

ADABASの分散データベース機能

石井義興 嶋田正裕
㈱ソフトウェア・エージェンシー

1. はじめに

ADABASに付加する分散データベース機能について述べる。ADABASの分散DB機能の特徴は、すでに分散処理されているものを統合するというよりはむしろ、現在集中処理されているものを分割/分散化し、メリットを得ようという所にある。この際、ユーザ・アプリケーション・プログラムはこの分割/分散化に際してまったく変更する必要はない。すなわちデータの分布が変化しても利用者が意識しなくてもすむということ、プログラムの組み直しをしないですむということである。このようなフレキシビリティをDBMSに持たせることは必ずしも容易ではないが、今後の発展を考えれば必須の機能である。

分割/分散化しているマシン(ノード)を接続する手段は単に通信回線によるのみではなく、チャンネル結合でも良く、また仮想マシン間でも良い。各ノードを接続するためのソフトウェアも開発済みである。各ノードとなり得るマシンはIBM370、308X、43XX、日立、富士通のMシリーズで、OSは、MVS、MVS/XA、VS1、DOS/VSE、VOS1、VOS1/ES、VOS2、VOS3、およびOSIV/F4である。異なるメーカーのハードを利用し、異なるOSを用いても良い。ADABASの分散データベース機能、ノード間接続ソフトウェアの総称はNET-WORKという名で呼ばれており、NET-WORKは多くのコンポーネントから構成されており、各コンポーネントには、それぞれ名前が付いている。

2. 集中処理の問題点

コンピュータの利用はバッチ処理から始まった。この事は、コンピュータの歴史を振り返ってみればすぐわかる。昔はコンピュータの価格が高く、したがって多くのコンピュータを設置することなど不可能であった。そのためデータを集め、バッチ化して処理することとなった。集中処理が当然のように思われ、しかも、それが最良の方法であるかのごとく

考えられていた時代さえあった。しかし企業体の中での活動は、そもそも人手によって分散処理されており、全社の処理が1ヵ所で集中的に処理されているわけではない。

人手によって分散処理されている作業をコンピュータライズするわけであるから、分散処理していた時の問題点を集中化によってカバーする事が当然要求され、システム・デザイン時にそれを加味するように考えられた。しかしその際、分散処理の利点も同時に失われるケースが多かった。現在は集中処理による問題点をむしろ分析し、分散処理されるべきものは、分散処理することによって、集中と分散のバランスを取り、最良のシステムとなるように考え直す時期にきている。

集中処理による問題点は次の通り。

◇バッチ処理による集中化はタイムラグが大きくなる。すなわち、ある部分では準備完了の状態となっても、どこかの部分のデータ準備がととのわないと、それによって処理が待たされる。

◇バッチ処理による集中化では、現場でのデータ・ハンドリング作業と中央での作業と2度手間になる。

この問題を解決するためには、バッチ処理をやめ、オンライン方式に切り換えれば一応かたずく。しかしオンライン方式による集中化は一般に非常にコンピュータ資源をくいコスト高となる。

現時点は、このオンラインによる集中処理方式から、一部分散処理へもどすことによって、コストを削減し、スピード・アップをはかろうと考え始めた段階といえる。

3. 中型複数コンピュータによる システムの有効性

「コンピュータのコスト・パフォーマンスは大型機になればなる程良い」という神話はIBM4300シリーズの出現と共に崩壊し

た。IBMを初めとする各コンピュータ・メーカーは常に新しいコンピュータ・シリーズを世に出し、しかも大型機を利用すればする程コスト・パフォーマンスは良かった。しかし表1に示す通り、4300シリーズ(中型機)のコスト・パフォーマンスはIBMの超大型コンピュータより安い。このことは中型コンピュータを複数台組合せることによって大型コンピュータの代りに用いることがコスト的にも見合うことを示している。しかもマシン・ダウン対策を含めて考えると、複数台の組合せによる方式の方がはるかにリライアブルである。

コストがかえって安くなる点は「大量生産によるコンピュータ(小型機)の方が大量生産でないコンピュータ(超大型機)より安い」ということから考えれば当然のことと見ることができ、今後もこの傾向は続くと考えられる。これによりコンピュータを複数台連結し、大型コンピュータの代りに使用するという考え方が成立する。

コンピュータの処理能力がもし限界に達したら今までのように単にマシンをグレード・アップするのではなく、「小型のCPUをHOSTマシンに付加することにより処理量の増大に対応しよう」という考え方もなりたつのである。

