して、高速、高信頼であることに加えて、物理的に1対多(放送)通信機能をサポート出来る特徴を有している。DDBSにおける通信処理は、冗長コピーの同時更新問題[BERNJP81]に見られる様に本質的に1対多通信に基づいている。この為、LISとしてのDDBSを共有媒体型LANで実現することは有効である。

本論文では、開発中のJDDDBS[TAKIM 78, 79, 80, 82a, b]で実現を試みている放送通信機能に基づいた通信処理方式について論じる。異種DBSのサポートについては、[TAKIM 79, 80, 82b]を参照されたい。

## 2. 分散データベースモデル

分散データベースシステム(DDBS)は、図2.1に示す様に通信網を結合させた複数のデータベース(DM<sub>1</sub>, ..., DM<sub>n</sub>)の集合としてモデル化する。各データベース(DM)は、データベース(DB)とDBプロセッサ(DP)とから成る。DB



DM : data module  
 DB : database      DP : DB processor  
 CN : communication network

図2.1 DDBSモデル

には、関係の集合が格納されている。DPは、自分のDB内の関係に対する関係演算と、組単位の変更が出来るプロセッサである。又、DPは、通信網CNを介して、他のDMへ/から関係を転送出来る。DPが受信した関係、DPが処理した結果関係は、自分のDB内に格納される。CNは、DM間での高信頼な基本通信機能を有している。DDBSに対する要求の解は、DM間での通信と、DM内での処理とによって得られる。

## 3. 放送網のモデル

放送網では、あるデータベース(DM)から発せられたメッセージを、他の全てのDMは、同時に受信出来る。Ethernet、無線網の様に、同軸ケーブル、無線といった通信媒体を共有した形態の網では、あるDMから出された信号は、網に結合されている全てのDMにまで受信出来る。従って、これ等の網では、放送通信機能が物理的にサポートされていると言える。

利用者から見た通信網は、届け先と届けたい情報量に対するコスト(全通信時間)によってモデル化出来る。ここで、 $C_i\{j_1, \dots, j_n\}(x)$  ( $n \geq 1$ ) を、DM<sub>i</sub> から、 $n$ 個のDM<sub>j<sub>1</sub></sub>, ..., DM<sub>j<sub>n</sub></sub> に情報量  $x$  を転送する時の全通信コスト(時間)とする。広域パケット交換網では、一対の通信実体間の高信頼な通信機能が、通信距離に依存せず通信量のみに応じてサポートされている。従って、(1)式の様になる。

$$C_i\{j_1, \dots, j_n\}(x) \cong n(a + b \cdot x) \quad \dots (1)$$

ここで、 $a$  と  $b$  とは定数である。 $a$  は通信リニアの初期化、切断のコストである。一対、 $R_i\{j_1, \dots, j_n\}(x)$  を、 $x$  の時の応答時間とする。同時に、 $n$ 本の通信リンクが設定出来る場合には、最良の応答時間  $a + b \cdot x$  をもたす。よって

$$R_i\{j_1, \dots, j_n\}(x) \geq a + b \cdot x \quad \dots (2)$$

これに対して、放送網では網上の信号を全てのDMが受信出来る為に、全通信

コストは、転送失の数にも独立に有り、(3)式の様になる。

$$C_i \{j_1, \dots, j_n\} (x) \cong a + b \cdot x \quad \dots (3)$$

又、放送網では、通信媒体が全てのDMで共有されている為に、同時に1つのDMのみが信号を出すことが出来る。よって、応答時間は(3)と同一になる。

$$R_i \{j_1, \dots, j_n\} (x) \cong a + b \cdot x \quad \dots (4)$$

実際の通信網では、経路選択、待行列、競合等による統計的遅延が加わると(4)式はより複雑なものとなる。ここでは、網は軽負荷であると仮定する。この仮定のもとで、(3)及び(4)式が成り立つ網を放送網と定義する。

#### 4. 分散問合せ処理

複数のデータモジュール(DM)内の関係を参照する問合せを、DM間での通信と、各DM内でのローカル処理と(i.e. 通信処理[TAJAK80])に分けて解くことを分散問合せ処理と呼ぶ。本章では、3.で定義した放送網を用いた分散問合せ処理アルゴリズムについて論じる。

