

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5021677号
(P5021677)

(45) 発行日 平成24年9月12日(2012.9.12)

(24) 登録日 平成24年6月22日(2012.6.22)

(51) Int.Cl.

F I

G 0 6 F 12/00 (2006.01)

G 0 6 F 12/00 5 3 5 P

請求項の数 9 (全 41 頁)

(21) 出願番号 特願2008-548512 (P2008-548512)
 (86) (22) 出願日 平成18年11月17日(2006.11.17)
 (65) 公表番号 特表2009-522638 (P2009-522638A)
 (43) 公表日 平成21年6月11日(2009.6.11)
 (86) 国際出願番号 PCT/US2006/044735
 (87) 国際公開番号 W02007/078444
 (87) 国際公開日 平成19年7月12日(2007.7.12)
 審査請求日 平成21年11月10日(2009.11.10)
 (31) 優先権主張番号 11/275,434
 (32) 優先日 平成17年12月30日(2005.12.30)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 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 ジョナサン アール、ハウエル
 アメリカ合衆国 98052 ワシントン
 州 レッドモンド ワン マイクロソフト
 ウェイ マイクロソフト コーポレーシ
 ョン インターナショナル パテント内

最終頁に続く

(54) 【発明の名称】 デルタページャを使用した状態の管理

(57) 【特許請求の範囲】

【請求項 1】

複数のレコードを有するデータベースであって、各レコードはある値へのアドレスの割り当てを含み、前記データベースは前記データベースの現行状態を追跡するための現行状態ポインタを有する、データベースを格納する記憶装置を有するコンピュータが、

前記データベースに対する複数の命令を含むトランザクションを受け取るステップと、
 第1の書き込みバッファを作成し、当該第1の書き込みバッファから、前記現行状態ポインタによって指示されるオブジェクトへのポインタを確立するステップと、

前記複数の命令の各命令について、当該命令が前記データベースに変更を適用する場合、追加の書き込みバッファを作成し、当該追加の書き込みバッファから、前記トランザクションの実行により作成された書き込みバッファのうち、直前に作成された書き込みバッファへのポインタを確立し、前記命令が前記データベースに適用しようとする変更を前記追加の書き込みバッファに格納するステップと、

前記現行状態ポインタに、前記トランザクションの実行により作成された書き込みバッファのうち、最後に作成された書き込みバッファを指示させることにより、前記トランザクションをコミットするステップと、

合体開始ポイントを識別するステップと、

前記合体開始ポイントと次の合体されていないオブジェクトとの間に合体済みオブジェクトを作成するステップと、

合体する可能性があるオブジェクトが存在しなくなるまで、前記合体済みオブジェクト

10

20

の次のオブジェクトに進み、当該次のオブジェクトに格納された変更を、前記合体済みオブジェクトに格納し、当該次のオブジェクトを解放するステップと、

前記現行状態ポインタに、前記合体済みオブジェクトを指示させるステップと
を含むことを特徴とする方法。

【請求項 2】

前記データベースにアクセスし、特定のアドレスの値を求めるステップであって、
前記特定のアドレスの割り当てが見つかるまで、または前記データベースに到達するま
で、前記現行状態ポインタによって指示されるオブジェクトから開始し、当該オブジェク
トからのポインタによって指示される次のオブジェクトに進む、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。

10

【請求項 3】

前記コミットするステップは、
前記第 1 の書き込みバッファからのポインタ、および前記現行状態ポインタが同じオブ
ジェクトを指示する場合、または
介入トランザクションが、前記トランザクションがアクセスしたアドレスの値を変更し
ていない場合に、実行されることを特徴とする請求項 1 に記載の方法。

【請求項 4】

前記トランザクションがコミットされない場合、前記トランザクションの実行により作
成された書き込みバッファを破棄することによって、前記トランザクションを打ち切るス
テップをさらに含むことを特徴とする請求項 3 に記載の方法。

20

【請求項 5】

現行状態を保持する命令に応答して、当該保持する命令を受けた時点で前記現行常状態
ポインタによって指示されるオブジェクトを指示するスナップショットポインタを作成す
るステップを含み、

前記次のオブジェクトがスナップショットポインタによって指示される場合、前記スナ
ップショットに前記合体済みオブジェクトを指示させ、前記合体済みオブジェクトと次の
合体されていないオブジェクトとの間に追加の合体済みオブジェクトを作成することを特
徴とする請求項 1 に記載の方法。

【請求項 6】

前記記憶装置が低速ストアであり、前記トランザクションの実行により作成されたバッ
ファおよび合体済みオブジェクトが高速ストアに格納されている場合に、

30

前記低速ストアに格納されたデータベースのキャッシュオブジェクトを前記高速ストア
に作成するステップと、

前記キャッシュオブジェクトから前記データベースへのポインタを確立するステップと
、

前記データベースを指示する、前記高速ストア内のオブジェクトからのポインタに、前
記キャッシュオブジェクトを指示させるステップと

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 7】

前記高速ストアはメモリを含み、前記低速ストアはディスクストレージを含むことを特
徴とする請求項 6 に記載の方法。

40

【請求項 8】

前記データベース内の少なくとも 1 つのアドレスの割り当てを前記キャッシュオブジェ
クトにコピーするステップと、

前記キャッシュオブジェクト内の少なくとも 1 つのアドレスの割り当てを前記データベ
ースにコピーするステップと、

のうちの少なくとも 1 つをさらに含むことを特徴とする請求項 6 に記載の方法。

【請求項 9】

新しいアドレスの割り当てを前記キャッシュオブジェクトに書き込むと同時に、前記新
しい割り当てを前記データベースに書き込むステップをさらに含むことを特徴とする請求

50

項 6 に記載の方法。

【発明の詳細な説明】

【背景技術】

【0001】

データベーストランザクションを実行する際の重要な問題は、データベースの整合性である。データベースの整合性を保持するためには、トランザクションが原子性 (a t o m i c i t y) および分離性などの諸要件を確実に監視することが重要である。原子性要件は、トランザクションのタスクのすべてが実行されるか、またはそれらがまったく実行されないかの、いずれかであることを要求するため、不完全なトランザクションはデータベースに適用されない。分離性要件は、トランザクションが他のトランザクションとは別々に処理されることを要求するため、トランザクションが、他のトランザクションの中間状態にアクセスすること、またはこれに干渉することはない。原子性要件および分離性要件を監視することによって、トランザクションが互いに破損し合うこと、およびデータベース全体を破損することを防ぐ。