このことは分割/分散処理がコスト面で優位に立つことを示している。特にオンライン・リアルタイム・システムを構成する場合はバックアップ・マシンを必要とする。バックアップ・マシンの能力はメイン・マシンとほぼ同程度の能力を持つ必要がある。メイン・マシン1台のコストで中規模マシンを複数台買入すれば、バックアップ・マシンを買入する必要がない。(図1参照)なぜなら3台のマシンが同時にマシンダウンを起すことは考えられないからである。

ただしこのシステムが正しく稼動するためには、3ヵ所に同一ユーザ・プログラムを入れることになるし、1つのDBを3台のマシンで同時にUPDATEできなければならない。つまりDBMS側にこれを可能にする機能が必要となる。

	4331-2 1M	4341-2 4M	3033N 8M	3081 16M
MIPS	4.0	1.2	3.7	10.4
価格(千ドル)	150	416	1,695	3,720
3081 1台の値段で何台買えるか	25台	9台	2台	1台
合計メモリ・サイズ	25M	36M	16M	16M
合計MIPS	10	10.8	7.4	10.4

表1

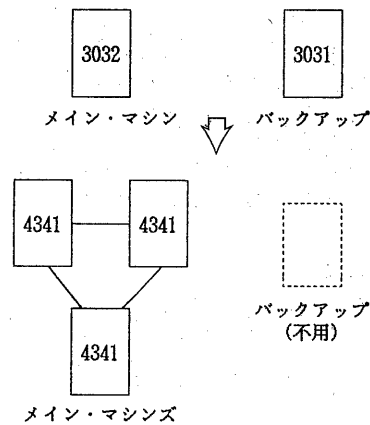


図1

4. 分散処理、分散DBの実例

今後の分散処理環境を予想し、(補)ソフトウェア・エージェンシーでは図2で示すようなマシン構成を作り、NET-WORKの開発テストを行なっている。わざわざ異なるメーカーのマシンを導入し、異なるOSを入れている。現時点ではVAX-11とIBMはまだ結合されていない。

この構成の特徴はデータベースのあるノードとそれをアクセスするプログラムがランするノードが異なってもかまわず、それぞれがユーザ・ニーズに対応して移動しても良い事を前提に作られたものである。

