

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5343399号  
(P5343399)

(45) 発行日 平成25年11月13日(2013.11.13)

(24) 登録日 平成25年8月23日(2013.8.23)

(51) Int.Cl. F I  
**G 0 6 F 1 2 / 0 0 ( 2 0 0 6 . 0 1 )**  
 G 0 6 F 1 2 / 0 0 5 3 5 R  
 G 0 6 F 1 2 / 0 0 5 1 8 A  
 G 0 6 F 1 2 / 0 0 5 3 5 C

請求項の数 7 (全 27 頁)

(21) 出願番号	特願2008-134031 (P2008-134031)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成20年5月22日(2008.5.22)		神奈川県川崎市中原区上小田中4丁目1番1号
(65) 公開番号	特開2009-282746 (P2009-282746A)	(74) 代理人	100070150 弁理士 伊東 忠彦
(43) 公開日	平成21年12月3日(2009.12.3)	(72) 発明者	石本 大明 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
審査請求日	平成23年1月18日(2011.1.18)	(72) 発明者	味沢 丞 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	伊藤 尚洋 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

最終頁に続く

(54) 【発明の名称】 管理プログラム、管理方法、及び管理装置

(57) 【特許請求の範囲】

【請求項1】

データベースにおける第一のデータの参照の開始時に時間の前後関係を識別可能な第一の時系列データを前記第一のデータを識別する識別情報と対応付けて記憶部に記憶して、該第一のデータを参照する処理を開始し、

前記第一のデータの更新要求に応じ、前記第一のデータに対応する第二のデータを前記データベースに生成し、前記更新要求に基づいて該第二のデータの更新を開始し、

前記第二のデータの更新の完了時に第二の時系列データを前記第二のデータを識別する識別情報と対応付けて前記記憶部に記憶し、

前記参照する処理の終了時に前記第一の時系列データと、当該第一の時系列データに対応付けて前記記憶部に記憶された前記第一のデータの識別情報とを前記記憶部から削除し

10

、  
前記参照する処理の終了より後の前記更新の完了に応じて、前記第二の時系列データより前の時間を示す前記第一の時系列データが前記記憶部に記憶されているかを判定し、

前記第二の時系列データより前の時間を示す第一の時系列データが前記記憶部に記憶されていない場合は前記第一のデータを削除する

ことをコンピュータに実行させることを特徴とする管理プログラム。

【請求項2】

前記削除する処理は、前記第二の時系列データより前の時間を示す前記第一の時系列データであって、当該第二の時系列データと同一の前記識別情報に対応付けられている前記

20

第一の時系列データが前記記憶部に記憶されていない場合は前記第一のデータを削除する請求項 1 記載の管理プログラム。

【請求項 3】

前記更新を開始する処理は、前記第一のデータ及び前記第二のデータに対してロックを適用する処理を含み、

前記管理プログラムは、

前記データの更新の完了に応じて、前記第二の時系列データより前の時間を示す前記第一の時系列データが前記記憶部に記憶されているかを判定した結果、当該第一の時系列データが前記記憶部に記憶されている場合は、前記ロックの内容を示す占有情報を前記記憶部に保持させ、

10

前記記憶部に保持されている占有情報に基づいて前記データベースの排他制御を行い、前記参照する処理の終了に応じて前記第二の時系列データより前の時間を示す前記第一の時系列データが前記記憶部に記憶されていない場合は前記占有情報を削除する

ことを前記コンピュータに実行させる請求項 1 又は 2 記載の管理プログラム。

【請求項 4】

前記第二の時系列データより前の時間を示す前記第一の時系列データに係る参照については前記第一のデータを参照対象とし、前記第二の時系列データより後の時間を示す前記第一の時系列データに係る参照については前記第二のデータを参照対象とする処理を前記コンピュータに実行させる請求項 1 乃至 3 いずれか一項記載の管理プログラム。

【請求項 5】

20

前記第一及び第二の時系列データには、複数ノード環境における各ノードを識別するノード識別子が関連付けられ、

前記参照する処理を開始する処理は、他ノードに対して前記第一の時系列データの生成を要求し該要求に応じて該他ノードで生成された前記第一の時系列データを取得し、

前記更新を開始する処理は、前記他ノードからの更新要求に応じ、前記第一のデータに対応する第二のデータを前記データベースに生成し、該第二のデータを更新し、

前記他ノードにおける前記更新の完了に応じて、該他ノードにおいて生成された前記第二の時系列データを取得し、

前記第二の時系列データの取得に応じて、前記他ノードより取得された前記第二の時系列データより前の時間を示す前記他ノードより取得された前記第一の時系列データが前記記憶部に記憶されていない場合は前記第一のデータを削除する

30

処理を前記コンピュータに実行させる請求項 1 乃至 4 いずれか一項記載の管理プログラム。

【請求項 6】

データベースにおける第一のデータの参照の開始時に時間の前後関係を識別可能な第一の時系列データを前記第一のデータを識別する識別情報と対応付けて記憶部に記憶して、該第一のデータを参照する処理を開始し、

前記第一のデータの更新要求に応じ、前記第一のデータに対応する第二のデータを前記データベースに生成し、前記更新要求に基づいて該第二のデータの更新を開始し、

前記第二のデータの更新の完了時に第二の時系列データを前記第二のデータを識別する識別情報と対応付けて前記記憶部に記憶し、

40

前記参照する処理の終了時に前記第一の時系列データと、当該第一の時系列データに対応付けて前記記憶部に記憶された前記第一のデータの識別情報とを前記記憶部から削除し

、前記参照する処理の終了より後の前記更新の完了に応じて、前記第二の時系列データより前の時間を示す前記第一の時系列データが前記記憶部に記憶されているかを判定し、

前記第二の時系列データより前の時間を示す第一の時系列データが前記記憶部に記憶されていない場合は前記第一のデータを削除する

ことをコンピュータが実行することを特徴とするデータ管理方法。

【請求項 7】

50

データベースにおける第一のデータの参照の開始時に時間の前後関係を識別可能な第一の時系列データを前記第一のデータを識別する識別情報と対応付けて記憶部に記憶して、該第一のデータを参照する処理を開始する参照開始手段と、

前記第一のデータの更新要求に応じ、前記第一のデータに対応する第二のデータを前記データベースに生成し、前記更新要求に基づいて該第二のデータの更新を開始する更新手段と、

前記第二のデータの更新の完了時に第二の時系列データを前記第二のデータを識別する識別情報と対応付けて前記記憶部に記憶する完了手段と、

前記参照する処理の終了時に前記第一の時系列データと、当該第一の時系列データに対応付けて前記記憶部に記憶された前記第一のデータの識別情報とを前記記憶部から削除する参照終了手段と、

10

前記参照する処理の終了より後の前記更新の完了に応じて、前記第二の時系列データより前の時間を示す前記第一の時系列データが前記記憶部に記憶されているかを判定する判定手段と、

前記第二の時系列データより前の時間を示す第一の時系列データが前記記憶部に記憶されていない場合は前記第一のデータを削除する削除手段とを有する管理装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、管理プログラム、管理方法、及び管理装置に関する。

20

【背景技術】

【0002】

従来、データベースでは、トランザクション（参照処理・更新処理）の多重実行において、同時実行性を確保しつつ、データの整合性や一貫性を保証するために、下記のような技術が利用されている。

（1）参照処理・更新処理をデータベース資源（行など）のロックにより制御する技術

（2）行をマルチバージョンで管理し、更新中にも更新前のデータを参照させる技術

（3）行を2バージョン以上で管理し、不要レコードの回収を行う技術

（1）では参照処理の際に対象資源に参照ロックを、更新処理の際に更新ロックをとることでデータの一貫性を保持している。参照ロックがとられている資源に対しては、参照ロックを取ろうとする処理のみ許可し、それ以外を排他する。更新ロックが取られている資源に対しては、どのロックの処理も排他する。このようにすることで、参照中や更新中の資源が更新されることを防ぎ、データの一貫性が保たれる。但し、参照処理と更新処理とをシリアライズ（排他待ち）することによるコンカレンシ（並列性）の向上に課題がある。また、データベース資源のロックによる結果として排他待ちが発生するため、業務構築において排他設計が困難となる。

30

【0003】

そこで、（2）では各レコードに複数のバージョンを作成できるようにしている。すなわち、更新の際に新しいバージョンが作成される。参照処理は参照可能なバージョンを参照する（すなわち、参照処理と更新処理において対象資源のバージョンを変える）ことで、参照処理と更新処理とにおいて排他待ちが発生しないようにしている。しかし、過去バージョンのデータの破棄タイミングが不明であるため、過去バージョンを破棄できず、データ領域が肥大化してしまうという問題がある。

40

【0004】

そこで、（3）ではレコードを2バージョン+ で管理することで、更新と参照とを同時に実行できるようにすると共に、各バージョンのレコードに参照カウンタを付加している。具体的には、参照開始後の対象物理レコードのアクセス時に当該物理レコードの参照カウンタに1を加算し、参照の終了時に1を減算する。参照カウンタに基づいて参照されなくなった（参照カウンタが0になった）削除可能な過去バージョンのレコード（過去レコード）を判別及び削除し、データ領域の肥大化を抑止している。但し、バージョン数を

50

固定した場合、当該バージョン数以上のバージョンを作成しようとする最古のバージョンが削除されてしまう。そうすると、最古のバージョンを用いる処理が実行できなくなるという不具合が発生する。そこで、更新トランザクションの完了時に、参照されている過去レコード（参照カウンタが1以上である過去レコード）がある場合は、遅延削除と呼ばれる偽の削除を行う。遅延削除とは、物理レコードに付加された削除ビットと呼ばれる偽の削除が行われたことを表すフラグを立てる処理をいう。遅延削除により、その過去レコードを参照する参照処理以外は当該過去レコードにアクセスできなくなる。全ての参照が完了すると（参照カウンタが0になると）、遅延削除が行われたレコードは削除される。

【特許文献1】特表2007-501468号公報

【非特許文献1】“Postgresql Documentation Manuals PostgreSQL 7.1 Multi-Version Concurrency Control”、[online]、[平成20年4月25日検索]、<<http://www.postgresql.org/docs/7.1/static/mvcc.html>>

【発明の開示】

【発明が解決しようとする課題】

【0005】

しかしながら、(3)では、参照カウンタを物理レコードに付加し、参照の際に更新している。したがって、参照カウンタの更新のたびにデータベースへの書き戻しのためのI/Oが発生する可能性があり、参照処理が低速になるという問題がある。