10

【0002】

ネットワークシステムでは、多くのユーザおよび多くのアプリケーションが同時に同じデータ上でトランザクションを実行しようとする場合がある。たとえスタンドアロン型システムであっても、複数のアプリケーションが同時に同じデータ上でトランザクションを実行しようとする場合があり、各アプリケーションの複数の処理スレッドが競合し合う可能性がある。トランザクションが競合または矛盾する変更をデータベースに適用する可能性があるすべての状況を防止することはもちろん、予測することも困難であるため、データベースの整合性を維持することは、重要な課題である。

20

【0003】

データベースの整合性を保持する方法の1つが、現行トランザクションによって使用されているデータベースの一部分をロックすることである。データベースの一部分をロックすることにより、現行トランザクションが完了するまで、他のトランザクションによるデータの読み取りまたは上書きを防止する。他のトランザクションによるデータベースの一部分の読み取りを防止することで、他のトランザクションが、現行トランザクションが変更する可能性のある値に基づく結果に達しないことが保証される。同様に、他のトランザクションがデータベースのその部分を上書きするのを防止することで、他のトランザクションが現行トランザクションに影響を与えないことが保証される。

30

【0004】

しかしながら、データベースの一部分をロックすることは、データの整合性を保持するのを助ける可能性はあるが、それに対抗するコストが生じる可能性もある。いくつかのトランザクションがアクセスしようとしているデータベースの一部分をロックすることで、トランザクションのバックログが大量に生じる可能性がある。待機中のトランザクションの一部またはすべてが、データベースの現行状態を変更しない可能性があるか、または、現行トランザクションが使用している特定の値からの読み取り、この値への書き込み、またはこの値の変更を行わない可能性がある。それでもやはり、データベースの一部分をロックすることは、たとえデータベースへのアクセスが遅くなっても、データベースの整合性を保持することになる。

40

【0005】

ロックの1つの形がリースである。リースとは、アパート、ビル、車、および設備のリースの場合とまったく同様の、限定期間にわたる独占許可である。したがって、あるトランザクションが選択されたデータに関するリースを認められた場合、そのトランザクションには、ある限定期間にわたる選択されたデータへの独占的なアクセス権が与えられる。データへの独占的アクセス権をある期間に限定することによって、トランザクションが完了した際にトランザクションがデータの解放に失敗した場合、またはトランザクションが実行されているシステムがクラッシュした場合、データはロック解除され、他のトランザクションはデータへのアクセスを無期限に待つ必要がなくなる。通常、選択されたデータ

50

を待っている各トランザクションは、事前にデータに関するキューに入れられたすべてのトランザクションがそのデータの使用を完了するまで、または各トランザクションに割り振られたリースが満了するまで、待つ必要がなくなる。

【0006】

リースされたデータへのアクセスをいくつかのトランザクションが待っている場合、潜在的にかなりの遅延が存在する可能性がある。加えて、データのリースを制御しているシステムは、データにアクセスしようとするすると障害を引き起こす可能性がある。このリースを制御しているシステムは、データアクセスに関する単一の制御ポイントとして、リース要求に応答しようとするするとトランザクションの遅延を発生させる可能性がある。さらにこのシステムは、選択された同じデータに関する複数の要求を受信している可能性があるため、これら多数の潜在的に繰り返される要求を処理することは、計算サイクルの無駄となり、結果として、同じデータまたは任意の他のデータを求めるトランザクションにさらなる遅延を発生させる。

【発明の開示】

【課題を解決するための手段】

【0007】

デルタページャ(delta pager)とは、原子的な分離されたトランザクションによってデータベースを維持するものである。トランザクションがデータベースを変更しようとする際、デルタページャは変更を書き込みバッファに格納し、介入するトランザクションが、事実上または実質上、トランザクションに依拠してデータベースの状態を変更しない場合、この変更を適用する。デルタページャは変更を適用して、データベースの状態を表す新しいデータ構造を形成するように、書き込みバッファとデータベースの現行状態とを結合させることによって、トランザクションをコミットする。デルタページャは、選択されたデータベースの状態を保持するために、デルタページャが遵守するスナップショットに従って、効率を維持するために書き込みバッファを合体(coalesce)する。デルタページャは、選択されたデータを永続ストアへ移動することによって、データベースの選択されたセクションを永続的にする。デルタページャは、データへの効率的なアクセスを促進するために、永続ストアと現行トランザクションとの間にキャッシュオブジェクトも提供する。

【0008】

この課題を解決するための手段は、以下の発明を実施するための最良の形態でより詳細に説明する概念の選択を簡略化された形で紹介するために提供される。この課題を解決するための手段は、主張された主題の主要な特徴または不可欠な特徴を識別することを意図するものでなく、また、主張された主題の範囲を決定する際の一助として使用されることを意図するものでもない。

【0009】

添付の図面を参照しながら、詳細な記述について説明する。図面では、3桁の参照番号の一番左の数字または4桁の参照番号の左側2つの数字が、その参照番号が最初に示された図面を識別するものである。同じ参照番号が異なる図面で使用される場合、同様または同一のアイテムであることを示す。

【発明を実施するための最良の形態】

【0010】

概要

「デルタページャ」という用語は、データベースの整合性を維持し、データベースの選択された状態を保持しながら、データベースへの変更を処理するための、方法およびシステムの諸実施形態を記述する。認識されたオブジェクト間でのポイントの不変性を遵守することによって、デルタページャは、効率性も維持しながらデータベースの状態の整合性を保持する。

