

# 米国における分散データベースシステム研究開発の現状

増永良文  
(図書館情報大学)

筆者は米国カリフォルニア州サンホセ (San Jose) 市にあるIBMサンホセ研究所の客員研究員として、昨年4月より今年3月までの一年間、分散型関係データベース管理システム R\* (アールスターと発音) のプロジェクト員として、その研究開発に従事する機会に恵まれた。本報告はこの体験を通して見聞した、米国における分散データベースシステムの研究開発の現状を記したものである。

## 1. 緒言.

分散型データベース管理システム(以下簡単に分散データベースシステムと称す)は1970年以後半に入り、急激に関心が高まると共に、幾つかの研究開発例をみるようになった。その原因としては次の項目を挙げられる。

- 分散向きとされた関係データベース管理システムがシングルサイトで実動した。分散データベースシステム開発の技術的裏付けを与えた。
- ARPANET に代表されるような計算機ネットワーク技術が十分に固まり、分散データベースシステム開発の基盤を固めた。
- データを意図的に複数サイトに分散させ、システムの availability を向上せよとするコンピュータアーキテクチャ側からの要求が高まった。
- 分散データベースシステムを構築して、生産性を向上せよという企業側の要求が高まった。

現在迄、米国では約10システム程、分散データベースシステムが研究開発されている。中には商品化されているものもある。それらを大体年代順にリストアップすると次のとおりである。

- SDD-1 (CCA社, マサチューセッツ州)
- 分散型 INGRESS (カリフォルニア大学バークレー分校)
- System R\* (IBM サンホセ研究所, カリフォルニア州)
- ENCOMPASS (タンデム コンピュータ社, カリフォルニア州)
- DDTs (ハネウエル 共同計算機科学センター, ミネソタ州)
- プライム コンピュータ社のシステム (マサチューセッツ州)
- Multibase (CCA社, マサチューセッツ州)
- DDM (CCA社, マサチューセッツ州)

本稿では、次章でこれらのシステムの研究開発状況、特徴等を要約する。System R\* は報告者がそのプロジェクトの一員として研究開発に従事したこともあり、章を改め、若干詳しくその設計方針、アーキテクチャを述べる。

## 2. サーベイ

表-1に前出のシステムの特徴を概略する。システム開発において、設計インテグレーションは時間の経過と共に変わるのが普通である。なるべく表

に記載の事項が事実と異なり、データはシステムの研究開発の最終あるいは最新のレポートからピックアップするよう努めた。以下表-1を補うべくシステム毎に特徴等を記す。

● SDD-1<sup>(1)</sup>: 汎用分散型関係データベースシステムのプロトタイプであった。同時制御 (concurrency control), 分散型値向処理, 信頼性等の技術課題で貢献した。同時制御のためコンフリクト (conflict) グラフ解析のメカニズムとタイムスタンプに基づいた同期プロトコルを導入した。コンフリクトグラフ解析のロックিং (locking) よりもトランザクションの直列化可能性 (serializability) が優れている。分散型値向処理では準結合 (semi-join) 演算を導入した。SDD-1では値向処理の最適化は、それを行なうために必要なサイト間のデータ転送量 (バイト) と定義したので、アクセスプランニングの問題はトランザクションとデータハースの状態が互いに排他的でデータ転送量最小の準結合演算列をプログラムすることと等しい。トランザクション管理は2フェイズコミットの一様を用いている。このプロジェクトは1977年秋に終了した。

● 分散型 / NGRES<sup>(2)</sup>: 非分散型関係データベース管理システム / NGRESの分散版である。汎用システムを目指している。同時制御, 分散型値向処理, パフォーマンスの測定, 解析等を主要技術課題としている。分散型値向処理の最適化は応答時間最小と通信量最小である。値向処理には / NGRESで用いられた傾向分解 (query decomposition) の手法が拡張されている。トランザクション管理は2フェイズコミット方式, 重複データの制御のために主コピーサイトという概念を用いた2フェイズロック方式を採用している。現在パフォーマンスの評価のため種々のベンチマークテストを実施中である。分散型 / NGRESをシングルサイトで保つと、非分散型 / NGRESに比べて、ローカルデータベース処理で20%程度速くすると報告されている。今後非対称なワークロード (workload) がパフォーマンスの劣化に与える影響の測定, 種々の同時制御法の評価, 等々多くのベンチマークテストを計画している模様である。

