

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3593366号
(P3593366)

(45) 発行日 平成16年11月24日(2004.11.24)

(24) 登録日 平成16年9月3日(2004.9.3)

(51) Int. Cl.⁷

F I

G06F 17/30

G06F 17/30 419B

G06F 12/00

G06F 12/00 531R

請求項の数 3 (全 23 頁)

| | |
|---|--|
| <p>(21) 出願番号 特願平6-222930 (22) 出願日 平成6年9月19日(1994.9.19) (65) 公開番号 特開平8-87510 (43) 公開日 平成8年4月2日(1996.4.2) 審査請求日 平成13年3月15日(2001.3.15)</p> | <p>(73) 特許権者 000005108 株式会社日立製作所 東京都千代田区神田駿河台四丁目6番地 (73) 特許権者 000233055 日立ソフトウェアエンジニアリング株式会社 神奈川県横浜市鶴見区末広町一丁目1番4 3 (74) 代理人 100075096 弁理士 作田 康夫 (74) 代理人 100068504 弁理士 小川 勝男</p> |
|---|--|

最終頁に続く

(54) 【発明の名称】 データベース管理方法

(57) 【特許請求の範囲】

【請求項1】

障害発生時の再開始回復処理を前提としたトランザクション処理を可能とするデータベース管理システムにおけるシンクポイントを取得するデータベース管理方法であって、前記シンクポイント取得開始時に外部記憶装置上のデータベースに書き込みされていない主記憶装置上にマッピングされたバッファ中のすべての更新されたページを管理するテーブルにシンクポイント取得処理中であるマークをつけ、前記シンクポイント取得処理中は前記マークのついたバッファ中の更新されたページを前記データベースに書き込み、前記シンクポイント取得処理中に前記マークのついたバッファ中の更新されたページを更新するトランザクションは前記更新されたページをアクセスする前に前記データベースに書き込みを行い、前記マークをはずした後アクセスを許可し、前記マークのついた更新ページをすべて前記データベースに書き込みが完了した時点シンクポイント取得完了とし、前記シンクポイントを取得する間隔がある一定量に達した時点より、前記シンクポイント取得開始時点まで前記バッファ中の更新ページを非同期に前記データベースに書き込みを行うプレシンクポイント取得処理を行うことにより、前記シンクポイント処理中であってもトランザクション処理続行可能とするとともに前記シンクポイント取得処理における前記バッファ中の更新ページ数を削減することを特徴とするデータベース管理方法。

10

20

【請求項 2】

請求項 1 記載のデータベース管理方法において、
前記プレシンクポイント取得処理において、前記シンクポイント取得開始時点までに前記データベースに書き込めなかった前記バッファ中の更新ページは、前記シンクポイント取得開始時にバッファ中のすべての更新されたページを管理するテーブルにシンクポイント取得処理中であるマークをつけることによりシンクポイントを取得することを特徴とするデータベース管理方法。

【請求項 3】

請求項 2 記載のデータベース管理方法において、
前記プレシンクポイント取得処理において、前記シンクポイント取得開始時にプレシンクポイント取得が完了したか否かを記憶しておき、前記プレシンクポイント取得が完了していない場合は前記データベースに書き込めなかった前記バッファ中の更新ページ数から前記プレシンクポイント取得開始点を早くするように変更し、前記プレシンクポイント取得がすでに完了している場合は、前記シンクポイント取得時に存在した更新ページ数から前記プレシンクポイント取得開始点を遅くするように変更し、前記プレシンクポイント取得開始点を補正することを特徴とするデータベース管理方法。

10

【発明の詳細な説明】**【0001】****【産業上の利用分野】**

本発明は、複数のトランザクションからアクセスされるデータを記憶するデータベースを管理するデータベース管理システムにおける、データベースのシンクポイント処理を行う方法に関する。

20

【0002】**【従来の技術】**

多重処理環境におけるデータベース管理システムでは、一定間隔おきにシンクポイントあるいはチェックポイントという確認点を設けている。これは、障害が発生しデータベースの再開回復処理を実行しなければならない場合、障害が発生した時点の直前のシンクポイントの状態にデータベースを復元し、そこから処理の再開をはかることによってデータベースの再開回復処理を行う。シンクポイントは、更新を含むトランザクション処理の進行とは非同期に取得することができる。

30

【0003】

データベース管理システム、特にリレーショナル・データベース管理システムは、データベース、スキーマ、表、行といったように論理的に階層的な資源で構成される。物理的に行を格納するための不揮発性記憶装置である磁気ディスクなどの場合、行を物理的に固定長のブロックあるいはページという単位に格納する。トランザクション処理において行われるデータベースアクセスは、このページを単位として主記憶装置上に確保されたバッファプールを介して行われる。

【0004】

あるトランザクションにおいて、必要なページがバッファプールに読み込まれ、バッファプール内で検索、更新処理が行われる。バッファプールに読み込まれたページは、できるだけ長くバッファプール上に存在するように管理されることによって、磁気ディスクなどの外部記憶装置との入出力処理オーバーヘッドを削減している。あるトランザクションにおいて行われたデータベースの更新処理については、その更新履歴情報をログとして取得する。ログは外部記憶装置のような不揮発性の記憶装置に永久的に記録される。ログレコードはプロセッサ内の主記憶装置上の揮発性のログバッファにまず書き込まれ、トランザクションが確定されるときに外部記憶装置上のログファイルに移される。この場合、トランザクションにおいて行われたバッファプール上のページの更新は、いずれ外部記憶装置上のデータベースに反映されなければならない。

40

【0005】

トランザクション確定時に、当該トランザクションによって更新されたバッファプール上

50

の全てのページを外部記憶装置上のデータベースに反映する方法がある。この方法では、トランザクション確定時の外部記憶装置への書き込みオーバーヘッドがかかる。そこで、もう一つの方法では、トランザクション確定時とは遅延させてバッファプールより更新されたページを外部記憶装置上のデータベースに反映することによって、同一ページへの複数のトランザクションの更新をまとめて外部記憶装置への書き込みを行うことで、トランザクション確定時の外部記憶装置への書き込みオーバーヘッドを削減するものである。

【 0 0 0 6 】

トランザクション確定時とよりも遅延させて外部記憶装置への書き込みを行う場合、外部記憶装置に更新されたページが書き込まれる前に障害等によりデータベース管理システムが終了してしまうとそのページの情報は失われる。この場合、システムの再開始が行われ、システムが終了した時点までに行われたトランザクションの回復を保証する必要がある。トランザクションが確定される前に障害が生じた場合は、そのトランザクションで行われたすべての更新を確実に元に戻す。また、トランザクション確定後に障害が生じた場合は、そのトランザクションにて行われた更新をすべて保証されなければならない。

10

【 0 0 0 7 】

システムの再開始時には、前回のシンクポイントより行われた全てのトランザクションを外部記憶装置上に保存されたログファイルからログレコードを読み出して回復する。ログレコードによるトランザクション回復時には、確定されたトランザクションであっても実際に更新したページについては、バッファプール上から外部記憶装置へ反映されていないことがある。その場合は、外部記憶装置上のデータベースから回復対象のページを再度バッファプール上に読み込み、ログレコードを基に更新処理を再実行 (R E D O) する。

20

【 0 0 0 8 】

逆に、確定されたトランザクションおよび確定されていないトランザクションで行われた更新がすでにバッファプールから外部記憶装置へ反映されてしまっている場合もある。これに対応する対策としては、ページ中に更新が発生したときのログのログシーケンス番号 (L S N) を持つておくことによって解決される。ページに対してあるトランザクションによって更新が行われると、そのときのログを取得するがそのログレコードの L S N を合わせてページ内に書き込んでおく。そうすると、バッファプールから更新されたページが外部記憶装置に書き込まれることによって、そのトランザクションで行った更新がどのログに対応する更新まで外部記憶装置に反映されているかがわかる。これにより、システムの再開始時に回復対象の更新ページすべてを回復する必要がなくなる。

30

【 0 0 0 9 】

更新対象ページをバッファプールに読み込んだ場合、読み込んだページの L S N とログレコードの L S N を比較し、ページの L S N の方がログレコードの L S N より大きいときは、そのログレコードに記録された更新情報はすでに反映されていることになる。ということは、ページの L S N の方がログレコードの L S N より小さいときだけ、ログを使用して R E D O を行えばよいということになる。データベース管理システムでは、こうした再開始時の処理時間をできる限り短時間で行うことが望ましい。そのためには、再開始処理時の R E D O 処理を極力削減することが望まれる。それには、通常、バッファプール上で更新されたページを適切な時期にシステムプロセスを使用して、バックグラウンドで外部記憶装置上のデータベースに書き込む。

40

【 0 0 1 0 】

ここで、再開始処理において直前のシンクポイントから回復処理を行う場合に前提とされるのは、直前のシンクポイント時点のデータベースの状態からログを用いて回復が行えることを保証することである。そのために、シンクポイントを取得する処理では、まず、ログレコードとしてシンクポイント開始ログを取得することにより開始される。この後、トランザクションの進行を停止させ、ログバッファ上のログレコードを外部記憶装置上のログファイルに書き込む。

【 0 0 1 1 】

ログの書き込みが終了すると、次はバッファプール上の全ての更新されたページを外部記

50

憶装置上のデータベースに書き込む。これが終わると、シンクポイント取得完了ログをログバッファに書き込んで、停止したトランザクション処理の続行を開始してシンクポイント取得を完了する。これによって、次のシンクポイントが取得されるまでに障害等によりシステムが終了しても、今、取得したシンクポイントのデータベースの状態が保証されたので、取得したシンクポイントを起点にして再開始処理を行うことができる。

【0012】

このシンクポイント取得処理では、説明したようにシンクポイント取得が開始されるとトランザクションの進行が一時的に妨げられるので、シンクポイント取得処理中はトランザクションのスループットが一時的に沈み込むという問題がある。シンクポイント取得処理中であっても、トランザクション処理を続行できる方法としては、シンクポイント取得処理においてバッファプール中のすべての更新ページ番号と各ページのLSNをログに出力しておくことによって、シンクポイント取得を完了する方法がある。この方法では、システム再開始時に直前のシンクポイント時に取得した更新ページ番号のログから最も古いLSN探索する。そして、最も古いLSNを持つログレコードを再開始の回復起点とするものである。しかし、シンクポイント取得時にバッファプールから長い間、外部記憶装置上のデータベースに書き込まれていないページがあると、REDO対象となるログレコード数が増加することになるため、再開始処理時間を最小にするための工夫が必要となる。

