

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2005-25750

(P2005-25750A)

(43) 公開日 平成17年1月27日(2005.1.27)

(51) Int. Cl.⁷
G06F 12/00

F I
G06F 12/00 518A

テーマコード(参考)
5B082

審査請求 未請求 請求項の数 44 O L (全 22 頁)

(21) 出願番号 特願2004-190397(P2004-190397)
(22) 出願日 平成16年6月28日(2004.6.28)
(31) 優先権主張番号 10/611,774
(32) 優先日 平成15年6月30日(2003.6.30)
(33) 優先権主張国 米国(US)

(71) 出願人 500046438
マイクロソフト コーポレーション
アメリカ合衆国 ワシントン州 9805
2-6399 レッドモンド ワン マイ
クロソフト ウェイ
(74) 代理人 100077481
弁理士 谷 義一
(74) 代理人 100088915
弁理士 阿部 和夫
(72) 発明者 マイケル ジェイ. ツウィリング
アメリカ合衆国 98052 ワシントン
州 レッドモンド ノースイースト 68
ストリート 15215

最終頁に続く

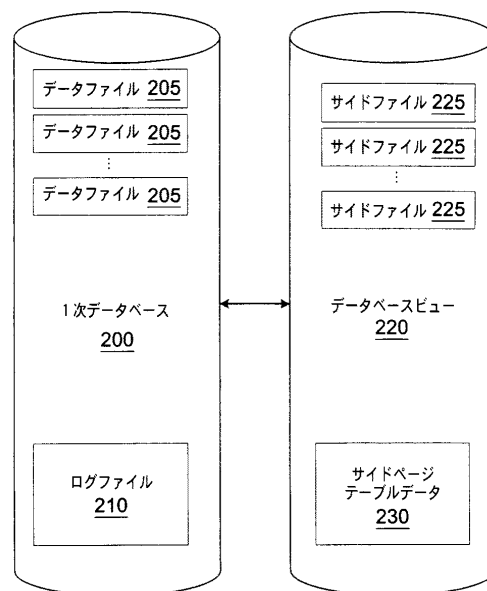
(54) 【発明の名称】 トランザクションの整合性を保つ書き込み時コピーのデータベース

(57) 【要約】

【課題】 前のある時点の既存のデータベースのトランザクションの整合性が保たれたビューを提供するデータベースのデータベースビューを作成すること。

【解決手段】 各データベースビューは、1次データベースとともに、前のある時点の1次データベースの内容を決定するのに必要なすべての情報を含む。データベースビューは、1次データベース内の各データファイルに対応するサイドファイルから成る。サイドファイルは、データベースビューが作成されてから変更された対応するデータファイルからのすべてのデータのコピーを含む。スペースを保つために、スパーズファイルをサイドファイルのために使用することができる。1次データベースからのページが修正され、古いバージョンがデータベースビューのサイドファイルに格納されているかどうかについて迅速に決定できるようにするために、ページテーブルデータが保持される。

【選択図】 図3



【特許請求の範囲】**【請求項 1】**

特定の時点のデータベースの内容を反映するトランザクションの整合性を保つデータを含むデータベースビューを提供する方法であって、前記データベースは、データ要素を含み、トランザクションログに関連付けられ、前記トランザクションログは、アクティブトランザクションおよび非アクティブトランザクションを含み、前記データベースビューは、前記データベースからの前記データ要素のうちの少なくとも1つの前のバージョンを格納する差分ストレージを含み、前記方法は、

前記時点に対応する前記トランザクションログにおける分割点を決定するステップと、
前記データベースに対する修正を行う前記分割点より前の各トランザクションを前記トランザクションログ上で見つけるステップと、

前記差分ストレージ内に前記修正のそれぞれを格納するステップと、

前記分割点より前の各アクティブトランザクションを前記トランザクションログ上で見つけるステップと、

前記差分ストレージ内の対応する任意の修正を元に戻すステップと
を備えたことを特徴とする方法。

【請求項 2】

前記差分ストレージを初期化するステップをさらに備えたことを特徴とする請求項 1 に記載の方法。

【請求項 3】

前記時点に対応する前記トランザクションログにおける分割点を決定する前記ステップは、ログの切り捨てを使用不可にするステップをさらに含むことを特徴とする請求項 1 に記載の方法。

【請求項 4】

前記データベースビューは、ページテーブルをさらに含み、前記方法は、前記ページテーブルを初期化するステップをさらに備えたことを特徴とする請求項 1 に記載の方法。

【請求項 5】

前記差分ストレージ内の対応する任意の修正を元に戻す前記ステップは、前記対応する修正を削除するステップを含むことを特徴とする請求項 1 に記載の方法。

【請求項 6】

前記差分ストレージ内の対応する任意の修正を元に戻す前記ステップは、

前記データベース内の対応する修正されていないデータを読み取るステップと、

前記対応する修正されていないデータを前記差分ストレージに書き込むステップと

を含むことを特徴とする請求項 1 に記載の方法。

【請求項 7】

前記データ要素のそれぞれは、1 ページのデータを含むことを特徴とする請求項 1 に記載の方法。

【請求項 8】

前記差分ストレージは、少なくとも1つのスパーファイルを含むことを特徴とする請求項 7 に記載の方法。

【請求項 9】

前記差分ストレージ内に前記修正のそれぞれを格納する前記ステップは、前記スパーファイルのうちの1つにおいてメモリの領域を割り振るステップを含むことを特徴とする請求項 8 に記載の方法。

【請求項 10】

前記データ要素のそれぞれは、1 ページのデータを含み、前記データベースビューは、ページテーブルをさらに含み、前記ページテーブルは、ページごとに、

前記ページが前記差分ストレージに格納されているかどうかを示す第1の格納データと

、
前記差分ストレージにおいて前記領域が割り振られているかどうかを示す第2の格納デ

10

20

30

40

50

ータと

を含むことを特徴とする請求項 9 に記載の方法。

【請求項 1 1】

前記方法は、

前記ページテーブルが無効であることを検出するステップと、

前記スパーファイル内の領域ごとに、前記領域が割り振られているかどうかを決定するステップと、

前記スパーファイル内の領域ごとに、前記領域が割り振られているかどうかに基づいて前記第 2 の格納データを設定するステップと

をさらに備えたことを特徴とする請求項 1 0 に記載の方法。

10

【請求項 1 2】

前記差分ストレージ内の特定のページにデータが格納されているかどうかの決定は、

前記第 1 の格納データをチェックし、前記特定のページが前記差分ストレージに格納されていることを、前記第 1 の格納データが示す場合、前記データは、前記差分ストレージ内の前記特定のページに格納されていることを決定するステップと、

前記第 2 の格納データをチェックし、前記領域が前記差分ストレージにおいて割り振られていないことを、前記第 2 の格納データが示す場合、前記データは、前記差分ストレージ内の前記特定のページに格納されていないことを決定するステップと、

前記ページが前記差分ストレージ内に格納されていることを、前記第 1 の格納データが示さず、前記領域が前記差分ストレージにおいて割り振られていないことを、前記第 2 の格納データが示さない場合、前記特定のページについての前記差分ストレージの対応するエリアからページデータを読み取り、前記対応するエリアからの前記ページデータが有効であるかどうかを決定するステップと

20

を含むことを特徴とする請求項 1 0 に記載の方法。

【請求項 1 3】

前記データベースビューは、ページテーブルをさらに含み、前記ページテーブルは、ページごとに、

前記ページが前記差分ストレージ内に格納されているかどうかを示す第 1 の格納データを含むことを特徴とする請求項 7 に記載の方法。

【請求項 1 4】

30

データが前記差分ストレージ内の特定のページ内に格納されているかどうかの決定は、

前記第 1 の格納データをチェックし、前記特定のページが前記差分ストレージに格納されていることを、前記第 1 の格納データが示す場合、前記データは、前記差分ストレージ内の前記特定のページに格納されていることを決定するステップと、

前記特定のページが前記差分ストレージ内に格納されていることを、前記第 1 の格納データが示さない場合、前記特定のページについての前記差分ストレージの対応するエリアからページデータを読み取り、前記対応するエリアからの前記ページデータが有効であるかどうかを決定するステップと