【0006】

また、DBMS (DataBase Management System) がダウンした場合、そのリカバリ時において物理レコードに書き込まれた参照カウンタなどに何らかの補正（例えば、全レコードの参照カウンタを0クリアする等）が必要とされる。補正が行われないと、リカバリ後に参照カウンタが0にならず、いつまでも不要なレコードが存在してしまう可能性があるからである。そのため、リカバリ時において、CPU、メモリ、2次記憶装置等へのI/Oの発生といったコストがかかるという問題がある。

【0007】

また、(3)では、更新完了後に参照される古いバージョンのレコードが不要となった際のレコードの回収、クリーンアップ等を、参照カウンタを0にした参照処理に行わせている。その結果、当該参照処理が低速となるという問題もある。

【0008】

更に、(3)では、通常の2バージョンに加え、遅延削除により参照中の参照処理以外にアクセスできないバージョンが存在する。そのため、データベース資源の容量を設計する際に、この遅延削除のための容量についても考慮しなければならない。例えば、参照中の或るレコードに対し2回更新を行うと、2回目の更新中の際に、遅延削除による最も古いバージョンのレコード、参照可能なレコード、及び更新中のレコードの3バージョン分の容量が必要となる。遅延削除はいつ、どの程度行われるのかを事前に把握するのは難しいため、データベース資源の容量を見積もるのが難しいという問題がある。

【0009】

本発明は、上記の点に鑑みてなされたものであって、データベースのアクセスの効率性とデータの一貫性とを適切に確保することのできる管理プログラム、管理方法、及び管理装置の提供を目的とする。

【課題を解決するための手段】

【0010】

そこで上記課題を解決するため、管理プログラムは、データベースにおける第一のデータの参照の開始時に時間の前後関係を識別可能な第一の時系列データを前記第一のデータを識別する識別情報と対応付けて記憶部に記憶して、該第一のデータを参照する処理を開始し、前記第一のデータの更新要求に応じ、前記第一のデータに対応する第二のデータを前記データベースに生成し、前記更新要求に基づいて該第二のデータの更新を開始し、前記第二のデータの更新の完了時に第二の時系列データを前記第二のデータを識別する識別情報と対応付けて前記記憶部に記憶し、前記参照する処理の終了時に前記第一の時系列デ

10

20

30

40

50

データを前記記憶部から削除し、前記更新の完了に応じて、前記第二の時系列データより前の時間を示す前記第一の時系列データが前記記憶部に記憶されているかを判定し、前記第二の時系列データより前の時間を示す第一の時系列データが前記記憶部に記憶されていない場合は前記第一のデータを削除することをコンピュータに実行させることを特徴とする。

【0011】

このような管理プログラムでは、データベースのアクセスの効率性とデータの一貫性とを適切に確保することができる。

【発明の効果】

【0012】

データベースのアクセスの効率性とデータの一貫性とを適切に確保することのできる管理プログラム、管理方法、及び管理装置を提供することができる。

【発明を実施するための最良の形態】

【0013】

以下、図面に基づいて本発明の実施の形態を説明する。図1は、第一の実施の形態におけるデータ管理システムの構成例を示す図である。同図において、データ管理システムは、データベースサーバ10及びアプリケーションサーバ20を含む。データベースサーバ10及びアプリケーションサーバ20は、LAN(Local Area Network)又はインターネット等のネットワーク(有線又は無線の別は問わない。)を介して接続されている。

【0014】

データベースサーバ10は、データベース及びデータベースに対するアクセス手段(操作手段)等が実装されたコンピュータであり、データベース11、SQL実行制御部12、アクセス部13、排他制御部14、及びゴーストデーモン15等を有する。

【0015】

データベース11は、記憶装置に体系的に記録されたデータ(情報)の集合である。本実施の形態では、データを表形式によって管理するRDB(Relational Database)を例とする。SQL実行制御部12は、SQL(Structured Query Language)によるデータベース11の操作指示を受信し、当該操作指示に応じた処理を制御する。アクセス制御部13は、SQL実行制御部12からの指示に応じ、データベース11に対する操作(更新、参照、削除等)を実行する。排他制御部14は、データベースに対する操作時における排他制御を行うことによりデータの一貫性を保証する。ゴーストデーモン15は、ゴーストトランザクションの管理等を行うデーモンプロセスである。ゴーストトランザクションとは、更新トランザクション又は削除トランザクションが所定の条件が満たされる場合に移行する(又は開始される)トランザクションをいう。ゴーストトランザクションは、トランザクション自体は実行せずに、移行元のトランザクションの占有情報(ロック情報)を引き継ぐ(保持する)。したがって、移行元のトランザクションがゴーストトランザクションに移行した場合、当該移行元のトランザクションが終了後も操作対象のデータに対するロックは維持される。なお、所定の条件とは、該当するトランザクションのコミット時において操作対象とされていたデータと同一のデータを参照している他のトランザクションが存在することである。

【0016】

アプリケーションサーバ20は、例えば、ユーザによる指示入力に応じてデータベースサーバ10に対してSQLによるデータベース11の操作指示を送信する。

【0017】

図2は、本発明の実施の形態におけるデータベースサーバのハードウェア構成例を示す図である。図2のデータベースサーバ10は、それぞれバスBで相互に接続されているドライブ装置100と、補助記憶装置102と、メモリ装置103と、CPU104と、インタフェース装置105とを有する。

【0018】

データベースサーバ10での処理を実現するプログラムは、CD-ROM等の記録媒体

10

20

30

40

50

101によって提供される。プログラムを記録した記録媒体101がドライブ装置100にセットされると、プログラムが記録媒体101からドライブ装置100を介して補助記憶装置102にインストールされる。補助記憶装置102は、インストールされたプログラムを格納すると共に、データベース11等を格納する。

【0019】

メモリ装置103は、プログラムの起動指示があった場合に、補助記憶装置102からプログラムを読み出して格納する。CPU104は、メモリ装置103に格納されたプログラムに従ってデータベースサーバ10に係る機能を実行する。インタフェース装置105は、ネットワークに接続するためのインタフェースとして用いられる。

【0020】

続いて、データベース11におけるデータ構造について説明する。図3は、本発明の実施の形態におけるデータベースのデータ構造の例を示す図である。なお、本実施の形態において「データ」とは、論理的なテーブルにおける一つの行（論理行）に相当するデータベース資源をいう。

【0021】

同図に示されるように、一つのデータは、物理的なレコードとして管理レコード、削除レコード、及び標準レコードの3つのレコードを含みうる。削除レコードとは、更新処理又は削除処理の際に一時的に生成されるレコードであり、対応するデータが更新中の際に参照対象とされる。削除レコードが生成されることにより、更新処理又は削除処理と参照処理との同時実行性が確保される。標準レコードは、対応するデータが削除されない限り常に存在するレコードであり、更新処理の際は更新対象とされる。また、削除レコードが有る場合、標準レコードは削除レコードの新バージョンのレコードであるともいえる。管理レコードは、二つの物理レコードアドレスとカレントレコードアドレスを含む。二つの物理レコードアドレスのうち一方は、削除レコードに対するアドレス情報（データベース11における物理的な位置情報）である。他方は、標準レコードに対するアドレス情報である。カレントレコードアドレスは、二つの物理レコードアドレスのうち、最新バージョンのレコード（カレントレコード）のアドレス情報が格納されている領域に対する位置情報である。なお、管理レコードはデータが存在する限り常に存在する。すなわち、管理レコード及び標準レコードは常に存在し、削除レコードは、更新処理が行われる際に生成される。

【0022】

以下、データベースサーバ10の処理手順について説明する。図4は、挿入トランザクションの処理手順を説明するためのフローチャートである。

【0023】

アプリケーションサーバ20よりデータの挿入指示を受信したSQL実行制御部12からの指示に応じ、アクセス部13は、挿入対象のデータ（以下、「挿入データ」という。）を挿入するための空きレコード（管理レコード及び標準レコードを生成可能な空き領域）をデータベース11より検索する（S101）。続いて、排他制御部14は、検索された空きレコードに対して更新ロックを適用する（更新ロックをかける）（S102）。本実施の形態において、更新ロックには2種類存在する。一つはEXロックといい、もう一つはSXロックという。管理レコードに対するEXロックは、対応するデータに対して変更（挿入、更新、又は削除）が行われている最中であることを示すロックである。また、管理レコードに対するSXロックは、対応するデータに対して変更は行われている最中ではないことを示す。一方、標準レコード又は削除レコードに対するEXロックは、ある特定の場合（ゴーストトランザクション後に開始した参照の場合）を除き参照対象とすることができないことを示す。また、標準レコード又は削除レコードに対するSXロックは参照対象とすることができることを示す。なお、EXロック、SXロックのそれぞれは更新ロックであることを除いてそれ自体に特段の絶対的な意味はない。すなわち、一方を他方と区別するための相対的な意味を有するに過ぎない。したがって、二つのロックは更新ロックとしての仕組みは同じであり、例えば、EXロック又はSXロックを識別するた

10

20

30

40

50

めの識別子が付加されていると考えてもよい。ここでは、対応するデータがこれから挿入されるため管理レコードに対してEXロックが適用される。また、挿入中のデータを参照させることはできないため、標準レコードに対してEXロックが適用される。

【0024】

続いて、アクセス部13は、空きレコードに対して管理レコード及び標準レコードを生成する(S103)。標準レコードの生成において、挿入データが標準レコードに登録される。続いて、アクセス部13は、生成された管理レコードに対して標準レコードのアドレス情報を物理レコードアドレスとして登録する(S104)。なお、ここでは、削除レコードは生成されないため、もう一つの物理レコードアドレスは空のままとされる。続いて、アクセス部13は、管理レコードのカレントレコードアドレスに、標準レコードに係る物理レコードアドレスの位置情報を登録する(S105)。続いて、コミット処理が実行される(S106)。なお、コミット処理の詳細については後述する。

10

【0025】

続いて、更新トランザクションについて説明する。図5は、更新トランザクションの処理手順を説明するためのフローチャートである。

【0026】

アプリケーションサーバ20よりデータの更新指示を受信したSQL実行制御部12からの指示に応じ、アクセス部13は、更新対象のデータ(以下、「更新データ」という。)に係るレコード(管理レコード及び標準レコード)をデータベース11より検索する(S121)。続いて、排他制御部14は、検索された管理レコード及び標準レコードに対してSXロックを適用する(S122)。続いて、アクセス部13は、標準レコードについて、更新後の状態を登録するための空きレコード(標準レコードに対応する1レコード分)をデータベースより検索する(S123)。続いて、排他制御部14は、更新処理を開始するため、検索された空きレコードと更新データに係る管理レコードとにEXロックを適用する(S124)。