● System R<sup>\*</sup>: サイトの自治権の確立を最大の目標として、システムの開発を行っている。非分散型関係データベース管理システムとして開発された System R (のちに SQL/DS として商品化) の分散版である。汎用システムを目指している。現在3台のIBMプロセッサ上で実動しており、インプレメンテーションは最終段階にきている。システムアーキテクチャの説明は次章に譲る。

● ENCOMPASS<sup>(3)</sup>: タンデム NonStop コンピュータシステムのデータベース部として開発された分散型データベース管理システムである。レポート生成レータ機能と関係データベース値向処理機能の融合を計っているが特徴がある。ENCOMPASSは ENSCRIBE, TMF, DDL, ENFORM, PATHWAY 等のサブシステムからなる。ENSCRIBEはファイルシステムであり、データへのアクセスや分散時のローテーション、スパーレンシー等を提供する。TMFはトランザクション管理システムである。DDLはデータディクショナリシステムである。ENFORMは上に述べたOAとDBの融合機能を提供している。これはコンパイラ/レポートライター (フロントエンド) と値向プロセッサ (バックエンド) の二つのプロセッサから

た。コンパイル/レポートライターは処理要求の発生したサイトで起動されるが、質問プロセサは質問処理のためのデータアクセスが最も重いサイトで起動され、ネットワーク上のデータ転送量を最小化しようとしている。ENFORMはユーザがその目的のために使用するユーザ言語の名前でもある。

● DDT<sup>(4)</sup>: DDTは Distributed Database Testbed Systemの略である。命名から判るように "testbed" として開発されている。プロトタイプは効率とシンボリック性を強調するのに対し、テストベッドはモジュラリティとフレキシビリティを強調する。種々の設計パラメータ(たとえば同時制御法)がパフォーマンスにどのような影響を及ぼすかを "量的" につかみとることを目的としている。6層のアーキテクチャをとり、ローカルデータベースはCO DASYL社のZDS/II, その上に関連モデルのローカルとグローバルの表現(シンボリック)スキーマのり, その上は ECR (Entity-Category-Relationship) モデルで記述されたグローバル概念(セマンティック)スキーマがのる。このアクセス言語は GORDAS (Graph Oriented Data Selection Language) と呼ばれており、現在この段階までしかインプリメントされているようである。

● プライムコンピュータ社の分散データベースシステム<sup>(5)</sup>: プライム社はCO DASYLタイプの非分散型データベース管理システムを所有しており、それを Primenet で結合、分散データベースシステムを開発しようとするものである。特にマルチバージョン(multi-version)分散読み出し専用トランザクションの管理技術の確立を課題とした。現在、インプリメンテーションはほぼ終了したと報告されている。分散質問処理の最適化等はまだ研究されているようである。

● Multibase<sup>(6)</sup>: 異種(heterogeneous)分散データベースシステムである。一般にデータモデルの異なるローカルスキーマを共通データモデルにスキーマ変換し、それらを統合するというアーキテクチャをとっている。共通モデルに Functional Data Model (FDM) を採用している。統合スキーマ上にグローバルスキーマを、質問言語 DAPLEX のビュー定義機能を用いて定義し、ユーザはそれとインタラクトする。したがって INGRES が導入した質問変形(query modification)の手法が質問処理に適用可能で、実際そうしている。1982年初めに一次のインプリメンテーションが終了したと報告されている。

● DDM<sup>(7)</sup>: データベース操作言語 DAPLEX をプログラミング言語 Ada に埋込んだ ADAPLEX をサポートできる汎用の分散型データベース管理システム DDM (Distributed Database Manager) として開発されている。DDM の機能的には SDD-1 や System R\* を基本にして設計されたが、ADAPLEX のサポート、重複データの取扱、読み出し専用トランザクションの最適化、障害時回復のアルゴリズム等で特色もあるという。システムを Ada でインプリメントしており、システムのポータビリティが向上している。現在進行中のプロジェクトである。