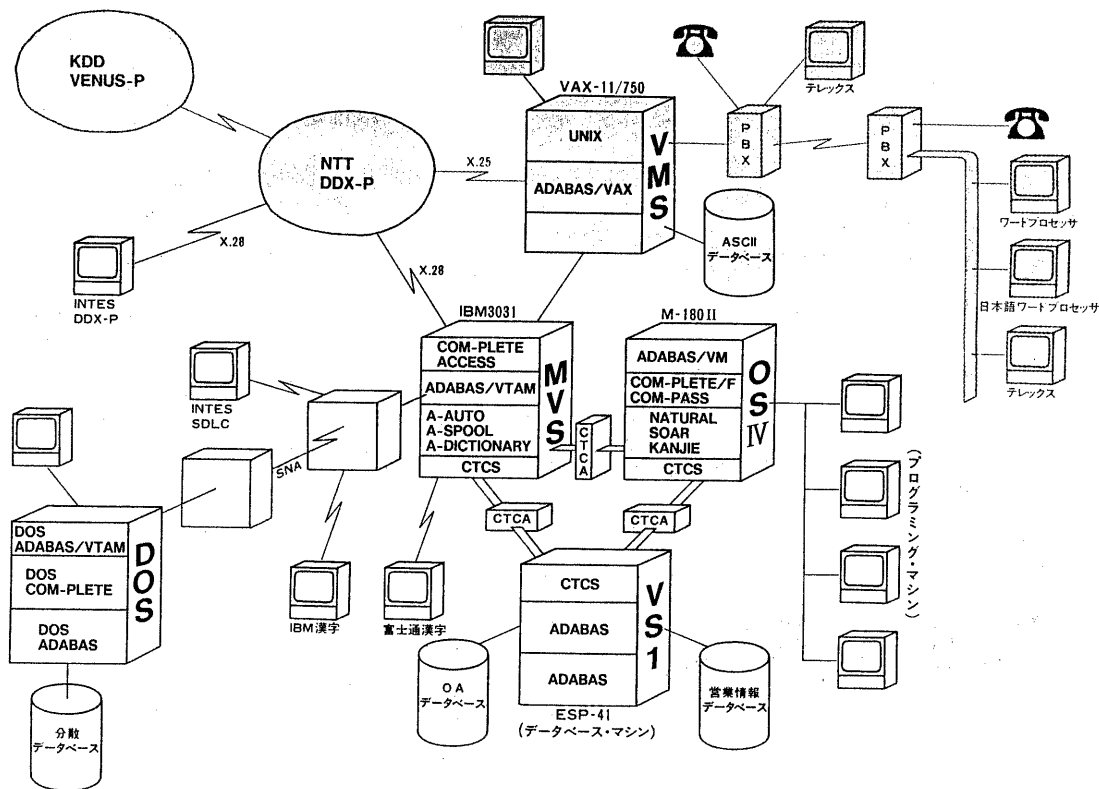


図2

*UNIXは、ベル研究所の登録商標です

図3に示す例は野村証券のCAPITALである。このシステムは集中処理をする際処理量が多く、メーカーが提供する最大機器構成を用いても処理不能なため、処理を分割し、3台のCPUを結合し、1つのDBをアクセス/更新する分散処理システムを構成した例である。CAPITALの場合は分散DBとはなっていないが、DBに対するアクセス/更新がさらに多量になった場合はDBを分割し、分散DB化し、それぞれのDBに1CPUを対応させることによって処理負荷を軽減させ、システムのニーズに対応させ成長させることができる。

5. 分散処理における問題点

1) 一般的な考慮点

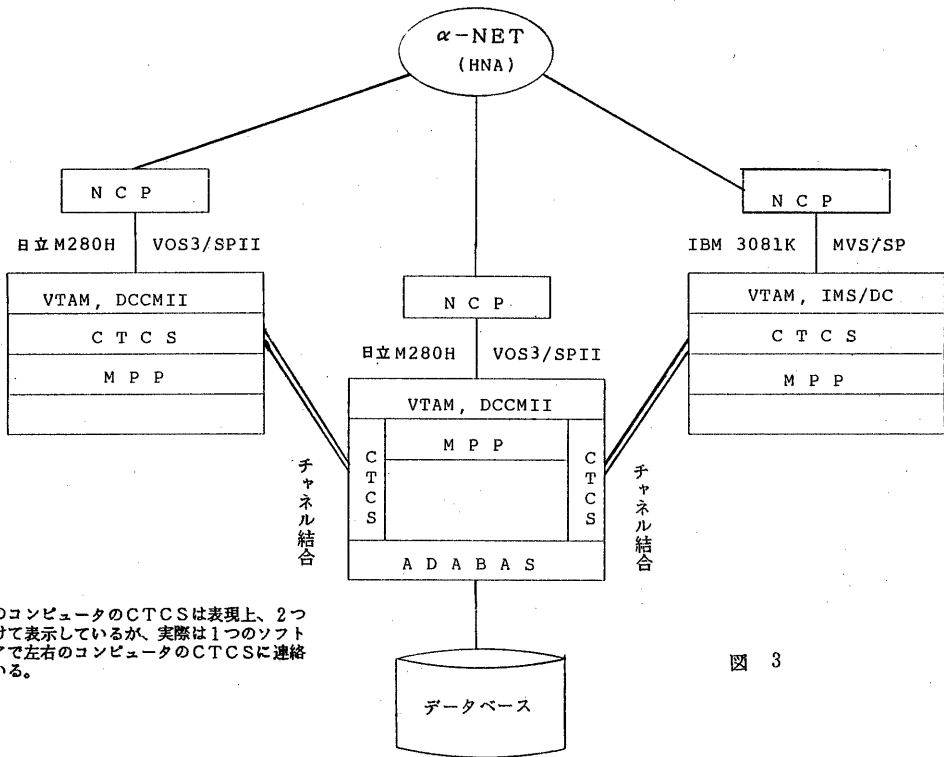
◇分散処理を構成する各コンポーネント、

たとえば、OS、データベース管理システム、コミュニケーション・ソフトウェア、エンドユーザ言語等は非同期に各々バージョン・アップされ、しかもそれらは統一的考え方にそって変更されるとはかぎらない。これにどう対処するか。