を含むことを特徴とする請求項 1 3 に記載の方法。

【請求項 1 5】

40

前記データベースビュー内の特定のデータ要素の要求を受け付けるステップと、

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されているかどうかを決定するステップと、

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されている場合、前記差分ストレージを読み取ることによって前記要求に応答するステップと、

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されている場合、前記データベースを読み取ることによって前記要求に応答するステップと

をさらに含むことを特徴とする請求項 1 に記載の方法。

【請求項 1 6】

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されてい

50

るかどうかを決定する前記ステップは、前記差分ストレージが前記場所に有効なデータを含んでいるかどうかを決定するステップを含むことを特徴とする請求項 15 に記載の方法。

【請求項 17】

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されているかどうかを決定する前記ステップは、ページテーブルを調べるステップを含むことを特徴とする請求項 15 に記載の方法。

【請求項 18】

前記方法は、

第 2 の特定のデータ要素の代わりに前記データベース内のある場所に第 1 の特定の値を格納する前記データベースに加えられた修正を検出するステップと、

前記データベースビュー内の対応する場所が有効なデータを含んでいるかどうかを決定するステップと、

前記データベースビュー内の前記対応する場所が有効なデータを含んでいない場合、前記第 2 の特定のデータ要素を前記対応する場所に書き込むステップと

をさらに含むことを特徴とする請求項 1 に記載の方法。

【請求項 19】

請求項 1 に記載の方法を実行することを特徴とする、オペレーティングシステム、複数のコンピュータ実行可能命令を格納するコンピュータ読取り可能媒体、コプロセッシング装置、コンピューティング装置、およびコンピュータ実行可能命令を含む変調されたデータ信号のうち少なくとも 1 つ。

【請求項 20】

特定の時点のデータベースの内容を反映するトランザクションの整合性を保つデータを含むデータベースビューを提供するシステムであって、前記データベースは、データ要素を含み、トランザクションログに関連付けられ、前記トランザクションログは、アクティブトランザクションおよび非アクティブトランザクションを含み、前記システムは、

前記時点に対応する前記トランザクションログにおける分割点を決定する分割点決定器と、

前記データベースに対する修正を行う前記分割点より前の各トランザクションを前記トランザクションログ上で見つける第 1 のトランザクションログアナライザと、

前記差分ストレージ内に前記修正のそれぞれを格納する差分ストレージと、

前記分割点より前の各アクティブトランザクションを前記トランザクションログ上で見つける第 2 のトランザクションログアナライザと、

前記差分ストレージ内の対応する任意の修正を元に戻す差分ストレージ変更器とを備えたことを特徴とするシステム。

【請求項 21】

前記分割点決定器は、ログの切り捨てを使用不可にするログ切り捨て使用不可器をさらに含むことを特徴とする請求項 20 に記載のシステム。

【請求項 22】

前記システムは、特定のデータ要素が前記差分ストレージに格納されているかどうかを示すデータを含むページテーブルをさらに含むことを特徴とする請求項 20 に記載のシステム。

【請求項 23】

前記ページテーブルは、前記特定のデータ要素が前記差分ストレージに格納されているかどうかを示す第 1 の格納データを含むことを特徴とする請求項 22 に記載のシステム。

【請求項 24】

前記差分ストレージは、スパーズファイルを含み、

前記ページテーブルは、前記差分ストレージにおいて前記特定のデータ要素に対応する領域が割り振られているかどうかを示す第 2 の格納データをさらに含むことを特徴とする請求項 23 に記載のシステム。

【請求項 25】

前記システムは、

前記スパーファイル内の領域ごとに、前記領域が割り振られているかどうかを決定する領域割り振り決定器と、

前記スパーファイル内の領域ごとに、前記領域が割り振られているかどうかに基づいて前記第2の格納データを設定する第2の格納データ設定機器と

をさらに備えたことを特徴とする請求項24に記載のシステム。

【請求項 26】

前記データベース内の特定のデータ要素の要求を受け付け、データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されているかどうかを決定し、データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されている場合、前記差分ストレージを読み取ることによって前記要求に応答し、データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されている場合、前記データベースを読み取ることによって前記要求に応答する要求応答器をさらに含むことを特徴とする請求項20に記載のシステム。

10

【請求項 27】

特定の時点のデータベースの内容を反映するトランザクションの整合性を保つデータを含むデータベースを提供するコンピュータ読取り可能媒体であって、前記データベースは、データ要素を含み、トランザクションログに関連付けられ、前記トランザクションログは、アクティブトランザクションおよび非アクティブトランザクションを含み、前記データベースは、前記データベースからの前記データベース要素のうちの少なくとも1つの前のバージョンを格納する差分ストレージを含み、前記コンピュータ読取り可能媒体は、

20

前記時点に対応する前記トランザクションログにおける分割点を決定するステップと、

前記データベースに対する修正を行う前記分割点より前の各トランザクションを前記トランザクションログ上で見つけるステップと、

前記差分ストレージ内に前記修正のそれぞれを格納するステップと、

前記分割点より前の各アクティブトランザクションを前記トランザクションログ上で見つけるステップと、

前記差分ストレージ内の対応する任意の修正を元に戻すステップと

30

を含む動作を実行する命令を備えたことを特徴とするコンピュータ読取り可能媒体。

【請求項 28】

前記動作は、前記差分ストレージを初期化するステップをさらに含むことを特徴とする請求項27に記載のコンピュータ読取り可能媒体。

【請求項 29】

前記時点に対応する前記トランザクションログにおける分割点を前記決定するステップは、ログの切り捨てを使用不可にするステップをさらに含むことを特徴とする請求項27に記載のコンピュータ読取り可能媒体。

【請求項 30】

前記データベースは、ページテーブルをさらに含み、前記動作は、前記ページテーブルを初期化するステップをさらに含むことを特徴とする請求項27に記載のコンピュータ読取り可能媒体。

40

【請求項 31】

前記差分ストレージ内の対応する任意の修正を元に戻す前記ステップは、前記対応する修正を削除するステップを含むことを特徴とする請求項27に記載のコンピュータ読取り可能媒体。

【請求項 32】

前記差分ストレージ内の対応する任意の修正を元に戻す前記ステップは、

前記データベース内の対応する修正されていないデータを読み取るステップと、

前記対応する修正されていないデータを前記差分ストレージに書き込むステップと

50

を含むことを特徴とする請求項 27 に記載のコンピュータ読取り可能媒体。

【請求項 33】

前記データ要素のそれぞれは、1 ページのデータを含むことを特徴とする請求項 27 に記載のコンピュータ読取り可能媒体。

【請求項 34】

前記差分ストレージは、少なくとも 1 つのスパースファイルを含むことを特徴とする請求項 33 に記載のコンピュータ読取り可能媒体。

【請求項 35】

前記差分ストレージ内に前記修正のそれぞれを格納する前記動作は、前記スパースファイルのうち 1 つにおいてメモリの領域を割り振るステップを含むことを特徴とする請求項 34 に記載のコンピュータ読取り可能媒体。

10

【請求項 36】

前記データ要素のそれぞれは、1 ページのデータを含み、前記データベースは、ページテーブルをさらに含み、前記ページテーブルは、ページごとに、

前記ページが前記差分ストレージに格納されているかどうかを示す第 1 の格納データと

、前記差分ストレージにおいて前記領域が割り振られているかどうかを示す第 2 の格納データと

を含むことを特徴とする請求項 35 に記載のコンピュータ読取り可能媒体。

【請求項 37】

20

前記動作は、

前記ページテーブルが無効であることを検出するステップと、

前記スパースファイル内の領域ごとに、前記領域が割り振られているかどうかを決定するステップと、

前記スパースファイル内の領域ごとに、前記領域が割り振られているかどうかに基づいて前記第 2 の格納データを設定するステップと

をさらに備えたことを特徴とする請求項 36 に記載のコンピュータ読取り可能媒体。

【請求項 38】

前記差分ストレージ内の特定のページにデータが格納されているかどうかの決定は、

前記第 1 の格納データをチェックし、前記特定のページが前記差分ストレージに格納されていることを、前記第 1 の格納データが示す場合、前記データは、前記差分ストレージ内の前記特定のページに格納されていることを決定するステップと、