システムは開発しているが、例えば分散型質問処理法、同時制御法等の理論を研究している研究教育機関は数多くあり、例えばスタンフォード大学、ベル研究所

所, ミシガン大学, メリーランド大学, ハーバード大学, MIT等を挙げられ。上に報告したシステムの研究開発状況を大別すれば次の二つに分けられる。

- 汎用分散型データベース管理システムの開発: CCA社のSDD1 → Multibase → DDMの流れ, IBMのSystem R → System R\*の流れ, UC BerkeleyのINGRES → 分散型INGRESの流れ, 筆で代表とする。
- OA (Office Automation) システムに組込型分散型データベース管理システムの開発: このタイプの管理システムのアーキテクチャが十分研究されているとは思えないが, タンデム社のENCOMPASSはOAのアプリケーションを主たる目的の一つとしている。シリコンバレーの小規模コンピュータ会社はこのあたりの市場を十分意識しているだろう。

分散データベースシステムが非分散のデータベースシステムと異なった取扱をせねばならぬ技術的問題は少なくとも次の二点を含めよう。

- 分散型箇所処理法
- 分散型のカタログ (DD/D) 管理法, 同時制御法, 障害時回復法

現在このシステム開発を通して, 決定的な解答が出ている訳ではない。典型的なシステムアーキテクチャを確立するために, 単に実験室レベルでのパフォーマンス向上のみならず, そのプロトタイプを実際企業業務に適用して, システムアーキテクチャの改善と評価を行わなくてはゆく必要がある。

区分 システム名	年代 <sup>*1</sup>	マシン	通信網	プログラミング言語	ユーザ言語	状況 <sup>*2</sup>	タイプ	インベリメンタ
SDD-1	1976~78(テ) 1979(経)	DEC-10と DEC-20 合計4台	Arpanet	BCPL 25,000行	Datalanguage (非線形言語)	imp/inact	institutional	CCA マサチューセッツ州
分散型 INGRES	1976~79(テ) 1979~(イ)	DEC VAX 11/80 DEC VAX 11/750 各5台	ETHERNET (3Mbit)	C 非分散/INGRES の1.7倍	QUEL EQUEL	imp/act	↑	カリフォルニア大学 バークレー分校
System R*	1979~(イ) 1984(三)	IBM 4341 Model x3台	SNA Virtual Circuit	PLS/III	SQL PLI/SQL	↑	↑	IBMサンホセ研究所 カリフォルニア州
ENCOMPASS	1981 商品発表	Tandem Non Stop J=2台	long-haul net(56KB)	-	ENFORM	↑	Commercial	タンデムコンピュータ社 カリフォルニア州
DDTS	1979~(イ)	Honeywell Level 6 mini- comp x3台	Bus 結合 (coax cable)	C 20,000行	GORDAS (on ECR)	↑	institutional	トネウエール CCSC ミネソタ州
特別居るし	~1982~ (イ)	7.5インチの2台 5.25インチの2台	Primernet	-	-	↑	↑	フェイスコンピュータ社 マサチューセッツ州
Multibase	1980~84 (三)	-	-	Ada	DAPLEX (on FDM)	↑	↑	CCA マサチューセッツ州
DDM	1982~84 (三)	DEC VAX マシン (on VMS)	-	Ada	DAPLEX ADAPLEX	des/act	↑	↑

表-1. 米国における分散データベースシステム概観

\*-1. (テ)は開発中, (経)はプロジェクト終了, (イ)はインベリメント, (三)はプロジェクト終了予定。  
\*-2 imp はインベリメント完了かインベリメント中, des は開発中, inact はプロジェクト終了, act は活動中。

### 3. システム R\* (7~17)

#### 3.1 システム R\* の設計指針

システム R\* が研究開発の技術的ゴールとしたのは次の三美である (11, 15)