##### 4.1 仮定

分散問合せ処理問題を考えるうえでの基本仮定について述べる。第1に、コストについて考える。放送網の通信コストは、(3)及び(4)式が成り立つとする。各DMでの処理コストは、広域網と比較して、放送網の高圧性とDMとして小型システムも考えることから、通信コストと比肩し得るものとする。\$R\_1, R\_2\$を関係、\$a, b\$を各々の属性とする。\$j(R\_1, R\_2), p(R\_1), r(R\_1)\$を各々、結合\$R\_1[A=b]R\_2\$, 制限\$R\_1[A=a]\$, 射影\$R\_1[A]\$とする。\$|R\_1|\$と\$|R\_2|\$を各々\$R\_1, R\_2\$のサイズとし、遅算\$\alpha\$に於いて\$P(\alpha)\$を凡の処理コストとする。又、\$st(a)\$を、属性\$a\$の遅延度[HEONAZ8]とする。この時、関係遅算コストを次の様に仮定する。

$$\begin{aligned} P(j(R_1, R_2)) &= |R_1| + |R_2| \\ P(r(R_1)) &= |R_1| \cdot st(a) \\ P(p(R_1)) &= |R_1| \end{aligned}$$

各DMには、互いに独立な関係が格納されていて、冗長コピーの存在は考えないものとする。関係の存在が不可視な視野(i.e. 全体概念スキーム[TAKIM78, 81])に対する問合せは、問合せ修正手法[STONM76]を用いて、各DM内の関係を参照する問合せ(全体問合せと呼ぶ)に既に交換されている[TAKIM82]とする。

問題を簡単にする為に、全体問合せは、sum, maxといった統計関数(aggregate)を含まず、結合として等結合のみから成るものとする。全体問合せは、関係計算言語QUEL[STONM75]によって記述されるものとする。

最後に、分散問合せ処理の最適化の目標について考える。通信処理スリジュール(又は単にスリジュール)とは、DM間でのデータ通信と、DM内での処理(関係遅算)との実行順序を表す順序集合とする。全体問合せに対して複数の可能なスリジュールが存在するが、この中から以下の目標を達成する最も速いスリジュールが選ばれ

- u) 全通信処理コストの最小化 [HEONAZ8].

#### (ii) 応答時間の最小化

### 4.2 放送網による特徴

分散問合せ処理の実現において、以下の点を検討する必要がある。

#### (i) 通信形態

##### (ii) スケジューリング決定手法

##### (iii) スケジューリング実行の管理系 (controller) の存在

(i)の通信形態としては、既存広域網の1対1通信と、共有媒体型のローカルネットワーク等の1対多通信とがある。我々は、後者を優先している。1対多通信を用いる手法としては、既に [TAKIMORI], [TOANN81], [KAMBERI] 等がある。[TOANN81]では制御情報だけが放送される。

(ii)のスケジューリングの決定手法としては、実行前に全てのスケジューリングを決定してしまふ静的な手法と、実行中にその都度、次に何をするか決定する動的決定法とがある。[TOANN81]では両者を混合した手法が取られている。

(iii)では、各DM間内の通信処理の実行の制御を、1つのDMが行う(集中制御)か、各DMが独自に判断する(分散制御)かの選択がある。[VINBERG]は、データフロー方式の分散制御を提案している。

我々の手法は、以下の特徴を有している。

#### (i) データ及び制御情報の放送 (ii) 動的決定 (iii) 完全分散制御

各DMから出された情報は、放送網で結合された他の全てのDMによって受信出来る。これによつて、1種の属性についての  $m$  個 ( $m \geq 2$ ) の異なったDM上の問合せ間の結合を行う時、1つの問合せを放送すれば、他の  $m-1$  個の問合せとの結合を  $m-1$  個のDM上で行える [TAKIMORI]。1回の通信で、 $m-1$  個の処理が出来る。又、各DMの状態情報 (e.g. 問合せのサイズ) を放送する存すれば、1回の通信によつて、他の全てのDMはこれを知ることが出来る。即ち、RDBS内の全てのDM<sub>1</sub>, ..., DM<sub>n</sub>は、RDBS全体についての同一状態情報を知ることが出来る。よつて、各DMは、自分の自ずるRDBS状態情報によつて、他とは独立に、かつ同一の決定を行える(完全分散制御)。又、スケジューリングの各ステップ毎に、各DMがその処理結果を放送することによつて、その時点での最適なスケジューリングを動的に決定出来る。これ等の特徴は、分散問合せ処理アルゴリズムを簡単なものとし、実現と運用とを容易なものとする。

### 4.3 全体問合せ

利用者から問合せを入力されたデータベース (DM) を、結果DM (RDM) と呼ぶ。結果DMとは、分散性が不可復元視野に基づいた問合せから、各DM内の問合せを参照する全体問合せを生成する。全体問合せは、QUELによつて、次の様に記述される。

$$\text{range } (r_1, R_1 : d_1) (r_2, R_2 : d_2) \dots (r_m, R_m : d_m);$$
$$\text{retrieve into } R (a_1 = \text{exp}_1, \dots, a_k = \text{exp}_k) \text{ where } \text{qual}_i \} \dots (5)$$