30

前記第 2 の格納データをチェックし、前記領域が前記差分ストレージにおいて割り振られていないことを、前記第 2 の格納データが示す場合、前記データは、前記差分ストレージ内の前記特定のページに格納されていないことを決定するステップと、

前記ページが前記差分ストレージ内に格納されていることを、前記第 1 の格納データが示さず、前記領域が前記差分ストレージにおいて割り振られていないことを、前記第 2 の格納データが示さない場合、前記特定のページについての前記差分ストレージの対応するエリアからページデータを読み取り、前記対応するエリアからの前記ページデータが有効であるかどうかを決定するステップと

40

を含むことを特徴とする請求項 36 に記載のコンピュータ読取り可能媒体。

【請求項 39】

前記データベースは、ページテーブルをさらに含み、前記ページテーブルは、ページごとに、前記ページが前記差分ストレージ内に格納されているかどうかを示す第 1 の格納データを含むことを特徴とする請求項 33 に記載のコンピュータ読取り可能媒体。

【請求項 40】

データが前記差分ストレージ内の特定のページ内に格納されているかどうかの決定は、

前記第 1 の格納データをチェックし、前記特定のページが前記差分ストレージに格納されていることを、前記第 1 の格納データが示す場合、前記データは、前記差分ストレージ内の前記特定のページに格納されることを決定するステップと、

50

前記ページが前記差分ストレージ内に格納されていることを、前記第1の格納データが示さない場合、前記特定のページについての前記差分ストレージの対応するエリアからページデータを読み取り、前記対応するエリアからの前記ページデータが有効であるかどうかを決定するステップと

を含むことを特徴とする請求項39に記載のコンピュータ読取り可能媒体。

【請求項41】

前記動作は、

前記データベース内の特定のデータ要素の要求を受け付けるステップと、

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されているかどうかを決定するステップと、

10

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されている場合、前記差分ストレージを読み取ることによって前記要求に応答するステップと、

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されている場合、前記データベースを読み取ることによって前記要求に応答するステップと

をさらに含むことを特徴とする請求項27に記載のコンピュータ読取り可能媒体。

【請求項42】

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されているかどうかを決定する前記ステップは、前記差分ストレージが前記場所に有効なデータを含んでいるかどうかを決定するステップを含むことを特徴とする請求項41に記載のコンピュータ読取り可能媒体。

20

【請求項43】

データが前記差分ストレージ内の前記特定のデータ要素に対応する場所に格納されているかどうかを決定する前記動作は、ページテーブルを調べるステップを含むことを特徴とする請求項41に記載のコンピュータ読取り可能媒体。

【請求項44】

前記コンピュータ読取り可能媒体は、

第2の特定のデータ要素の代わりに前記データベース内のある場所に第1の特定の値を格納する前記データベースに加えられた修正を検出するステップと、

前記データベース内の対応する場所が有効なデータを含んでいるかどうかを決定するステップと、

30

前記データベース内の前記対応する場所が有効なデータを含んでいない場合、前記第2の特定のデータ要素を前記対応する場所に書き込むステップと

をさらに含むことを特徴とする請求項27に記載のコンピュータ読取り可能媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般的に、データベースシステムの分野に関する。より詳細には、本発明はトランザクションの整合性を保つデータベースのコピーの作成に関する。

【背景技術】

【0002】

40

データベースおよびデータベース製品は、今日よく使用されている。データベースは、レコード、テーブル、インデックスなどのオブジェクトに情報を格納する。データベースに格納されている現在のバージョンの情報に加えて、前バージョンの情報もユーザにとって有用となり得る。

【0003】

前バージョンのデータに関する情報を提供する1つの方法は、ユーザが関心を持つある時点のデータベースの完全なコピーを作成することである。これが行なわれるときは、データベースに関連するすべてのファイルが格納される。しかし、ある量の記憶スペースに格納されているデータベースでは、各コピーもその量の記憶スペースを必要とするため、この技術ではスペースを使いすぎる。またこの手法では、コピー自体が大量のデータの移

50

動に關与するため、時間もかかりすぎる。

【0004】

変更されたデータのコピーのみを格納する書き込み時コピー (copy-on-write) 機構を使用して、ディスクボリュームスナップショット (disk volume shapshot) を提供することができる。元のデータベースに変更が加えられるたびに、前に格納されていたデータが変更されたデータのコピーに書き込まれる。しかし、このボリュームスナップショットでは、トランザクションの整合性が保たれない。つまり、格納された変更は、進行中のトランザクションの一部である可能性があり、したがってボリュームスナップショットは、あるトランザクションに関する中途半端な情報を含んでいる可能性がある。こうしたトランザクションの不整合の可能性があるため、アプリケーションがボリュームスナップショットを使用するには特殊なフックが必要である。このように作成されたボリュームスナップショットは、格納されているインデックスおよびテーブルの構造上の修正が不完全であるため、構造的にも整合性のないものとなり得る。さらに、ボリュームスナップショットでは、コピーの粒度はボリュームレベルである。非データベースデータもコピーすることがあるが、これは不要であり、時間およびリソースの無駄につながる。

10

【0005】

また、特定のトランザクションによってアクセスされたデータを使用可能にすることによってデータベースのバージョンング (versioning) が提供されている。所与のトランザクションについてこのデータのみが格納され、したがって所与のトランザクションによってアクセスされたものに関するデータが必要な場合、(トランザクションによって任意の変更が加えられる前に) アクセスされたデータをそのトランザクションを開始したユーザに提供することはできる。しかしこのデータは、複数のトランザクションまたはユーザからはアクセスできない。複数のトランザクションまたはユーザが同じデータにアクセスを試みた場合、これらの異なるアクセスごとにデータが異なり得る。さらに、こうしたバージョンングは、データベースサーバの再起動後は永続的ではない。

20

【発明の開示】

【発明が解決しようとする課題】

【0006】

したがって、迅速に作成することができ、トランザクションの整合性を保ち、整合性のある情報を複数のトランザクションまたはユーザに提供し、データベースサーバの再起動後も永続性のあるデータベースを表示する方法が必要となる。

30

【課題を解決するための手段】

【0007】

データベースの完全なコピーを作成することなく、前のある時点 (データベースビューが作成されたとき) の既存のデータベースのトランザクションの整合性が保たれたビューを提供するデータベースのデータベースビューを作成する。1次データベースとは使用中のデータベースであり、1または複数のデータベースビューが作成される。

【0008】

各データベースビューは、1次データベースとともに、前のある時点の1次データベースの内容を決定するのに必要なすべての情報を含む。データベースビューは、1次データベース内の各データファイルに対応するサイドファイルから成る。サイドファイルは、データベースビューが作成されてから変更された対応するデータファイルからのすべてのデータのコピーを含む。

40

【0009】

一実装形態では、データベースのデータベースビューが実装される。データベースは、データ要素を含み、トランザクションログに關連付けられる。トランザクションログは、アクティブトランザクションおよび非アクティブトランザクションの両方を含む。データベースビューは、現在のバージョンのデータベースと前のバージョンのものとの差のストレージ (「差分ストレージ (difference storage)」) を含む。これは、サイドファイルとして知られているファイルに格納される。

50

【0010】

新しいデータがデータベースに書き込まれた場合は必ず、前バージョンのデータのコピーが差分ストレージに格納される。別の変更が加えられた場合、差分ストレージはすでに前バージョンのデータのコピーを含んでいるため、新しいデータは格納されない。

【0011】

一実施形態で、トランザクションログが分析され、前記トランザクションログにおける分割点(split point)が識別される。修正を行う、ログ内の分割点より前の各トランザクションが見つけれ、修正の結果が差分ストレージに書き込まれる。次いでアクティブトランザクションによってもたらされるトランザクションログに対する修正は、データベースビューを更新してトランザクションを元に戻すことによって、データベースビュー内で元に戻される。このように、分割点の時点でのデータベースのトランザクションの整合性を保つビューが提供される。