20

【0027】

続いて、アクセス部14は、更新前の標準レコードに対して更新した結果を空きレコードに登録する(S125)。したがって、この時点において標準レコードには更新前の状態が、空レコードには更新後の状態が登録されていることになる。なお、図5における以下の説明において、更新前の状態が登録されたレコード(これまでの標準レコード)を「旧標準レコード」と呼ぶ。また、更新後の状態が登録されたレコードを(これまでの空きレコード)を「標準レコード」と呼ぶ。

30

【0028】

続いて、アクセス部13は、管理レコードに対して標準レコードのアドレス情報を物理レコードアドレスとして登録する(S126)。続いて、アクセス部13は、管理レコードのカレントレコードアドレスに、標準レコードに係る物理レコードアドレスの位置情報を登録する(S127)。続いて、アクセス部13は、旧標準レコードを削除レコードに変更する(S128)。削除レコードへの変更とは、該当するレコードが削除レコードであることを示す情報(削除フラグ)を記録することをいう。削除フラグは、該当するレコード(物理レコード)に記録してもよいし、該当するレコードに関連付けてメモリ装置103上に記録してもよい。続いて、コミット処理が実行される(S129)。

40

【0029】

続いて、削除トランザクションについて説明する。図6は、削除トランザクションの処理手順を説明するためのフローチャートである。

【0030】

アプリケーションサーバ20よりデータの削除指示を受信したSQL実行制御部12からの指示に応じ、アクセス部13は、削除対象のデータ(以下、「削除データ」という。)に係るレコード(管理レコード及び標準レコード)をデータベース11より検索する(S141)。続いて、排他制御部14は、検索された管理レコード及び標準レコードに対してSXロックを適用する(S142)。続いて、排他制御部14は、削除処理を開始す

50

るため、管理レコードにEXロックを適用する(S143)。続いて、アクセス部14は、標準レコードを削除レコードに変更する(S144)。標準レコードから削除レコードへの変更については、図5のS128において説明した通りである(すなわち、該当レコードについて削除フラグを記録する。)。続いて、コミット処理が実行される(S145)。

**【0031】**

図4、図5、又は図6において説明した処理手順より明らかなように、挿入処理中、更新処理中、又は削除処理中における、各レコードのロック状態は次の通りである。なお、「処理中」とは、コミット前の状態をいう。

**【0032】**

図7は、操作対象とされたレコードのロック状態を示す図である。同図において、(A)は挿入処理中のロック状態を示す。挿入処理中には、管理レコード及び挿入レコードのそれぞれにEXロックが適用される(図4参照)。(B)は更新処理中のロック状態を示す。更新処理中には、管理レコードと標準レコードに対してEXロックが適用され、削除レコード(又は削除レコードに変更される前の旧標準レコード)にはSXロックが適用される(図5参照)。したがって、更新処理中において、削除レコードは参照可能とされる。(C)は削除処理中のロック状態を示す。削除処理中には、管理レコードに対してEXロックが適用され、削除レコード(又は削除レコードに変更される前の標準レコード)にはSXロックが適用される。したがって、削除処理中において、削除レコードは参照可能とされる。なお、各レコードのロックの内容(状態)を示す情報(占有情報)は、トランザクション情報の一部としてトランザクションごとにメモリ装置103に記録(保持)される。

**【0033】**

続いて、ステップS106、S129、及びS145におけるコミット処理について説明する。図8は、コミット処理の処理手順を説明するためのフローチャートである。同図の説明におけるレコードは、コミットの対象とされたトランザクション(以下、「カレントトランザクション」という。)において操作対象とされたデータに係るレコードをいう。

**【0034】**

コミットの開始に際し、アクセス部13は、カレントトランザクションに対して発行データを発行し、当該発行データをカレントトランザクションに関連付けて発行データ格納領域に登録する(S201)。発行データ格納領域は、メモリ装置103内における所定のメモリ領域である。

**【0035】**

図9は、発行データの構成例を示す図である。同図において、発行データ150は、ノード番号、対象データID、時系列番号、及び発行時期種別等の情報を含む。ノード番号は、ロードシェア(負荷分散)環境(複数ノード環境)のようにデータベースサーバ10が複数存在する構成において各データベースサーバ10(ノード)を識別するための番号(識別子)である。したがって、図1に示されるように単一ノードの環境の場合は必須ではない。対象データIDは、発行データ150の発行対象とされたトランザクションにおいて操作対象とされたデータの識別子である。複数のデータが操作対象とされた場合は複数の対象データIDが登録される。時系列番号は、発行データ150の発行順序(時間の前後関係)を識別するための番号である。時系列番号は、例えば、通番で各発行データ150に割り当てられる。但し、時系列番号は、時間の前後関係を識別できるデータであれば所定のものに限定されない。例えば、時刻であってもよいし、他の符号であってもよい。発行時期種別は、発行データ150が発行された時期(タイミング)の種別(参照の開始又はコミットの開始時)を示す情報である。すなわち、発行データ150は、ステップS201のタイミング(コミットの開始時)のみならず、参照の開始時にも発行される。なお、参照の開始時に発行された発行データ150は、当該参照の終了に応じて削除される。

10

20

30

40

50



## 【 0 0 3 6 】

続いて、アクセス部 1 3 は、当該コミットより前に開始され、カレントトランザクションにおいて操作対象とされたデータを参照する参照処理であって、現時点においても実行中のもの（以下、「並列参照処理」という。）の存否を判定する（S 2 0 2）。この判定は、発行データ 1 5 0 に基づいて行われる。すなわち、ステップ S 2 0 1 において発行された発行データ 1 5 0 の対象データ ID と少なくとも一部が重複する対象データ ID を含み、かつ、発行時期種別が参照の開始を示し、かつステップ S 2 0 1 にて発行された発行データの時系列番号より小さい発行データ 1 5 0 の存否が判定される。斯かる発行データ 1 5 0 が存在する場合、並列参照処理は存在すると判定する。斯かる発行データ 1 5 0 が存在しない場合、並列参照処理は存在しないと判定する。

10

## 【 0 0 3 7 】

例えば、図 1 0 は、発行データに基づく並列参照処理の存否の判定内容を説明するための図である。同図では、発行データ格納領域 1 6 0 に記録されている発行データ 1 5 0 が発行順に配列されている。

## 【 0 0 3 8 】

同図において、発行データ 1 5 0 a は、参照の開始時に発行され、時系列番号が「1 0 0」の発行データ 1 5 0 である。一方、発行データ 1 5 0 b は、更新のコミット時に発行され、時系列番号が「1 0 1」の発行データある。ここで、発行データ 1 5 0 a の時系列番号（1 0 0）は、発行データ 1 5 0 b の時系列番号（1 0 1）より小さい。したがって、発行データ 1 5 0 a に係る参照は発行データ 1 5 0 b に係る更新のコミットよりも先に実行されたものであることが分かる。また、発行データ 1 5 0 a が存在していることにより、当該参照はまだ終了していないことが分かる。よって、この場合、並列参照処理は存在すると判定される。

20

## 【 0 0 3 9 】

並列参照処理が存在しない場合（S 2 0 2 で N O）、アクセス部 1 3 は、当該通知に応じ、カレントトランザクションをコミットする（S 2 0 3）。続いてアクセス部 1 3 は、削除レコードが有る場合は削除レコードを削除する（S 2 0 4）。図 7 の（B）、（C）より、更新又は削除が行われた場合は削除レコードが存在する。続いて、アクセス部 1 3 は、管理レコードに登録されている物理データアドレスのうち、削除レコードに対応する物理データアドレスをクリアする（S 2 0 5）。続いて、当該コミットが削除処理を含むトランザクションのコミットである場合（S 2 0 6 で Y E S）、アクセス部 1 3 は、削除されたデータに係る管理レコードを削除する（S 2 0 7）。続いて、アクセス部 1 3 は、ステップ S 2 0 1 において生成された発行データ 1 5 0 を削除する（S 2 0 8）。続いて、排他制御部 1 4 は、カレントトランザクションによってロックされていたレコードに対してアンロック（ロックの解放）を実行する（S 2 0 9）。より具体的には、カレントトランザクションに係る占有情報をクリア（削除）する。

30

## 【 0 0 4 0 】

一方、並列参照処理が存在する場合（S 2 0 2 で Y E S）、アクセス部 1 3 は並列参照処理が存在することをゴーストデーモン 1 5 に通知する。当該通知に応じ、ゴーストデーモン 1 5 は、並列参照処理によって参照されているレコードが削除されてしまうのを防止するためゴーストトランザクションを生成する（S 2 1 0）。この際、ゴーストデーモン 1 5 は、当該ゴーストトランザクションにカレントトランザクションの占有情報及び発行データ 1 5 0 を引き継がせる。占有情報がゴーストトランザクションに引き継がれることにより（すなわち、カレントトランザクションのコミットに応じてその占有情報が削除されることが回避されることにより）、並列参照処理によって参照されているレコードの削除が防止される。すなわち、排他制御部 1 4 は、ゴーストトランザクションが有する占有情報にも基づいて（をも参照して）、データベース 1 1 に対する排他制御を行うからである。ゴーストトランザクションの実体は、例えば、占有情報を管理するためにメモリ装置 1 0 3 内に生成されるデータでよい。したがって、占有情報のゴーストトランザクションへの引継ぎとは当該占有情報をメモリ装置 1 0 3 内に保持させることに相当する。なお、

40

50

当該データが、他の通常のトランザクションの情報と同様の形式の場合は、ゴーストトランザクションであるか否かを示す情報を当該データに記録すればよい。続いて、アクセス部13は、カレントトランザクションをコミットする(S211)。これにより、カレントトランザクションは、ゴーストトランザクションに移行したことになる。

#### 【0041】

図8より明らかなように、ゴーストトランザクションは、例えば、次のような契機によって生成される(発生する)。図11は、ゴーストトランザクションの発生の契機の具体例を示す図である。図中、左から右へ時間が経過している。また、太線の矢印は、参照処理又は更新トランザクションを示す。太い破線の矢印はゴーストトランザクションを示す。また、参照処理の対象データと更新トランザクションの対象データの少なくとも一部は重複する。