$R_i$  は、DM <sub>$i$</sub>  に存在する問合せであり、range 又は  $R_i$  に対する組変数  $r_i$  を定義している ( $i=1, \dots, m, m \geq 1$ )。Rは結果問合せ名で、 $a_1, \dots, a_k$  はその属性である。exp <sub>$i$</sub>  は、変数  $a_i, b_i$  ( $r_i$  は問合せ  $R_i$  の組変数で、 $b_i$  は  $R_i$  の属性) とする算術式である

( $l=1, \dots, n, n \geq 1$ ).  $qual$  は 結合述語  $r_i.b_i = r_j.b_j$  ( $r_i \neq r_j$ ). 制限述語  $r_i.b_i \neq r_j.b_j$  ( $b_i \neq b_j$ ) 又は  $r_i.b_i \theta v$  ( $\theta \in \{<, =, >, \neq\}$ ,  $v$  は定数) との又種の述語の積として表される論理式と仮定する.

例として, 図4.1に示す関係を考えよう.  $DM_1$  には部と部員情報が格納され,  $DM_2$  には従業員の会議予定(日時(date), 場所(place), テーマ(theme))が, 又  $DM_3$  には社内のプロジェクト構成情報が格納されている. この時, 「"A"氏と"DBBS"の会議に出席する従業員の部と所属プロジェクト名」を知りたいとすると, 全体問合せはQUELで次の様に表示する.

```
range (e, EHP:1) (de, DE:2) (s, SCDL:2)
      (p, PROJ:3) (s1, SCDL:2);
retrieve into R(d.dname, p.pname)
where s1.ename = "A" and de.eno = e.eno and
      e.eno = s.eno and s1.thema = s.thema and
      s1.thema = "DBBS" and s.eno = p.pno and
      e.eno = p.pno ; --- (6)
```

$DM_1$  EHP (eno, ename, tel)  
DE (dno, dname, eno)  
 $DM_2$  SCDL (eno, ename, date, place, thema)  
 $DM_3$  PROJ (pno, pname, eno, position)

図4.1 例)

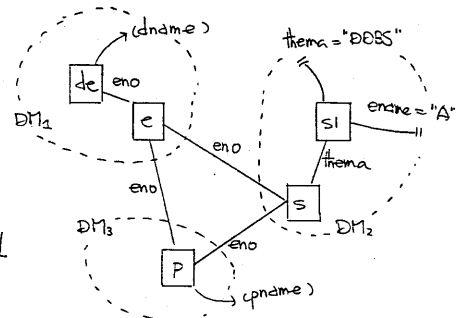


図4.2 (6)の問合せグラフ

(6)の問合せは, 図4.2の様により表現出来る. 節点は相変数を表し, 節点間の辺は結合式を表す. 節点に付加されている  $\rightarrow$  と  $\rightarrow$  は, 名々制限式と目標属性を表している.

#### 4.4 初期ローカル問合せ処理

全体問合せの中で, 各DMで独立に処理出来る部分と, 各DMで処理し, DM間の通信処理に必要な部分のみから成る関係を生成することを初期ローカル問合せ処理(ILQP)と呼ぶ. 図4.2の例では, 点線で囲まれた部分が, 各DMで処理され, 図4.3の様な結果となる. 即ち, DM内の結合, 制限式が処理され, GCS問合せの目標属性と共に, DM間の結合属性だけから成る結果関係が生成される.

ILQPの目的は, 後述するDM間との通信処理の為の通信量の縮小と共に, 異種DBSから条件を満足するデータを関係形式に導き出す目的を有している. CODASYLモデルの様に構造を有したデータモデルでは, 関係モデルの様な閉包性を有していない. 通信処理では中間結果の動的な生成, 消滅操作が必要となり, 閉包性が求められず, 従って, まず異種DBSから, 関係インタフェース LDP [AKIN80, 82b] によって結果を関係として導き出す必要がある. ILQPで生成された関係は, 各DMのDBに格納される.