10

【0012】

本発明の他の特徴については後述する。

上記の概略および以下の好ましい実施形態の詳細な説明は、添付の図面と併せ読めば、よりよく理解できる。本発明を説明するために本発明の構造例を図面に示しているが、本発明は、開示されている特定の方法および手段に限定されるものではない。

【発明を実施するための最良の形態】

【0013】

(コンピューティング環境の例)

図1は、本発明の態様を実施できるコンピューティング環境の例を示している。コンピューティングシステム環境100は、適したコンピューティング環境の一例にすぎず、本発明の使用または機能の範囲に関する限定を示唆するものではない。また、コンピューティング環境100を、動作環境100の例に示した構成要素のいずれか1つ、またはその組合せに関連する依存性または必要条件を有しているものと解釈すべきではない。

20

【0014】

本発明は、他の多くの汎用または専用コンピューティングシステム環境または構成で動作可能である。本発明との使用に適したよく知られているコンピューティングシステム、環境、および/または構成の例には、それだけには限定されないが、パーソナルコンピュータ、サーバコンピュータ、ハンドヘルドまたはラップトップ装置、マルチプロセッサシステム、マイクロプロセッサベースのシステム、セットトップボックス、プログラム可能家庭用電化製品、ネットワークPC、ミニコンピュータ、メインフレームコンピュータ、組み込みシステム、上記の任意のシステムまたは装置を含む分散コンピューティング環境などがある。

30

【0015】

本発明は、コンピュータによって実行されるプログラムモジュールなどのコンピュータ実行可能命令の一般的な文脈で説明することができる。一般にプログラムモジュールは、特定のタスクを実行する、または特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、構成要素、データ構造などを含む。また、本発明は、タスクが通信ネットワークまたは他のデータ通信媒体によってリンクされているリモート処理装置によって実行される分散コンピューティング環境でも実施することができる。分散コンピューティング環境では、プログラムモジュールおよび他のデータを、メモリ記憶装置を含むローカルおよびリモートのコンピュータ記憶媒体に置くことができる。

40

【0016】

図1を参照すると、本発明を実施するシステムの例は、汎用コンピューティング装置をコンピュータ110の形で含んでいる。コンピュータ110の構成要素は、それだけには限定されないが、処理ユニット120、システムメモリ130、およびシステムメモリを含む様々なシステム構成要素を処理ユニット120に結合するシステムバス121を含む。処理ユニット120は、マルチスレッドプロセッサ上でサポートされるものなど、複数の論理処理ユニットを表し得る。システムバス121は、様々なバスアーキテクチャのう

50

ちの任意のものを使用するメモリバスまたはメモリコントローラ、周辺バス、およびローカルバスを含むいくつかのタイプのバス構造のうちどんなものでもよい。こうしたアーキテクチャには、それだけには限定されないが一例として、ISA (Industry Standard Architecture) バス、MCA (Micro Channel Architecture) バス、EISA (Enhanced ISA) バス、VESA (Video Electronics Standards Association) ローカルバス、および(メザンバスとしても知られている) PCI (Peripheral Component Interconnect) バスなどがある。システムバス121は、通信装置間のポイントツーポイント接続、スイッチングファブリックなどとして実装することもできる。

【0017】

コンピュータ110は、一般に様々なコンピュータ読取り可能媒体を含む。コンピュータ読取り可能媒体は、コンピュータ110からアクセスできる使用可能な任意の媒体とすることができ、揮発性および不揮発性媒体、リムーバブルおよび非リムーバブル媒体を含む。コンピュータ読取り可能媒体は、それだけには限定されないが一例として、コンピュータ記憶媒体および通信媒体を含み得る。コンピュータ記憶媒体には、コンピュータ可読命令、データ構造、プログラムモジュール、他のデータなど、情報を記憶するための任意の方法または技術で実施される揮発性および不揮発性のリムーバブルおよび非リムーバブル媒体がある。コンピュータ記憶媒体には、それだけには限定されないが、RAM、ROM、EEPROM、フラッシュメモリまたは他のメモリ技術、CD-ROM、デジタル多用途ディスク(DVD)または他の光ディスク記憶装置、磁気カセット、磁気テープ、磁気ディスク記憶装置または他の磁気記憶装置、または所望の情報の格納に使用でき、コンピュータ110からアクセスできる他の任意の媒体などがある。通信媒体は一般に、コンピュータ可読命令、データ構造、プログラムモジュール、または他のデータを搬送波または他の移送機構などの変調されたデータ信号に組み込む。これには任意の情報配送媒体がある。「変調されたデータ信号」という用語は、信号内の情報を符号化するように設定または変更された1または複数のその特徴を有する信号を意味する。通信媒体には、それだけには限定されないが一例として、有線ネットワーク、直接配線された接続などの有線媒体、および音響、RF、赤外線、その他の無線媒体などの無線媒体がある。また、上記のどんな組合せでもコンピュータ読取り可能媒体の範囲内に含まれるものとする。

【0018】

システムメモリ130は、読取り専用メモリ(ROM)131やランダムアクセスメモリ(RAM)132など、揮発性および/または不揮発性メモリの形のコンピュータ記憶媒体を含む。基本入出力システム133(BIOS)は、例えば起動中など、コンピュータ110内の要素間での情報の転送を助ける基本ルーチンを含み、一般にROM131に格納されている。RAM132は一般に、処理ユニット120から直接アクセス可能な、かつ/または処理ユニット120が現在処理しているデータおよび/またはプログラムモジュールを含む。図1は、それだけには限定されないが一例として、オペレーティングシステム134、アプリケーションプログラム135、他のプログラムモジュール136、およびプログラムデータ137を示している。

【0019】

コンピュータ110は、他のリムーバブル/非リムーバブル、揮発性/不揮発性コンピュータ記憶媒体を含むこともできる。一例にすぎないが、図1は、非リムーバブル不揮発性磁気媒体から読み取り、あるいはそこに書き込むハードディスクドライブ140、リムーバブル不揮発性磁気ディスク152から読み取り、あるいはそこに書き込む磁気ディスクドライブ151、およびCD-ROMや他の光媒体など、リムーバブル不揮発性光ディスク156から読み取り、あるいはそこに書き込む光ディスクドライブ155を示している。動作環境の例で使用できる他のリムーバブル/非リムーバブル、揮発性/不揮発性コンピュータ記憶媒体には、それだけには限定されないが、磁気テープカセット、フラッシュメモリカード、デジタル多用途ディスク、デジタルビデオテープ、半導体RAM、半導体ROMなどがある。ハードディスクドライブ141は一般に、インタフェース140などの非リムーバブルメモリインタフェースを介してシステムバス121に接続され、磁気

10

20

30

40

50

ディスクドライブ 1 5 1 および光ディスクドライブ 1 5 5 は一般に、インタフェース 1 5 0 などのリムーバブルメモリインタフェースによってシステムバス 1 2 1 に接続される。

【 0 0 2 0 】

上述し、図 1 に示したドライブおよびその関連のコンピュータ記憶媒体は、コンピュータ可読命令、データ構造、プログラムモジュール、およびコンピュータ 1 1 0 の他のデータの記憶を提供する。図 1 では例えば、ハードディスクドライブ 1 4 1 は、オペレーティングシステム 1 4 4、アプリケーションプログラム 1 4 5、他のプログラムモジュール 1 4 6、およびプログラムデータ 1 4 7 を格納するものとして示されている。これらの構成要素は、オペレーティングシステム 1 3 4、アプリケーションプログラム 1 3 5、他のプログラムモジュール 1 3 6、およびプログラムデータ 1 3 7 と同じであっても、異なっている。ユーザは、キーボード 1 6 2、および一般にマウス、トラックボール、またはタッチパッドと呼ばれるポインティング装置 1 6 1 などの入力装置を介してコマンドおよび情報をコンピュータ 1 1 0 に入力することができる。他の入力装置（図示せず）には、マイクロフォン、ジョイスティック、ゲームパッド、衛星パラボラアンテナ、スキャナなどがある。これらおよび他の入力装置は、しばしばシステムバスに結合されているユーザ入力インタフェース 1 6 0 を介して処理ユニット 1 2 0 に接続されるが、パラレルポート、ゲームポート、ユニバーサルシリアルバス（USB）など他のインタフェースおよびバス構造で接続してもよい。モニタ 1 9 1 または他のタイプの表示装置もまた、ビデオインタフェース 1 9 0 などのインタフェースを介してシステムバス 1 2 1 に接続される。モニタに加えて、コンピュータは、出力周辺インタフェース 1 9 5 を介して接続できるスピーカ 1 9 7、プリンタ 1 9 6 などの他の周辺出力装置を含むこともできる。