- サイトの自治権の確立
- 使い易さ
- 高いパフォーマンス

サイトの自治権の確立はシステムのモジュラリティを上げ、システムの増強をたやすくし、誤りに対する回復力を強める。しかし、このためには一切の集中制御の概念を捨て去った、完全に分散型のアーキテクチャを開発しなければならぬ。使い易さは新しく R\* に加入したユーザに、加入のための特別なロード (例えばカタログの書換、アプリケーションプログラムの変更) を一切やり取りするもので、他にシングルサイトイメージの保持、高レベル非手続き的言語 SQL のサポートを意味する。高いパフォーマンスは勿論サイトの自治権の確立や使い易さといった課題とは相矛盾する設計ゴールであるが、種々の最適化技法をとり入れて、それと達成しようとするものである。

#### 3.2 システム R\* の基本アーキテクチャ

システム R\* は均質 (homogeneous) な関係データベース管理システムのあり方である。二に均質とは各システムが (version は異なっても) 同じ外部インターフェイスを持ち、共通の通信プロトコルに従って動作することという。図-1に R\* のシステムの基本構造を示す。R\* は大別して、• データベースマネージャ、• トランザクションマネージャ、• 通信マネージャからなる。データベースマネージャは単に RDS (Research Data System, 分散データベース均相マネージャ) と RSS (Research Storage System, ローカルデータベースマネージャ)

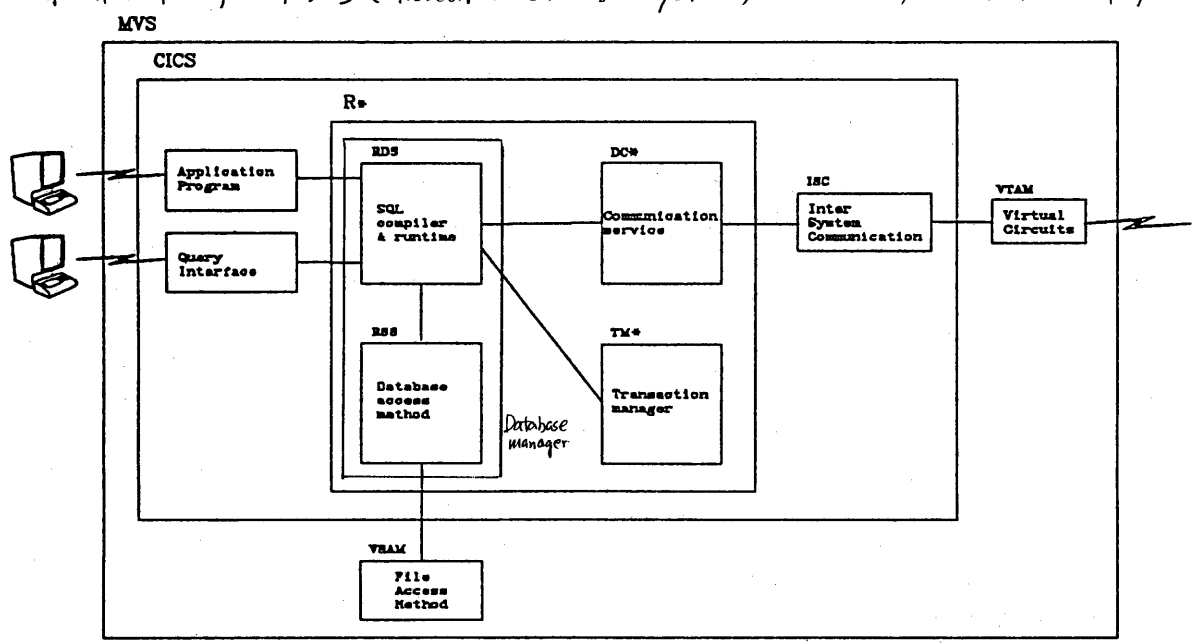


Figure 1: The R\* System Structure (文献 16 頁)