10

#### 【0042】

同図では、まず、参照処理が開始されている(S21)。参照処理の開始に応じ、当該参照処理に対する発行データ150rが発行され、発行データ格納領域160に登録される(S22)。その後更新トランザクションが開始される(S23)。本実施の形態では、参照処理の際に参照ロックは行われない。したがって、既に参照中のデータに対して更新を行おうとしても排他されない。更新トランザクションがコミットされると更新トランザクションに対する発行データ150uが発行され、発行データ格納領域160に登録される(S24)。このタイミングで、更新トランザクションがゴーストトランザクションへ移行すべきか否かが判定される。ここでは、発行データ150uの時系列番号よりも小さな時系列番号を有する発行データ150rが存在する(S25)。したがって、更新トランザクションはゴーストトランザクションへ移行する(S26)。その後、ゴーストトランザクションへの移行の要因となった参照処理が終了すると(S27)、当該参照処理に対する発行データ150rが削除される(S28)。その結果、発行データ150uより時系列番号が小さい発行データ150は存在しないため、発行データ150uが削除され(S29)、ゴーストトランザクションは終了する(S30)。

20

#### 【0043】

続いて、データの参照処理について説明する。図12は、参照処理の処理手順を説明するためのフローチャートである。

#### 【0044】

アプリケーションサーバ20よりデータの参照指示を受信したSQL実行制御部12からの指示に応じ、アクセス部13は、当該参照処理を実行するトランザクション(カレントトランザクション)に対して発行データ150(図9参照)を発行し、発行データ格納領域160に登録する(S301)。なお、ここで生成される発行データ150の発行時期種別には参照の開始を示す値が登録される。

30

#### 【0045】

続いて、データベースサーバ10は、参照対象のデータ(以下、「参照データ」)の参照を開始する(S302)。なお、上記したように、本実施の形態では参照時において参照ロックは行われない。

#### 【0046】

まず、排他制御部14は、参照データに係るレコードは既にロック(占有)されているか否かを確認(判定)する(S303)。ロックされているか否かは、現在存続(実行)中のトランザクション(ゴーストトランザクションも含む。)の占有情報に基づいて判定される。なお、以下の図12の説明におけるレコードは、参照データに係るレコードをいう。

40

#### 【0047】

ロックされていない場合(S303でNO)、又はロックされていたとしてもロックの主体がカレントトランザクションである場合(S304でYES)、アクセス部13は、管理レコードにおけるカレントレコードアドレスが指す物理レコードアドレスに係る標準レコードを参照対象とする(S305)。

50

## 【 0 0 4 8 】

一方、他のトランザクションによってロックされている場合（S 3 0 4でNO）、排他制御部 1 4は、ロックの主体がゴーストトランザクションであるか否かを判定する（S 3 0 6）。ゴーストトランザクションであるか否かは、占有情報に関連付けて記録されている情報（ゴーストトランザクションであるか否かを示す情報）に基づいて判定すればよい。

## 【 0 0 4 9 】

ロックの主体がゴーストトランザクションで無い場合（これは、ロックの主体のトランザクションがまだ存続している場合に相当する。）（S 3 0 6でNO）、アクセス部 1 3は、標準レコード及び削除レコードのうち、SXロックが適用されているレコードを参照対象とする（S 3 0 7）。例えば、図 7（B）又は（C）の状態であれば、削除レコードが参照される。但し、例えば、（C）において、標準レコードが削除レコードへ変更される前であれば（図 6 参照）、標準レコードが参照される。一方、ロックの主体がゴーストトランザクションである場合（これは、ロックの主体のトランザクションが既にコミットされている場合に相当する。）（S 3 0 6でYES）、アクセス部 1 3は、生成されている発行データ 1 5 0に基づいて参照対象を決定（判定）する（S 3 0 8）。当該ステップの詳細については後述する。

## 【 0 0 5 0 】

続いて、アクセス部 1 3は、参照対象とされたレコードを参照し（S 3 0 9）、参照を終了させる（S 3 1 0）。続いて、アクセス部 1 3は、ステップ S 3 0 1において発行した発行データ 1 5 0を削除する（S 3 1 1）。

## 【 0 0 5 1 】

続いて、発行データ 1 5 0を削除することにより終了するゴーストトランザクションの有無の判定、及び終了が可能なゴーストトランザクションが存在する場合は、終了可能なすべてのゴーストトランザクションの終了処理が実行される。したがって、ゴーストトランザクションが存在しない場合、以降の図 1 2中の破線で囲まれた処理P4は実行されない。なお、P4の処理はゴーストデーモン 1 5にて実施される処理であり、参照処理とは非同期に実施される。

## 【 0 0 5 2 】

まず、アクセス部 1 3は、削除した発行データ 1 5 0がゴーストトランザクションの最後の要因の参照処理に関するものであるか否かを判定する（S 3 1 2）。ゴーストトランザクションの要因となった参照処理とは、ゴーストトランザクションの移行元のトランザクションのコミット前に同一データに対して行われた参照処理をいう。

## 【 0 0 5 3 】

当該判定内容の具体例を図 1 3に示す。図 1 3は、ゴーストトランザクションの最後の要因となった参照処理の終了の判定内容を説明するための図である。図 1 3中、図 1 0と同一部分には同一符号を付している。但し、同図において発行データ 1 5 0 bはゴーストトランザクションに係る発行データ 1 5 0であるとする。したがって、発行データ 1 5 0 aは、ゴーストトランザクションの要因となった参照処理に係る発行データ 1 5 0である。

## 【 0 0 5 4 】

同図では、発行データ 1 5 0 bより時系列番号が小さい発行データ 1 5 0 aは削除されている（図中×は削除を示す。）。また、発行データ 1 5 0 a以外に、発行データ 1 5 0 bより時系列番号が小さい発行データ 1 5 0は存在しない。したがって、この場合、削除された発行データ 1 5 0 aがゴーストトランザクションの最後の要因であった（ゴーストトランザクションの要因となった参照処理は全て終了している）と判定する。一方、図 1 0のような状態の場合、ゴーストトランザクションの要因となった参照処理の全ては終了していない（終了した参照処理以外にまだゴーストトランザクションの要因となる参照処理が存在する）と判定する。

## 【 0 0 5 5 】

10

20

30

40

50

終了可能と判断されたゴーストトランザクションが存在する場合（S 3 1 2でYES）、アクセス部 1 3はその旨をゴーストデーモン 1 5に通知する。当該通知に応じ、ゴーストデーモン 1 5は、削除レコードが有る場合は削除レコードを削除する（S 3 1 3）。続いて、ゴーストデーモン 1 5は、管理レコードに登録されている物理データアドレスのうち、削除レコードに対応する物理データアドレスをクリアする（S 3 1 4）。続いて、ゴーストデーモン 1 5は、当該ゴーストトランザクションが開始されたタイミングが削除処理を含むトランザクションのコミットであるか否かを判定する（S 3 1 5）。当該判定は、当該ゴーストトランザクションに係る発行データ 1 5 0の発行時期種別に基づいて行えばよい。削除処理を含むトランザクションのコミットである場合（S 3 1 5でYES）、ゴーストデーモン 1 5は、削除されたデータに係る管理レコードを削除する（S 3 1 6）。続いて、ゴーストデーモン 1 5は、当該ゴーストトランザクションに係る発行データ 1 5 0を削除する（S 3 1 7）。続いて、ゴーストデーモン 1 5は、当該ゴーストトランザクションを終了させる（S 3 1 8）。具体的には、当該ゴーストトランザクションの実体としてのデータを削除する。したがって、当該データに含まれる占有情報もクリア（削除）され、ゴーストトランザクションによるロックがアンロックされる。

【 0 0 5 6 】

このように、ゴーストトランザクションの要因となった参照処理が全て終了することにより、ゴーストトランザクションは終了し、ゴーストトランザクションによるロックは解放される。

【 0 0 5 7 】

続いて、ステップ S 3 0 8の詳細について説明する。図 1 4は、ゴーストトランザクションによって占有されているデータの参照対象レコードの判定処理を説明するためのフローチャートである。

【 0 0 5 8 】

ステップ S 3 0 8 1において、アクセス部 1 3は、ステップ S 3 0 1において参照処理に対して発行された発行データ 1 5 0の時系列番号（参照番号）と、ロックの主体のゴーストトランザクションに係る発行データ 1 5 0の時系列番号（コミット番号）とを比較する。参照番号の時系列番号がコミット番号の時系列番号未満の場合（S 3 0 8 2でYES）、アクセス部 1 3は、S Xロックが適用されているレコードを参照対象とする（S 3 0 8 3）。参照番号がコミット番号より大きい場合（S 3 0 8 2でNO）、アクセス部 1 3は、管理レコードのカレントレコードアドレスが指す物理レコードアドレスに係るE Xロックが適用されたレコードを参照対象とする（S 3 0 8 4）。

【 0 0 5 9 】

図 1 4の処理の具体例を示す。図 1 5は、ゴーストトランザクションによって占有されているデータの参照対象レコードの判定処理の具体例を示す図である。図中、左から右へ時間が経過している。また、太線の矢印は参照処理又は更新トランザクションを示す。太い破線の矢印はゴーストトランザクションを示す。

【 0 0 6 0 】

同図では、まず、参照処理 Aが開始されて、時系列番号が 9 9の発行データ 1 5 0 Aが発行されている。続いて、更新トランザクションがコミットされ、時系列番号が 1 0 0の発行データ 1 5 0 Uが発行されている。更新トランザクションのコミット時において、参照処理 Aは並列参照処理であるため、更新トランザクションはゴーストトランザクションに移行する。続いて、参照処理 Bが開始され、時系列番号 1 0 1の発行データ 1 5 0 Bが発行されている。

【 0 0 6 1 】

その後、時間 t 1のタイミングで参照処理 A及びBの双方が、ゴーストトランザクションによって占有されているデータを参照しようとする。この際、参照処理 Aは、ゴーストトランザクションの時系列番号（1 0 0）より小さい時系列番号（9 9）を有するため（すなわち、参照処理 Aは更新トランザクションのコミットより前に開始しているため）、S Xロックのレコード（更新前のレコード（削除レコード））を参照対象とする（S 3 0

10

20

30

40

50

83に対応)。一方、参照処理Bは、ゴーストランザクションの時系列番号(100)より大きい時系列番号(101)を有するため(すなわち、参照処理Bは更新ランザクションのコミットより後に開始しているため)、EXロックのレコード(更新後のレコード(標準レコード))を参照対象とする(S3084に対応)。このように、時系列番号は、参照対象とする物理レコードの判定にも用いられる。

【0062】

上述したように、第一の実施の形態によれば、更新中のデータについても、参照処理の開始時点のコミット済みの整合性の有るレコード(削除レコード)を、排他待ちすることなく参照させることができる。したがって、データベースへのアクセスの効率化を図ることができる。

10

【0063】

また、発行データ150により更新処理等のコミットの時期や参照処理の開始時期の前後関係を把握することができる。したがって、コミット時において同一データに対する参照処理の有無を適切に判定することができる。したがって、更新処理等のコミットに応じて参照中のレコード(削除レコード)が削除されることを防止することができる。また、参照の終了時において、発行データ150が削除され、それに応じて更新処理等のコミットより前に発行された発行データ150の存在の有無が判定される。その結果、斯かる発行データ150が存在しない場合は削除レコードが削除される。したがって、適切なタイミングで削除レコードを削除することができる。