【0011】

1つのモードでは、デルタページャは、データベースに格納されたすべてのデータの完

10

20

30

40

50

全マッピングを含む、データベースの現行状態を追跡するために、現行状態ポインタを使用することによって、トランザクションを管理する。完全マッピングでは、データベースに関して定義された各アドレスまたはアドレスの表現について、データ値またはヌル (null) 値が存在することになる。部分マッピングは、トランザクションが開始された時点で存在するデータベースの現行状態に対して、トランザクションが適用しようとする変更を格納するために、デルタページャが使用する、1つまたは複数のオブジェクトを含むことができる。トランザクションがデータベースにコミットされた場合、新しい現行状態を形成するために、データベースのオリジナルの現行状態に部分マッピングが付加される。新しく適用された部分マッピングを介して新しい現行状態にアクセスされた場合、部分マッピングにおけるアドレスへの値の割り当ては、オリジナル状態におけるそれらよりも優先されるため、データベースのオリジナル状態におけるデータの値は効果的に上書きまたは変更される。

10

【 0 0 1 2 】

本説明で使用される用語を明確にするためおよび区別するために、データベースのマッピングとは、アドレス値の、データ値への割り当てを言い表す。データベースの完全マッピングとは、データベース内のすべてのアドレスの、それらのアドレスで格納されたデータ値 (またはヌル値) への割り当ての、完全なセットを言い表す。部分マッピングとは、トランザクションがデータを書き込もうとするアドレス、およびトランザクションがアドレスに書き込もうとするデータ値を含む、トランザクションがデータベースに対して実行するかまたは実行しようとする変更または上書きのセットを言い表す。トランザクションがデータベースにコミットされた場合、部分マッピングを以前の完全マッピングに追加することが可能であり、結果としてデータベースの現在のマッピングが更新される。したがって、データベースの現在の完全マッピングは一連の部分マッピング (a series of partial mappings) を含むことが可能であり、後で追加された部分マッピングは、以前の部分マッピングに含まれるアドレスの割り当てを上書きまたは変更することができる。

20

【 0 0 1 3 】

データベースの状態とは、データベースのマッピングが格納される、メモリ、ディスク、または大容量記憶装置に格納されたオブジェクトの集まりを言い表す。状態内のオブジェクトは、トランザクションがデータベースに書き込もうとする変更、ならびにバッファを内部に集めることができる他のオブジェクトを格納するために、デルタページャがトランザクションに対して作成するバッファを含む。バッファ、他のオブジェクト、ならびに、これらのオブジェクトが結合、合体、および他の方法で処理される方法について、以下で詳細に説明する。データベースの現行状態は、データベースの現在の完全マッピングを提示するそれらのオブジェクトを含む。以下で説明するように、デルタページャは、トランザクションが開始された時点で存在するデータベースの状態などの、他の状態も維持する。デルタページャは、トランザクションをデータベースにコミットすべきであるかどうかを判断するために、この状態を追跡する。デルタページャは、それらの状態、ならびに後で使用するためにそれらの状態が表すデータベースのマッピングを保持するために、選択された状態も維持する。

30

40

【 0 0 1 4 】

各トランザクションが開始されると、デルタページャは、トランザクションが開始された時点で存在していたデータベースの状態に基づいて、トランザクションを実行させることができる。その後、トランザクションは、書き込みバッファなどの少なくとも1つのオブジェクトを作成し、トランザクションが開始された時点で存在していたデータベースの状態に書き込みバッファをポイントさせる。デルタページャは、一連のバッファを使用することが可能であり、ここで各バッファは、トランザクションがデータベースに適用しようとする単一の変更を格納し、各バッファは、データベースの状態を保持するために先行バッファをポイントする。別の方法として、デルタページャは、トランザクションが実行しようとするすべての変更を、トランザクションが開始された時点のデータベースの現行

50

状態をポイントする累積書き込みバッファに格納することもできる。トランザクションが開始された時点の状態をポイントするバッファまたは他のオブジェクトは、トランザクションがデータベースに適用しようとする変更の部分マッピングを表す。

【 0 0 1 5 】

1つまたは複数の介在トランザクションが現行トランザクションの適用を妨げているというインディケーションを、デルタページが見つけなかった場合、デルタページは、データベースの現行状態の追跡にデルタポインタが使用する現行状態ポインタを、トランザクション用に作成された最後のバッファまたは唯一のバッファをポイントするように変更することによって、トランザクションをコミットする。このバッファをポイントするように現行状態ポインタを変更することによって、データベースの新しい更新された状態を生成し、データベースの完全マッピングにおけるデータ値を変更することができる。

10

【 0 0 1 6 】

デルタページは、読み取りバッファおよび書き込みバッファを使用してトランザクションを処理することができる。読み取りバッファは、トランザクションがアクセスするデータを追跡する。現行トランザクションが完了した場合、デルタページは読み取りバッファを使用して、トランザクションをコミットするかどうかを決定する。読み取りバッファが、介入トランザクションによって上書きまたは変更されたデータを現行トランザクションが読み取ったことを示す場合、デルタページはトランザクションを打ち切ることができる。

【 0 0 1 7 】

20

デルタページがデータベースへの複数の変更をコミットした後、デルタページは効率性のために、データベースの現行状態においてバッファまたは他のオブジェクトを合体することができる。デルタページは、1つの部分マッピングにおける複数のオブジェクトを、他の部分マッピングにおけるオブジェクトと合体することができる。2つの部分マッピングのうちの新しい方におけるアドレスの割り当ては、古い方の部分マッピングに含められた割り当てに追加されるか、またはこれと置き換えられることになる。合体の結果、メモリオブジェクトまたは永続オブジェクトを含む1つまたは複数の合体オブジェクトが生じる可能性があるか、あるいはデルタページは、バッファに格納された変更をデータベースのオリジナル状態に適用することによってバッファを合体することができる。デルタページは、選択された状態を保持できるようにするためにスナップショットポインタを提供する。スナップショットポインタに関連付けられたいかなる状態も、スナップショットポインタによって保持された状態に対して実行された変更を表すバッファと合体されることはない。