からなる。RDPの仕事はネーム解説, 権限チェック, アクセスパス選択, ビュー生成, コード生成である。RSPはローカルトランザクション管理, テーブル情報管理, 同時制御, 障害時回復, 物理的ディスク管理の仕事をする。トランザクションマネージャの仕事は分散された(複数サイトの)トランザクションがそのトランザクションのプロセスが実行された全てのサイトでコミット(commit)されたか、棄却(abort)されたかを保証することである。2相コミットプロトコルの一変形である入れ子型のものが使われている。通信マネージャはバッチルサーキットを保持, 入力要求をデータベース, トランザクション各マネージャに適切にルーティングする。

### 3.3 システムR\*のカatalog管理体系<sup>(9)</sup>

分散システムではデータ共有化実現のためのオブジェクト(関係テーブル, ビュー等)の命名法が必要である。システムR\*ではサイトの自治権確立が大きな技術的課題に存っているので, 次のような命名法をとった。オブジェクトは生成者名@生成者サイト名.オブジェクト名@オブジェクト誕生サイト名の4項からなるSWN(System Wide Name)をもつ。したがって, YOSHI@SJ. SAL@SJ と YOSHI@SJ. SAL@NY とはたゞと同一のSALテーブルでも異なるオブジェクトとして認識される。デフォルトはYOSHIがSJでlogonしSALという関係テーブルを作れば, 自動的にYOSHI@SJ. SAL@SJ というSWNを生じる。オブジェクトが格納されているサイトを格納サイトという。格納サイトがSWNの構成要素に入っているのは, 元来オブジェクトは移動(migrate)するものであり, それをSWNに組入れるとオブジェクトのローテーション・トランスパーレンシーが損なわれるからである。オブジェクトの格納サイトはその生成時, SQL文のIN句(例として DEFINE TABLE SAL@NY IN SJ;)で指定できる。オブジェクトの移動はMIGRATE文で指定できる。シノニム(別名)はシステムRと同様定義できる。

分散データベースシステムのどのサイトでどういったオブジェクトが誕生, 格納されているか, またどのようなアクセスパスを有しているのか, 等の情報を格納する。ため, システムR\*は複数サイトの目的のための関係テーブルを用意しており, システムカATALOGと呼んでいる。別の言葉ではDD/D(Database Dictionary/Directory)である。システムR\*のサイトの自治権確立という観点から, システムカATALOGは完全に分散している。つまりR\*にはローカルカATALOG, グローバルカATALOGといった区別は一切ない。システムカATALOGのうちSYSCATALOGテーブルは, オブジェクトの生成, 格納等の情報を保持している。オブジェクトがあるサイトのSYSCATALOGテーブルにエントリされるのは ●それがそのサイトで誕生したか, ●そのサイトに格納されているか, あるいは●そのサイトで使われ, キャッシュ情報として保持される, かのいづれかのことである。オブジェクトの誕生サイトのSYSCATALOGテーブルにはそのオブジェクトが消滅したの限り, それを記録する。オブジェクトの統計量等はその格納サイトのSYSCATALOGテーブルに記入される。他に, シノニム, ビュー, インデックス, 権限付与, 等の目的のためにテーブルが定義されている。このカATALOG体系は次のことを実現できる。

- 新規に網に加入してもシステムカタログ、応用プログラムを変更しない。  
 • オブジェクトの移動(migration)に応用プログラムは無関係である。  
 • 高々2回のリモートテーブル参照で、オブジェクトの格納サイトがわかる。  
 (つまり、質問処理を受付けたサイトはまず自身のSYSCATALOGテーブル  
 をみる。 各サイトが誕生サイト存続に格納サイトが判る。 もし格  
 納サイト存続もれも判る。 もしキャッシュされていて、その情報が  
 out-of-dateでなれば、直ちに格納サイトが判る。 もしout-of-date  
 なら、指示されたサイトに要求してもそのオブジェクトは無いのだから、  
 オブジェクトのSWNからその誕生サイトを識別し、誕生サイトに現在の  
 格納サイトを問合せる。)
- キャッシュ情報を格納するためパフォーマンスが向上する。(キャッシュさ  
 れたオブジェクトが他のトランザクションにより変更をうけた場合に、  
 out-of-date化されたが、カタログにバージョン数を保持することにより、  
 発見、訂正される。 中3.4節参照。)