10

【0013】

これらの従来技術は、C. Mohan, Don Haderle等の"ARIES: A Transaction Recovery Method Supporting Fine-Granularity Locking and Partial Rollbacks Using Write-Ahead Logging", ACM Transactions on Database Systems、第17巻、第1号、pp. 94-162(1992年3月)に詳細に説明されている。さらに、特開平5-6297号公報に開示されているトランザクション処理方法およびシステムにも詳細に説明されている。

20

【0014】

【発明が解決しようとする課題】

本発明の目的は、シンクポイント取得処理が発生してもトランザクション処理のスループットを低下させないようにすることにある。システム再開始時における再開始処理時間の削減するための方法を提供することが関連した目的である。特に、シンクポイント取得処理の開始時点においてバッファプール中に存在した更新ページとそれ以外の非更新ページのアクセスをシンクポイント取得中でも同時にアクセスするための方法を提供することである。

30

【0015】

【課題を解決するための手段】

上記目的を達成するために、本発明のデータベース管理方法では、障害発生時の再開始回復処理を前提としたトランザクション処理を可能とするデータベース管理システムにおけるシンクポイントを取得する方法において、シンクポイント取得開始時に外部記憶装置上のデータベースに書き込みされていない主記憶装置上にマッピングされたバッファ中のすべての更新されたページを管理するテーブルにシンクポイント取得処理中であるマークをつける。

40

【0016】

シンクポイント取得処理中はマークのついたバッファ中の更新されたページをデータベースに書き込み、シンクポイント取得処理中にマークのついたバッファ中の更新されたページを更新するトランザクションは更新されたページをアクセスする前にデータベースに書き込みを行い、マークをはずした後アクセスを許可し、マークのついた更新ページをすべてデータベースに書き込みが完了した時点シンクポイント取得完了とすることにより、シンクポイント処理中であってもトランザクション処理続行可能とすることを特徴とするデータベース管理方法によって解決される。

50

【 0 0 1 7 】

【作用】

次に、第一の目的を達成するため、本発明のデータベース管理方法は、複数のトランザクションが、各々データベースをアクセスしてバッファプール上に読み込んだページに対して更新処理を行っている場合、ある時点でデータベースのシンクポイント取得時点に達すると、バッファプール上の更新されたページに対応するバッファを管理するテーブルにシンクポイント取得中であるマークをつけることにより、バックグラウンドでデータベースへの書き込みを行うプロセスは、該バッファプール上でシンクポイント取得中であるマークのついたバッファに対応するページを対象に外部記憶装置上のデータベースへ書き込みを行うようにする。これによって、シンクポイント取得時にログファイルへの書き込みとログに対応した更新ページのデータベースへの書き込みを同期することができる。

10

【 0 0 1 8 】

シンクポイント取得中に、該マークのついたバッファに格納されたページに対するアクセス要求は、一旦バッファプール上から外部記憶装置上のデータベースに対して書き込みを行ってから、アクセスを許可するようにする。これによって、シンクポイント取得中であっても、トランザクションを停止させることなく、システムとしてのスループットを低下させることがない。

【 0 0 1 9 】

【実施例】

以下、本発明の一実施例を図面を用いて具体的に説明する。

20

図2は、本発明の一実施例にかかるコンピュータシステム10の構成を示す。コンピュータシステム10は、CPU12、主記憶装置14、磁気ディスク等の外部記憶装置16及び多数の端末18で構成される。主記憶装置2上には、データベース管理システム20が置かれ、外部記憶装置16上にはデータベース管理システム20が管理するデータベース36a、36b及びデータベース更新に関する更新履歴情報であるログファイル37が格納される。

【 0 0 2 0 】

データベース管理システム20は、ユーザからのデータベース問い合わせ要求であるSQLを受け取り、構文解析意味解析処理を通してデータベースアクセスの最適なアクセス経路を決定する最適化処理を行ない、決定したアクセス経路に基づいてデータベース処理用の内部処理コードを生成する質問解析処理部21、生成された内部処理コードを基にデータベース・アクセスを行うデータベース演算処理部22、外部記憶装置16に格納されたデータベース16からデータを主記憶装置14上に確保したバッファプールとの間でデータの入出力を行うバッファプール管理部23、前記端末18から投入されたトランザクションを管理するトランザクション管理部24、バッファプール34の上で行われたデータの更新をトランザクションとは非同期に外部記憶装置16上のデータベース36に書き込みを行う遅延書き込み処理部25、定期的にデータベースを整合性が取れた状態にすることを保証する処理を行うシンクポイント取得処理部26及びトランザクションで行われたデータベースへの更新履歴情報であるログ(以下、ログと略す)を管理し、各ログレコードのログシーケンス番号(以下、LSNと略す)を割り振り、トランザクションの完了時に主記憶装置14上のログバッファ35から外部記憶装置16上のログファイル37に書き込みを行うログ書き込み処理部27で構成される。前記バッファプール管理部23は、主記憶装置14上に確保したバッファプール34をデータベースのデータを蓄積する単位である物理適に固定長のページと対応させて管理するためのバッファプール管理テーブル31、ページ管理テーブル32及びハッシュ管理テーブル33を使用する。

30

40

【 0 0 2 1 】

図3は、先に述べたデータベースの入出力単位であるページ40の構造を示す。データベース管理システム、特にリレーショナルデータベース管理システムではデータは複数の行からなるリレーションと呼ばれる表形式を成している。リレーションは、データベース管理システム20の入出力単位である物理的に固定長の複数のページに分割され、外部記憶

50

装置 16 のデータベース 36 に格納されている。ページ 40 には行が格納され、ページ内の状態はページ内制御情報 42 で管理され、格納された行 1 ~ 行 6 は各々ページ内の相対アドレスを持っているスロット領域 44 によりポイントされている。

【 0022 】

図 4 は、図 3 に示したページ内制御情報 42 に含まれる管理情報を示す。ページ内制御情報 42 には、ページ内に割り当てられたスロット数である割当スロット数 421、実際に格納されている行に使用されたスロット数である使用スロット数 422、ページ内の空き領域長を管理する空き領域長 423、ページ内の未使用領域の先頭を相対アドレスで管理する未使用領域先頭オフセット 424 及びページ内で行われた更新を反映したログの LSN 425 で構成される。この LSN 425 は、データベース 36 からバッファプール 34 に読み込まれた時に設定されている LSN 425 が実際にその LSN 425 までのログがすべて反映されていることを示す。逆に、バッファプール 34 上で行われたページ上の行に対する更新が行われると、その更新に対応したログの LSN が書き込まれ、バッファプール 34 からデータベース 36 に書き込まれた時には、その LSN が示す更新はデータベースに反映したことが保証される。もちろん、バッファプール 34 からデータベース 36 に書き込まれる前に、一般的に良く知られる WAL (Write Ahead Logging) プロトコルによりログがログファイル 37 に書き込まれることは言うまでもない。

10

【 0023 】

図 5 にバッファプール管理部 23 がデータベース 36 との入出力を行うためのバッファプール 34 を管理する各テーブルの関連を示す。バッファプール管理テーブル 31 は、ハッシュ管理テーブル 33 及びページ管理テーブル 32 を管理している。バッファプール管理テーブル 31 は、図 6 に示すような情報を持つ。バッファプール 34 が持つページを格納するバッファ面数 311、ハッシュ管理テーブル 33 へのポインタ 312、できるだけバッファプール上にページを留めるように管理するための手法である LRU (Least Recently Used) アルゴリズムに基づいた LRU チェインの最も新しいページを管理するページ管理テーブル 32 へのポインタ 313、逆に最も使用されていない古いページを管理するページ管理テーブルへのポインタ 314、前記 LRU チェインとは別に更新されたページについての LRU チェインを管理するための最も新しいページを管理するページ管理テーブル 32 へのポインタ 315、逆に、最も使用されていない古いページを管理するページ管理テーブルへのポインタ 316、ページ管理テーブル 32 が管理するページに対応したバッファが無効状態である。すなわち、FREE 状態であるページ管理テーブルのポインタ 317、バッファプール 34 上に存在する更新ページ数 318 及び遅延書き込み状態を示すフラグ 319 から成る。

20

30

【 0024 】

次に、図 6 にページ管理テーブル 32 の情報について示す。ページ管理テーブル 32 には、記憶したページのオブジェクト識別子 321、バッファ番号 322、LRU チェイン又は FREE チェインに接続された場合の前方ポインタ 323 と後方ポインタ 324、ハッシュチェインに接続された場合の前方ポインタ 325 と後方ポインタ 326、ページ番号 327、ロックカウンタ 328、排他モードでページを入力した状態を示す XGET フラグ 329、出力要求フラグ 330、データベースからの読み込み中である状態を示す入力中フラグ 331、データベースへの書き込み中である状態を示す出力中フラグ 332、データベースからの入出力待ち状態である場合の入出力待ちチェイン 333、ロック待ちチェイン 334、及び本発明の特長であるシンクポイント取得時の出力対象ページとなったことを示すシンクポイント取得中フラグ 335 で構成される。

40

【 0025 】

図 5 を用いて、各テーブルの関連について説明する。バッファプール管理テーブル 31 からハッシュ管理テーブル 33 へは、ハッシュ管理テーブルポインタ 312 でポイントされ、ハッシュ管理テーブルはデータベースを構成するページ番号からハッシュングすることによって該当するハッシュエントリにそのページを管理するページ管理テーブル 32 が接続

50

される。例えば、ハッシュエントリ 1 には同じハッシュ値を持つページ管理テーブル 3 2 a 及び 3 2 b が接続される。バッファプール管理テーブル 3 1 の L R U チェインにはページ管理テーブル 3 2 上の出力要求フラグ 3 3 0 が出力要求状態でないものがチェーンされ、ページ管理テーブル 3 2 c、及び 3 2 f がチェーンされている。バッファプール管理テーブル 3 1 の更新 L R U チェインにはページ管理テーブル 3 2 上の出力要求フラグ 3 3 0 が出力要求状態のものがチェーンされ、ページ管理テーブル 3 2 d、及び 3 2 h がチェーンされている。また、F R E E チェインにはページ管理テーブル 3 2 i、及び 3 2 j がチェーンされており、バッファプール管理部 2 3 はバッファプール上に存在しないページを入力する場合、この F R E E チェイン 3 1 7 のページ管理テーブル 3 2 を割り当て、当該ページ管理テーブルに対応したバッファにページを入力する。