【0064】

20

また、更新処理等のコミット時において、同一データに対する参照処理が存在する場合は当該更新処理等に係るランザクションの占有情報はコミットに応じて削除されずにゴーストランザクションに引き継がれて保持される。そして、ゴーストランザクションは、当該参照処理が終了するまで存続する。したがって、参照対象とされているデータに対して適切に排他制御を行うことができる。

【0065】

また、各データについて管理されるバージョンは最大で2バージョン(標準レコードと削除レコード)であるため、データベース格納領域の容量設計を容易化することができる。

【0066】

30

また、従来技術における参照カウンタと異なり、発行データ150は高速にアクセス可能なメモリ装置103内に記録される(データベース11の物理レコードには記録されない)。したがって、参照処理の性能の劣化を防止することができる。仮に、従来技術の参照カウンタの全てをメモリ上で管理したとしても、本実施の形態の方が有利である。すなわち、参照処理が1度に複数のレコード、例えば1000件のレコードを参照する場合を考える。この場合、従来技術であれば、1000件分の参照カウンタを更新することになり、その分多くのCPUやメモリを消費することになる。しかし、本実施の形態では、1000件分の対象データIDが一つの発行データ150によって管理される。したがって、一つの発行データ150を生成又は削除すればよい。このように、本実施の形態では、1度に参照するデータ数が増えるほど、従来技術に対して有利となる。

40

【0067】

また、発行データ150がメモリ装置103において管理されることにより、データベース11において障害が発生した場合であっても、参照カウンタの補正等の処理は不要である。したがって、高速にリカバリ処理を行うことができる。

【0068】

上記より、本実施の形態は、参照処理が多く、参照処理の負荷が高いデータベースに対して特に効果的であると言える。

【0069】

また、本実施の形態では、不要となった(参照されなくなった)削除データの削除は、ゴーストデーモン15によってバックグラウンドで行われるため、参照処理自体に対する

50

負荷の増大が抑止される。

【 0 0 7 0 】

なお、発行データ150における対象データIDは、必須項目ではない。対象データIDが無い場合、コミット時(図8)や参照終了時(図12)等において、参照対象又は更新対象等が異なる発行データ150も比較対象とされる。しかし、この場合であっても、更新と参照との並列性を維持しつつ、更新等のコミットによる参照中のレコードの削除を防止することができる。但し、この場合、ゴーストデーモンの終了時期が遅れ、それによりゴーストデーモンによるロックが必要以上に長く維持される可能性がある。その結果、他の更新処理等が排他される可能性がある。このような可能性を考慮すれば、本実施の形態のように、発行データ150において対象データIDを管理することが望ましい。

10

【 0 0 7 1 】

次に、第二の実施の形態について説明する。第二の実施の形態では第一の実施の形態と異なる点について説明する。したがって、特に言及しない点については第一の実施の形態と同様でよい。

【 0 0 7 2 】

図16は、第二の実施の形態におけるデータ管理システムの構成例を示す図である。図16中、図1と同一部分には同一符号を付し、その説明は省略する。

【 0 0 7 3 】

同図に示されるように、第二の実施の形態では2つ以上のノード(データベースサーバ10(データベースサーバ10a、10b))がアプリケーションサーバ20に接続されている。すなわち、第二の実施の形態におけるデータ管理システムではロードシェア環境が実現されている。したがって、アプリケーションサーバ20は任意のデータベースサーバ10に接続(コネクション)してSQL文を実行させることができる。なお、各データベースサーバ10の構成例は、図1に示されるものと同様でよい。二つのデータベース10の構成要素を区別する場合、データベースサーバ10aの構成要素については参照番号の末尾に「a」を付加する。データベースサーバ10bの構成要素については参照番号の末尾に「b」を付加する。

20

【 0 0 7 4 】

図17は、第二の実施の形態におけるデータ管理システムの動作概要を説明するための図である。なお、同図では、アクセス部13、排他制御部14、及びゴーストデーモン15は、便宜上省略されている。

30

【 0 0 7 5 】

アプリケーションサーバ20といずれか一方のデータベースサーバ10(同図の例では、データベースサーバ10a)とのコネクションが確立された後、他方のデータベースサーバ10(データベースサーバ10b)によって管理されているデータに対するSQLの実行要求が行われるとする(S21)。そうすると、データベースサーバ10aのSQL実行制御部12aは、他ノードのSQL実行制御部12bにSQLの実行制御を依頼する(S22)。この場合、アプリケーションサーバ20とのコネクションを確立したデータベースサーバ10aを「チーフノード」という。また、チーフノードからSQLの実行制御を依頼されたデータベースサーバ10bを「アシスタントノード」という。

40

【 0 0 7 6 】

ところで、ロードシェア環境では、各ノードのシステム時刻が必ずしも一致するとは限らない。したがって、各ノードのシステム時刻に基づいてトランザクション等の前後関係を判定した場合、正しい前後関係が得られない可能性がある。そこで、第二の実施の形態では、ロードシェア環境における斯かる問題点を考慮して、第一の実施の形態を利用してデータの一貫性及びトランザクションの同時実行性を適切に確保した例を説明する。

【 0 0 7 7 】

図18は、第二の実施の形態の処理概要を説明するための図である。また、図19は、第二の実施の形態の処理概要において各ノードで管理される発行データの遷移を示す図である。図中、ノード0はデータベースサーバ10aを示す。また、ノード1はデータベー

50

スサーバ10bを示す。なお、図18において参照処理の対象データと更新トランザクションの対象データは同じであるとする。

【0078】

例えば、ノード1において参照処理が開始されると(S31)、ノード1において発行データ150a(時系列番号:200)が発行される(S32)。なお、ノード1は当該参照処理のチーフである。続いて、ノード1において参照処理が開始されたことがノード0に通知され、ノード0において当該参照処理に対する発行データ150b(時系列番号:100)が発行される(S33)。なお、時系列番号はノードごとに通番である。続いて、ノード0において発行された発行データ150bがノード1の参照に送信される(S34)。この段階において、各ノードに管理されている発行データ150は、図19(A)の通りとなる。すなわち、ノード0には発行データ150bが管理されている。また、ノード1には発行データ150aが、ノード1の参照には150a及び150bが管理されている。

10

【0079】

続いて、ノード0において、ノード1に属するデータに対する更新を含む要求に応じ更新トランザクションが開始されたとする(S35)。この場合、ノード0は更新トランザクションのマスタとして、ノード1に対して当該更新トランザクションを依頼する(S36)。したがって、ノード1はアシスタントとして更新トランザクションを開始する(S37)。続いて、ノード0において更新トランザクションのコミットが開始されると、ノード0において発行データ150c(時系列番号:101)が発行される(S38)。

20

【0080】

続いて、ノード0において更新トランザクションのゴーストトランザクションへの移行の可否が判定される。図19(B)を参照すると、ノード0には更新のコミットに係る発行データ150cより時系列番号が小さい参照に係る発行データ150bが存在する。したがって、更新トランザクションをゴーストトランザクションに移行させる(S40)。

【0081】

ノード0は、発行された発行データ150cをコミットの通知及びゴーストトランザクションへの移行の可否の判定結果と共にノード1に送信する(S39)。この段階において、各ノードに管理されている発行データ150は、図19(B)の通りとなる。すなわち、(A)の発行データ150に加え、発行データ150cがノード0及びそれぞれのノードの更新トランザクションに追加されている。

30

【0082】

また、ノード1はコミットの通知及びゴーストトランザクションへの移行の可否の判定結果に応じ、アシスタントとしての更新トランザクションのコミット及びゴーストトランザクションへの移行を行う(S41)。なお、図中において、ノード番号が記載されていない発行データ150は、自ノードで発行された発行データ150であることを示す。

【0083】

続いて、ノード1において参照処理が終了すると、ノード1は、当該参照処理に対応する発行データ150aを削除する(S42)。また、ノード1は、当該参照処理に対してノード0で発行された発行データ150bの削除をノード0に要求する(S43)。当該要求に応じ、ノード0は、ノード0において管理されている発行データ150bを削除する。この段階において、各ノードに管理されている発行データ150は、図19(C)の通りとなる。発行データ150a及び150bの削除により、ノード0及び各ノードのゴーストトランザクションが、発行データ150cを保持している。

40

【0084】

続いて、ノード1及びノード0のそれぞれは、参照処理の発行データ150a及び150bの削除に応じ、ゴーストトランザクションの終了判定を行う。図19(C)より、ノード0には、発行データ150cより時系列番号が小さい参照の発行データ150は存在しない。したがって、ノード0が保持する発行データ150cを削除し、それぞれのゴーストトランザクションを終了させる。この段階において各ノード及び各種処理において管

50

理されている発行データ150は、図19(D)の通りとなる。

【0085】

図18及び図19を用いて説明した処理内容を一般化してフローチャートを用いて説明する。図20は、第二の実施の形態における参照処理の処理手順を説明するためのフローチャートである。同図は、第一の実施の形態における図12に対応する。なお、同図において、「チーフ」は、アプリケーションサーバ20より参照処理の要求を受け付けたノード(データベースサーバ10)を示す。「他ノード」は、チーフ以外の全ノードを示す。

【0086】

アプリケーションサーバ20よりデータの参照指示を受信したチーフのSQL実行制御部12からの指示に応じ、チーフのアクセス部13は、当該参照処理を実行するトランザクション(カレントトランザクション)に対して発行データ150(図9参照)を発行(メモリ装置103内に生成)する(S501)。なお、ここで生成される発行データ150の発行時期種別には参照の開始を示す値が登録される。続いて、チーフのアクセス部13は、当該参照処理に対応する発行データ150の発行を他ノードの各アクセス部13に対して要求する(S502)。他ノードの各アクセス部13は、当該要求を受けて、発行データ150を発行し、各ノードの発行データ格納領域に登録する(S503)。続いて、他ノードの各アクセス部13は、それぞれが発行した発行データ150をチーフのアクセス部13に対して送信する(S504)。チーフのアクセス部13は、他ノードからの発行データ150を受信する(S505)。これにより、他ノードの各発行データ150がチーフに収集される。