30

【 0 0 1 8 】

デルタページは、データベースおよびデータベースに対して実行された変更を、永続的にすることができる。デルタページは、データベースの現行状態を維持するために使用された1つまたは複数のオブジェクトを、永続的な不揮発性ストレージにコピーすることができる。デルタページは、デルタページが永続的にするために不揮発性ストレージに移動した現行状態の一部を含む、大容量記憶媒体に格納されたデータベースの現行状態の一部へのより高速なアクセスを提供するための、キャッシュオブジェクトも備える。デルタページは、領域照会を容易にするために、このキャッシュオブジェクトを使用してヌルデータの領域をデータベースの状態に格納することができる。

40

【 0 0 1 9 】

書き込みバッファを使用したトランザクションの処理

図1は、データベースの現行状態100を示す。現行状態100は、データベースの例示的なオリジナル状態110を含む。データベースは、揮発性メモリまたは不揮発性メモリに常駐することができる。デルタページは、現行状態ポインタ120を使用して、図1ではオリジナル状態110のみを含む現行状態をポイントすることができる。

【 0 0 2 0 】

オリジナル状態110は、レコード0 130、レコード1 140、レコード2 1

50

50、レコード3 160、およびレコード4 170の、それぞれが現在はヌル値を格納している5つのレコードを維持するオブジェクトである。ヌル値は、実際にオリジナル状態に格納されている必要がないことに留意されたい。またオリジナル状態110は、トランザクションが変更を適用する可能性のある、ヌルおよび非ヌルの値を含む、任意の初期状態を含むことができる。

【0021】

図2は、複数の命令210~270を含む例示的トランザクションである、トランザクション1 200を示す。命令1 220、「x = Read(3)」などのいくつかの命令は、データベースからデータを読み取る。命令3 240、「Write(2, x)」などの他の命令は、結果としてデータベース内のアドレスに値を割り当てるかまたは再割り当てすることになる、データベースへの変更を適用しようとする。

10

【0022】

図3は、データベースに対してトランザクション1 200などのトランザクションを開始する、デルタページの1つのモードを示す。トランザクション1 200は、トランザクション開始命令、trans_start 210で開始される。トランザクション1 200が開始された後、現行状態ポインタ120は、トランザクションが開始された時点のデータベースの状態、すなわちオリジナル状態110を、引き続きポイントする。デルタページは第1の書き込みバッファ310を作成し、デルタページは、トランザクション1 200が適用しようとする変更をここに格納することになる。デルタページは、現行状態ポインタ120によって示された、トランザクション2 200が開始された時点の状態に、書き込みバッファ310をポイントさせる。デルタページはトランザクションポインタ320も作成し、このトランザクションポインタ320を第1の書き込みバッファ300にポイントさせ、この第1の書き込みバッファは現行状態をポイントする。

20

【0023】

図4~6は、デルタページがトランザクション1 200を処理する方法を示す。図2を参照すると、命令1 220、「x = Read(3)」は、ヌルである現行値xをアドレス3から読み取る。命令2 230、「x = x + 1」は、xの値を1だけ増分し、このケースではxの値を1に変更する。命令3 240、「Write(2, x)」は、コンピューティングシステムにxの値をアドレス2に書き込むように命じる。

30

【0024】

デルタページは値1をオリジナル状態110のアドレス2には書き込まない。その代わりに図4に示されるように、デルタページは第2の書き込みバッファ410を作成する。第2の書き込みバッファ410は第1の書き込みバッファ310をポイントし、デルタページは、トランザクション2 200が、トランザクションが開始された時点のデータベースの状態を修正する方法を示す、データベースの部分マッピングを維持するために、トランザクションポインタ320が第2の書き込みバッファ410をポイントするように切り替える。デルタページは、このアドレス2への値1の書き込みの変更を、第2の書き込みバッファ410に格納する。

【0025】

40

図5は、トランザクション1 200がデータベースに実行しようとするデルタページの第2の変更への応答を示す。命令4 250、「Write(4, 「DOG」)」は、文字列「DOG」をアドレス4に格納するようにコンピューティングシステムに命じる。ここでも、変更をオリジナル状態110に書き込む代わりに、デルタページは、第2の書き込みバッファ410をポイントする第3の書き込みバッファ510を作成する。さらにデルタページは、トランザクション1 200がこれを修正しようとした場合にデータベースの状態を保持するために、トランザクションポインタ320が第3の書き込みバッファ510をポイントするように変更し、文字列「DOG」を第3の書き込みバッファ510に格納する。したがって、トランザクションポインタ320から読み取る、トランザクション21200に関するデータベースの部分マッピングは、アドレス2が値1を

50

格納し、アドレス4が文字列「DOG」を格納するという修正を含む。しかしながら、トランザクションは、オリジナル状態110も、現行状態ポインタ120によって示される現行状態も変更していない。

【0026】

図6は、デルタページャがトランザクション1200をデータベースにコミットする方法を示す。トランザクションポインタ320は、書き込みバッファ510、410、および310を通るか、または現行状態ポインタ120がポイントするのと同じ状態をポイントするため、デルタページャは、トランザクション1200がデータベースにコミットされないようにするために、いかなる介入トランザクションもデータベースの現行状態を変更していないものと判別する。したがってデルタページャは、書き込みバッファ310、410、および510によって修正された場合にオリジナル状態110を含む、トランザクション1200に対して作成された部分マッピングをポイントするように現行状態ポインタ120を変更することによって、トランザクション1200をデータベースにコミットする。したがってデルタページャは、トランザクション1200によって追加された最新の書き込みバッファである、第3の書き込みバッファ510をポイントするように、現行状態ポインタ120を変更する。次にデルタページャは、トランザクションポインタ320を削除する。これで現行状態ポインタ120は、書き込みバッファ310、410、および510に格納された変更によって更新されたオリジナル状態110を含む、データベースの更新された状態をポイントすることになる。