10

【 0 0 2 6 】

ここで、本発明に関わるトランザクション処理について図 8 から図 1 1 を用いて説明する。図 2 における端末 1 8 からトランザクションが投入されるとトランザクション管理部 2 4 はトランザクション開始処理を行う。トランザクション開始処理の概略処理フローを図 8 に示す。トランザクションが開始されると、まず、トランザクション管理テーブルを当該トランザクション用に割り当て、初期化する(ステップ 5 2)。当該テーブルにはトランザクション番号、トランザクション開始ログの L S N 及びカレントな当該トランザクションの L S N などが格納される。

【 0 0 2 7 】

次に、トランザクション開始ログを作成し、ログ書き込み処理部 2 7 にログの出力要求を発行し、主記憶装置 1 4 上のログバッファ 3 5 へ出力し、データベース管理システム内の一貫した L S N を割り振り、L S N を取得する(ステップ 5 4)。取得した L S N を当該トランザクションのトランザクション管理テーブルのトランザクション開始ログとして登録する(ステップ 5 6)。こうして、トランザクション開始処理が完了すると、データベースに対する質問処理が実行される。データベースに対する質問処理のうち、更新処理が行われた場合の概略処理フローを図 9 に示す。当該データベース管理システムが行レベルの排他制御を行うデータベース管理システムの場合、更新対象となった行に排他モードでロックを取得する(ステップ 6 1)。更新対象行にロックを確保できたならば、バッファプール管理部 2 3 に対して更新対象行が記憶されているページを排他モードでページ入力要求を発行する(ステップ 6 2)。

20

30

【 0 0 2 8 】

そして、更新対象行への更新操作を行う前に更新操作を回復するための U N D O ログを作成し、ログ書き込み処理部 2 7 に対し書き込み要求を発行するとともに U N D O ログの L S N を取得する(ステップ 6 3)。U N D O ログが取得できたら、行の更新操作を行う(ステップ 6 4)。そして、更新した行に対する再実行による回復のための R E D O ログを出力し、R E D O ログの L S N を取得する(ステップ 6 5)。取得した R E D O ログの L S N をページ内制御情報 4 2 の L S N 4 2 5 に書き込み、更新する(ステップ 6 6)。R E D O ログの L S N を更新した後、当該ページに対する更新をデータベースに反映するため、バッファプール管理部 2 3 に対してページ出力要求を行う(ステップ 6 8)。

【 0 0 2 9 】

ただし、当該要求は直ちにデータベース 1 6 に書き込まれるのではなく、当該ページを管理するページ管理テーブル 3 2 の出力要求フラグ 3 3 0 を出力要求状態に設定するのみとする。出力要求状態に設定することで、当該バッファプールの置換アルゴリズムにしたがって、バッファスチール(他のページを入力するためにデータベースに強制的に書きだす)されるか、遅延書き込み処理部 2 5 によって非同期にデータベース 1 6 に書き込まれるまで、主記憶装置 1 4 上のバッファプール 3 4 上に置かれる。次に、トランザクション内で行ったデータベースへのアクセス処理が終わるとトランザクションを有効なものとするために 2 フェーズに分けてトランザクションを完結させる。

40

【 0 0 3 0 】

まず、図 1 0 にトランザクションの P R E P A R E 処理フェーズについて説明する。まず

50

、トランザクションのPREPAREログをログバッファに書き込む(ステップ72)。次いで、当該PREPAREログをログファイル37に強制出力する(ステップ74)。PREPAREフェーズが終わるとトランザクションを正常に終了させるためにCOMMITフェーズに入る。COMMITフェーズの処理概略フローを図11に示す。トランザクションCOMMITログをログバッファに出力する(ステップ82)。COMMITログをログファイル37に強制出力する。これにより、トランザクションが確約できたので当該トランザクションで取得したすべての排他資源に対するロックを解放する(ステップ84)。COMMITされたトランザクションは、トランザクション実行時に割り当てられたトランザクション管理テーブルが不要となるため、システム内のアクティブトランザクション表から当該COMMITされたトランザクションを取り除くとともにトランザクション管理テーブルを解放する(ステップ86)。

10

【0031】

次に、本発明に関わるシンクポイント取得処理について、概略処理フローを図1を用いて説明する。シンクポイントの設定の目的は、システムが障害により停止してしまっても確実にシステムすなわちデータベースを回復でき、システムを存続できることを確かかなものにする事である。また、シンクポイントを設定することによってデータベースの回復時のREDOの必要なログの使用を小さくさせることにもある。シンクポイントを取得するタイミングは、システムインプリメンテーションによって様々であるが、定期的な何らかのアクションによって実施される。

【0032】

例えば、ログファイルへのログの出力回数が一定の回数に達した時点によってシンクポイントを取得する方法がある。この場合、シンクポイント取得のトリガを与えるのはログ書き込み処理部27である。まず、シンクポイント取得処理部26にシンクポイント取得処理要求が渡されると、シンクポイント取得開始ログを作成し、ログバッファに書き込む(ステップ251)。次にシンクポイント取得開始時点でのバッファプールの物理的整合性を一時的に保証するためにロックを確保する(ステップ252)。当該ロックは、ラッチと呼ぶこともあり、データベースの論理的な整合性を保証するために用いるロックとは異なり、一般のロックよりも保持時間が短い。バッファプールのロックを確保できたならば、当該ロックが解放されるまでの間はトランザクションによるバッファプールへのアクセスは一時的に待たされる。

20

30

【0033】

次にバッファプール34中に存在するすべてのページのページ管理テーブルを走査し、出力要求フラグ330が出力要求状態になっているページについてシンクポイント取得中フラグ335をシンクポイント取得状態に設定する(ステップ253)。ステップ253と同時にシンクポイント取得においてデータベース16へ書き込む対象となるページを管理するページ管理テーブルリスト38を作成する(ステップ254)。

【0034】

図12に出力対象ページ管理テーブルリスト38の構成を示す。出力対象ページ管理テーブルリスト38には、当該リスト作成時にリストに登録されたバッファプール34中の対象ページのページ管理テーブル32のアドレス382aから382dとその出力対象ページ数であるシンクポイント出力ページ数カウンタ381で構成される。登録されたリストの最後には、次の出力対象ページ管理テーブルがないことを示す情報328eとして0を設定する。そして、作成された出力対象ページ管理テーブルリスト38の当該リストが作成された時点でシンクポイント出力ページ数カウンタ381が決定される(ステップ255)。当該カウンタは、セマフォにより実現されていてもよい。

40

【0035】

当該シンクポイントでデータベースへの書きだしを行うページが確約されたので直ちにバッファプールのロックを解放する(ステップ256)。当該バッファプールのロック期間中には、一切の外部記憶装置との入出力処理を行わないのでCPU処理だけの短いロック時間となる。これにより、バッファプールへのロック待ちとなっていたトランザクション

50

の処理が続行可能となる。そして、作成された出力対象のページ管理テーブルリストに基づいて当該シンクポイントで出力すべきページを遅延書き込み処理部 25 に対して要求しデータベースへの出力を開始する(ステップ 257)。

【0036】

シンクポイントにおいて出力対象となったすべてのページの出力が完了したか否かは、前記シンクポイント出力ページ数カウンタ 381 が 0 になったかどうかによって決定される(ステップ 258)。したがって、シンクポイント出力ページ数カウンタ 381 が 0 になるまで当該処理は待ち状態となる。これは、シンクポイント出力対象ページの書き込み処理中に他のトランザクションによって書き込み処理が行われることがあるため、他のトランザクションによって書き込み処理が行われ、シンクポイント出力ページ数カウンタ 381 をカウントダウンして 0 になるまで同期を取る必要があるからである。シンクポイントにおいて出力対象となったすべてのページの出力が完了すると、すなわちシンクポイント出力ページ数カウンタ 381 が 0 になると待ち状態が解除されるので、シンクポイント取得完了ログを作成し、ログバッファに書き込むと同時にログファイルへの書き込みが強制される(ステップ 259)。

10

【0037】

シンクポイント取得処理が完了した時点で、当該シンクポイントをカレントな有効シンクポイントとするため当該シンクポイント開始ログの LSN をカレント有効シンクポイントとする(ステップ 260)。この有効シンクポイントの LSN は外部記憶装置 16 のような不揮発な記憶装置に格納されることが望ましい。つまり、システムがこの後、障害や故障によりシステムが停止しても、システムの再開始時にこの有効シンクポイントの LSN を読み出して、当該 LSN を起点としてデータベースの回復が行えるようになる。次に、図 1 で示したステップ 257 のシンクポイント出力対象ページの書き込み処理の処理が概略フローを図 13 を用いて説明する。

20

【0038】

まず、出力対象となった出力対象ページ管理テーブルリスト 38 の先頭アドレスにカレントポジション(以下、CP と略す)を設定する(ステップ 2561)。次に、当該 CP の指すページ管理テーブルのアドレスが 0 か否かを判定し、0 すなわちすべての出力対象ページのデータベースへの書き込みが完了したか出力対象ページがない場合は当該出力処理を終了する(ステップ 2562)。以降、出力対象ページが存在する間、以下の処理を繰り返す。まず、ページ管理テーブル 32 のロックカウンタ 328 をチェックし、ロックカウンタ 328 が 0 以上ならば当該出力処理の対象からはずし、当該ページを使用中のトランザクションによって出力を強制するようにする。ロックカウンタ 328 が 0 の場合、バッファプールにロックを確保する(ステップ 2564)。そして、出力中フラグ 332 を出力中状態に設定することによってデータベース 16 への書き込みが完了するまで他のトランザクションによる参照を待たせる(ステップ 2565)。