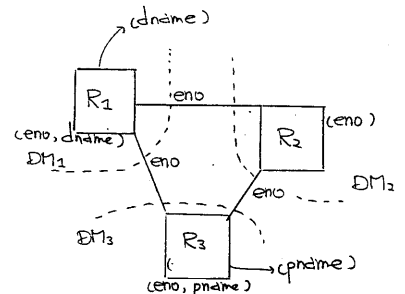


図4.3 ILQP後の問合せグラフ

#### 4.5 単純問合せのDM間通信処理

ILQP後の全併問合せは、DM間の結合だけから成る。まず、図4.4に示す様に、各DM間、 $m$ 個の問合せ  $R_1, \dots, R_n$  間に、単一結合属性  $a$  がある全併問合せを単純問合せと呼び、これについて考える。[TAKI81]では、動的決定と分散制御に基づき次のアルゴリズムを提案している。

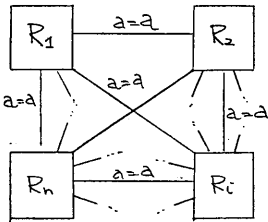


図4.4 単純問合せ

- (1)  $\Omega = \{R_1, \dots, R_n\}$  とする。  $\Omega' \leftarrow \Omega$ ;
- (2)  $R \in \Omega$  と  $|R[a]|$  が最大のものを見つけ、これを放送し、 $\Omega$  から  $R$  を除く。
- (3)  $\forall R' \in \Omega \quad R' \leftarrow R'[a=a](R[a])$ ;
- (4) (2), (3) を  $\Omega = \emptyset$  となるまで繰り返す。最後の問合せの放送時に、(3)では  $\forall R' \in (\Omega' - \Omega)$  に対して行う。
- (5) の終了した時点で、各DMは正の acknowledgement (ACK) に得られた結果のパフォーマンス情報(サイズカーボリタ)を乗せて放送する。全てのDMは、他の全てのDM内の問合せのパフォーマンス情報を知ることが出来るので、(2)に

ついて同一の決定を行うことが出来る。[TAKI81]では、更に[HEVNA80]の1対1通信に基づいたアルゴリズムとの比較から、大幅な通信コストの減少をもたらすことを示している。これを、アルゴリズム BA (broadcast algorithm) と呼ぶ。

ここでは、更に、[KAMBY81]にみる圧縮技法を加えることにし、より通信量を減少出来るアルゴリズムを示す。先々の方式と同様な方式が、上林[ KAMBY82]に、独立に存在している。アルゴリズム BBA (bit-map BBA) と呼ぶ。単に結合結果のカーボリタのみを返している。この為、 $m$ 個の問合せに対して、 $m-1$ 回かの問合せの放送と処理 Ack 転送が必要になる。これに対して、問合せの放送の ACK に、カーボリタに加えて、生成された問合せ  $R'[a]$  の  $R[a]$  に対するビットマップを乗せてしまう方式である。各DMは他のDMからの ACK を受信し、ビットマップの値を取っていき、最終結果を得ることが出来る。この方式にせよ、1回の問合せと、各DMでの処理と Ack 放送だけで処理させてしまう。これを、アルゴリズム BBA (bit-map BBA) と呼ぶ。

##### アルゴリズム BBA

- (1)  $\Sigma = \{R_1, \dots, R_n\}$  とする。
- (2)  $R \in \Sigma$  と  $|R[a]|$  が最大のものを見つけ、これを放送する。  $\Sigma \leftarrow \Sigma - \{R\}$ ;
- (3)  $\forall R' \in \Sigma$  と、  $R' \leftarrow R'[a=a](R[a])$ ;  
 $R'[a]$  の  $R[a]$  に対する bit-map  $BM'_a$  を生成する。
- (4)  $BM'_a$  の放送を行い、他の全てのDMからの bit-map  $BM_a$  を待つ。
- (5)  $BM_a$  を受信した後は  $BM'_a \leftarrow BM'_a \wedge BM_a$ ;
- (6) 全てのDMからの  $BM_a$  を受信した後は、 $BM'_a$  を基にして、最終結果  $R'[a]$  を生成する。

図4.5には、アルゴリズム BA の例を示してある。

次に通信処理コストについて考える。BAにおいて、問合せ  $R_1, R_2, \dots, R_n$  の順に放送されるとする。又、属性  $a$  の選択度を  $\alpha$  とし、各問合せの  $a$  の値は均一分散しているとする。この時、通信コスト  $C_{BA}$  と処理コスト  $P_{BA}$  は (7), (8) 式の最