【0027】

図6の現行状態を使用して、トランザクションがデータベースにアクセスする場合、デルタページャは、現行状態ポインタ120によって識別される、トランザクションが開始される時点で存在するデータベースの状態を使用する。この現行状態を使用して、トランザクションがバッファ510内でアドレス4のコンテンツを識別しようとした場合、文字列「DOG」を見つけることになる。他方で、トランザクションがアドレス2のコンテンツを求めた場合、バッファ510内には見つからないことになる。したがって、トランザクションはバッファ410に進み、ここでトランザクションは値1を見つけることになる。図6の例で、トランザクションが他の値を求めた場合、オリジナル状態110が他のアドレスについてヌル値のみを含むことを見つけるまで、トランザクションはすべてのバッファにアクセスすることになる。このようにしてデルタページャは、データベースの現行状態にバッファまたは他のオブジェクトを追加することによって、データを変更する。

【0028】

デルタページャの不変条件(invariant)とは、デルタページャが作成したそれぞれのポインタの不変性を遵守することである。不変性の不変条件とは、各オブジェクトが依存する状態、したがって各オブジェクトがポイントする状態が、ポインタがそのオブジェクトをポイントする限りは変更されないことを維持するものである。たとえば、第1の書き込みバッファ310は、データベースのオリジナル状態110をポイントする。その後の変更にかかわらず、トランザクション1200はデータベースに適用し、第1の書き込みバッファ310は、オリジナル状態110のみを含むデータベースの状態を常にポイントする。同様に、第2の書き込みバッファ410は第1の書き込みバッファ310をポイントする。第2の書き込みバッファ410は、オリジナル状態110およびオリジナル状態110を修正しない書き込みバッファ310を含む、データベースの不変状態もポイントする。第3の書き込みバッファ510は、第2の書き込みバッファ410をポイントするため、第2の書き込みバッファ410によって修正された場合のオリジナル状態110を含む、データベースのそれ独自の部分マッピングをポイントする。これらのポインタの不変性は、以下で説明する利点を有する。

【0029】

累積書き込みバッファおよび単一書き込みバッファの再書き込み

今説明したモードでは、デルタページャは、トランザクションを開始するため、およびトランザクションがデータベースに適用しようとする各変更を格納するために、新しい書

10

20

30

40

50

き込みバッファを作成する。別の方法として、デルタページは、累積バッファなどの異なるタイプのオブジェクトを採用することができるか、またはデルタページは、1つまたは複数の変更を格納するために既存のバッファを上書きすることができる。

【0030】

図7は、累積書き込みバッファを使用してトランザクションを処理するデルタページを示す。図3の例と同様に、デルタページは、現行状態ポインタ120によって示されるように、トランザクションが開始されたときに存在した現行状態をポイントする第1の累積バッファ710を作成することによって、トランザクション1200を開始する。デルタページは、トランザクションポインタ320を再度作成し、これを、最も新しく作成された書き込みバッファにポイントさせて、トランザクション1200によって修正されるようにデータベースの部分マッピングを維持する。

10

【0031】

図8は、第1の累積バッファ710をポイントする第2の累積バッファ810を追加することによってトランザクションの処理を続行する、デルタページを示す。デルタページは、トランザクションポインタ320を、最も新しく作成された書き込みバッファである第2の累積バッファ810にポイントさせる。デルタページは、トランザクション1200がデータベースに適用しようとする、アドレス2の値を1に設定する変更を、第2の累積バッファ810に格納する。トランザクション1200はこれまで1つの変更しか試みてこなかったため、第2の累積バッファ810は累積的には見えない。

【0032】

20

図9は、第2の累積バッファ810をポイントする第3の累積バッファ910を追加することによってトランザクション1200の処理を続行する、デルタページを示す。デルタページは、トランザクションポインタ320を、最も新しく作成された書き込みバッファである第3の累積バッファ910にポイントさせる。ここでデルタページは、トランザクション1200がデータベースに適用しようとする、アドレス2の値を1に設定する変更と、アドレス4で文字列「DOG」を格納する変更の両方を格納する。

【0033】

図10および図11は、デルタページが累積書き込みバッファを使用するモードの利点を示す。図10は、第3の累積バッファ910に格納されたトランザクション1200によって求められたすべての変更と共に、トランザクション1200がデータベースに適用することになるすべての変更の部分マッピングが、単一のオブジェクトである累積書き込みバッファ910に含められることを示す。図10の例に示されるように、部分マッピングは、たとえ累積バッファ710および810が除外された場合であっても保持される。

30

【0034】

不変性の不変条件により、いかなるトランザクションもトランザクションの中間状態に干渉またはアクセスすることができないため、トランザクションポインタ320によって示される状態の整合性を維持するために書き込みバッファ710および810を保存する必要はない。結果として、累積書き込みバッファ710および810が格納した変更が、最も新しく作成された累積バッファ910にも格納された後で、累積書き込みバッファ710および810を維持する必要はない。

40

【0035】

したがって、図11が示すように、デルタページが新しい累積バッファを作成し、そのそれぞれに書き込む場合、デルタページは先行するバッファを破棄することができる。第1の段階1100は、アドレス2の値を1に設定するという第1の変更が第2の累積バッファ810に書き込まれた後の、トランザクション1200の現行状態を示す。第2の累積バッファ810は、第1の書き込みバッファ710に格納されたいずれの変更も含むため、第2の段階1110でデルタページは、第2の累積バッファ810を、第1の累積バッファ710がポイントした状態にポイントさせ、トランザクションポインタ320によって示された状態を変更せずに、第1の累積バッファ710およびそのポインタ

50

を解放することができる。トランザクションポインタ320がポイントする状態は、第2の累積バッファ810に格納された同じデータのサブセットのみを含むバッファを通過しないため、より短い。第3の段階1120で、デルタページは、第1の累積バッファ710を格納するために使用されたメモリを解放し、他に使用するためにメモリを確保することができる。デルタページが、トランザクション1200の次の変更を格納するために、図11に示されたものと同様の段階を適用した場合、結果は図10に示されたものと同じになる。