【0087】

続いて、チーフは、参照対象のデータ(以下、「参照データ」)の参照を開始する(S506)。まず、チーフは、図12の破線P1で囲まれた処理(S303~S308)と同様の処理手順によって参照対象のレコードを判定する(S507)。参照データがチーフのみに属する場合(S508でNO)、チーフのアクセス部13は、参照対象とされたレコードを参照する(S509)。一方、参照データが他ノード(アシスタント)に属し、参照処理をアシスタントに派生させる(依頼する)必要がある場合(S508でYES)、チーフのアクセス部13は、アシスタントのアクセス部13に対して参照処理を派生させる(S510)。続いて、アシスタントにおいて図12の破線P1で囲まれた部分(S303~S308)と同様の処理手順が実行され、参照対象のレコードが判定される(S511)。続いて、アシスタントのアクセス部13は、参照対象とされたレコードを参照する(S512)。

【0088】

ステップS509及びS512に続いて、チーフのアクセス部13は参照処理を終了させる(S513)。続いて、チーフのアクセス部13は、ステップS501において発行した発行データ150を削除する(S514)。続いて、チーフのアクセス部13は、他ノードのアクセス部13に対して発行データ150の削除を要求する(S515)。続いて、チーフは、図12の破線P2で囲まれた処理(S312~S318)と同様の処理手順によってチーフにおけるゴーストトランザクションの終了化処理を実行する(S516)。

【0089】

一方、発行データ150の削除要求を受信した他ノードの各アクセス部13は、ステップS503においてそれぞれのノードで発行された発行データ150を削除する(S517)。続いて、他ノードは、図12の破線P2で囲まれた処理(S312~S318)と同様の処理手順によって各他ノードにおけるゴーストトランザクションの終了化処理を実行する(S518)。

【0090】

続いて、図21は、第二の実施の形態におけるコミット処理の処理手順を説明するためのフローチャートである。同図は、第一の実施の形態における図8に対応する。

【0091】

10

20

30

40

50



ステップS601～S603では、チーフにおいて、図8のステップS201～S203と同様の処理手順が実行される。続いて、チーフは、図8の破線P3で囲まれた処理(S204～S208)と同様の処理手順によってレコードの削除処理を実行する(S604)。続いて、チーフの排他制御部14は、コミット対象のトランザクション(カレントトランザクション)によってロックされていたレコードに対してアンロック(ロックの解放)を実行する(S605)。

【0092】

一方、並列参照処理が存在する場合であって(S602でYES)、カレントトランザクションがアシスタントに依頼されている場合(S606でYES)、チーフのアクセス部13は、アシスタントにステップS601で発行された発行データ150とコミットの  
10 実行及びゴーストトランザクションへと移行することを通知する(S607、S608)。

【0093】

一方、並列参照処理が存在する場合であって(S602でYES)、カレントトランザクションがアシスタントに依頼されていない場合(S606でNO)、及びステップS608に続いて、チーフは、図8のステップS210、S211と同様の処理手順によって、ゴーストトランザクションの生成、カレントトランザクションのコミットを実行する(S609、S610)。

【0094】

上述したように、第二の実施の形態によれば、自ノードと他ノードとのシステム時間が異なる場合であっても他ノードから依頼された更新と自ノードの参照との前後関係を適切に把握することができる。  
20

【0095】

以上、本発明の実施例について詳述したが、本発明は斯かる特定の実施形態に限定されるものではなく、特許請求の範囲に記載された本発明の要旨の範囲内において、種々の変形・変更が可能である。

【0096】

以上の説明に関し、更に以下の項を開示する。

(付記1)

データベースにおける第一のレコードの参照の開始時に時間の前後関係を識別可能な第一の時系列データをメモリ装置に生成後、該第一のレコードを参照する参照開始手順と、  
30 前記第一のレコードの更新要求に応じ、前記第一のレコードに対応する第二のレコードを前記データベースに生成し、該第二のレコードを更新する更新手順と、  
前記更新のコミット時に第二の時系列データを前記メモリ装置に生成するコミット手順と、

前記参照の終了時に前記第一の時系列データを削除する参照終了手順と、

前記コミット手順又は前記参照終了手順に応じて前記第二の時系列データより前に生成された前記第一の時系列データが無くなった場合は前記第一のレコードを削除する削除手順とをコンピュータに実行させるためのデータ管理プログラム。  
40

(付記2)

前記第一及び第二の時系列データには、前記参照又は更新の対象とされたレコードの識別子が関連付けられ、

前記削除手順は、前記第二の時系列データより前に生成された前記第一の時系列データであって、当該第二の時系列データと同一の前記レコードの識別子に関連付けられている前記第一の時系列データが無くなった場合は前記第一のレコードを削除する付記1記載のデータ管理プログラム。

(付記3)

前記更新手順は、前記第一のレコード及び前記第二のレコードに対してロックを適用し、

前記コミット手順に応じて、前記第二の時系列データより前に生成された前記第一の時  
50

系列データが有る場合は前記ロックの内容を示す占有情報を前記メモリ装置に保持させる占有情報保持手順と、

前記メモリ装置に保持されている占有情報に基づいて前記データベースの排他制御を行う排他制御手順と、

前記参照終了手順に応じて前記第二の時系列データより前に生成された前記第一の時系列データが無くなった場合は前記占有情報を削除する占有情報削除手順とを有する付記 1 又は 2 記載のデータ管理プログラム。

(付記 4)

前記第二の時系列データより前に生成された前記第一の時系列データに係る参照については前記第一のレコードを参照対象とし、前記第二の時系列データより後に生成された前記第一の時系列データに係る参照については前記第二のレコードを参照対象とする参照手順を有する付記 1 乃至 3 いずれか一項記載のデータ管理プログラム。

10

(付記 5)

前記第一及び第二の時系列データには、複数ノード環境における各ノードを識別するノード識別子が関連付けられ、

前記参照開始手順は、他ノードに対して前記第一の時系列データの生成を要求し該要求に応じて該他ノードで生成された前記第一の時系列データを取得し、

前記更新手順は、前記他ノードからの更新要求に応じ、前記第一のレコードに対応する第二のレコードを前記データベースに生成し、該第二のレコードを更新し、

前記他ノードにおける前記更新に対するコミットに応じて、該他ノードにおいて生成された前記第二の時系列データを取得する他ノードコミット手順と、

20

前記他ノードコミット手順に応じて、前記他ノードより取得された前記第二の時系列データより前に生成された前記他ノードより取得された前記第一の時系列データが無くなった場合は前記第一のレコードを削除する第二の削除手順とを有する付記 1 乃至 4 いずれか一項記載のデータ管理プログラム。

(付記 6)

コンピュータが実行するデータ管理方法であって、

データベースにおける第一のレコードの参照の開始時に時間の前後関係を識別可能な第一の時系列データをメモリ装置に生成後、該第一のレコードを参照する参照開始手順と、

前記第一のレコードの更新要求に応じ、前記第一のレコードに対応する第二のレコードを前記データベースに生成し、該第二のレコードを更新する更新手順と、

30

前記更新のコミット時に第二の時系列データを前記メモリ装置に生成するコミット手順と、

前記参照の終了時に前記第一の時系列データを削除する参照終了手順と、

前記コミット手順又は前記参照終了手順に応じて前記第二の時系列データより前に生成された前記第一の時系列データが無くなった場合は前記第一のレコードを削除する削除手順とを有するデータ管理方法。

(付記 7)

前記第一及び第二の時系列データには、前記参照又は更新の対象とされたレコードの識別子が関連付けられ、

40

前記削除手順は、前記第二の時系列データより前に生成された前記第一の時系列データであって、当該第二の時系列データと同一の前記レコードの識別子に関連付けられている前記第一の時系列データが無くなった場合は前記第一のレコードを削除する付記 6 記載のデータ管理方法。

(付記 8)

前記更新手順は、前記第一のレコード及び前記第二のレコードに対してロックを適用し、

前記コミット手順に応じて、前記第二の時系列データより前に生成された前記第一の時系列データが有る場合は前記ロックの内容を示す占有情報を前記メモリ装置に保持させる占有情報保持手順と、

50

前記メモリ装置に保持されている占有情報に基づいて前記データベースの排他制御を行う排他制御手順と、

前記参照終了手順に応じて前記第二の時系列データより前に生成された前記第一の時系列データが無くなった場合は前記占有情報を削除する占有情報削除手順とを有する付記 6 又は 7 記載のデータ管理方法。

(付記 9)

前記第二の時系列データより前に生成された前記第一の時系列データに係る参照については前記第一のレコードを参照対象とし、前記第二の時系列データより後に生成された前記第一の時系列データに係る参照については前記第二のレコードを参照対象とする参照手順を有する付記 6 乃至 8 いずれか一項記載のデータ管理方法。

10

(付記 10)

前記第一及び第二の時系列データには、複数ノード環境における各ノードを識別するノード識別子が関連付けられ、

前記参照開始手順は、他ノードに対して前記第一の時系列データの生成を要求し該要求に応じて該他ノードで生成された前記第一の時系列データを取得し、

前記更新手順は、前記他ノードからの更新要求に応じ、前記第一のレコードに対応する第二のレコードを前記データベースに生成し、該第二のレコードを更新し、

前記他ノードにおける前記更新に対するコミットに応じて、該他ノードにおいて生成された前記第二の時系列データを取得する他ノードコミット手順と、

前記他ノードコミット手順に応じて、前記他ノードより取得された前記第二の時系列データより前に生成された前記他ノードより取得された前記第一の時系列データが無くなった場合は前記第一のレコードを削除する第二の削除手順とを有する付記 6 乃至 9 いずれか一項記載のデータ管理方法。

20

(付記 11)

データベースにおける第一のレコードの参照の開始時に時間の前後関係を識別可能な第一の時系列データをメモリ装置に生成後、該第一のレコードを参照する参照開始手段と、

前記第一のレコードの更新要求に応じ、前記第一のレコードに対応する第二のレコードを前記データベースに生成し、該第二のレコードを更新する更新手段と、

前記更新のコミット時に第二の時系列データを前記メモリ装置に生成するコミット手段と、

30

前記参照の終了時に前記第一の時系列データを削除する参照終了手段と、

前記コミット手段又は前記参照終了手段に応じて前記第二の時系列データより前に生成された前記第一の時系列データが無くなった場合は前記第一のレコードを削除する削除手段とを有するデータ管理装置。

(付記 12)

前記第一及び第二の時系列データには、前記参照又は更新の対象とされたレコードの識別子が関連付けられ、

前記削除手段は、前記第二の時系列データより前に生成された前記第一の時系列データであって、当該第二の時系列データと同一の前記レコードの識別子に関連付けられている前記第一の時系列データが無くなった場合は前記第一のレコードを削除する付記 11 記載のデータ管理装置。

40

(付記 13)