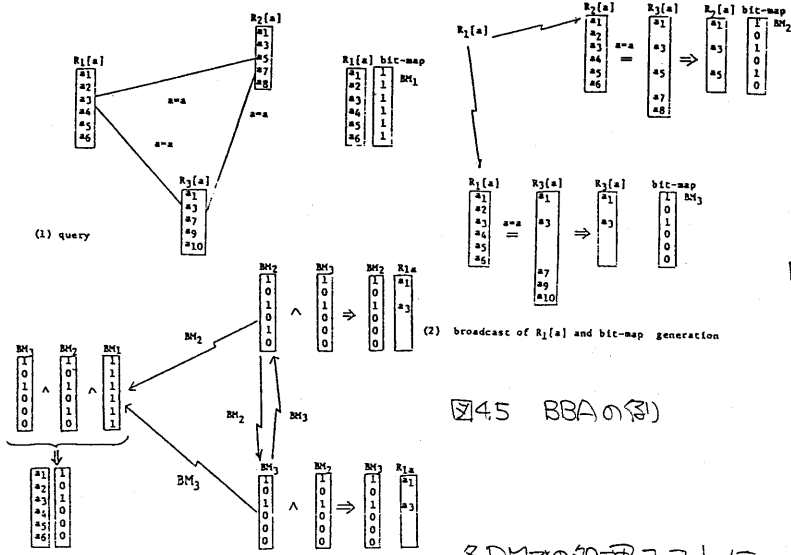


図4.5 BBAの例

の値になる。

$$C_{BA} = |R_1(a)| + \alpha |R_2(a)| + \dots + \alpha^{m-1} |R_m(a)|$$

$$= \sum_{i=1}^m \alpha^{i-1} |R_i(a)| \quad \dots (7)$$

$$P_{BA} = (\alpha^{m-1} |R_1(a)| + |R_2(a)| + \dots + |R_m(a)|) + \alpha(\alpha^{m-2} |R_1(a)| + |R_2(a)| + \dots + |R_m(a)|) + \dots + (\alpha |R_1(a)| + \alpha^2 |R_2(a)| + \dots + \alpha^{m-1} |R_m(a)|) \quad \dots (8)$$

(3) broadcast of bit-map response  
 各BM中の処理コストは、4.1の仮定に基づいている。一つの問合せ  $R_i(a)$  の返送は、 $m-1$  の結合処理  $J(R_1(a), R_2(a), \dots, J(R_{i-1}(a), R_i(a)))$  をもたすので、処理コストは、 $(m-1) \cdot |R_i(a)| + \sum_{j=2}^m |R_j(a)|$  となる。この結合処理で  $R_i(a)$  は、 $\alpha \cdot |R_i(a)|$  のサイズになる。又、(7)では ACK メッセージの通信コストを無視してある。

BBA の通信コスト  $C_{BBA}$ 、処理コスト  $P_{BBA}$  は、 $R_i(a)$  が放送された関係とすると (9)、(10) 式の値になる。

$$C_{BBA} = |R_1(a)| + (m-1) \cdot \text{card}(R_1(a)) \quad \dots (9)$$

$$P_{BBA} = \sum_{i=2}^m (|R_1(a)| + |R_i(a)|) = (m-1) |R_1(a)| + \sum_{i=2}^m |R_i(a)| \quad \dots (10)$$

(9) 式の右2項は、ビット2つの転送コストである。処理コストでは、ビット2つの演算が主要場内で出来るとして、無視してある。  $P_{BBA} < P_{BA}$  は明らかである。  $|R_1(a)|$  は最小であるので

$C_{BA} \geq |R_1(a)| (1 + \alpha + \dots + \alpha^{m-1})$   $C_{BBA} = |R_1(a)| (1 + \frac{m-1}{|\text{card}|})$ 。よって、 $|a| \geq (m-1) \frac{1-\alpha}{\alpha-\alpha^m}$  であれば、 $C_{BBA} \leq C_{BA}$  である。  $|a|$  は2ビット (16 bits)、 $m=2$  の時  $\alpha \geq 0.06$  であれば、BBAの劣が有利になる。

応答時間  $R_{BA}$ 、 $R_{BBA}$  は、次の値になる。

$$R_{BA} = C_{BA} + \sum_{i=1}^m \max_{j \in \{1, \dots, m\}} (\alpha^{i-1} [|R_j(a)| + |R_i(a)|]) \quad \dots (11)$$

$$R_{BBA} = C_{BBA} + \max_{i \in \{2, \dots, m\}} (|R_1(a)| + |R_i(a)|) \quad \dots (12)$$

図4.6には、放送網における、1対1通信を用いた [HEVNAR8]、或2の BA, BBA の各々に於ける通信コストをまとめてある。

(1), (2) 式で  $\alpha=0$ 、 $b=1$  としてある。データは、[HEVNAR8] を用いた。図から、放送通信機能を用いた BA, BBA の優位性が解かるとともに、データ圧縮技法 (ビット2つ) の利用効果も明らかである。