【0036】

図12は、デルタページがトランザクションを処理するために使用するモードを示す。流れ図1200は、ブロック1202がトランザクションを受け取ると開始される。ブロック1204は、トランザクションがデータベースに適用しようとする任意の変更をポイントすることになる、トランザクションポインタを作成する。ブロック1206は、第1の書き込みバッファを作成し、そのバッファを、トランザクションが受け取られた時点で存在する現行状態にポイントさせる。ブロック1208は、トランザクションポインタを第1の書き込みバッファにポイントさせ、第1の書き込みバッファはトランザクションが開始された時点の現行状態をポイントする。

【0037】

ブロック1210は、トランザクションが、たとえばデータベースからデータを読み取るだけでなく、データベースに変更を適用しようとするかどうかを判別する。変更を適用しようとする場合、ブロック1212は追加の書き込みバッファを作成し、これを、以前作成されたバッファにポイントさせる。追加の書き込みバッファは、単一書き込みバッファまたは累積書き込みバッファとすることができる。1つのモードで累積書き込みバッファが使用される場合、ブロック1212は、追加の累積バッファを、以前のバッファがポイントした状態にポイントさせ、図11を参照しながら説明したように以前のバッファを解放する。

【0038】

ブロック1214は、トランザクションがデータベースに適用しようとする変更を、追加の書き込みバッファに格納する。ブロック1216は、トランザクションポインタを追加の書き込みバッファにポイントさせる。ブロック1218は、トランザクションが任意の追加の変更をデータベースに適用しようとするかどうかを判別する。適用しようとする場合、流れ図1200はブロック1212にループし、ブロック1212は追加のバッファを作成する。他方で、ブロック1218が、トランザクションは追加の変更を適用しようとしなないものと判別した場合、後で説明するように、ブロック1220は変更のコミットを試行することになる。ブロック1224は、他のトランザクションの受け取りを待つ。

【0039】

他方で、ブロック1210が、トランザクションがいかなる変更もデータベースに適用しようとしなかったものと判別した場合、第1の書き込みバッファおよびトランザクションポインタは不要であるため、ブロック1222はこれらを解放する。ブロック1224は、再度、他のトランザクションの受け取りを待つ。

【0040】

デルタページは、一連の累積バッファを使用する代わりに、図13に示されたような単一の再書き込み可能書き込みバッファを使用することもできる。第1の段階1300で、デルタページはトランザクションを受け取り、単一書き込みバッファ1310を作成して、これを現行状態にポイントさせる。次にデルタページは切り替えトランザクションポインタ1320を作成し、これを、単一書き込みバッファ1310にポイントさせる。切り替えポインタ1320は、以下で説明するように、デルタページ内のポインタの不変性を保持する。このようにバッファに上書きまたは再書き込みすることにより、それぞれの変更を格納するために複数のバッファを使用する代わりに、トランザクションが適用しようとする変更を格納するために単一の累積バッファを使用して、デルタページに

10

20

30

40

50

よって認識される不変性の不変条件を保持する。

【 0 0 4 1 】

単一書き込みバッファ 1 3 1 0 の使用の第 2 の段階 1 3 3 0 で、デルタページャは、切り替えポインタ 1 3 2 0 を単一書き込みバッファ 1 3 1 0 を避けてポイントさせ、トランザクションによって示された変更または追加の変更で単一書き込みバッファを更新する。前述のように、デルタページャはすべてのポインタの不変性を保持する。したがって技術的には、切り替えポインタが引き続き単一書き込みバッファ 1 3 1 0 をポイントする場合、バッファ 1 3 1 0 に変更を書き込むと、不変性の不変条件が破られることになり、切り替えポインタによって示される状態は変化するように見える。したがって 1 つのモードでは、デルタページャは、単一書き込みバッファ 1 3 1 0 を避けてポイントすることが可能な切り替えポインタ 1 3 2 0 を使用する。

10

【 0 0 4 2 】

第 3 の段階 1 3 4 0 で、デルタページャは単一書き込みバッファを更新すると、トランザクションがデータベースに適用しようとする変更の部分マッピングを維持するために、再度切り替えポインタ 1 3 2 0 を単一書き込みバッファ 1 3 1 0 にポイントさせる。第 4 の段階 1 3 5 0 で、デルタページャは、再度、切り替えポインタ 1 3 2 0 を単一書き込みバッファ 1 3 1 0 を避けてポイントさせ、トランザクションによって示された他の変更を追加するように単一書き込みバッファを更新する。第 5 の段階 1 3 6 0 で、デルタページャが追加の変更を単一書き込みバッファに追加すると、デルタページャは部分マッピングを維持するために、再度切り替えポインタ 1 3 2 0 を単一書き込みバッファ 1 3 1 0 にポイントさせる。

20

【 0 0 4 3 】

デルタページャは、実際にはバッファ 1 3 1 0 を避けてポイントされる切り替えポインタ 1 3 2 0 を、事実上使用する必要がないことに留意されたい。不変性の不変条件が論理的に保持されるのを示すために、デルタページャは、切り替えポインタ 1 3 2 0 によって示された再書き込み可能バッファを使用するものとして説明される。別の方法として、システムのどの部分も再書き込みが発生した間ポインタを逆参照 (r e f e r e n c e) しないうに、再書き込みが原子的に発生した場合、不変条件が依然として有効に保持されている間、通常のトランザクションポインタ 3 2 0 を代わりに使用することが可能であり、ポインタは、だれもその値を監視していない場合、ポインタではない。したがってこの限られたケースでは、デルタページャは、不変性の不変条件を保持することからインプレース再書き込みを使用する。

30

【 0 0 4 4 】

図 1 4 は、デルタページャが単一書き込みバッファを使用することによってトランザクションを処理するために使用する、図 1 3 に示されたモードを示す。流れ図 1 4 0 0 は、ブロック 1 4 0 2 がトランザクションを受け取ると開始される。ブロック 1 4 0 4 は、トランザクションがデータベースに適用しようとする変更の部分マッピングを示すことになる、切り替えポインタを作成する。ブロック 1 4 0 6 は、単一書き込みバッファを作成し、そのバッファを、トランザクションが受け取られた時点の現行状態にポイントさせる。ブロック 1 4 0 8 は、トランザクションがデータベースに適用しようとする部分マッピングを示すために、切り替えポインタを単一書き込みバッファにポイントさせる。