10

20

30

40

【 0 0 2 1 】

コンピュータ 1 1 0 は、リモートコンピュータ 1 8 0 など 1 または複数のリモートコンピュータへの論理接続を使用してネットワーク環境で動作することができる。リモートコンピュータ 1 8 0 は、パーソナルコンピュータ、サーバ、ルータ、ネットワーク PC、ピア装置、または他の一般のネットワークノードでよく、一般にコンピュータ 1 1 0 に関連して上述した多くまたはすべての要素を含むが、図 1 にはメモリ記憶装置 1 8 1 のみを示している。図 1 に示した論理接続は、ローカルエリアネットワーク（LAN）1 7 1 および広域エリアネットワーク（WAN）1 7 3 を含むが、他のネットワークを含んでいてもよい。こうしたネットワーキング環境は、オフィス、全社規模のコンピュータネットワーク、イントラネット、およびインターネットではごく一般的である。

【 0 0 2 2 】

LAN ネットワーキング環境で使用する場合、コンピュータ 1 1 0 は、ネットワークインタフェースまたはアダプタ 1 7 0 を介して LAN 1 7 1 に接続される。WAN ネットワーキング環境で使用する場合、コンピュータ 1 1 0 は一般に、モデム 1 7 2、またはインターネットなど WAN 1 7 3 を介して通信を確立する他の手段を含む。モデム 1 7 2 は、内蔵のものでも外付けのものでもよく、ユーザ入力インタフェース 1 6 0 または他の適切な機構を介してシステムバス 1 2 1 に接続することができる。ネットワーク環境では、コンピュータ 1 1 0 に関連して示したプログラムモジュール、またはその一部をリモートメモリ記憶装置に格納することができる。図 1 は、それだけには限定されないが一例として、リモートアプリケーションプログラム 1 8 5 をメモリ装置 1 8 1 上に存在するものとして示している。図示したネットワーク接続は例であり、コンピュータ間の通信リンクを確立する他の手段を使用してもよいことは理解されよう。

【 0 0 2 3 】

（データベースおよびデータベースビュー）

一般にデータベースは、データベースファイルおよびログファイルの 2 つの種類から成る。ログファイルは、ある期間にわたってデータベースファイルに追加された

50

変更を記載する一連のログレコードを含む。ログレコードは、ログシーケンス番号（LSN）で識別することができる。図2に示すように、一実施形態では、1次データベース200は、1組のデータベースファイル205およびログファイル210から成る。データファイルは、ページと呼ばれる複数のストレージの塊に分割される。

【0024】

データベースの完全なコピーを作成することなく、前のある時点の既存のデータベースのトランザクションの整合性が保たれたビューを提供するデータベースのデータベースビューが作成される。データベースビューは、データベースとともに、前のある時点のデータベースのコピーを生成するのに必要なすべての情報を含む。しかしデータベースビューは、それ自体情報のすべてを含んでいるわけではなく、したがって全コピーより小さいサイズになり得る。さらに、修正がデータベースに加えられると、その場でビューが作成され、それによってコスト（時間および処理）を時間にわたって分散することができる。前の時点でデータベースビューからコピーを作成した場合、時間および処理のコストは、一時に集中することになる。さらに、データベースに対する更新アクティビティが実行している間にデータベースビューを作成することができる。1次データベースとは使用中のデータベースであり、1または複数のデータベースビューが作成される。

10

【0025】

上述したように、データベースビューは、1次データベースとともに、前のある時点の1次データベースの内容を決定するのに必要なすべての情報を含む。データベースビューは、1次データベース内の各データファイルに対応するサイドファイルから成る。サイドファイルは、データベースビューが作成されてから変更された対応するデータファイルからのすべてのデータのコピーを含む。一実施形態では、サイドファイル内のページから1次ファイル内のページにテーブルをマッピングする必要を回避するために、サイドファイルがスパースファイル（sparse file）に格納される。スパースファイル内の実際に書き込まれるファイルの部分のみが記憶スペースを必要とする。ファイルの他のすべての領域は割り振りされていない。他の実施形態では、サイドファイルのストレージは、スパースファイル内にはない。

20

【0026】

一実施形態では、スパースファイル機構は、標準領域サイズで動作する。1つの領域内のデータがスパースファイルに書き込まれると、たとえデータが領域全体を埋めなくても、領域全体のためのスペースが割り振られる。このスペースが割り振られ、そこからの読み取りが可能であるため、有効な値で埋められた領域のエリアと、その領域内の任意のストレージが必要な場合にスパースファイルの粒度によってあるサイズの領域が割り振られることが求められるために存在する領域のエリアとの間に区別をつける必要がある。

30

【0027】

データベースビューは、データベースビューが作成されてから1次データベースで変更されたすべてのデータの元の値を含んでいるため、データベースビューの作成時点のデータベースのデータは、データベースビューから読み取ることができる。データベースビューからのデータの要求に応答するため、サイドファイルが要求のデータを含んでいる場合、データは、データベースビューのサイドファイルから読み取られる。読み取るべきデータであってサイドファイルに存在しないものは、データベースビューが作成されてから変更されておらず、1次データベースから読み取られる。

40

【0028】

一実施形態で、サイドファイルは、1次データベースからのデータのページを含む。1次データベースの任意のページの任意のデータが変更されると、そのデータのページがサイドファイルに格納される。本発明では、1次データベース内のデータの単位としてページについて説明しているが、1次データベースのデータの他の単位を使用することもできることを企図している。

【0029】

どのデータがサイドファイルに書き込まれたか、またどのデータを1次データベースか

50

ら読み取るべきかを決定するために、サイドファイル内の有効なデータの存在を確認する必要がある。一実施形態では、有効なデータが存在するかどうかを確認するためにサイドファイルを直接読み取る。別の実施形態では、サイドページテーブルが作成され、所与のページが存在するかどうか、また有効かどうかに関するデータが格納される。

【0030】

一実施形態では、1次データベース内のページごとに、サイドページテーブルは、ページは変更されておらず、それを1次データベースから読み取るべきか、またはページが変更されているのでそれをサイドファイルから読み取るべきかに関する情報を格納する。サイドページテーブルによって、所与のページがサイドファイルに存在するかどうかの迅速な決定が可能になる。

10

【0031】

(1ビットおよび2ビットのページテーブル)

一実施形態では、サイドファイルおよびスパーズファイル機構は、いずれも同じページ/領域サイズを使用する。つまり、サイドファイルが1次データベースから格納したページは、任意のメモリがスパーズファイルに書き込まれたときにスパーズファイルが格納した領域と同じサイズである。例えば、スパーズファイル領域が8KBであり、1次データベースから格納されたページも8KBである場合、ページサイズおよび領域サイズは等しい。この場合、埋められる任意の領域は、1次データベースから読み取られたページによって完全に埋められ、無効なデータがその領域に格納される可能性はない。

【0032】

20

別の実施形態では、いくつかのサイドファイル領域が各ページに正確に対応する。スパーズファイル領域は8KB(キロバイト)であり、1次データベースから格納されたページは16KBである場合、サイドファイルに格納される各ページによって2つの領域が埋められる。この場合も、埋められる任意の領域は、1次データベースから読み取られたページからの内容によって完全に埋められる。この場合もまた、無効なデータがその領域に格納される可能性はない。

【0033】

これらの実施形態の場合、サイドページテーブルは、サイドファイル内のページごとに1ビットの情報を保持するメモリ内ビットマップを含む。サイドファイル内のページごとに、対応するビットは、ページがサイドファイル内にあるかどうかを示す。