#### 4.6 一般問合せのDM間伝送処理

次に、異ったDMに存在する  $m$  個の関係  $R_1, \dots, R_m$  ( $m \geq 2$ ) が、 $m$  ( $\geq 1$ ) 個の結合属性  $a_1, \dots, a_m$  に、2等結合されている一般問合せについて考える。木型問合せ [BKNP 対] でない問合せに対しては、各属性の属性の組の値が転送される必要がある。以下に、一般問合せに対するアルゴリズム AGA を示す。ここで、 $\Omega_{R_i}$  を関係  $R_i$  の結合属性の集合とする。 $\Omega$  を全ての結合属性の集合とする。各結合属性  $a \in \Omega$  に対して、 $\mathcal{Q}_a$  を  $a$  を有する関係集合とする。 $\Sigma$  を  $R_1, \dots, R_m$  の集合とする。

##### アルゴリズム AGA

- (1)  $STA \leftarrow \emptyset; \forall R_i \in \Sigma \quad STA_i \leftarrow \emptyset; i \leftarrow 1;$
- (2)  $\forall a \in \Omega \quad \forall R \in \mathcal{Q}_a$   $R_i$  が最小のものを見つける。これを  $R_i[a_i]$  とする。 $a_i$  を  $i$  番目のスーパ属性と呼び、 $R_i$  を  $i$ -スレ関係と呼ぶ。 $R_i[a_i]$  の転送を行う。 $STA \leftarrow STA \cup \{a_i\}; \Omega \leftarrow \Omega - \{a_i\};$
- (3)  $a_i \in \Omega_{R_j}$  有る全ての DM で、 $\forall a \in STA_j \quad R_j \leftarrow R_j[a_i]; \quad STA_j \leftarrow STA_j \cup \{a_i\};$   
 結合処理  $R_j \leftarrow R_j[a_i = a_j](R_i[a_i]); \quad R_{ja} \leftarrow R_i[a_i];$   
 $\forall a \in STA_j$  に対して、 $R_j \cap R_j[a_i]$  であれば、 $R_j[a_i]$  の  $R_{ja}$  に対するビットマップ  $BM_{ja}$  を作り、同時に  $\forall b \in (\Omega_{R_j} - STA_j)$  に対して、 $R_j[b]$  のサイズ  $s_{jb}$  をもとめる。
- (4) これらの  $BM_{ja}$  と  $s_{jb} \in ACK$  を  $i$ -ジとして転送する。他の DM から ACK を待つ。
- (5) 他の DM から ACK を受信したならば、 $BM_{ja} \leftarrow BM_{ja} \wedge BM_{ka};$
- (6) 全ての DM から ACK を受信したならば、 $\forall a \in STA_j \quad BM_{ja}$  に基づいて  $R_j$  を再構成する。
- (7)  $i \leftarrow i + 1;$   $\Omega \neq \emptyset$  ならば (2) ~ (7) を繰り返す。

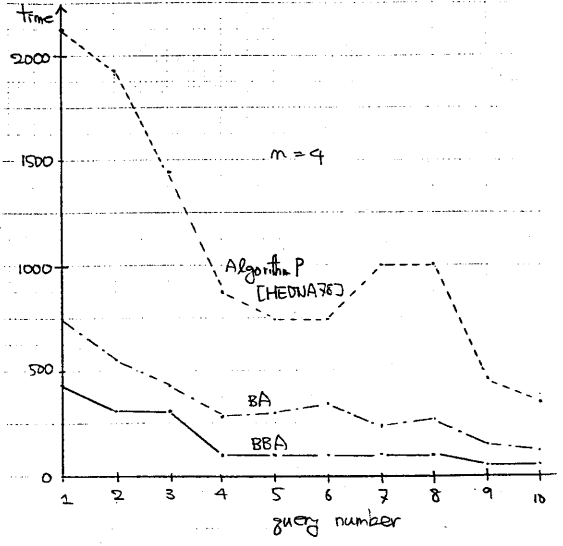


図4.6 全通信時間

最終的に、スーパ属性は、 $a_1, \dots, a_m$  の順に取られたとする。 $i$  番目の  $i$ -スレ関係  $R_i[a_i]$  のサイズを  $|R_i[a_i]|$  とする。この時、 $i$  番目の  $a_i$  を処理する為の通信コスト  $C_i$  は、(13)式の様になる。

$$C_i = |R_i[a_i]| + \sum_{\substack{j=1 \\ j \neq i}}^m \delta_{ij} \cdot \left[ \text{card}(R_i[a_i]) + \sum_{l=1}^{i-1} \gamma_{jl} \cdot \text{card}(R_j[a_l]) \right] \quad \dots (13)$$