30

【0039】

次に、ロックカウンタ 328 が 0 ということは更新 LRU チェインに接続されているということであるので、更新 LRU チェインから切り離し、参照 LRU チェインの先頭に接続する。さらに、出力要求フラグ 330 を OFF にする(ステップ 2567)。本処理が完了するとバッファプールのロックを解放する(ステップ 2568)。そして、データベース 16 への書き込み処理を行う(ステップ 2569)。データベース 16 への書き込みが完了した時点で出力中フラグ 332 を OFF にする(ステップ 2570)。出力が完了するまでの間にトランザクションから参照された場合には入出力待ちチェイン 333 に参照したトランザクションが登録されているので、登録されたすべてのトランザクションをアクティブにするための同期制御処理を行う(ステップ 2571、2572)。こうして、一つのページの処理が終了すると CP を次のページ管理テーブルリストのアドレスに更新する。当該出力処理によって、シンクポイント取得開始ログ以前のすべてのログに対応した更新がデータベースに反映されることになる。

40

【0040】

50

シンクポイント取得処理によって選択された出力対象ページが出力処理の対象からはずされたページについては、バッファプール管理部23がトランザクションの処理によって確実にデータベース16に書きだされることを保証している。つまり、出力対象となったページを管理するページ管理テーブル32にはシンクポイント取得中フラグ335が設定されたままの状態になっており、トランザクションによって当該シンクポイント取得中フラグ335が設定されているページを使用する場合あるいは使用を終了する場合にデータベース16への出力を強制することによってトランザクション処理を止めることなく確実にシンクポイント取得処理が行われる。バッファプール管理部23においてどのようにシンクポイント出力対象ページがデータベース16に書きだされているかを図14から図19を用いて説明する。

10

【0041】

図14は、トランザクション内で行われたデータベースに対する更新処理において図9のステップ67で示した更新対象行格納ページ出力要求処理のバッファプール管理部23における処理フローを示す。まず、バッファプールにロックを確保する(ステップ401)。次にページ管理テーブル32の出力要求フラグ330がONであるかどうかチェックし(ステップ402)、ONならばさらにシンクポイント取得中フラグ335がONか否かをチェックし(ステップ403)、ONということは当該ページがシンクポイント取得中であるので当該トランザクションによってデータベース16への書き込みを強制される必要がある。

【0042】

まず、出力要求フラグ330をOFFにし(ステップ404)、参照LRUチェーンの先頭に接続する(ステップ405)。そして、出力中フラグ332をONにし(ステップ406)、ロックカウンタ328を0にし(ステップ407)、この時点でバッファプールのロックを解放する(ステップ408)。バッファプールのロックを解放したら、直ちにデータベース16への書き込みを行う(ステップ409)。データベース16への書き込みが完了した時点で出力中フラグ332をOFFにする(ステップ410)。シンクポイント取得時の対象となったページをデータベース36に書き込みを行ったのでシンクポイント出力ページ数カウンタ381から1を減ずる(ステップ411)。出力が完了するまでの間に他のトランザクションから参照された場合にはロック待ちチェーン334に参照したトランザクションが登録されているので、登録された先頭のトランザクションをアク

20

30

【0043】

図15は、トランザクション内で行われたデータベースに対するアクセスにおいてページの入力要求処理のバッファプール管理部23における処理フローを示す。基本的に要求されたページがバッファプール上に存在すればデータベース16からの読み込み無しで入力処理が完了し、バッファプール上に存在しない場合はデータベース16からの入力処理を伴う。その際、バッファプール上にページが存在し、ページ入力要求モードが排他モードの場合でかつシンクポイント取得中フラグ335がONになっているページについてはあらかじめデータベース16への書き込みを強制した後、入力要求処理を満たすようにする。

40

【0044】

まず、バッファプール34にロックを確保し(ステップ501)、バッファプール34上に目的のページが存在するか否かをハッシュ管理テーブル33を使用してサーチする(ステップ502)。目的のページがバッファプール34上に存在したならばページロックモードの種別によって処理を振り分ける(ステップ503、504)。ページロックモード

50

が共用モード（Sモード）ならば、共用モード入力処理を行い（ステップ505）、ページロックモードが排他モード（Xモード）ならば排他モード入力処理を行う（ステップ506）。目的のページが存在しない場合は空きバッファが有るか否かをバッファプール管理テーブル31のFREEチェインポインタ317によりチェックし（ステップ507）、FREEチェインポインタ317に接続されたバッファが存在する場合は当該FREEチェインに接続されたページ管理テーブル32を割り当てる（ステップ508）。FREEチェインポインタ317に該当するページ管理テーブル32が接続されていない場合は、参照LRUチェイン及び更新LRUチェインより最も使用されていないページ管理テーブルをスチールの対象とし、目的ページ入力用のバッファとして割り当てる（ステップ509）。

10

【0045】

こうして、割り当てられたバッファのページ管理テーブル32の入力中フラグ331を入力中状態に設定するためONにし（ステップ510）、ロックカウンタ328に1を設定し（ステップ511）、バッファプールのロックを解放する（ステップ512）。この場合、ページロックモードが排他モードの場合はXGET中フラグ329もONに設定する。そして、目的のページをデータベース16から読み込み（ステップ513）、読み込みが完了すると直ちに入力中フラグ331をOFFにするとともに入出力待ちチェインあるいはロック待ちチェインに登録されているトランザクションがあれば同期制御処理を行う（ステップ514、515、516）。

【0046】

20

次に、ページロックモードが共用モードの場合の入力処理の処理フローを図16に示す。まず、ロックカウンタ328が0か否かをチェックし（ステップ601）、0ならば参照LRUチェイン又は更新LRUチェインに存在することになる。出力要求フラグ330がONか否かをチェックし（ステップ602）、ONならば更新LRUチェインに存在するので更新LRUチェインから切り離す（ステップ603）。出力要求フラグ330がOFFならば参照LRUチェインに存在するので参照LRUチェインから切り離す（ステップ604）。そして、ロックカウンタ328に1を設定し（ステップ605）、バッファプールのロックを解放する（ステップ607）。

【0047】

また、ステップ601によりロックカウンタ328が0でない場合、すなわちすでに他のトランザクションが使用中の場合は、ロック待ちチェイン334に登録されているトランザクションが存在するか、又は当該ページのロックモードが排他モードで使用中であることを示すXGETフラグがONになっているか否かをチェックし（ステップ607）、いずれかの条件を満たす場合、ロック待ちチェイン334の最後尾に当該入力要求トランザクションを登録し（ステップ609）、一旦バッファプールのロックを解放し（ステップ609）、当該ロック待ち状態から解放されるのを同期制御処理によって待つ（ステップ610）。ステップ607の条件をいずれも満たさない場合、すなわち、共用モードで使用中のページの場合は、ロックカウンタ328を1増やし（ステップ611）、入出力中フラグがONか否かをチェックし（ステップ612）、データベース16からの入力中の場合は、入出力待ちチェイン333に当該トランザクションを登録し（ステップ613）、一旦バッファプールのロックを解放した後（ステップ614）、当該入出力が完了し、当該トランザクションをアクティブ状態にしてくれるまで待つ同期制御処理を行う（ステップ615）。ステップ612で入出力状態でない場合は、直ちにバッファプールのロックを解放することによって処理を終了する（ステップ616）。

30

40

【0048】

次に、ページロックモードが排他モードの場合の入力処理の処理フローを図17に示す。まず、ロックカウンタ328が0か否かをチェックし（ステップ701）、0ならば参照LRUチェイン又は更新LRUチェインに存在することになる。出力要求フラグ330がONか否かをチェックし（ステップ702）、ONならば更新LRUチェインに存在するので更新LRUチェインから切り離す（ステップ703）。出力要求フラグ330がOF

50

Fならば参照LRUチェーンに存在するので参照LRUチェーンから切り離す(ステップ704)。そして、ロックカウンタ328に1を設定する(ステップ705)。

【0049】

そして、ロックモードが排他モードであるのでXGETフラグ329をONにし、入力中フラグ331がONか否かをチェックし(ステップ707)、データベース16からの入力中の場合は入出力待ちチェーン333に当該トランザクションを登録し(ステップ708)、一旦バッファプールのロックを解放した後(ステップ709)、当該入出力が完了し、当該トランザクションをアクティブ状態にしてくれるまで待つ同期制御処理を行う(ステップ710)。ステップ707で入出力状態でない場合は、シンクポイント取得中フラグ335がONか否かをチェックし(ステップ714)、ONの場合は出力中フラグ332をONにし(ステップ715)、一旦バッファプールのロックを解放する(ステップ716)。

10

【0050】

そして、目的のページをデータベース16へ書き込みを強制し(ステップ717)、書き込み完了後、出力中フラグ332をOFFにし(ステップ718)、さらにシンクポイント取得中フラグ335及び出力要求フラグ330をOFFにする。また、シンクポイント取得の出力対象ページをデータベース36に出力したのでシンクポイント出力ページ数カウンタ381を1減ずる(ステップ719)。入出力待ちチェーンに登録されているトランザクションがあれば同期制御処理を行う(ステップ720、721)。ステップ714によりシンクポイント取得中フラグ335がOFFの場合は直ちにバッファプールのロックを解放する(ステップ722)。

20

【0051】

図18は、トランザクションにおいてアクセスしたページの使用が終わった場合にページのロックを解放する処理の処理フローを示す。まず、バッファプールのロックを確保する(ステップ801)。次にシンクポイント取得中フラグ335がONで有るか否かをチェックし(ステップ802)、ONの場合は出力中フラグ332をONにし(ステップ803)、一旦バッファプールのロックを解放する(ステップ804)。そして、当該ページをデータベース16へ書き込む(ステップ805)。データベース16への書き込みが完了するとページ管理テーブル32内の出力中フラグ332、シンクポイント取得中フラグ335、及び出力要求フラグ330をOFFにする。また、シンクポイント取得の出力対象ページをデータベース36に出力したのでシンクポイント出力ページ数カウンタ381を1減ずる(ステップ806)。

30

【0052】

次に、もしこのデータベースへの書き込み最中に当該ページへのアクセスが行われた場合には入出力待ちチェーン333に登録されているトランザクションをアクティブにするための同期制御処理を行った後(ステップ807、808)、バッファプールのロックを再度確保する(ステップ809)。次に、ステップ802においてシンクポイント取得中フラグ335がOFFの場合、あるいはONの場合にデータベース16への書き込みが終了するとロックカウンタ328を1減らし(ステップ810)、ロックカウンタ328が0でかつXGETフラグ329がONの場合はXGETフラグ329をOFFにする(ステップ811、812、813)。