30

【0034】

別の実施形態では、サイドファイル領域の粒度は、1次データベースから格納されたページの粒度より大きい。例えば、サイドファイルの各領域は64KBであり、ページのサイズは8KBである場合、サイドファイル内の領域の存在は、必ずしもその領域内のすべての情報が1次データベースからの有効なデータであることを示すとは限らない。1つのページのみがサイドファイルにコピーされる場合、この例では、割り振られた領域内の64KBのうち8KBのみが有効なデータを含む。別の実施形態では、一部のサイドファイルページが領域にわたって分散される。

【0035】

これらの実施形態では、サイドページテーブルは、サイドファイル内のページごとに2ビットの情報を保持する2つのメモリ内ビットマップを含む。これらをビット1およびビット2と呼ぶ。サイドファイル内のページごとに、対応するビットは、(ビット1)ページが確かにサイドファイル内にあるかどうか、(ビット2)ページがサイドファイル内にある可能性があるかどうかを示す。ビット2は、ページがサイドファイル内に格納される領域が割り振られていることを示すものとも考えられる。しかし、後述するように、一実施形態では、このビット2は、サイドページテーブルが再構築されたときにのみ設定される。

40

【0036】

ビットマップは、メモリ内で維持され、したがって永続的ではない可能性がある。ビットマップは、消去されると、スパーズファイル情報から再構築される。スパーズファイル

50

を調べ、ページごとに、サイドファイルでページが配置される領域にメモリが割り振られている場合、ビット2は、そのページがサイドファイル内にある可能性があることを示すように設定される。ページごとに、ビット1は、ページがサイドファイル内にあることを確定しないことを示すように設定される。

【0037】

サイドページテーブルが永続的となるように維持されると、領域およびページの粒度を無視することができ、1ビットのサイドページテーブルを使用することができる。しかし一実施形態では、データベースサーバの再起動後、永続的なデータベースビューをサポートするために、2ビットのページテーブルを使用する。

【0038】

一実施形態では、サイドファイルのためのページテーブルは作成されない。この場合、コピーがデータベースビュー内のページから作成されているかどうかを決定することが必要なときはいつでも、データベースビューを調べる。本発明では、1ビットまたは2ビットのページテーブルが存在する実施形態に関して後述するが、ページテーブルがなく、データベースビューを検査して、1次データベースからコピーされたページが含まれているかどうかを決定する必要がある実施形態についても検討する。

【0039】

図3に示すように、1次データベース200のデータベースビュー220は、サイドファイル225から成る。1次データベース200内のデータファイル205のそれぞれは、データベースビュー220内に対応するサイドファイル225を有する。さらに、サイドページテーブルデータ230は、データベースビュー220のメモリ内に格納される。一実施形態では、サイドページテーブルデータ230は、サイドファイル225のすべてをカバーする1つのサイドページテーブルである。別の実施形態では、個別のサイドページテーブルがサイドファイル225ごとに存在する。

【0040】

(トランザクションログ)

データベースでは、トランザクションログとは、トランザクションログを最後にバックアップしてからデータベースに対して実行されたすべてのトランザクションのシリアルレコードである。トランザクションログを使用して、障害時にデータベースを回復する。一実施形態では、トランザクションログは循環キューとしてモデル化される。トランザクションログは、ログの非アクティブな部分を削除することによって切り捨てることができる。この非アクティブ部分は、回復する必要のない完了したトランザクションを含む。逆に、トランザクションログのアクティブ部分は、完了したトランザクション、および依然として稼動しており、まだ完了していないトランザクション(アクティブトランザクション)を含む。切り捨ては、トランザクションログが拡大し続け、より多くのスペースを使用できるようにする代わりに、トランザクションログ内の非アクティブなスペースを最小限に抑えるために行われる。

【0041】

アクティブトランザクションは、トランザクションの不整合を引き起こす可能性がある。アクティブトランザクションでは、データファイルの一部の修正がバッファキャッシュからデータファイルに書き込まれなかった可能性があり、データファイル内に未完了のトランザクションからの一部の修正がある可能性がある。ログファイルを使用して、データベースの回復によって確実にトランザクションの整合性が保たれるようにする。これは、A R I E S (Algorithms for Recovery and Isolation Exploiting Semantics) 型の回復を使用して行われる。データファイルに書き込まれていない可能性のある、ログに記録されているすべての修正は、データベースへの修正を行うことによってロールフォワードされる。データベースの保全性を確保するために、トランザクションログ内にある未完了のすべてのトランザクションは、データベースに対する修正を元に戻すことによってロールバックされる。

【0042】

10

20

30

40

50

(データベースビューの作成)

データベースビューを作成するために、データベースビューの物理的構造(サイドファイルおよびページテーブル)を初期化する必要がある。まず、サイドファイル225が1次データベース200内のデータファイル205ごとに作成される。上述したように、サイドファイルはスパーファイルでよい。あるいは別の実施形態では、データファイル205と同じサイズの非スパーファイルとすることができる。サイドファイル225は、1次データベース200内のデータファイル205に関連付けられる。

【0043】

トランザクションは継続的に行われており、データベースビューではトランザクションの整合性が保たれるため、データベースビューの作成中、トランザクションログを使用する必要がある。データベースビューに使用すべきトランザクションに関する情報が確実に破棄されないようにするために、1次データベース200において、(存在する場合は)ログの切り捨てを使用不可にする。

10

【0044】

一実施形態では、サイドページテーブル230は、データベースビューのために初期化される。最初に、サイドページテーブルは、サイドファイル225内にページが存在しないことを示すように設定され、2ビットのサイドページテーブルの場合、サイドファイル225にページが潜在的に、または確実に存在しないことを示すように設定される。

【0045】

初期化が完了すると、データベースビューは、「オンライン」になる用意ができてい

る。データベースビューは、この時点で1次データベース200と平行して稼動しており、修正が行われると、修正されたページの元の値のコピー(すなわち更新が行われる前のページの内容)がデータベースビューに格納される。図5は、本発明の一実施形態によるデータベースのトランザクションの整合性を保つビューを実装する方法のフロー図を示している。図5のステップ500に示すように、トランザクションログにおいて分割点が決定される。この分割点は、データベースビューが表す時点に対応する。データベースビューが作成されると、1次データベース200上のログの端部のLSNが取得される。このLSNが、1次データベース200およびデータベースビュー220が分岐を開始する「分割点」である。次いで1次データベース200は、データベースビュー処理が要求されるようにマーク付けされる。後述するように、1次データベース200内のデータベースビューサポートが開始する。

20

30

【0046】

データベースビューが整合性を保つようにするために、分割点より前の1次データベース200のログを分析して、分割点の時点でどのトランザクションがアクティブであったかを決定する必要がある。ログ内で(分割点の時点で)最も古いアクティブトランザクションが識別される。その最も古いアクティブトランザクションの前にログの切り捨てが使用可能となる。

【0047】

ARIES型回復と同じように、分割点より前の最も古いアクティブトランザクションからの1次データベース200のログ内のすべての操作がデータベースビューに対して行

われる。図4は、本発明の好ましい実施形態によるトランザクションログの例、ログファイル210を示すブロック図である。ログファイル210内のログエントリは、ログエントリ400、410、420、430、440、450、460、480、490、および499を含む。分割点475が決定される。トランザクションは、引き続きログに書き込まれ、しかし切り捨ては使用不可である。ログファイル210を検査し、最も古いアクティブトランザクションから分割点まで(図4の例ではログエントリn400からログエントリn+7まで)のトランザクションの結果、データベースに対する任意の修正がサイドファイル225に対して行われる。これらの各トランザクションにおける修正の結果がサイドファイル225に格納される。次いでこれらのトランザクションを検査する。ログ内の任意のアクティブトランザクションによってログファイルに書き込まれた修正、例え

40

50

ばログエントリ n 4 0 0、ログエントリ n + 2 4 2 0、ログエントリ n + 6 がサイドファイル 2 2 5 において元に戻される。

【 0 0 4 8 】

図 5 のステップ 5 0 0 からわかるように、トランザクションログにおいて分割点を選択する。次にステップ 5 1 0 で、データベースに対して修正を行う各トランザクションを見つめる。一実施形態では、トランザクションを分析し、トランザクションの結果としてある値がデータベースのある場所に書き込まれる場合、ステップ 5 2 0 からわかるように、以下で詳述するデータベースビューを修正する方法を使用して、その値がサイドファイル内の対応する場所に保存される。このように、データベースに書き込む必要のあるすべての変更（「ダーティページ」の変更など）がデータベースビューのために格納される。