40

【 0 0 4 5 】

ブロック 1 4 1 0 は、トランザクションがデータベースに変更を適用しようとするかどうかを判別する。適用しようとする場合、ブロック 1 4 1 2 は切り替えポインタを単一書き込みバッファを避けてポイントさせる。ブロック 1 4 1 4 は、変更を単一書き込みバッファに格納する。図 1 3 に示されるように、単一書き込みバッファは累積的であるため、複数の変更を単一書き込みバッファに追加することができる。ブロック 1 4 1 6 は、再度切り替えポインタを単一書き込みバッファにポイントさせる。

【 0 0 4 6 】

ブロック 1 4 1 8 は、トランザクションがデータベースに任意の追加の変更を適用しよ

50

うとすることがかを判別する。適用しようとする場合、流れ図 1 4 0 0 はブロック 1 4 1 2 にループし、ブロック 1 4 1 2 は追加のバッファを作成する。他方で、ブロック 1 4 1 8 が、トランザクションが追加の変更を適用しようとしなないものと判別した場合、後で説明するように、ブロック 1 4 2 0 は変更のコミットを試行することになる。ブロック 1 4 2 4 は、他のトランザクションの受け取りを待つ。

【 0 0 4 7 】

他方で、ブロック 1 4 1 0 が、トランザクションがデータベースにいかなる変更も適用しようとしなかったものと判別した場合、単一書き込みバッファおよび切り替えポインタは不要であるため、ブロック 1 4 2 2 はこれらを解放する。ブロック 1 4 2 4 は他のトランザクションの受け取りを待つ。

10

【 0 0 4 8 】

トランザクションのコミット

デルタページャは、トランザクションをコミットするタイミングを決定するための異なるモードをサポートする。デルタページャは、複数のトランザクションが、各トランザクションが開始された時点で存在するデータベースの状態に同時にアクセスできるようにする。しかしながら、複数のトランザクションが同時に実行された場合、最初のデータベースへのコミットにより、他の同時に実行しているトランザクションが依拠するデータベースの状態が変更される可能性がある。古くなったかまたは取って代わられたデータベースの状態に基づくトランザクションがデータベースにコミットされると、無効の結果につながる可能性がある。したがってデルタページャは、現行トランザクションが依拠したデータベースの状態に介入トランザクションが干渉しなかった場合にのみ、現行トランザクションをコミットする。

20

【 0 0 4 9 】

図 1 5 は、デルタページャがトランザクション 1 2 0 0 と同時に処理し、結果としてトランザクションの競合を発生させる可能性のある、第 2 のトランザクションである、トランザクション 2 (T x 2) 1 5 0 0 を示す。競合する可能性のあるトランザクションを考えた場合、デルタページャは、どのトランザクションをデータベースにコミットし、どのトランザクションを打ち切るかを決定する。トランザクション 2 1 5 0 0 は、複数の命令 1 5 1 0 ~ 1 5 7 0 を呼び出す。命令 1 1 5 2 0、「x = R e a d (0)」などのいくつかの命令は、データベースからデータを読み取ろうとするが、命令 4 1 5 5 0、「W r i t e (3 , 「 C A T 」)」などの他の命令は、変更をデータベースに適用しようとする。

30

【 0 0 5 0 】

図 1 6 は、デルタページャがトランザクション 1 2 0 0 およびトランザクション 2 1 5 0 0 の両方を受け取り、同時に処理する状況を示す。2 つのトランザクションのうちのどちらかを、どちらかのトランザクションがデータベースにコミットされる前に、他方よりも先に開始することができる。どちらのトランザクションも、現行状態ポインタ 1 2 0 が引き続きオリジナル状態 1 1 0、すなわちその時点での現行状態をポイントするものとして開始された。

【 0 0 5 1 】

40

デルタページャは、前述のように、現行状態ポインタ 1 2 0 によって示されたオリジナル状態 1 1 0 をポイントする第 1 の書き込みバッファ 3 1 0 を作成することによって、トランザクション 1 2 0 0 を開始し、そのトランザクションに関して最も新しく作成された書き込みバッファをポイントするトランザクションポインタ 3 2 0 を作成する。デルタページャは、トランザクション 1 2 0 0 がデータベースに対して実行しようとする変更を格納する、追加の書き込みバッファ 4 1 0 および 5 1 0 を作成する。

【 0 0 5 2 】

トランザクション 2 1 5 0 0 の場合、デルタページャは、オリジナル状態 1 1 0 をポイントする第 1 の書き込みバッファ 1 6 1 0 と、そのトランザクションに関して最も新しく作成された書き込みバッファを示すためのトランザクションポインタ 1 6 2 0 とを作成

50

する。同様にデルタページは、トランザクション2 1500に関する追加の書き込みバッファを作成する。図16では、命令1 1520である「x = Read(0)」、および命令2 1530である「y = Read(1)」は、どちらもオリジナル状態110からヌル値を読み取る。したがって条件付き命令3 1540、「if y = x then」は真であり、命令4 1550、「Write(3,「CAT」)」が実行される。デルタページは、第1の書き込みバッファ1610をポイントする追加の書き込みバッファ1630を作成し、命令4 1550によって示された変更をこの追加の書き込みバッファ1630に格納し、トランザクションポインタ1620を、この追加の書き込みバッファ1630にポイントさせる。トランザクションポインタ320およびトランザクションポインタ1620は、どちらも、各トランザクションがデータベースに適用することになる変更を含む部分マッピングをポイントする。

10

【0053】

トランザクション1 200およびトランザクション2 1500は、データベースに対して競合する変更を行う可能性があるか、あるいは、命令3 1540などの条件付き命令は、他のトランザクションが変更する可能性のあるデータに依拠する可能性がある。デルタページのモードは、これらの潜在的な競合を様々な方法で説明する。