40

【0053】

そして、出力要求フラグ330がONの場合は更新LRUチェーンへ接続し、出力要求フラグ330がOFFの場合は参照LRUチェーンへ接続する(ステップ814、815、816)。ステップ811においてロックカウンタ328が0でない場合、すなわち他にまだ使用中のトランザクションがいる場合とLRUチェーンに接続している最中にロック待ちチェーン334に登録されたトランザクションがある場合は、そのトランザクションのロック待ちを解除するための同期制御処理を行い、バッファプールのロックを解放する(ステップ817、818、819)。

【0054】

50

図19は、トランザクションにおいてアクセスしたページが不要になった場合に、当該ページを管理するページ管理テーブルをFREEチェインへ登録する処理の処理フローを示す。まず、バッファプールのロックを確保する(ステップ901)。次にシンクポイント取得中フラグ335がONで有るか否かをチェックし(ステップ902)、ONの場合は出力中フラグ332をONにし(ステップ903)、一旦バッファプールのロックを解放する(ステップ904)。そして、当該ページをデータベース16へ書き込む(ステップ905)。データベース16への書き込みが完了するとページ管理テーブル32内の出力中フラグ332、シンクポイント取得中フラグ335、及び出力要求フラグ330をOFFにする。また、シンクポイント取得の出力対象ページをデータベース36に出力したのでシンクポイント出力ページ数カウンタ381を1減ずる(ステップ906)。

10

【0055】

次に、もしこのデータベースへの書き込み最中に当該ページへのアクセスが行われた場合には入出力待ちチェイン333に登録されているトランザクションをアクティブにするための同期制御処理を行った後(ステップ907、908)、バッファプールのロックを再度確保する(ステップ909)。次に、ステップ802においてシンクポイント取得中フラグ335がOFFの場合、あるいはONの場合にデータベース16への書き込みが終了するとロックカウンタ328を0に設定する(ステップ910)。

【0056】

さらに、XGETフラグ329及び出力要求フラグ330をOFFにする(ステップ911、912)。FREEチェインに接続する前にロック待ちチェイン334に登録されたトランザクションがある場合は参照LRUチェインに接続するようにし、そのトランザクションのロック待ちを解除するための同期制御処理を行う(ステップ913、914、915)。ステップ913においてロック待ちチェイン334に登録されているトランザクションがない場合は、バッファプール管理テーブル31のFREEチェインに当該ページ管理テーブル32を接続し、ハッシュチェインからも切り離す(ステップ917、918)。これらの処理が終了するとバッファプールのロックを解放する(ステップ916)。

20

【0057】

以上、本発明を実施例に基づき具体的に説明したが、本発明により構成されたデータベース管理システムがトランザクション処理を始めた場合、システムのスループット、すなわち単位時間当たりのトランザクション処理量を時間の経過で示すと図20のような効果になる。トランザクションが一様に投入された場合のシステムのスループットを示したものであり、ある一定間隔においてシンクポイントが発生してもスループットの低下は見られない性能効果が得られる。

30

【0058】

次に、本発明の効果をさらに引き出すための実施例について説明する。前述した実施例では、シンクポイント取得開始時点のバッファプール上の更新ページをシンクポイント取得時にデータベースの整合性を保証するようにする処理について説明したが、通常シンクポイント間隔では図2に示した遅延書き込み処理部25がバッファプール上の更新ページ数を基に非同期にデータベースへの書き込みを行っているので、シンクポイント取得時点になるべくバッファプール上の更新ページ数が少なくなるように制御されている。

40

【0059】

しかし、図23の(a)を見てわかるようにシンクポイント取得時点の更新ページ数がバッファプール全体に閉める割合が大きい場合には、データベースへ書き出すページ数が増加するためシンクポイント取得に要する処理時間が長くなってしまふことがある。これは、通常シンクポイント間隔で行われる遅延書き込みの頻度を多くする、すなわちバッファプール全体に占める更新ページ数の比率を小さい値に設定することによって、その水準に達した時点で遅延書き込みを起動する方法が考えられる。ここでは、他の方法としてシンクポイント取得開始時点までに、あらかじめ遅延書き込み処理を利用してバッファプール中の更新ページをできるだけ多くデータベースに書き出すためのプレシンクポイントを行う処理の実施例を示す。

50

【 0 0 6 0 】

図 2 1 は、図 2 におけるログ書き込み処理部 2 7 において実施されるログファイル 3 7 への書き込み回数をシンクポイント間隔としてプレシンクポイント及びシンクポイントを取得する処理の制御を行う処理の概略フロ - を示す。トランザクションのコミットやロ - ルバックにおいてログバッファ 3 5 に書き込まれたログをログファイル 3 7 に強制出力する場合、まず、ログファイル 3 7 への書き込みを行う（ステップ 2 7 1 0）。ログファイル 3 7 への書き込みが終了すると、シンクポイント間隔を制御するためのログ出力回数を 1 増やす（ステップ 2 7 1 2）。

【 0 0 6 1 】

そして、ログ出力回数がプレシンクポイント取得点に到達したか否かを比較し（ステップ 2 7 1 2）、プレシンクポイント取得点に到達した場合は、バッファプールにロックを確保し（ステップ 2 7 1 6）、バッファプール管理テーブル 3 1 の遅延書き込み中フラグ 3 1 9 を ON にする（ステップ 2 7 1 8）。遅延書き込み中フラグ 3 1 9 を ON にすると、バッファプールのロックを解放し（ステップ 2 7 2 0）、プレシンクポイント取得処理に制御を渡す（ステップ 2 7 2 2）。ステップ 2 7 1 4 により、プレシンクポイント取得点に達していないか、あるいは通過した場合は、シンクポイント取得点に達したか否かを比較する（ステップ 2 7 2 4）。プレシンクポイントに達した場合は、さらに遅延書き込み中フラグ 3 1 9 が ON か否かをチェックする（ステップ 2 7 2 6）。遅延書き込み中フラグ 3 1 9 が ON の場合は、まだプレシンクポイント取得処理が完了していないので、プレシンクポイント取得処理を直ちに停止させる要求をプレシンクポイント取得処理を行っているプロセスに通達する（ステップ 2 7 2 8）。この方法は、シグナルやセマフォなどによって実現される。

【 0 0 6 2 】

プレシンクポイント取得処理が停止するか、あるいはすでに終了していた場合にはシンクポイント取得処理に制御を渡す（ステップ 2 7 3 0）。シンクポイント取得処理が終了すると、ログ出力回数を 0 にクリアする（ステップ 2 7 3 2）。次に、プレシンクポイント取得処理の処理フロ - を図 2 2 を用いて説明する。プレシンクポイント取得処理では基本的にバッファプール管理テーブル 3 1 の更新 L R U チェインに登録されたページを対象にデータベ - ス 3 6 への書き込みを行う。あらかじめ、プレシンクポイント取得点における更新 L R U チェイン中のページの L S N 4 2 5 を参照し、最も最新の L S N をプレシンク

【 0 0 6 3 】

まず、更新 L R U チェインの先頭アドレスすなわち最も古い更新ページを管理しているページ管理テーブルのアドレスをカレントポジション（以下、C P と略す）に設定する（ステップ 2 7 4 0）。C P の指すページ管理テーブル 3 2 が管理するページの L S N 4 2 5 をプレシンクポイント取得時の最大 L S N と比較し（ステップ 2 7 4 1）、C P の指すページ管理テーブル 3 2 が管理するページの L S N 4 2 5 の方が小さい間、プレシンクポイント取得処理のデータベ - ス 3 6 への出力対象ページとして処理を行う。

【 0 0 6 4 】

次に、プレシンクポイント取得処理中にシンクポイント取得点に達した場合は、プレシンクポイント取得停止要求が発行されていることがあるため、プレシンクポイント取得停止要求が発行されたか否かをチェックし（ステップ 2 7 4 2）、プレシンクポイント取得停止要求が発行されていない場合は、プレシンクポイント取得処理を継続する。まず、バッファプールにロックを確保する（ステップ 2 7 4 3）。そして、出力中フラグを出力中状態に設定することによってデータベ - ス 2 6 への書き込みが完了するまで他のトランザクションによる参照を待たせる（ステップ 2 7 4 4）。

【 0 0 6 5 】

次に、更新 L R U チェインから切離し、参照 L R U チェインに登録する（ステップ 2 7 4 5）。そして、出力要求フラグ 3 3 0 を OFF にする（ステップ 2 7 4 6）。本処理が完了するとバッファプールのロックを解放する（ステップ 2 7 4 7）。そして、データベ -

10

20

30

40

50

ス16への書き込みを行う(ステップ2748)。デ-タベ-ス16への書き込みが完了した時点で出力中フラグ332をOFFにする(ステップ2749)。デ-タベ-ス16への書き込みが完了するまでの間にトランザクションから参照された場合には入出力待ちチェーン333に参照したトランザクションが登録されているので、登録されたすべてのトランザクションをアクティブにするための同期制御処理を行う(ステップ2751、2752)。ステップ2741、2742によって、プレシクポイント取得処理が終了したか、プレシクポイント取得停止要求を受け付けた場合は、プレシクポイント取得処理の終了処理として遅延書き込み中フラグ319をOFFにする。

【0066】

プレシクポイントを設けた場合の性能効果について図23に示す。プレシクポイント取得時点は、シクポイント間隔の設定方法としてログの出力回数を使用する場合にはログ出力回数がある一定数に達した時点とする。シクポイント間隔の設定方法は、ログの出力回数のみではなく他の別の実施方法であってもよい。例えば、一定時間で設定されるようなシクポイント間隔ならばその一定時間に達する前のある時間に達した時点

10

【0067】

本実施例では、シクポイント間隔をログの出力回数によって設定される場合について述べる。この場合、プレシクポイント取得開始のトリガを与えるのは図2におけるログ書き込み処理部27である。ログ書き込み処理部27は、トランザクションからのログの書き込み要求に応じてログバッファ35にキャッシュし、ログバッファ35が満杯になると、トランザクションのコミット、ロールバック時点によってログファイル37への書き出しが強制される。その場合、ログファイル37への出力回数を主記憶装置上に記憶しておき、シクポイント間隔のある一定数に達した時点で遅延書き込み処理部25に制御を渡す。