◇分散処理は全体が完成することではなく、ステップ・バイ・ステップに変化しながら稼働しなければならない。このような柔軟性をどう保証するか。

◇分散処理環境を作り上げてゆくプロセスで、単独に稼働しているシステムに悪影響を与えないようにするにはどうするか。

◇将来提供される新しいコミュニケーション方式や、ディスクにも最少の労力で対応できることが望ましい。



●中央のコンピュータのCTCSは表現上、2つに分けて表示しているが、実際は1つのソフトウェアで左右のコンピュータのCTCSに接続している。

図 3

●CAPITALサービスで構築したマルチホスト単一DB方式の概念図。中央のM-280Hに証券情報を提供するためのADABASデータベースがあり、2台のホスト・コンピュータ (M-280HとI-3081K) がチャンネル結合する。図の四角で囲んだのは主要なソフトウェアで、NCP (Network Communication Program) はネットワーク通信プログラム、VTAMとDCCM-IIは通信プログラム、MPPはメッセージ処理プログラム、CTCSにチャンネル結合プログラムである。

2) 分散処理環境を実現するための機能

- ◇分散ネットワークのダイナミックな再構成たとえば、ノードの位置、アプリケーションの位置の再構成があっても簡単に対応できること。
- ◇孤立ファイルを取扱う機能のみならず、重複ファイル、分散ファイルを取扱うことも可能なこと。孤立、重複、分散ファイルの定義は次でふれる。
- ◇データベースの管理、ネットワークの管理は、どこのノードからでも可能なこと。
- ◇各ノードのすべてにシステム担当者や、アプリケーション担当者を置かなくても良いこと。

- ◇メッセージ・ジャーナリング、ファイルあるいはデータベース・リカバリ、パフォーマンス・チューニング、アカウントリングおよびセキュリティの定義/変更がローカルにもグローバルにも行なうことができること。
- ◇遠隔地のコンピュータ資源をローカルにも使用できること。
- ◇コンピュータ間のパスやコンピュータ自身がダウンしても全メッセージの送信を保証し、全トランザクションを完了させられること。
- ◇ネットワーク内の代理パスやバックアップ・パスをダイナミックに選択できること。
- ◇仮想マシン・ノードを置き、それによってテストが可能なこと。

◇非分散化アプリケーションに対しては上位方向の互換性があること。

6. 分散データ環境における ファイルのタイプ

分散処理環境内のデータ・ファイルには次の3つのファイルタイプを考える必要がある。

◇孤立ファイル（普通のファイル）
分散処理環境内に1個だけ存在し、それは1個のデータベース中にのみある。普通のファイルはこのタイプである。

◇重複ファイル
分散処理環境内に同一のファイルが重複して2個以上存在し、それぞれ別のデータベース内にある。

◇分散ファイル
同一タイプのレコードが物理的に複数のファイルに分割されており、ユーザは論理的には、1個のファイルと認識しているもの。

孤立ファイルはいわゆる普通のファイルのことで、重複ファイルは分散処理の際ひんぱんに使用されるマスタ・ファイル類のことであり、重複してデータをもつことによって遠くまでデータをリファラーしに行くことなくアクセスのスピード・アップをはかるうというタイプのファイルである。分散ファイルとはいわゆる伝票類を入れるファイルのことで、伝票は各部所に分散して発生するが、このタイプのデータは分散ファイルに入れることになる。

7. 論理DB、ファイルと 物理DB、ファイル

ユーザあるいはユーザ・プログラムが認知するDB、ファイルを論理DB、ファイルと名付ける。1論理DB、ファイルは複数物理DB、ファイルに分割され格納されても良い。

物理DBとは1 ADABAS DBMSによって管理されるDBを指し、その物理DB内の各ファイルを物理ファイルと呼ぶ。集中処理されている状態は論理DBと物理DBがまったく同一と考える。集中処理されているDB（論理）を分割（物理的）し、分散DBとしても、論理的には1DBとみる。論理的には今までと同一であるから、そのDBをアクセス/更新するユーザ・アプリケーションは変更する必要はない。集中処理を分割し分散DB化した簡単な例を図4に示す。