### 3.4 システムR\*の分散型質問処理体系<sup>(8, 12, 13, 14)</sup>

PLI/SQLプログラムは一般にSQL文を含み、プレコンパイル時にこのSQL文  
 の表わす質問が分散コンパイルされる。 システムR\*では、そのプログラム  
 が提示したサイト(主サイトと呼ぶ)をコーディネーター(coordinator)と  
 し、そのプログラムの質問処理に関係するサイトを従サイトとし、分散型コンパ  
 イルを行なうという方式を採用している。 この方式はシステムR\*が設計目標  
 にあがっているサイトの自治性の確立にもあつち、システムカタログの構造とも適合  
 する。 分散型コンパイルの結果、従サイトに低レベルのアクセスモジュールが格  
 納される。 これはランタイム時にアクセスインタプリタにより実行される。

主サイトでの質問コンパイルは、1. SQL文の構文解析とネーム解  
 読、2. パーズ木を基にカタログ参照、オーソリゼーション検査、ビュー合成、  
 3. グローバル(アクセス)プランニング(この結果、1.の時点で送られた従サイトへ  
 のSQL文と3.で送られるグローバルプランを従サイトに送る)、4. アクセス  
 ジェネレーションをして主サイトのサブセッションに格納、というステップをと  
 る。 従サイトでは送られてきたSQL文とグローバルプランからローカルア  
 クセスプランを生成する。(そのSQL文に保護ビューの存在するとこのグローバル  
 プランも生成する。) 主サイトがグローバルプランを作成するので、それが  
 使用したカタログデータがout-of-dateであれば、従サイトで与えられたバージョン  
 数を検査して、主サイトに報告、主サイトはグローバルプランニングをやり直す。

グローバルプランは次のコスト関数 $T_c$ が最小となるように行なう。

$$T_c = W_1 \times \text{送信されたメッセージ数} + W_2 \times \text{送信されたバイト数} \\ + W_3 \times \text{フェッチされたディスクページ数} + W_4 \times \text{要求されたデータ数}$$

関係テーブルの(自然)結合をとる場合のコストを考へてみると、それは●結合  
 順序、●結合法(システムR\*では入れ子型かマージ型)、●各テーブルのアク  
 セスパス、●結合場所(サイト)、●テーブルの送信法(全体か一部か)、等に依  
 る。 例へば、二つのテーブルのマルチサイト結合を考へてみると、R\*で

は次の場合を想定している。(入れ子型の正式には入れ子ルーブリック型)

- アウターテーブルをインナーサイトに送る。インナーサイトではアウターテーブルの全ての到着順に、入れ子型かマージ型で結合を開始する。
- インナーテーブルをアウターサイトに送る。アウターサイトではインナーテーブル全体を格納した後、マージ型結合を行う。
- インナーテーブルをアウターサイトに送る。アウターサイトでは必要のインナーテーブル部分が到着次第、入れ子型かマージ型で結合を開始する。
- アウターテーブルもインナーテーブルもオミサイトに送りそこで結合する。
- オミサイトにアウターテーブルを格納し、必要のインナーテーブル部分がオミサイトに到着次第、結合を行う。

グローバルプランの生成はアクセスパス  
 選択探索木を breadth-first に総あたり  
 でサーチして決める。グローバルプラン  
 はミニプランの系列として表わされる。  
 図-2はサイトCにあるテーブル1をサイ  
 トAに送り、それをアウターテーブルとし  
 て、入れ子型でテーブル3と結合し、次い  
 でテーブル2がサイトBからサイトAに送  
 られ、それが格納されたあと、テーブル1  
 と3の結合結果とマージ型で結合するプ  
 ランを表わしている。图中、ボックスが一  
 つのミニプランに対応する。VT (Virtual  
 table) はリモートサイトからテーブルをとり  
 こんでおく placeholder である。グ  
 ロバルプラン中、連続したミニプランの最大サブ  
 セットをサブセクシオンと呼び、VTは他のサイ  
 トのサブセクシオンへの call を表わしている。  
 サブセクシオン全体がアクセスモジュールである。<sup>(13)</sup>