ここで、 $\delta_{ij} = \begin{cases} 1 & \text{if } a_i \in \Omega_{R_j} \\ 0 & \text{otherwise} \end{cases}$        $\gamma_{jl} = \begin{cases} 1 & \text{if } a_l \in STA_j \\ 0 & \text{otherwise} \end{cases}$

(13)式のカ1項は  $R_i[a_i]$  の転送コスト、カ2項は  $BM_{ja}$ 、カ3項は  $BM_{ja}$  ( $a_l \in STA_j$ ) の転送コストである。4.5と同様に、ビットマップ処理コストを無視すると  $i$  番目の  $a_i$  の処理コスト  $P_i$  は、(14)式の様になる。

$$P_i = \sum_{j=1, j \neq i}^m \delta_{ij} \cdot [ |R_i[a_i]| + |R_j[a_i]| ] \quad \dots (14)$$

応答時間  $R_i$  は、 $R_i = C_i + \max_j ( \delta_{ij} [ |R_i[a_i]| + |R_j[a_i]| ] ) \quad \dots (15)$  とする。

全通信コスト、処理コスト、応答時間は、これらの和 ( $i=1, \dots, m$ ) とする。



- 一般問合せも、放送網を用いることにより、容易にかつ有効に処理出来る。処理コストと通信コストの重み付けの割合は、今後の課題である。

## 5. 分散更新処理

利用者から見た原子的な実行単位は、トランザクションと呼ばれる。更新処理とは、この原子性を守る必要がある。実行が原子的であるとは、正常終了した時のみ、この実行によるデータ変化を、他のトランザクションが参照出来、もし正常終了しなければ何の影響も与えないことである。トランザクション実行の原子性を保障する為の制御は、コミットメント制御と呼ばれる。即ち、トランザクションが、 $m$ 個のデータオブジェクト (e.g. 組)  $x_1, \dots, x_n$  ( $n \geq 2$ ) の更新を行おうとする時、 $x_1, \dots, x_n$  の全ての更新が行われたか、全く行われなかったかのどちらかになる様に更新実行は制御される。この為の手法としては、二相コミット法[2C]がある。これは、次の様である。

(i)  $x_1, \dots, x_n$  に対する更新データを結果DM (RDM) に生成する。 (ii)  $x_1, \dots, x_n$  の元データをモジュール (各 DM<sub>i</sub>,  $i=1, \dots, n$ ) に、更新データを含めた precommit  $pc_i$  ( $i=1, \dots, n$ ) 命令を送る。 (iii) 各 DM<sub>i</sub> は、precommit  $pc_i$  を受信したならば、更新データを自分の安全領域 (e.g. ログ) に待避する。格納出来たならば precommitted を RDM に送る。失敗すれば abort を送る。 (iv) RDM は、全ての DM<sub>i</sub> から precommitted を受信すれば、各 DM<sub>i</sub> に commit 命令を出す。全ての DM から precommit を受信出来なければ、各 DM<sub>i</sub> に abort を送る。 (v) 各 DM<sub>i</sub> は RDM から commit を受信すれば、安全領域への更新データにより、DB<sub>i</sub> 内の  $x_i$  の物理的な更新を行い、acknowledgement (ack) を RDM に送る。この様に、全部で  $4m$  回の通信が必要になる。

放送網を用いて、分散制御を行わせる時、二相コミット法を次の様に簡単化出来る。

- (i) 放送通信による分散問合せ処理において、 $x_1, \dots, x_n$  に対する更新データが生成される。生成された DM が、更新データを放送する [precommit に相当]。
- (ii) 各 DM<sub>i</sub> は、受信又は生成された更新データを安全領域に待避する。待避出来たならば、precommitted を放送する。
- (iii) 他の全ての DM から precommit を待つ。全てから precommit を受信したならば、DB<sub>i</sub> の  $x_i$  に対する更新を、安全領域のデータを用いて行う。
- (iv) 更新が出来たならば、Ack を放送する。

必要は通信回数は (i) で 1 回、(ii) で  $m$  回、(iii) で  $m$  回、都合  $2m + 1$  回となり、前者に対して通信遅延を半減できる。冗長コピーを保有する場合も、冗長性を意識することなく分散更新を行うことが出来る。

## 6. まとめと実現

本論文では、当初から現在開発中の分散型データベースにおいて実現を試みている分散問合せ処理と分散更新処理アルゴリズムについて論じた。Ethernet, 無線網の様に、通信媒体を共有した形態の通信網による物理的に損失される放送 (1対多) 通信機能を用いたアルゴリズムを考えた。放送網を用いることにより、各モジュール (DM) での分散制御と、自立的なスケジューリング決定が可能となり、データベースの分散問合せ、更新処理の実現と運用を簡単化することがで