20

【0068】

例えば、シクポイント間隔が1000であった場合、ログファイルへの出力回数が1000に達した時点がシクポイント取得時点となり、プレシクポイントはその80%の800に達した時点とするようにする。通常、遅延書き込み処理では、バッファプール中の更新ページ数がある一定数に達するとデータベースへの書き込みを開始するが、すべての更新ページがデータベースへ書き出されるのではなく、一定のページ数に留めることが多い。しかし、プレシクポイント取得時点では、その時点に存在するバッファプール上のすべての更新ページをデータベースに書き出す。そうすると、図20の(b)及び(c)のようにシクポイント取得時点までにバッファプール上のほとんどの更新ページが書き出されるようになる。

30

【0069】

理想的には、シクポイント取得時点にバッファプール上の更新ページ数が0になることが望ましいが、通常(b)のようにプレシクポイントで更新ページをデータベースに書き出す処理を行ってもシクポイント取得時点までに全てのページを書き出せないことがある。その場合は、残った更新ページをシクポイント取得処理にてデータベースへの書き出しを行うようになる。また、逆に(c)のようにプレシクポイントでバッファプール上の更新ページを全てデータベースに書き出せたが、シクポイント取得時点に達しない場合もある。その場合は、シクポイント取得時点に存在するバッファプール上の更新ページがシクポイント取得時点でのデータベースへの書き出し対象ページとなる。

40

【0070】

プレシクポイントの取得開始点は、図20の(b)、(c)のようにシクポイント取得開始時点での更新ページ数を残さないようにするため、シクポイント取得開始時にプレシクポイント取得が完了したか否かを記憶しておき、プレシクポイント取得が完了していない場合はデータベースへ書き込みができなかった更新ページ数からプレシクポイント取得開始点を早くするようにし、プレシクポイント取得がすでに完了している場合は、シクポイント取得時に存在した更新ページ数からプレシクポイント取得開始点

50

を遅くするようにプレシंकポイント取得開始点を補正し、次のシंकポイントのプレシंकポイント取得開始点とする処理を行う。

【 0 0 7 1 】

【 発 明 の 効 果 】

以上、説明したように、本発明によれば、データベース管理システムにおいてシंकポイント取得開始時にバッファプール上に存在するすべての更新ページにマークをつけるだけで、そのマークのついた更新ページをデータベースへの書き込み対象とするため、シंकポイント取得中であってもトランザクション処理を全く停止させることなく入出力処理を可能にするので、シंकポイント取得中のトランザクション処理量が0になることがない。また、プレシंकポイント処理を行うことにより、シंकポイント取得時点のバッファ 10
プール上に存在する更新ページ数を削減することができるので、シंकポイント取得処理時間を大幅に短縮できる。

【 図 面 の 簡 単 な 説 明 】

【 図 1 】 本 発 明 の 特 長 と な る シ ン ク ポ イ ン ト 取 得 処 理 の 概 略 処 理 フ ロ ー を 示 す。

【 図 2 】 本 発 明 を 実 施 す る た め の デ ー タ ベ ー ス 管 理 シ ス テ ム の 構 成 図 を 示 す。

【 図 3 】 デ - タ ベ - ス を 構 成 す る ペ - ジ の 構 造 を 示 す。

【 図 4 】 図 3 に お け る ペ - ジ 内 の 制 御 情 報 の 構 成 を 示 す。

【 図 5 】 図 2 に お け る バ ッ フ ァ プ ー ル を 管 理 す る テ - ブ ル の 構 造 を 示 す。

【 図 6 】 図 2 に お け る バ ッ フ ァ プ ー ル 管 理 テ - ブ ル の 構 成 情 報 を 示 す。

【 図 7 】 図 2 に お け る バ ッ フ ァ プ ー ル 中 の 各 ペ - ジ を 管 理 す る ペ - ジ 管 理 テ - ブ ル の 構 成 20
情 報 を 示 す。

【 図 8 】 ト ラ ン ザ ク シ ョ ン の 開 始 処 理 の 処 理 フ ロ ー を 示 す。

【 図 9 】 あ る ト ラ ン ザ ク シ ョ ン に お け る デ - タ ベ - ス の 更 新 処 理 の 処 理 フ ロ ー を 示 す。

【 図 1 0 】 ト ラ ン ザ ク シ ョ ン の P R E P A R E 処 理 の 処 理 フ ロ ー を 示 す。

【 図 1 1 】 ト ラ ン ザ ク シ ョ ン の C O M M I T 処 理 の 処 理 フ ロ ー を 示 す。

【 図 1 2 】 出 力 対 象 ペ - ジ 管 理 テ - ブ ル リ ス ト の 構 成 情 報 を 示 す。

【 図 1 3 】 シ ン ク ポ イ ン ト 出 力 対 象 ペ - ジ の 書 き 込 み 処 理 の 処 理 フ ロ ー を 示 す。

【 図 1 4 】 図 2 の バ ッ フ ァ プ ー ル 管 理 部 に お け る 出 力 要 求 処 理 の 処 理 フ ロ ー を 示 す。

【 図 1 5 】 図 2 の バ ッ フ ァ プ ー ル 管 理 部 に お け る 入 力 要 求 処 理 の 処 理 フ ロ ー を 示 す。

【 図 1 6 】 図 3 の 共 用 モ - ド の ペ - ジ 入 力 処 理 の 処 理 フ ロ ー の 詳 細 を 示 す。 30

【 図 1 7 】 図 3 の 排 他 モ - ド の ペ - ジ 入 力 処 理 の 処 理 フ ロ ー の 詳 細 を 示 す。

【 図 1 8 】 図 2 の バ ッ フ ァ プ ー ル 管 理 部 に お け る ペ - ジ ロ ッ ク 解 放 処 理 の 処 理 フ ロ ー を 示
す。

【 図 1 9 】 図 2 の バ ッ フ ァ プ ー ル 管 理 部 に お け る バ ッ フ ァ F R E E 処 理 の 処 理 フ ロ ー を 示
す。

【 図 2 0 】 本 発 明 実 施 時 の デ - タ ベ - ス 管 理 シ ス テ ム の 性 能 特 性 の グ ラ フ を 示 す。

【 図 2 1 】 プ レ シ ン ク ポ イ ン ト 及 び シ ン ク ポ イ ン ト 取 得 を 制 御 す る ロ グ 書 き 込 み 処 理 の 処
理 フ ロ ー を 示 す。

【 図 2 2 】 プ レ シ ン ク ポ イ ン ト 取 得 処 理 の 処 理 フ ロ ー を 示 す。

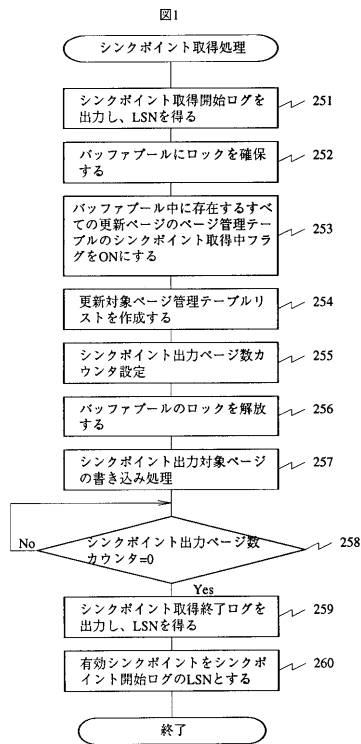
【 図 2 3 】 本 発 明 実 施 時 に 効 果 を 上 げ る た め の プ レ シ ン ク ポ イ ン ト 処 理 を 行 っ た 場 合 の バ
ッ フ ァ プ ー ル に お け る 更 新 ペ - ジ 数 の 時 間 経 過 に お け る 比 率 の グ ラ フ を 示 す。 40

【 符 号 の 説 明 】

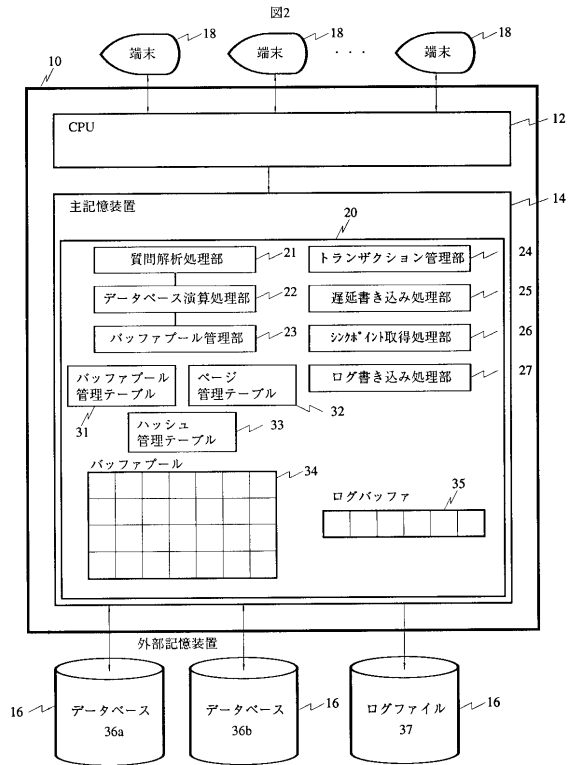
1 0 : コ ン ピ ュ - タ シ ス テ ム 、 1 2 : C P U 、 1 4 : 主 記 憶 装 置 、 1 6 : 外 部 記 憶 装 置 、
1 8 : 端 末 、 2 0 : デ - タ ベ - ス 管 理 シ ス テ ム 、 2 1 : 質 問 解 析 処 理 部 、 2 2 : デ - タ ベ
- ス 演 算 処 理 部 、 2 3 : バ ッ フ ァ プ ー ル 管 理 部 、 2 4 : ト ラ ン ザ ク シ ョ ン 管 理 部 、 2 5 :
遅 延 書 き 込 み 処 理 部 、 2 6 : シ ン ク ポ イ ン ト 取 得 処 理 部 、 2 7 : ロ グ 書 き 込 み 処 理 部 、 3
1 : バ ッ フ ァ プ ー ル 管 理 テ - ブ ル 、 3 2 : ペ - ジ 管 理 テ - ブ ル 、 3 3 : ハ ッ シ ュ 管 理 テ -
ブ ル 、 3 4 : バ ッ フ ァ プ ー ル 、 3 5 : ロ グ バ ッ フ ァ 、 3 6 : デ - タ ベ - ス 、 3 7 : ロ グ フ
ァ イ ル 、 3 8 : 出 力 対 象 ペ - ジ 管 理 テ - ブ ル リ ス ト 、 4 0 : ペ - ジ 、 4 2 : ペ - ジ 内 制 御
情 報 、 4 4 : ス ロ ッ ト 領 域 、 3 1 9 : 遅 延 書 き 込 み 中 フ ラ グ 、 3 2 8 : ロ ッ ク カ ウ ン タ 、 50