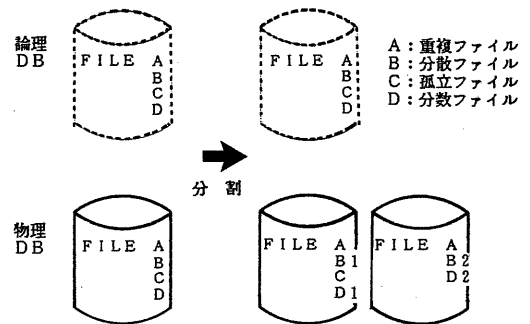


図 4

論理DB、ファイルと物理DBファイルとの対応はADABASの中のADANETと呼ばれるコンポーネントによってダイナミックに行なわれる（図5参照）。物理DBは同一ノード中であって良いし、異なるノード中にあっても良く、また仮想マシン上にあっても良い。異なるノード中にある場合、両者はチャネル結合（CTCA）あるいは通信回線（VTAM）で結合されていなければならない。仮想マシンの時はVMを利用し、VMTFにはCMS、DOS/VSE、VS1、MVS等多種類のOS下にそれぞれADABASを入れ、それらの中で自由にコミュニケーションできる。これらの機能を実現するためにCTCS、ADABAS/VTAM、ADABAS/VMと名付けられたソフトウェア・プロダクトが開発されずで利用されている。

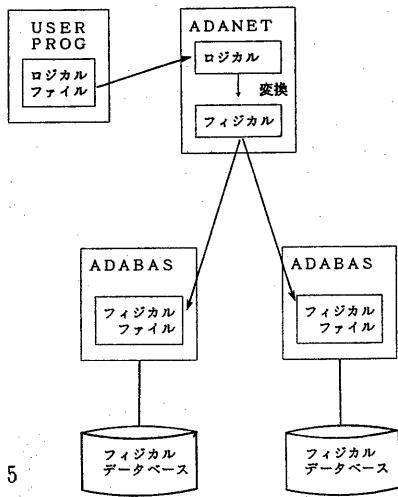


図 5

8. 処理の分割/分散化とデータベースの分割/分散化のプロセス

ADABASを利用し集中処理されているデータ処理を分割/分散化するプロセスをここで示す。図6はチャンネル結合 (CTCA) を利用し処理を分割した例である。STEP 1は集中処理の段階である。STEP 2はSTEP 1でCPU不足になった際、DBMSのCPU負荷をバックエンドCPUに移動させ、CPU負荷を軽減させる例であり、ADABASデータベース・マシンはこの考えを実現したものである。マシン間の結合のためのソフトウェアはCTCSと名付けられている。STEP 3は処理をさらに2分割した例である。

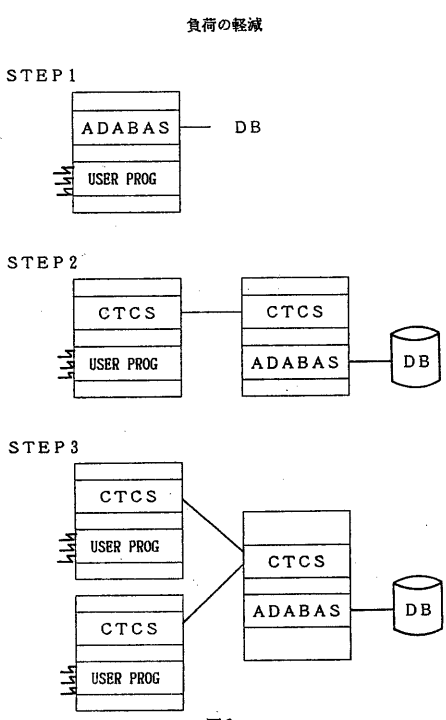


図 6

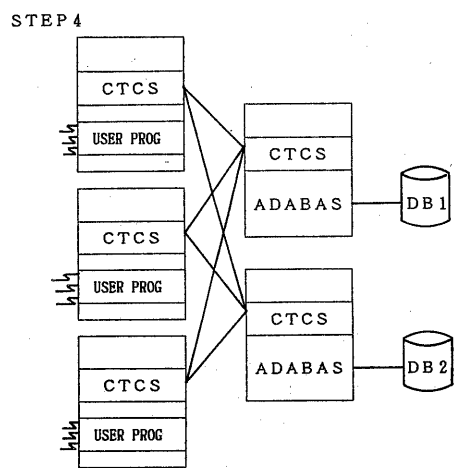
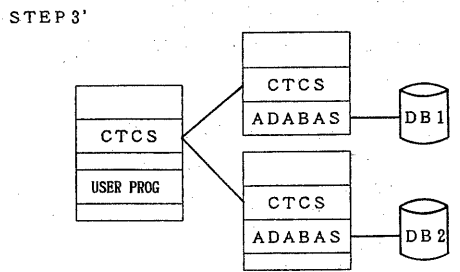