10

【 0 0 4 9 】

しかし、一部のトランザクションは、依然としてコミットされていない。したがって、ログ内の分割点までのこれらのアクティブトランザクションを探し出し（ステップ 5 3 0）、元に戻す（ステップ 5 4 0）必要がある。一実施形態で、未完了のトランザクションがデータベース内のある場所の値を変更する場合、上記でサイドファイルに追加された変更は、サイドファイルから削除される。代替実施形態では、後述するように、データベースビューを修正し、サイドファイル内のデータを分割点の時点でのデータベース内のデータと一致するように設定することによって、トランザクションが元に戻される。

【 0 0 5 0 】

このように、ログからの完了したトランザクションのみがデータベースビューに反映される。ログ上のトランザクションがデータベースビューにおいて反映されると、元に戻された、分割点が生じたときにアクティブであったトランザクションを除いて、1次データベース 2 0 0 上でログの切り捨てが使用可能になる。データベースビュー処理は使用可能であったため、1次データベース 2 0 0 に変更が加えられると、データベースビューが更新されるので、データベースビューを使用して分割点の時点での1次データベース 2 0 0 の内容を決定することができる。

20

【 0 0 5 1 】

（データベースビューの回復）

データベースサーバが（正常または異常に）シャットダウンした後で再起動すると、データベースビューを再度初期化する必要がある。そのために、メモリ内に格納されているサイドページテーブルを再度初期化する必要がある。

30

【 0 0 5 2 】

サイドページテーブルを再初期化するには、2ビットのサイドページテーブルの実装では、割り振られたサイドページテーブル内の領域ごとに、割り振られた領域内のページごとのサイドページテーブル内のデータ（ビット 2）は、ページがサイドファイル 2 2 5 に書き込まれた可能性があることを示すように設定される。他のすべてのページのサイドページテーブル内のデータは、ページがサイドファイル 2 2 5 に書き込まれた可能性がないことを示すように設定される。ただし、ページがサイドファイル 2 2 5 に書き込まれたことを確定するものではなく、したがって、ビット 1 は、最初は設定されない。

【 0 0 5 3 】

あるいは、2ビットのサイドページテーブルの実装、または1ビットのサイドページテーブルの実装では、上述したように、サイドファイル 2 2 5 を検査して、ページごとに、サイドファイル 2 2 5 内のページが有効であるかどうかを決定することができる。ページテーブルは、存在しているページごとに、ページが実際にサイドファイル 2 2 5 に存在することを示すように設定される。すべての他のページは、ページがサイドファイル 2 2 5 に存在しないことを示すように設定される。

40

【 0 0 5 4 】

（1次データベース内のデータベースビューサポート）

データが上書きされる前にデータベースビューが1次データベース 2 0 0 から情報を格納するようにするために、1次データベース 2 0 0 は、データベースビューの作成をサポ

50

ートする必要がある。1次データベース200が修正するページごとに、ページがデータベースビュー内にあるかどうかに関して決定を行う必要がある。データベースビュー内にページが存在する場合、それは正しいバージョンのページである。例えばこれは、前の修正が1次データベース200内のそのページに加えられたときであり得る。ページが1次200内で再度変更されると、データベースビューのバージョンを変更する必要がある。

【0055】

ページが変更されているという情報を1次データベース200から受信したとき、ページがサイドファイル225にある場合は何もする必要がない。ページがサイドファイル225内にない場合、ページをサイドファイル225に書き込む必要があり、サイドページテーブルで適切なビットを設定する必要がある。2ビットのページテーブルがある場合、次の表1で示すように、そのページでのビット1およびビット2の場合で3つの可能性がある。

10

【0056】

【表1】

表1 2ビットのページテーブルの様々なケース

	ビット1はページがサイドファイル内に確実にあることを示す	ビット1はページがサイドファイル内に確実にあることを示さない
ビット2はページがサイドファイル内にある可能性があることを示す	ケース1: ページはサイドファイル内にある	ケース2: ページはサイドファイル内にある可能性がある
ビット2はページがサイドファイル内に絶対になことを示す	ケース1: ページはサイドファイル内にある [あるいはケース4: 無効]	ケース3: ページは絶対にサイドファイル内にない

20

【0057】

一実施形態では、ページがサイドファイル225内に確実にあることを、ビット1が示すとき、ビット2は無視される。したがって表1に示すように、ページがサイドファイル225内に確実にあることを、ビット1が示す場合には、ビット2がどのように示していても、ページはサイドファイル225内にあるものとみなされる。代替実施形態では、ページがサイドファイル225に確実に存在することを、ビット1が示すように設定されると、ビット2は、ページがサイドファイル225内にある可能性があることを示すように設定される。またこの代替実施形態では、ページはサイドファイル225に確実にあることを、ビット1が示しているが、ビット2はサイドファイル225に絶対になことを示しているとき、このケースは無効であり、エラーが生じる。

30

【0058】

1次データベース200は、2ビットのページテーブルの場合にページが変更されていることを示しているとき、上述の場合にとるべき動作は次の通りである。

ケース1: 何も行わない

40

ケース2: ページがサイドファイル225内にあるかどうかを決定し、ない場合はページをサイドファイル225に書き込む

ケース3: ページをサイドファイル225に書き込む。

【0059】

ページがサイドファイル225に書き込まれるとき、ケース1またはケース2のいずれかで、1次データベース200内の古いバージョン(現在1次データベース200によって修正されているバージョン)のページがサイドファイル225に書き込まれる。さらに、ページテーブルは、ページが現在サイドファイル225内にあることを示すように設定され、その後のページへのすべての書き込みはケース1に従って処理され、データベースビューに適切なページは、サイドファイル225に格納されたままである。

50

【0060】

ケース2において、ページがサイドファイル225内にあるかどうかを決定するために、ページに対応するデータがサイドファイル225から読み取られる。データが有効であり、前のバージョンのページがサイドファイル225内にある場合、それを上書きする必要はない。一実施形態では、ページに対応するページテーブルのビット1は、ページがサイドファイル225に確実に存在することを示すように設定されるため、ページへの将来の書き込みはケース1に基づいて処理される。

【0061】

有効なデータはまだその領域に書き込まれていないことを示すために、新しく割り振られた領域に配置されたデータによってデータの無効を示すことができる。例えば、すべてゼロしかないデータベースのページがないことがわかっている場合、新しく割り振られた領域にすべてゼロを書き込むことができる。そうである場合、サイドファイル225内のページの存在は、割り振られた領域の一部であり、ゼロ以外の何らかのデータを含むサイドファイル225内の対応するページによって示される。

【0062】

(データベースビューの読み取り)

表1に詳しく示したケースは、データベースビューに格納されているデータの読み取りを行うのにも有用である。ページ内のデータをデータベースビューから読み取るとき、ページがサイドファイル225内に存在している場合は、そこから読み取る必要がある。そこにはない場合、ページは1次データベース200から読み取る必要がある。2ビットのページテーブルシステムでは、3つのケースの場合にとるべきアクションは、次の通りである。

ケース1：サイドファイル225からページを読み取る

ケース2：ページがサイドファイル225内にあるかどうかを決定し、そこにある場合はサイドファイル225からページを読み取り、そこにはない場合は1次データベース200からページを読み取る

ケース3：1次データベース200からページを読み取る。

【0063】

(データベースビューの修正)

データベースビューは、前のある時点でのデータベースの状態を表す。ユーザは、そのデータベースビューをデータベースとして使用することを選択することができる。例えばユーザは、前の時点のデータベースビューに対してそのアクションが行われたかのように、データベースビューに対してアクションを行い、データベースのデータベースビューを作成することができる。さらに、初期化中、上記で詳述したように、データベースビューに対してトランザクションを実行し、元に戻すこともできる。