329 : XGETフラグ、330 : 出力要求フラグ、335 : シンクポイント取得中フラグ、381 : シンクポイント出力ページ数カウンタ、421 : 割り当 slots 数、422 : 使用 slots 数、423 : 空き領域長、424 : 未使用領域先頭オフセット、425 : LSN

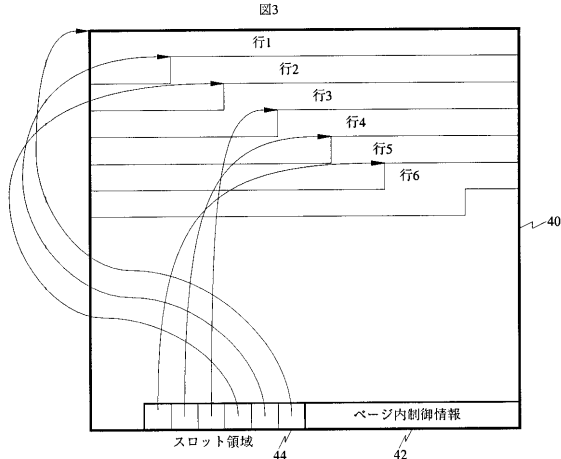
【 図 1 】



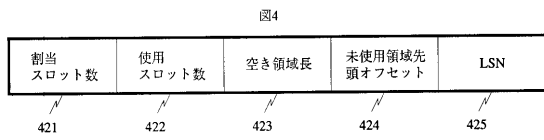
【 図 2 】



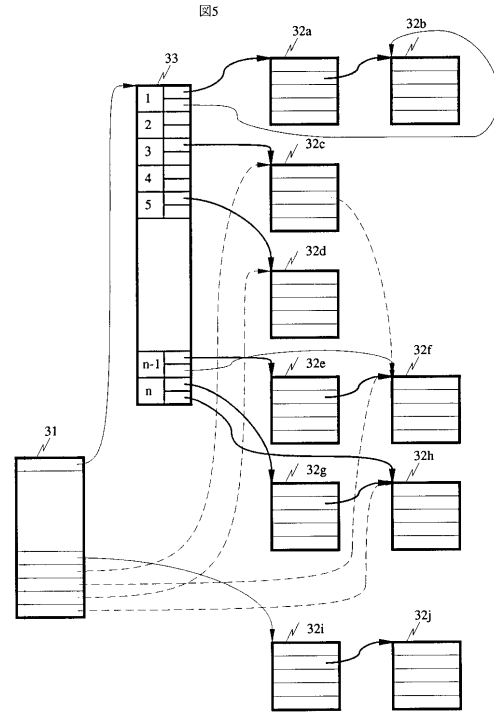
【 図 3 】



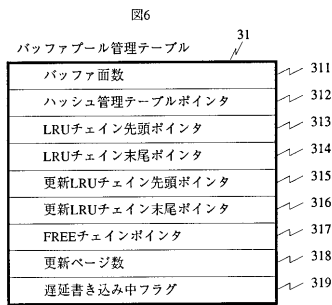
【 図 4 】



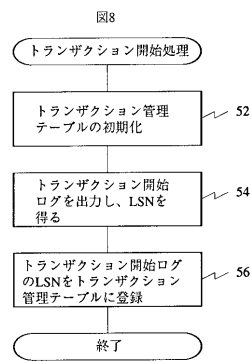
【 図 5 】



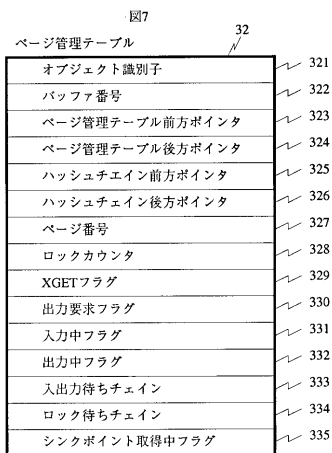
【 図 6 】



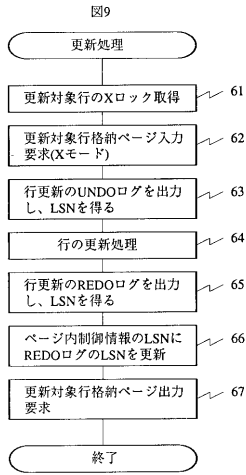
【 図 8 】



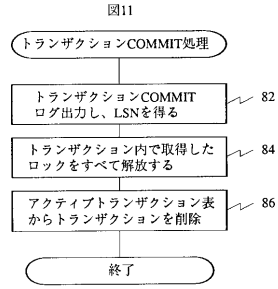
【 図 7 】



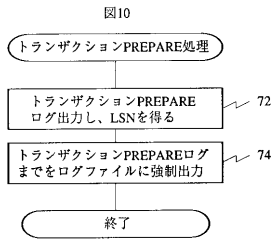
【 図 9 】



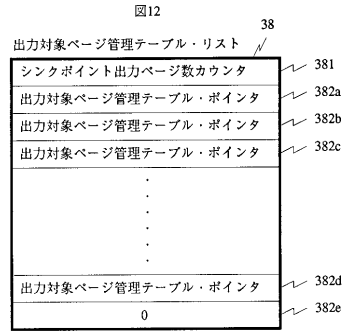
【 図 1 1 】



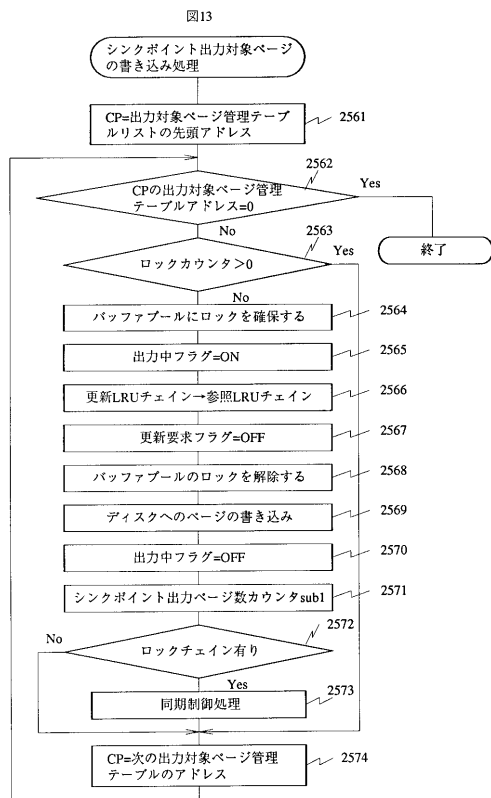
【 図 1 0 】



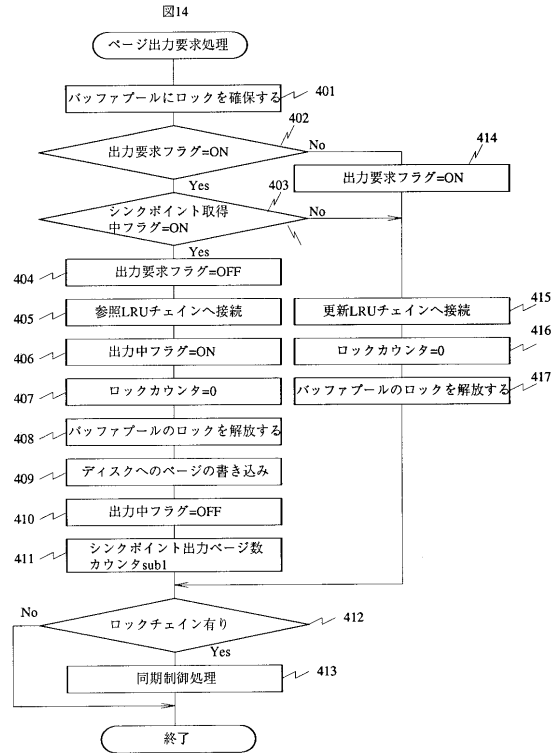
【 図 1 2 】



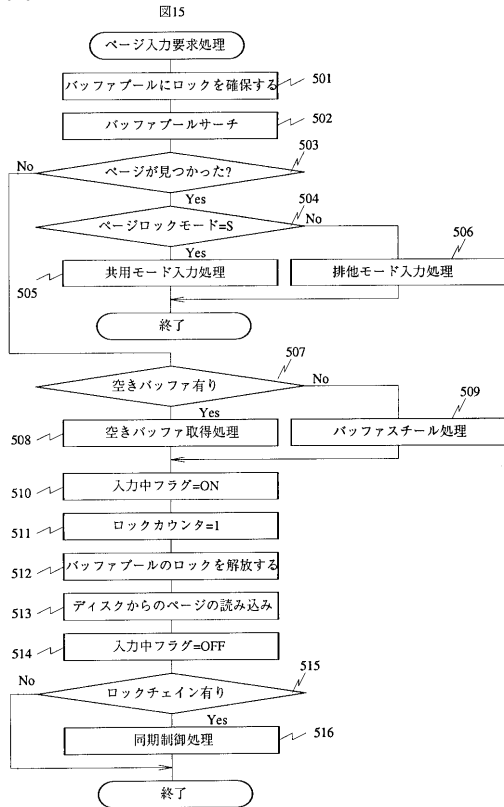
【 図 1 3 】



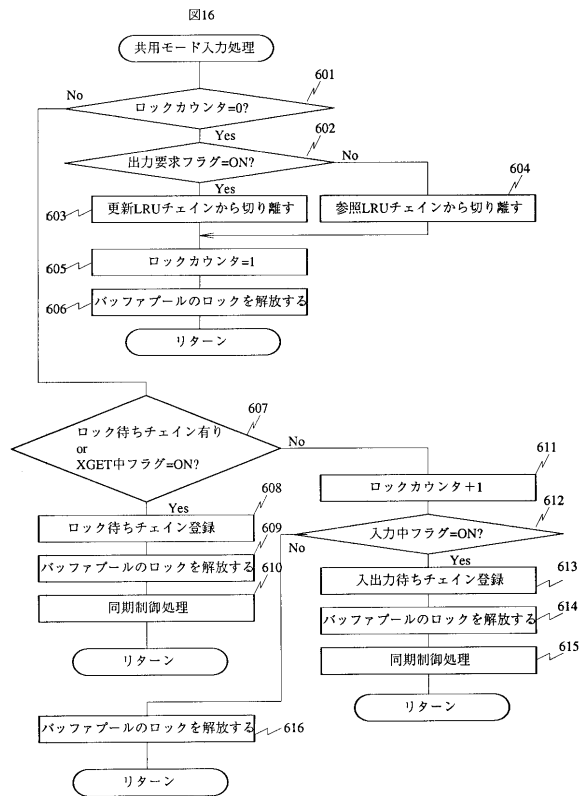
【 図 1 4 】



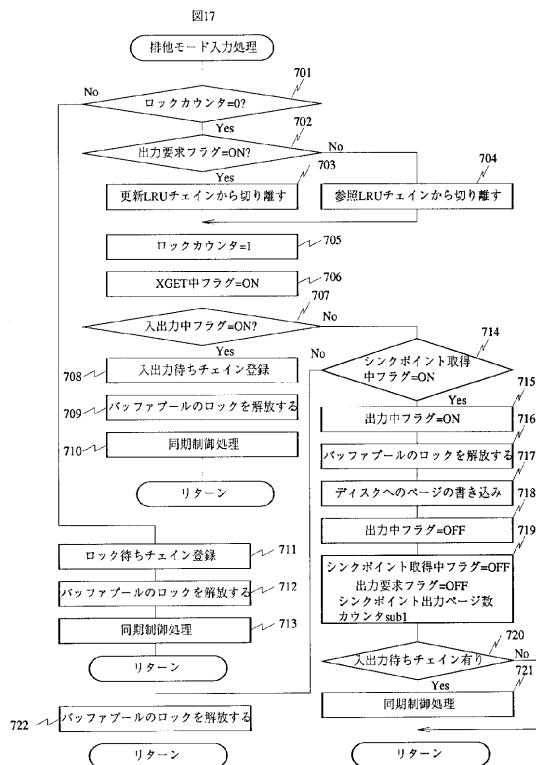