図 7

図3に示した野村證券のCAPITALはこの例である。STEP 2、3はデータベースが分割/分散化されていない。図7のSTEP 3'、STEP 4は処理/DBMS負荷がさらに増大した時にデータベースを分割/分散化した例である。いずれの場合もユーザ・アプリケーションを変更する必要はない。

図8は通信回線を利用し遠隔地からもデータベースへのアクセス/更新を行なう例である。STEP 1は集中処理である。STEP 2は遠隔地からもデータベースをアクセス/更新する例であり、マシン間のコミュニケーションはVTAMを利用し、ADABAS/VTAMと名付けられたソフトウェアがユーザ・プログラムから出されたADABASに対するリクエストをデータベースのあるマシンへ転送する。STEP 3はデータベースが分散化されている場合の例である。この場合もユーザ・アプリケーションはデータベースがある場所を意識する必要がない。

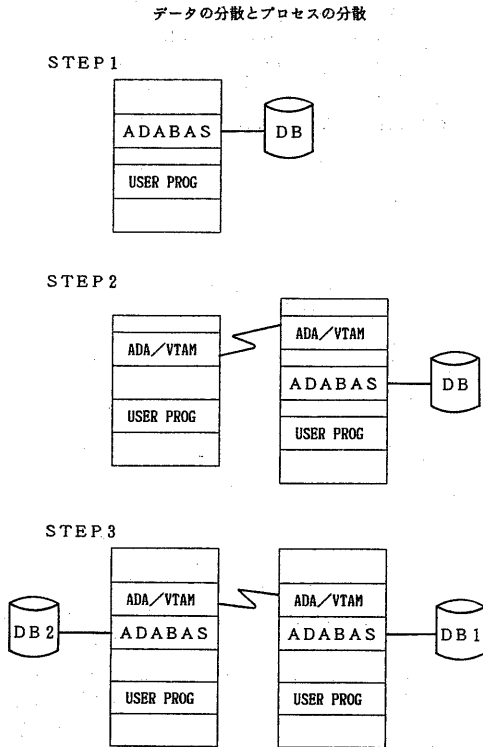


図8

9. 分散環境内の問い合わせ処理

孤立ファイル

孤立ファイルについては、ノードおよびレコードを選択したいファイルをDD (データディクショナリ) を利用しユニークに識別できる。したがって高レベル・エンド・ユーザ言語では論理ファイル名を用いてコーディングすればよい。したがって論理ファイルに対応するすべてのファイルが孤立ファイルの時は、必要に応じてリクエストを転送し、各ノードやDBMSから返されるレコードをユーザ・プログラムへリターンすれば良い。

分散ファイル

分散ファイルはいくつかのノード「論理的には1個のファイル」が分散しているものをいうが、DDには単に、あるファイルが分散ファイルであるという情報を持たせるのみならず、いかなるロジックでデータが分散しているかのロジックも同時に入っている。一番簡単なロジックは各物理ファイル内でユニークにつけられる番号 (一般にタプル・IDと呼ばれる分散ファイル全体でユニークな番号など) によって分けるロジックである。ユーザ・プログラムから分散ファイルに対して検索する場合はDDからどれとどれが各々の物理的分散ファイルであるかを判断し、各々の物理ファイルに対し検索する。ADANETは答を単にマージすれば良く、マージされたものをユーザ・プログラムへ返せば良い。各レコードのリードはADANETによって判断され必要に応じて他のノードへリクエストが転送され、対応レコードが読まれる。

重複ファイル

重複ファイルにたいする検索はユーザ・プログラムが稼働しているコンピュータ (オリジネーティング・ノード) にもしその重複ファイルがあれば、単にそれに検索すればよい。もしオリジネーティング・ノードに重複ファイルがないときは一番近い重複ファイルへリクエストが転送されアクセスされる。

10. 統合化された分散データベースの更新

データベースが物理的に分散した位置にあり、しかもそのいくつかの物理データベースにまたがったデータ更新を1論理トランザクションで行なう場合は一般にマシン・ダウン、DBMSアベンド、回線ダウン、TPモニタ・アベンド等の問題が起るとデータ間の不整合を起す可能性がでてくる。