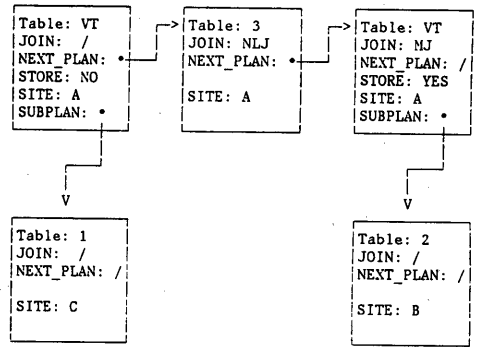


Figure 2. A Global Access Plan (文献13引)

システム R\* のアクセスモジュールは ● アクセスパスの欠落, ● データオブジェクト  
 の変化, ● アクセス権の撤回, に依り無効化 (invalidation) される。依存性  
 の損壊は一般にはアクセスモジュール全体に及ぶものでなく、その幾つかのサブ  
 セクシオンの範囲に収まるとはが異なることに注目する。無効化に自動的に対応  
 するため、コンパイル時にアクセスモジュールの依存性がカタログ (たゞしは、  
 usage カタログ) に記録される。

アクセスモジュールを再有効化 (re-validate) するため、再コンパイルが行な  
 われる。システム R\* ではサブセクシオン単位で再コンパイルをするグローバル再  
 コンパイルと、サブセクシオン単位で再コンパイルをするローカル再コンパイル  
 の二つの方式を用意している。入れ子型結合のインナーテーブルのインデック  
 スが欠落していた時にはグローバル再コンパイルを、マージ型で同様の欠落があ  
 った時とアウターテーブルのインデックスが欠落していた時はローカル再コンパ  
 イルが良い。アクセス権の撤回による無効化はローカル再コンパイルが良い。  
 ローカル再コンパイルが良いとしている時は、その結果得られたアクセスモジ  
 ュールがほぼ最適であるからである。<sup>(12)</sup>



### 3.5 システム R\* のトランザクション管理体系 <sup>(7, 11, 16, 17, 18)</sup>

Ad-hocなSQLターミナルユーザではSQL文が、PLI/SQLプログラムではBEGIN TRANSACTIONで始まりEND TRANSACTIONに終了一連のステートメントが一つのトランザクションを成すのはSystem Rと同じである。

システム R\* のトランザクションが、システム R のトランザクションと異なる点はそのマルチサイトトランザクションである点である。トランザクションはコミット (commit) されるかアボート (abort) される。マルチサイトトランザクションではそのトランザクションがデータベースアクセスと更新で関与した全てのサイトが一律に (uniform) コミットするかアボートしてほしい。そのため R\* は入札型 2-相コミットプロトコルを使用している。トランザクションのプロセス木を作る。根プロセスはそのトランザクションに対応し、根プロセスとオニレベルのプロセスは R\* のバーチャルサーキットで結合している。以下のレベルも同様である。トランザクションのコミットかアボートかを決定するために、コーディネータ (coordinator, フォロワーのトランザクションのサイト) はオニレベルに "prepare to commit" 要求を出す。(この結果サブセッションがアクセスモジュールに格納され実行される。) そして "vote yes" か "vote no" を送る。オニレベルがあれば、オニレベルから同様に入札型に付長型でオニレベルに "prepare to commit" を出す。これが "collecting" 相 (オニ相) である。オニ相でコーディネータは投票の結果をみてコミットかアボートかの指令を出す。プロセス木のノード数を  $N$  とすると、メッセージの総数は  $4(N-1)$  とある。R\* では読み出し専用の部分木に対してはオニ相の会話を省略し、最適化を計っている。<sup>(16)</sup>

トランザクションリカバリにはデータグラムを使う。障害のおきたサイトではバーチャルサーキットでの会話は行わないので、それが出来上がるまでデータグラムを使う。

トランザクションの同時制衡はシステム R と同様、2-相ロック方式を使っている。<sup>(18)</sup> トランザクションには R\* 全体で唯一識別可能なトランザクション番号 (サイト名とローカル時刻の組) をつけている。この結果、同一サイトの異なるプロセスでも同一のトランザクション番号を持っていれば、データオブジェクトのロックを共用させてやることができる。