前記更新手段は、前記第一のレコード及び前記第二のレコードに対してロックを適用し、

前記コミット手段に応じて、前記第二の時系列データより前に生成された前記第一の時系列データが有る場合は前記ロックの内容を示す占有情報を前記メモリ装置に保持させる占有情報保持手段と、

前記メモリ装置に保持されている占有情報に基づいて前記データベースの排他制御を行う排他制御手段と、

前記参照終了手段に応じて前記第二の時系列データより前に生成された前記第一の時系

50

列データが無くなった場合は前記占有情報を削除する占有情報削除手段とを有する付記 1 1 又は 1 2 記載のデータ管理装置。

(付記 1 4)

前記第二の時系列データより前に生成された前記第一の時系列データに係る参照については前記第一のレコードを参照対象とし、前記第二の時系列データより後に生成された前記第一の時系列データに係る参照については前記第二のレコードを参照対象とする参照手段を有する付記 1 1 乃至 1 3 いずれか一項記載のデータ管理装置。

(付記 1 5)

前記第一及び第二の時系列データには、複数ノード環境における各ノードを識別するノード識別子が関連付けられ、

前記参照開始手段は、他ノードに対して前記第一の時系列データの生成を要求し該要求に応じて該他ノードで生成された前記第一の時系列データを取得し、

前記更新手段は、前記他ノードからの更新要求に応じ、前記第一のレコードに対応する第二のレコードを前記データベースに生成し、該第二のレコードを更新し、

前記他ノードにおける前記更新に対するコミットに応じて、該他ノードにおいて生成された前記第二の時系列データを取得する他ノードコミット手段と、

前記他ノードコミット手段に応じて、前記他ノードより取得された前記第二の時系列データより前に生成された前記他ノードより取得された前記第一の時系列データが無くなった場合は前記第一のレコードを削除する第二の削除手段とを有する付記 1 1 乃至 1 4 いずれか一項記載のデータ管理装置。

【図面の簡単な説明】

【0097】

【図 1】第一の実施の形態におけるデータ管理システムの構成例を示す図である。

【図 2】本発明の実施の形態におけるデータベースサーバのハードウェア構成例を示す図である。

【図 3】本発明の実施の形態におけるデータベースのデータ構造の例を示す図である。

【図 4】挿入トランザクションの処理手順を説明するためのフローチャートである。

【図 5】更新トランザクションの処理手順を説明するためのフローチャートである。

【図 6】削除トランザクションの処理手順を説明するためのフローチャートである。

【図 7】操作対象とされたレコードのロック状態を示す図である。

【図 8】コミット処理の処理手順を説明するためのフローチャートである。

【図 9】発行データの構成例を示す図である。

【図 10】発行データに基づく並列参照処理の存否の判定内容を説明するための図である。

【図 11】ゴーストトランザクションの発生契機の実例を示す図である。

【図 12】参照処理の処理手順を説明するためのフローチャートである。

【図 13】ゴーストトランザクションの最後の要因となった全ての参照処理の終了の判定内容を説明するための図である。

【図 14】ゴーストトランザクションによって占有されているデータの参照対象レコードの判定処理を説明するためのフローチャートである。

【図 15】ゴーストトランザクションによって占有されているデータの参照対象レコードの判定処理の実例を示す図である。

【図 16】第二の実施の形態におけるデータ管理システムの構成例を示す図である。

【図 17】第二の実施の形態におけるデータ管理システムの動作概要を説明するための図である。

【図 18】第二の実施の形態の処理概要を説明するための図である。

【図 19】第二の実施の形態の処理概要において各ノードで管理される発行データの遷移を示す図である。

【図 20】第二の実施の形態における参照処理の処理手順を説明するためのフローチャートである。

10

20

30

40

50

【図21】第二の実施の形態におけるコミット処理の処理手順を説明するためのフローチャートである。

【符号の説明】

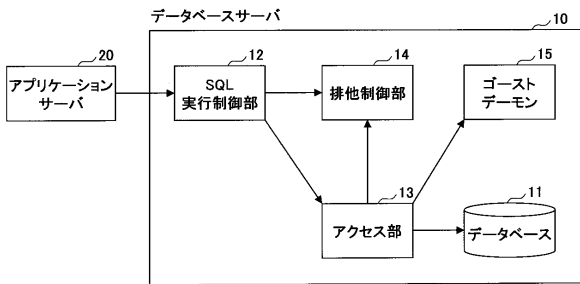
【0098】

- 10、10a、10b データベースサーバ
- 11、11a、11b データベース
- 12、12a、12b SQL実行制御部
- 13 アクセス部
- 14 排他制御部
- 15 ゴーストデーモン
- 20 アプリケーションサーバ
- 100 ドライブ装置
- 101 記録媒体
- 102 補助記憶装置
- 103 メモリ装置
- 104 CPU
- 105 インタフェース装置
- B バス