この問題を完全に解決するためには、まず各ノードのすべてのDBMSが、まず完全なリカバリ機能を持っている必要がある。その上でさらにマルチ・ノード間でのダウンによる不整合をオートマチックに修復する機能をADABAS、ADANETに持たせなければならない。マルチ・ノード間でのリカバリに関しては11. で述べる。

以下まず排他制御に関して述べよう。排他制御を実施するための機能としては普通行なわれているユーザ占有機能を用いる。

ユーザ占有（レコードのホールド）の単位は、すなわち1トランザクション内でアップデートのためにロック（ホールド）する最小単位はデータ・レコードとする。こうすれば容易に分散データベース環境にもユーザ占有概念を拡張できる。ホールド要求については更新のためのレコードを読む時のみ必要であり、ユーザ占有の設定にDBMS内で消費されるオーバーヘッドはたいしたものではない。

ホールド付の読み込み命令が、あるノードに対して発行される場合、個々のレコードのホールド状態はそのローカル・ノード内で設定すればよい。すなわち、個々のユーザに対して、ホールドしたレコードに関する複数のタブルIDリストがいくつかの別々のノードのDBMS内に存在することになる。トランザクション処理中に、もし、ユーザ・プログラムとあるノードとの間のコミュニケーション問題が生じた場合は、そのユーザがホールドした全レコードを、ADABASが自動的に開放すればよい。

ファイルのタイプごとに、ユーザ占有を設定する方法は、次の通り。

◇孤立ファイルについては、データ・レコードは常に1つのユニークに識別されるノードから選択されるので、そのレコードをホールドする。

◇重複ファイルについては、重複ファイルがあることをユーザは全く意識する必要はない。重複ファイルのあるレコードに対してホールド要求があると、全重複ファイルのレコードを自動的にホールドする。ユーザは重複ファイルの全てについて該当レコードがホールド状態になったとき、リクエストに対する答を受け取るだけである。この場合ユーザ占有は1個別レコードについて複数ノードにわたることになる。重複ファイル内でユーザ・レコードがホールドできないと、ユーザはレコードがホールド可能となるまで待たされる。もしデッドロック状態になるとトランザクションを排他せよという要求をADANETがユーザ・プログラムへ返せばよい。

そのトランザクション処理に関係するすべてのレコードがホールドされるとユーザ・プログラムではホールドされた各レコードに対し更新処理を行なうことができる。更新処理に関しても各ファイル・タイプにそってそれぞれは次の状態を扱わなければならない。

◇孤立ファイルおよび分散ファイルについては、分散データベース内で単一のレコードだけが更新される。

◇重複ファイルは2つの方法で更新される。重複ファイルはDD内に即更新かオフライン更新と定義できる。即更新とは、重複ファイルが存在するいずれかのノード内でファイルが更新されると、これら重複ファイルも同一トランザクション内で直ちに更新されることを示す。オフラインとは、オンラインで更新すべき分散データベース内のファイルは1マスタ・ファイルのみであり、ある特定時点でマスタ・ファイルをコピーするまで重複ファイルはそのままであることを示す。即更新については、ユーザ占有は前記の全重複ファイルに対して設定される。実際の更新操作では、ユーザが発行する更新コマンドは重複ファイルの全コピーに対し

てADANETからダブって発行される。更新動作に入る前に重複ファイルの全コピーに対してユーザ占有が設定されるので、全ファイルは同期がとれユーザ干渉（インタフェランス）は起こらないことが自動的に保証される。このような方法をとれば、ユーザにとってトランザクション処理は単一のデータベース環境の場合と同じである。

11. 分散環境における トランザクション処理と同期化

非分散環境と同様、自動的に論理トランザクション・レベルで自動的にリカバリ・リストを扱う方法を考えよう。

ユーザ・プログラムがDBMSと同一のマシン上で稼動する単一データベース・システムとは異なり、分散データベース環境では複数コンポーネントを含む。

◇各データ・ノード内のプロセッサと周辺機器

◇データ・ノード間の通信回線

単一データ、単一プロセッサ環境との主な相違は、コンポーネントの1つに障害があった場合、ネットワーク全体が操作不能になるわけではないことである。個々のノードは操作を続行できる。

ネットワーク内で、ユーザ・プログラムやリクエストを処理するノードをオリジネーティング・ノードといい、トランザクションに含まれるリモート・データベースを持つ全ノードはターゲット・ノードという。