【0064】

データベースビューを修正するには、修正は、データベースビュー内のデータに基づく必要がある、その結果得られたページをデータベースビューに格納する必要がある。ページのデータベースビューにデータが存在しない場合、修正は、1次データベース200内のデータに基づく必要がある、その結果得られたページをデータベースビューに格納する必要がある。

【0065】

2ビットのページテーブルシステムでは、3つのケースの場合にとるべきアクションは、以下の通りである。

ケース1：サイドファイル225からページを読み取り、修正を行い、ページをサイドファイル225に書き込む

ケース2：ページがサイドファイル225内に存在するかどうかを決定し、存在する場合はケース1と同様に処理し、存在しない場合はケース3と同様に処理する

ケース3：1次データベース200からページを読み取り、サイドファイル225にページを書き込み、ページがサイドファイル225にあることを示すようにページテーブ

10

20

30

40

50

ルを設定する。ページへの修正を行い、適切な場合はサイドファイル 2 2 5 に修正されたページを書き込む。

【 0 0 6 6 】

上記の例は、単に説明の目的で提供したものであり、決して本発明を限定するものと解釈されるものではないことに注意されたい。本発明を様々な実施形態に関連して説明してきたが、本明細書で使用した単語は、限定するものではなく説明および例示のためのものであることを理解されたい。さらに、本明細書において本発明を特定的手段、材料、および実施形態に関連して説明してきたが、本発明は、本明細書に開示した詳細に限定されるものではなく、添付の特許請求の範囲内に含まれるものなど機能的に等価なすべての構造、方法、用途にまで拡張される。本明細書の教示の恩恵を受ける当分野の技術者は、それに様々な修正を行うことができ、本発明の態様の範囲および意図から逸脱することなく変更を加えることができる。

10

【 図面の簡単な説明 】

【 0 0 6 7 】

【 図 1 】 本発明の態様を実施できるコンピューティング環境の例を示すブロック図である。

【 図 2 】 本発明の一実施形態によるデータベースを示すブロック図である。

【 図 3 】 本発明の一実施形態によるデータベースビューおよびデータベースを示すブロック図である。

【 図 4 】 本発明の一実施形態によるトランザクションログの例を示すブロック図である。

20

【 図 5 】 本発明の一実施形態によるトランザクションの整合性を保つデータベースビューを実装する方法を示すフロー図である。

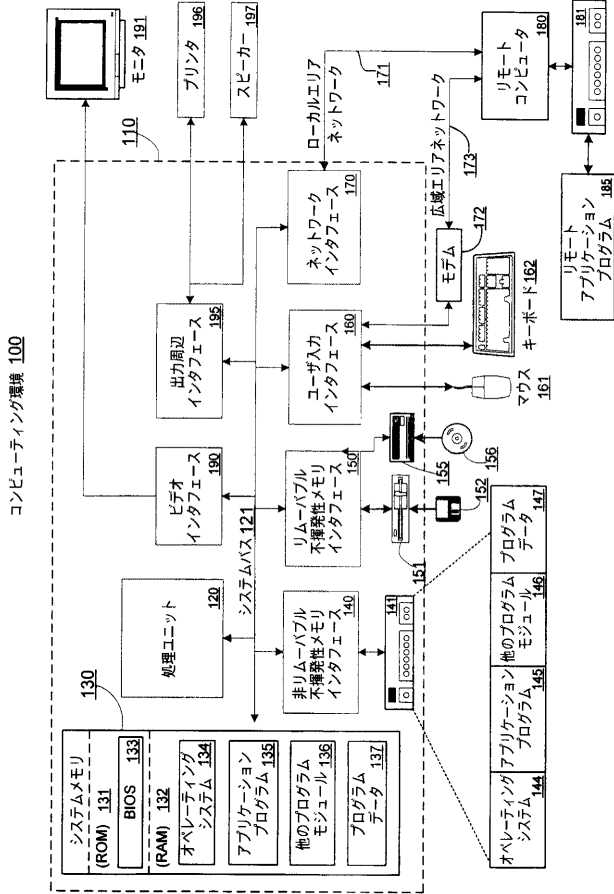
【 符号の説明 】

【 0 0 6 8 】

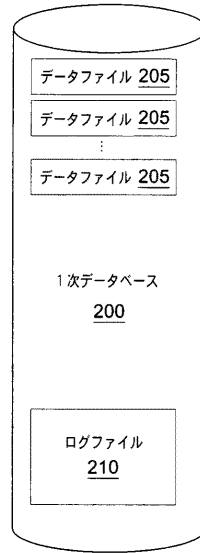
2 0 0 1 次データベース
 2 0 5 データファイル
 2 1 0 ログファイル
 2 2 0 データベースビュー
 2 2 5 サイドファイル
 2 3 0 サイドページテーブルデータ

30

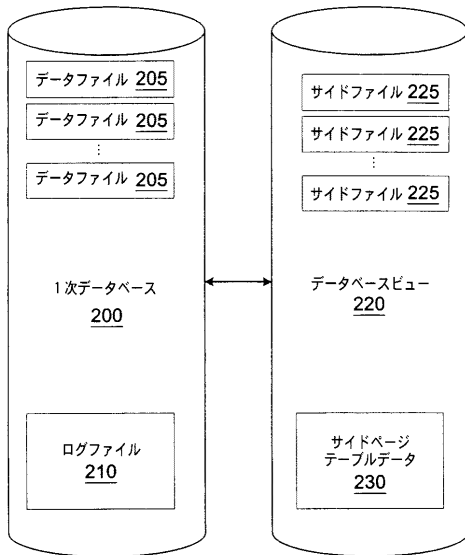
【図1】



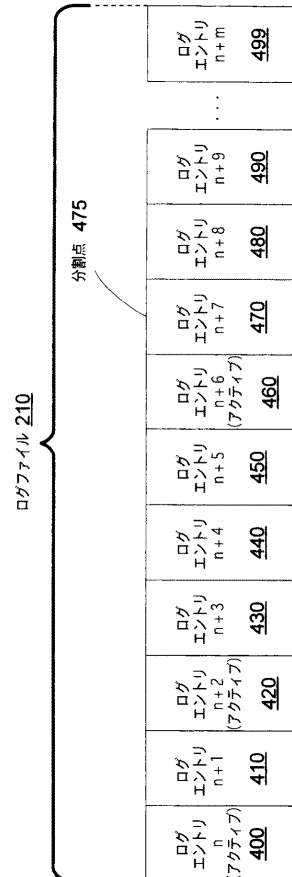
【図2】



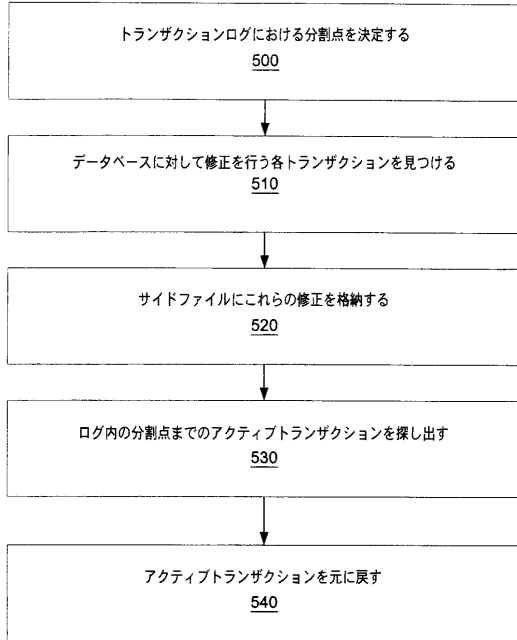
【図3】



【図4】



【 図 5 】



フロントページの続き

(72)発明者 ルイス エス・ブルック

アメリカ合衆国 98053 ワシントン州 レッドモンド ノースイースト 60 プレイス
21009

(72)発明者 サミート エイチ・アガーウォル

アメリカ合衆国 98052 ワシントン州 レッドモンド 149 プレイス ノースイースト
8127 ユニット ビー309

(72)発明者 ヤン カンロン

アメリカ合衆国 98029 ワシントン州 イサコア サウスイースト 41 ドライブ 25
001

Fターム(参考) 5B082 GB06