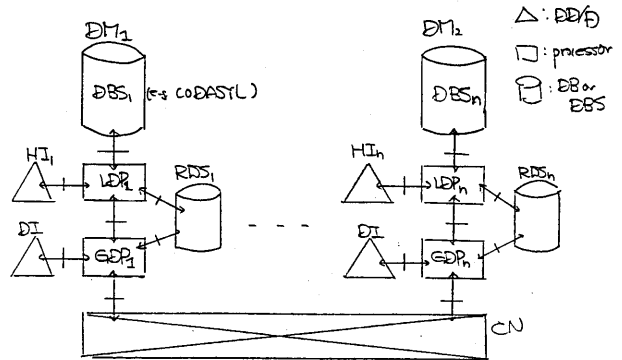
きる。今年度は、複数DMを、同一ホスト(M-170F)に実現し、JIPNETのインターネット又はAIM DB/DCによってDM間通信を行う予定である。既に、CODASTL DBS上の関係システム-システム LDP-V2 [TAKIZAWA, 82b, c] は AIM (M-170F) 及び ADBS (Acos-700) 上で実行している。

EDPは、本論文のDPに相当し、

RWSは、DBに相当する。

### 謝辞

JEDBSの実現に協力頂いているシステム社の鈴木信、当協会の菅塚実、塚本元子の各氏に感謝します。日産御指導頂いている当協会開発部の山本欣子、小関直美の両氏に感謝します。



LDP: local ES processor HI: heterogeneity information  
 GEP: global ES processor DI: distribution information  
 RWS: relational working storage CN: communication network

図6-1 JEDBSの構成

### 参考文献

[ADIB1780] Adiba, M., et al., "An Overview of the Polyphase Distributed Database Management System," Proc. of the IFIP'80, Oct. 1980, pp.475-479

[BERN781] Bernstein, P.A., et al., "Concurrency Control in DBDS," ACM Comp. Surv., Vol.13, No.2, June 1981.

[BERN782] Bernstein, P.A., et al., "Final Reducers of Relational Queries Using Multi-Attribute Sem.-Join," Proc. of IEEE Comp. Net. Dec. 1978

[GRAY78] Gray, J.N., "Notes on DB OS," in Operating Systems, Springer-Verlag, 1979

[HEVNA78] Hevner, A., et al., "Query Processing on a Distributed Database," Proc. of the 3rd Berkeley Workshop, Aug 1978

[HEVNA78b] Hevner, A., et al., "Optimization of Data Access in Distributed Systems," TR31, Purdue Univ., July 1978.

[KAMB781] Kamayoshi, Y., "データベースの分散処理の理論と実装," 電子情報学, Sep. 1981

[KAMB782] Kamayoshi, Y., "Theory of Databases," Comp. Sci. & Technology (Kijima, T. ed.), North-Holland ed OHM, 1982

[LEBI780] Lebitan, J. et al., "SIRIUS: A French Nationwide Project on DBDS," Proc. of the 6th VLDB, Oct. 1978

[STON785] Stonebraker, M. et al., "Design and Implementation of INGRES," ACM TODS, 1975

[SMIT781] Smith, J., et al., "Multibase - Integrating Heterogeneous Distributed Database Systems," ATIPS Conf. Proc., 1981

[ROTH780] Rothnie, J.B. et al., "Introduction to a System for DDB (SDO-1)," ACM TOBS, Vol.5, No.1, Mar. 1980

[TAKI786] Takizawa, M. et al., "Resource Integration and Data Sharing on Heterogeneous RDS," Proc. ICC'78, Sep. 1978

[TAKI785] Takizawa, M. et al., "The Four-Schema Concept as the Gross-Architecture of Distributed Database and Heterogeneity Problems," JIP, Vol.2, No.3, Nov. 1979

[TAKI780] Takizawa, M. et al., "Query Translation in Distributed Databases," Proc. of the IFIP'80, Oct. 1980.

[TAKI781] Takizawa, M., "分散型データベース-IIの通信処理," 電子情報学, June 1981.

[TAKI782a] Takizawa, M., "Distribution Problems in Distributed Databases - Integration and Query Decomposition," to appear in JIP, 1982

[TAKI782b] Takizawa, M. et al., "CODASTLとAIM-システムに対する関係システム-システム - LDP V1.5 - の設計と実装," to appear in IPSJ 論文誌, 1982

[TOAN781] Vine, D.H., "A Data Flow Solution for Implementing Distributed Queries," Proc. of the 5th Berkeley, Feb. 1981