サイトの自治権確立というシステム R\* のゴールに抵触しないデッドロックの検出と解消のアルゴリズムが開発されている。検出アルゴリズムでもトランザクション番号が必要である。デッドロックが起った場合、犠牲となるのはトランザクション番号が最も大ま (youngest) トランザクションである。このグローバルなデッドロックの検出と解消アルゴリズムの経路にインポートされている。<sup>(17)</sup>

## 4. 結言

システム R\* のプロジェクトは 10 人程のプロジェクトである。現在、プロジェクトリーダーは Bob Yost, メンバーは Bruce Lindsay, C. Mohan, Guy Lohman, George Lapis, Paul Wilms, Laura Haas, Elisa Bertino である。Pat Selinger は本年一月にプロジェクトを引かれた。本稿では、

ビューレポートについての記述を省略した。別の機会に譲る。テストプログラムもシステム R (及び SQL/DS) の遺産を卒直して、膨大な数ある。システム R\* がシステム R と同様、商品化されるかどうかの決定は未だなされてい  
ない。

**【謝辞】** 著者は IBM サンホセ研究行でシステム R\* の一員として、その研究開発に携わった好機を与えて下さった日本 IBM (株) に心から感謝の意を表す。著者がまた、気持ちよく雇員研究員として一年間を送らせていただいた R\* プロジェクトのメンバーの皆様、そして IBM サンホセ研究行の皆様にも心から感謝の意を表す。

## 【文献】

- (1) J. B. Rothnie 他, Introduction to a System for Distributed Databases (SDD-1), ACM TODS 5-1 ('80)
- (2) M. Stonebraker 他, Performance Analysis of Distributed Database Systems, IE<sup>3</sup>, Database Engineering 5-4 ('82)
- (3) J. Nauman, ENCOMPASS: Evolution of a Distributed Database/Transaction System, *ibid.*
- (4) S. K. Rahimi, A Structural View of Honeywell's Distributed Database Testbed System DDTs, *ibid.*
- (5) D. DuBourdieu, Survey of Current Research at Prime Computer Inc. in Distributed Database Management Systems, *ibid.*
- (6) A. Chan, Distributed Database Management Research at Computer Corporation of America, *ibid.*
- (7) B. Lindsay 他, Notes on Distributed Databases, IBM Res. Rep. RJ 2571 ('79)
- (8) P. G. Selinger 他, Access Path Selection in Distributed Database Management Systems, IBM Res. Rep. RJ 2883 ('80)
- (9) B. Lindsay, Object Naming and Catalog Management for a Distributed Database Manager, IBM Res. Rep. RJ 2914 ('80)
- (10) B. Lindsay 他, Site Autonomy Issues in R\*: A Distributed Database Management System, IBM Res. Rep. RJ 2927 ('80)
- (11) R. Williams 他, R\*: An Overview of the Architecture, IBM Res. Rep. RJ 3325 ('81)
- (12) P. Ng, Distributed Compilation and Recompilation of Database Queries, IBM Res. Rep. RJ 3375 ('82)
- (13) D. Daniels, Query Compilation in a Distributed Database System, IBM Res. Rep. RJ 3423 ('82)
- (14) D. Daniels 他, An Introduction to Distributed Query Compilation in R\*, IBM Res. Rep. RJ 3497 ('82)
- (15) L. Haas 他, R\*: A Research Project on Distributed Relational DBMS, IE<sup>3</sup>, Database Engineering 5-4 ('82)
- (16) B. Lindsay 他, Computation and Communication in R\*: A Distributed Database Manager, IBM Res. Rep. RJ 3740 ('83)
- (17) R. Obermark, Distributed Deadlock Detection Algorithm, ACM TODS 7-2 ('82)
- (18) P. A. Bernstein 他, Concurrency Control in Distributed Database Systems, Computing Surveys, 13-2 ('81)