オリジネーティング・ノード以外のターゲット・ノードのマシン・ダウン、オリジネーティング・ノードとターゲット・ノードとの間の回線ダウンなどいろいろのケースについても、トランザクション・リカバリ機能を持たねばならない。これは単に各ノード内のDBMSが単独でリカバリ機能を有しているのみでは問題を解決することができない。

この問題を解決するために次のアーキテクチャを提案する。

まず各DBMSはEND OF TRANSACTION COMMANDをうけるようにする。DBMSによっては内部的にGETによってターミナルからデータを受け取った時それ以前のトランザクション処理を「終了」とみなすものがあるが、ここで提案する方法ではEND OF TRANSACTIONコマンドを出してもらうようにする。

その上でADABAS、ADANETはそのEND OF TRANSACTION（以下ETと略す）を受取った時、次の手順で処理する。

①ユーザ・プログラムはETコマンドを発行する。

②PET（フレリミナリET）コマンドをADANETが自動的に生成し、そのトランザクションに関連する全ノードに対し発行する。オリジネーティング・ノード内のADANETは、どのターゲット・ノードに、PETコマンドを発行したかの全ノードのリストをディスクにセーブする。PETコマンドは、戻し処理できるようなログ・データに対する全I/Oを生成する。PETコマンドはユーザがホールドしているレコードは開放しない。

③ターゲット・ノードに対する全PETコマンドにOKの応答があると、ETコマンドがADANETによって並行して全ノードに対して送られる。当ETコマンドは全ノード内で、ユーザがホールドしていた全リソースを開放する。もしターゲット・ノードの1つから応答がないと、他の全ノードに対してトランザクションのバックアウト（戻し処理）指示がADANETから出され、オリジネーティング・ノードのユーザはバックアウトしたことを知らされる。

④全ノードから応答を受取ると、ユーザにトランザクション終了の旨知らされる。

当手法では、最小のオーバーヘッドで分散処理環境にトランザクション処理を設定できる。

このトランザクション処理に含まれるノードやバスに起こり得る障害や、障害後のリスタートは、1つのノードでトランザクション処理を行なうと同様を扱える。ネットワーク内の障害の最悪のケースは、ADANETが自動的に出すPETと最終のETコマンド間の時間である。

トランザクションに含まれる全ノードに並行してPETおよびETコマンドが送られたとき、この時間は最小となる。これは500マイクロ秒台であろう。なぜならPETではすべてI/Oがすでに出されているからこの間の時間は多くはない。しかしこの間に回線が操作不能となるときは、ターゲット・ノードはトランザクションを終了すべきかバックアウトすべきか決められない。そこでADANETがまずオリジネーティング・ノードとコミュニケーションしトランザクションの状態を確認する。

すなわちトランザクションの待ち時間中に回線が操作可能となるかもしれないし、この場合にはターゲット・ノードは停止中のトランザクションの終了方法を決めるのに必要な情報をADANETから自動的に受け取る。オリジネーティング・ノードとコミュニケーションが確立できないと、停止中のトランザクションがオペレータ・コンソールに表示される。

次の情報を表示する。

◇トランザクション番号

◇オリジネーティング・ノード

◇他のターゲット・ノード

この場合、オペレータは再動作するのに次のどれかの方法をとればよい。

◇回線が間もなく操作可能になると予想できればトランザクションを待ち状態にする。

◇電話でオリジネーティング・ノードのオペレータを呼び出し、特定トランザクションのステータスをきく。

◇トランザクションをバックアウトするか終了するかを任意に決定する。この場合、当トランザクションの全更新操作を特定のログ上にロギングしておけば、後で必要なときこの全情報を使用すればよい。あるノード内のプロセッサのリスタートについては、停止中のトランザクションのバックアウト/終了処理に対する全コミュニケーションも自動的に行なわれ、リスタートは単一プロセッサ環境のように扱える。

12. おわりに

ADABAS分散処理/分散データベース機能について述べた。ADABASの分散データベース機能は集中処理から処理と、データベースの両方を分割/分散化することによって全体的処理効率を向上させ、かつ、フレキシビリティを増大させようというねらいを持っている。日本においてもすでにこの方法は実用化され稼働している。今後大量生産に伴う中規模マシンのコストが大巾に安くなることが予想され、処理とデータベースの分割/分散化の方向は益々確乎たるものとなるであろう。