【 図 1 5 】



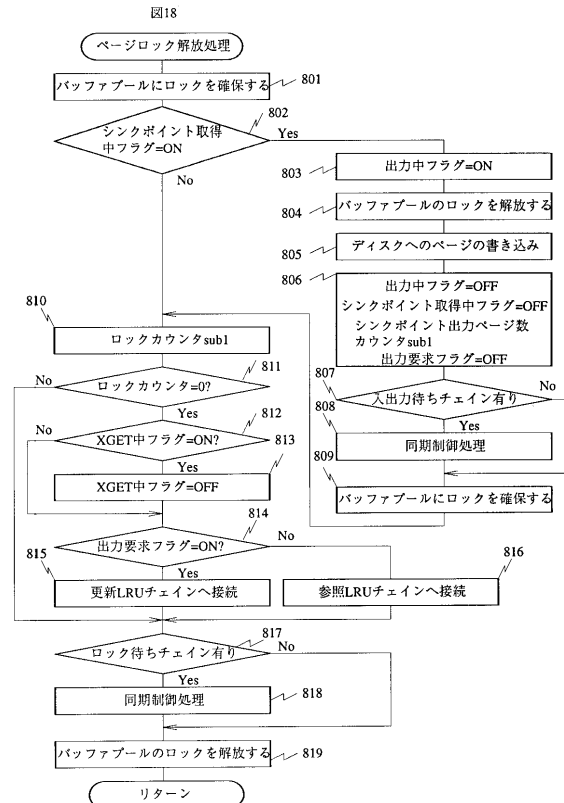
【 図 1 6 】



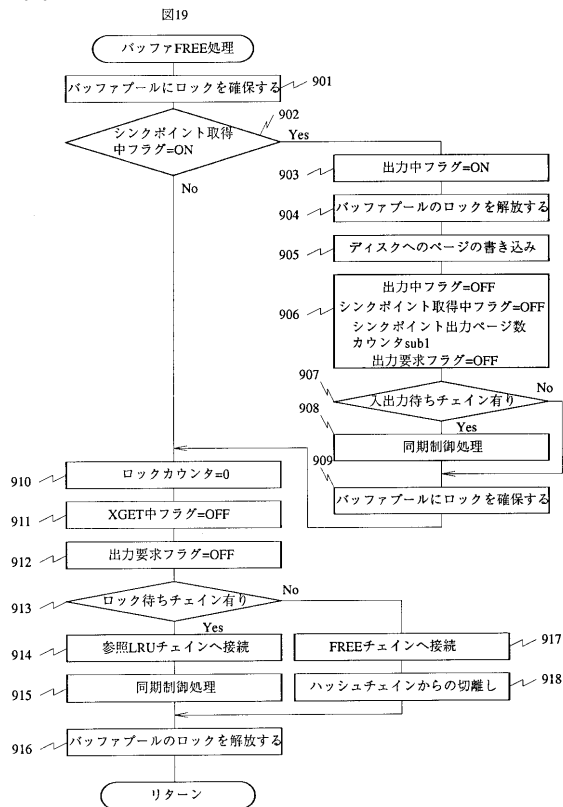
【 図 1 7 】



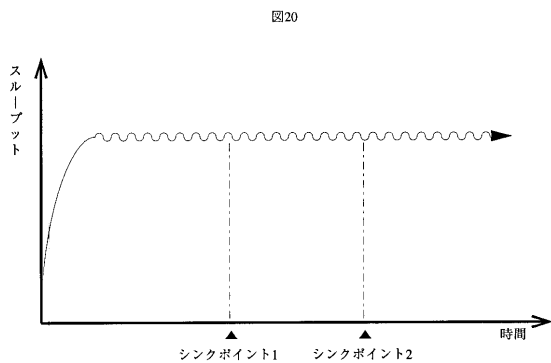
【 図 1 8 】



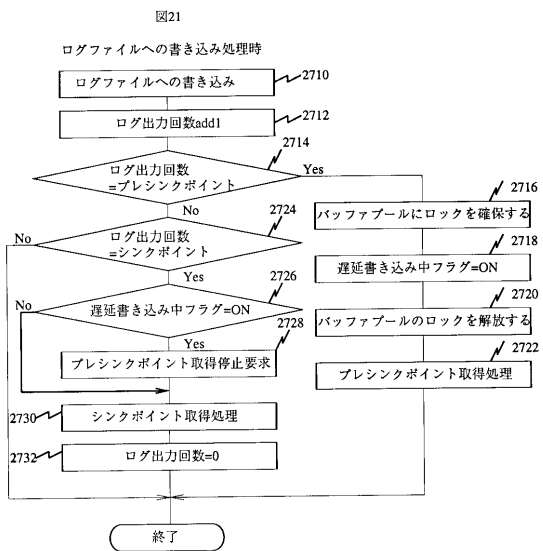
【 図 19 】



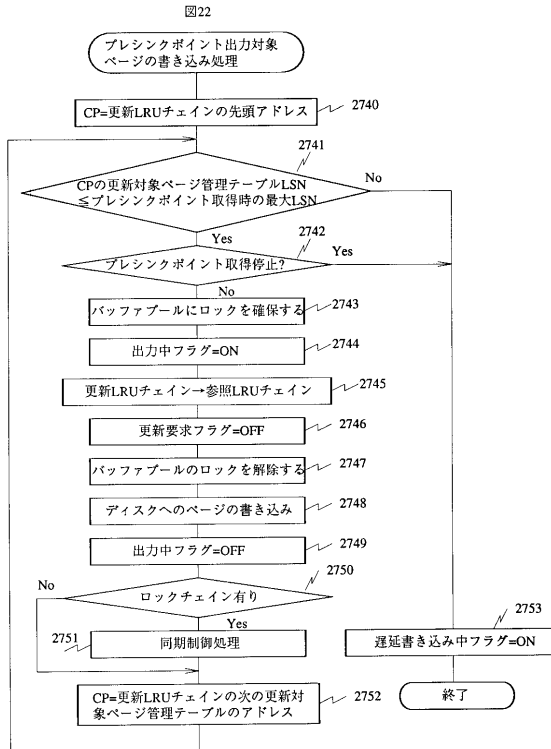
【 図 20 】



【 図 21 】

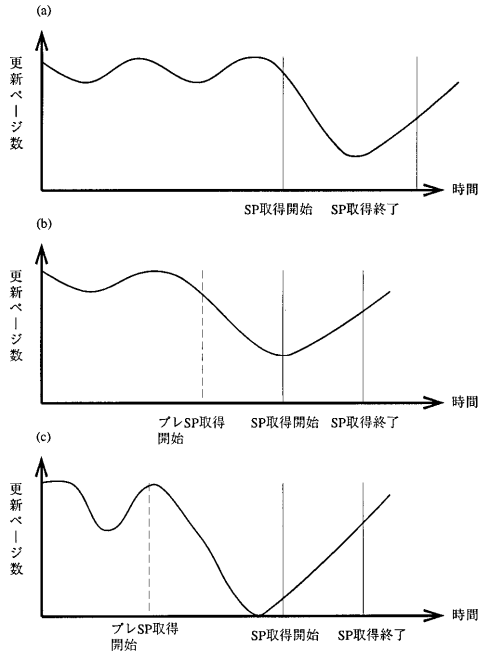


【 図 22 】



【 図 2 3 】

図23



フロントページの続き

- (72)発明者 河村 信男
神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所 システム開発研究所内
- (72)発明者 正井 一夫
神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所 ソフトウェア開発本部内
- (72)発明者 山下 信之
神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所 システム開発研究所内
- (72)発明者 永井 浩
神奈川県横浜市中区尾上町6丁目81番地 日立ソフトウェアエンジニアリング株式会社内

審査官 辻本 泰隆

- (56)参考文献 特開平04-282733(JP,A)
特開平06-214857(JP,A)
特開平06-214848(JP,A)
特開昭63-195755(JP,A)
特開平05-158778(JP,A)
特開平03-268146(JP,A)
特開昭63-133240(JP,A)

- (58)調査した分野(Int.Cl.⁷, DB名)
G06F 17/30,
G06F 12/00