10

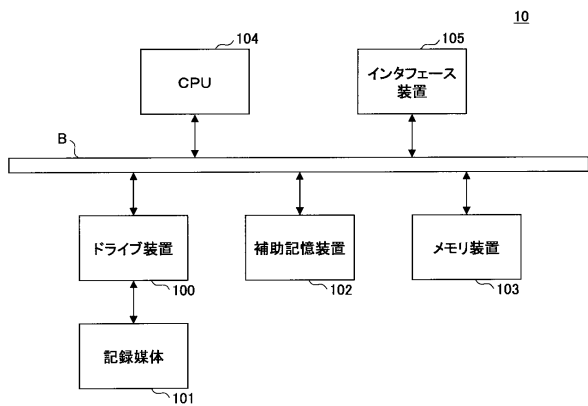
【図1】

第一の実施の形態におけるデータ管理システムの構成例を示す図



【図2】

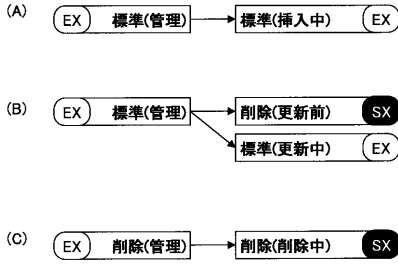
本発明の実施の形態におけるデータベースサーバのハードウェア構成例を示す図





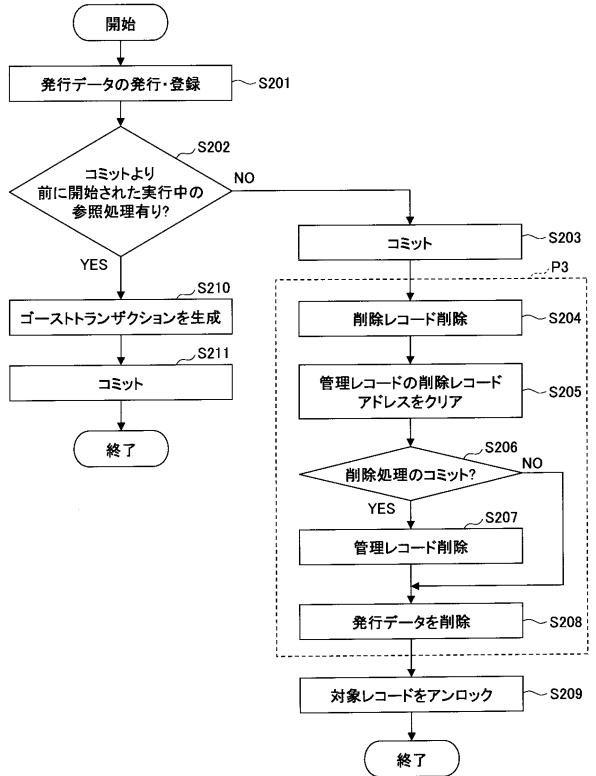
【図7】

操作対象とされたレコードのロック状態を示す図



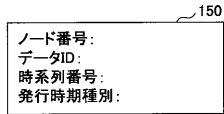
【図8】

コミット処理の処理手順を説明するためのフローチャート



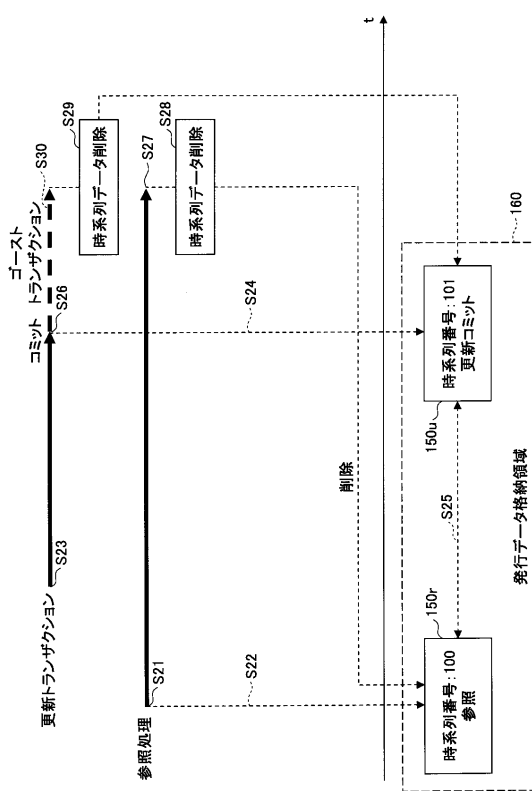
【図9】

発行データの構成例を示す図



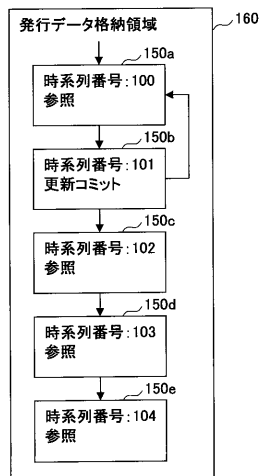
【図11】

ゴーストランザクションの発生の契機を具体例を示す図

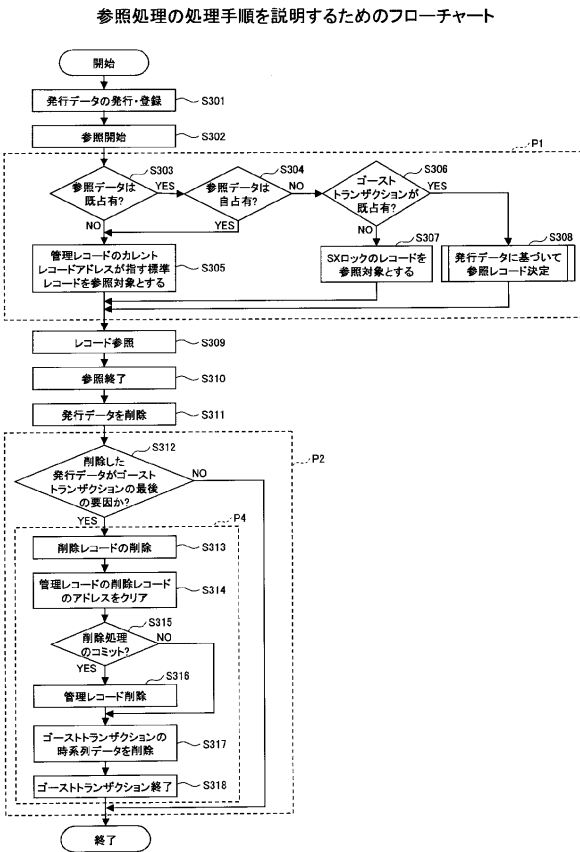


【図10】

発行データに基づく並列参照処理の存否の判定内容を説明するための図

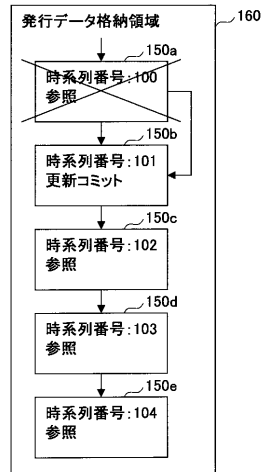


【図12】



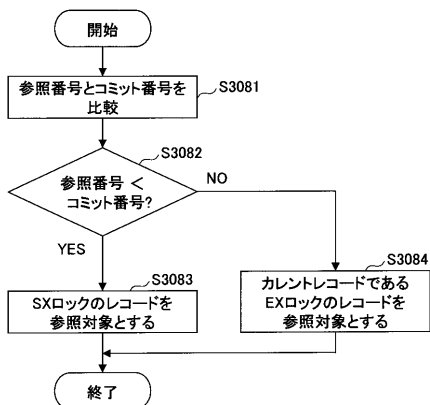
【図13】

ゴーストトランザクションの最後の要因となった全ての参照処理の終了の判定内容を説明するための図



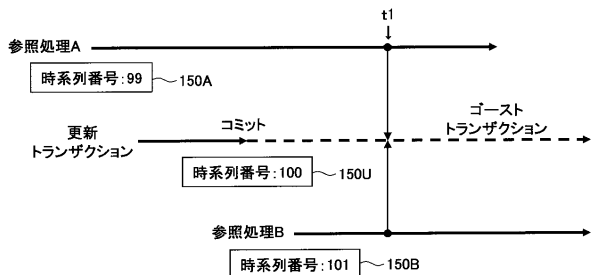
【図14】

ゴーストトランザクションによって占有されているデータの参照対象レコードの判定処理を説明するためのフローチャート



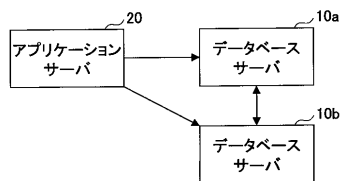
【図15】

ゴーストトランザクションによって占有されているデータの参照対象レコードの判定処理の具体例を示す図



【図16】

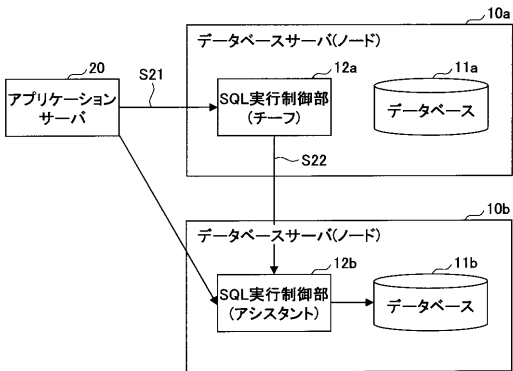
第二の実施の形態におけるデータ管理システムの構成例を示す図





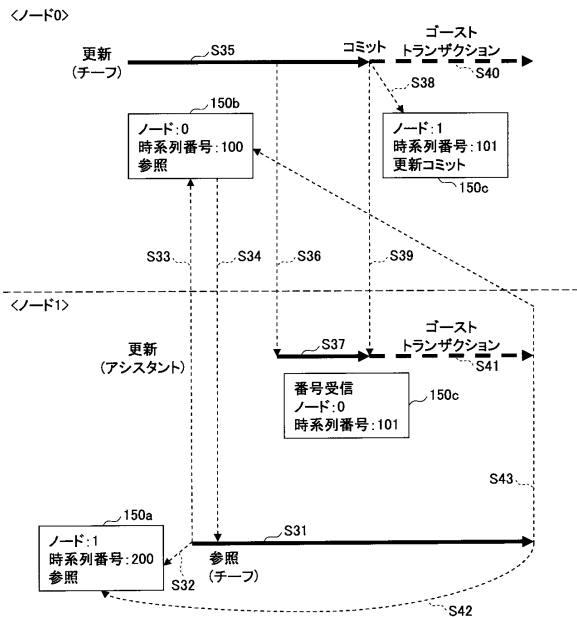
【図17】

第二の実施の形態におけるデータ管理システムの動作概要を説明するための図



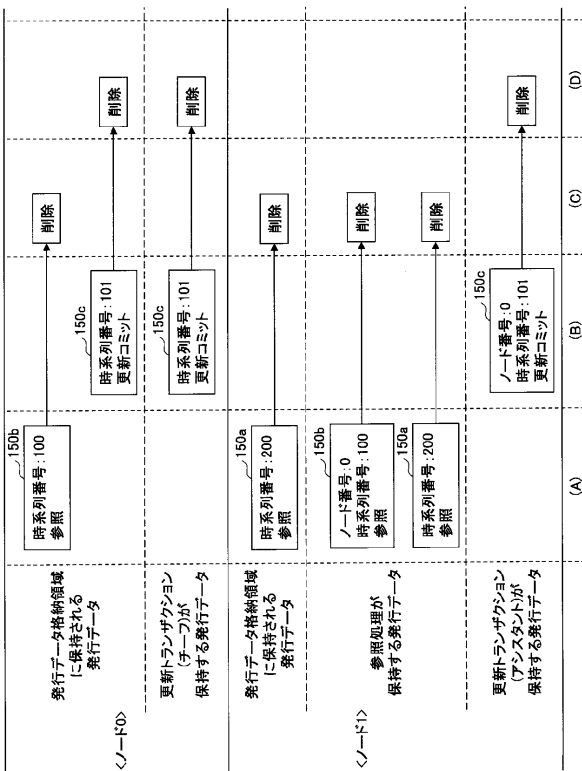
【図18】

第二の実施の形態の処理概要を説明するための図



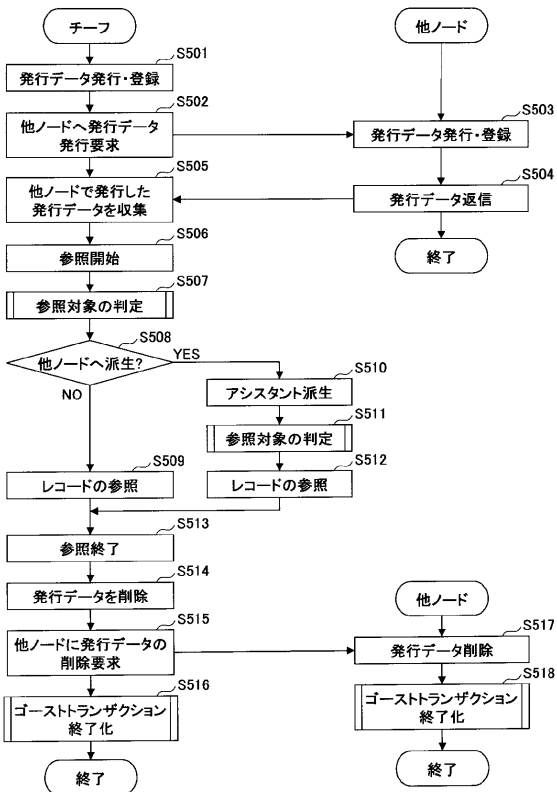
【図19】

第二の実施の形態の処理概要において各ノードで管理される発行データの遷移を示す図



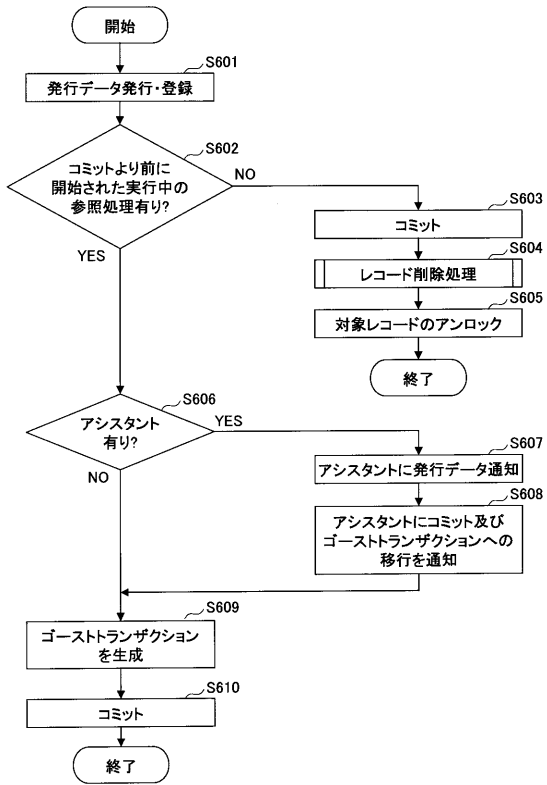
【図20】

第二の実施の形態における参照処理の処理手順を説明するためのフローチャート



【図 2 1】

第二の実施の形態における  
コミット処理の処理手順を説明するためのフローチャート



---

フロントページの続き

(72)発明者 丹羽 雅史

神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

審査官 田川 泰宏

(56)参考文献 特表2007-501468(JP,A)

特開平08-235046(JP,A)

特開2005-234945(JP,A)

特開2003-271436(JP,A)

片岡 良治, 部分更新と全数検索の混在処理に適した多版並行処理制御方式, 電子情報通信学会論文誌, 社団法人電子情報通信学会, 1991年 3月25日, 第J74-D-I巻 第3号, p. 224-231

國枝 和雄, 適応型時刻印方式に基づく同時実行制御方式, 情報処理学会論文誌, 社団法人情報処理学会, 1992年 6月15日, 第33巻 第6号, p. 802-811

鈴木 幸市, 第8回 トランザクション制御とMVCC, DB Magazine, 株式会社翔泳社, 2004年 7月 1日, 第14巻 第4号, p. 199-203

(58)調査した分野(Int.Cl., DB名)

G06F 12/00