【0054】

介入トランザクションが現行状態を変更した場合のトランザクションの打ち切り

図17および図18は、デルタページがデータベースの状態を変更する介入トランザクションをコミットした後にトランザクションを打ち切ることによって、デルタページが競合する可能性のあるトランザクションを処理するために使用する1つのモードを示す。図17では、図6に関して前述したように、デルタページは、トランザクション1 200をデータベースにコミットしており、トランザクションポインタ320は除去され、現行状態ポインタ120は、トランザクションに関して作成された最新の書き込みバッファ510をポイントするように切り替えられる。トランザクション1 200をデータベースにコミットした後、デルタページは、トランザクション2 1500が実行を完了したものと決定し、データベースへコミットしようとする。

20

【0055】

トランザクション2 1500をデータベースにコミットするかどうかを決定する場合、デルタページは、トランザクション1 200などの介入トランザクションが、現行のトランザクション、すなわちトランザクション2 1500が開始された時点で存在したデータベースの状態を変更したかどうかを判別する。この判別を行う場合、デルタページは、第1の書き込みバッファ1610がポイントする状態を決定することによって、トランザクション2 1500が開始された時点のデータベースの状態と比較する。次にデルタページは、これを、現行状態ポインタ120がポイントするデータベースの現行状態と比較する。言い換えれば、デルタページは、それらが両方とも同じオブジェクトをポイントしているかどうかを判別するために、現行状態ポインタ120と、トランザクションポインタ1620ならびに挿入された書き込みバッファ1610および1630とを比較する。トランザクション1 200をデータベースにコミットした後、デルタページは、現行状態ポインタ120が書き込みバッファ510をポイントする一方で、トランザクション2 1500は、その第1の書き込みバッファ1610を通してオリジナル状態110をポイントすることを決定する。

30

40

【0056】

トランザクション2 1500は、現行状態ポインタ120によって示される現行状態以外の状態をポイントするため、デルタページはトランザクション2 1500を打ち切る。図18は、デルタページが、データベースの現行状態の一部とはならないメモリを解放するために、書き込みバッファ1610および1630ならびにトランザクションポインタ1620を削除することによって、トランザクション2 1500を打ち切ることを示す。

【0057】

50

現行状態におけるいずれかの変更のためにトランザクションを打ち切るとは、競合するトランザクションのコミットを避けるためのリスクを嫌った手法である。図16～図18の例に示されるように、トランザクション1 200は、トランザクション2 1500が実行された時点で存在したデータベースの状態におけるいかなる値も変更しなかった。しかしながら、介入するトランザクション1 200が、トランザクション2 1500が依拠したデータを変更した可能性はある。トランザクション1 200はデータベースの現行状態を変更したため、トランザクション2 1500が依拠したデータを変更した可能性がある。

【0058】

トランザクション2 1500は打ち切られたが、図19および図20に示されるように、このトランザクションは再開することができる。図19で、デルタページは、トランザクション1 200がコミットした後に、トランザクション2 1500を開始する。したがって、トランザクション2 1500が開始された場合、デルタページは、トランザクション2 1500が再開された時点で存在する状態をポイントする、第1の書き込みバッファ1910を作成する。この状態は、トランザクション1 200によって状態に追加された書き込みバッファ510をポイントする、現行状態ポインタ120によって示される。次にデルタページは、元々は書き込みバッファ1910をポイントしていたトランザクションポインタ1920を作成する。デルタページは、トランザクション2 1500がデータベースに適用しようとする変更をデルタページが格納する追加の書き込みバッファ1930を追加し、その後、デルタページは、トランザクション2 1500がデータベースに適用しようとする変更を示す部分マッピングを示すために追加の書き込みバッファ1930をポイントするように、トランザクションポインタ1920を変更する。

【0059】

今回は、トランザクション2 1500が完了したものとデルタページが決定し、データベースにトランザクションをコミットしようとする、デルタページには、トランザクション2 1500が依拠した状態からデータベースの現行状態を変更した介入トランザクションがないことがわかる。トランザクションポインタ1920は、書き込みバッファ1930および1910を通して、書き込みバッファ510をポイントする。現行状態ポインタが依然として書き込みバッファ510をポイントしているため、トランザクションポインタ1920は依然として現行状態をポイントする。

【0060】

図20は、デルタページがトランザクション2 1500をデータベースにコミットしていることを示す。デルタページは、トランザクション2 1500によって追加された最新のバッファ、すなわち書き込みバッファ1930をポイントするように、現行状態ポインタ120を変更する。次にデルタページは、トランザクションポインタ1920を削除する。現行状態ポインタが現在ポイントしている現行状態は、トランザクション1 200によって適用された書き込みバッファ410および510と、トランザクション2 1500によって適用された書き込みバッファ1930とに格納された変更によって修正された場合、オリジナル状態110を含む。

【0061】

図21は、デルタページがトランザクションをコミットするかまたは打ち切るかどうかを決定するために使用する、第1のモードを示す。流れ図2100は、ブロック2102で開始される。ブロック2102は、完了されたトランザクションがデータベースへのコミットを求める場合を決定し、流れ図2100は、トランザクションがデータベースへのコミットを求めるまで、ブロック2102にループする。ブロック2102が、コミットしようとする完了されたトランザクションを検出した場合、ブロック2104は、現行状態が、たとえば現行トランザクションに関してトランザクションポインタによって示されるように、トランザクションが開始された時点で存在した状態と同じであるかどうかを判別する。言い換えれば、ブロック2104は、トランザクションのために作成された第

10

20

30

40

50

1の書き込みバッファが、現行状態ポインタがポイントするのと同じ状態をポイントするかどうかを判別する。ポイントする場合、ブロック2106は、現行状態ポインタに、トランザクションポインタもポイントするバッファである、トランザクションのために作成された最新のバッファをポイントさせ、トランザクションがコミットされる結果としてデータベースの更新を完了する。ブロック2108はトランザクションポインタを削除し、関連するメモリを解放する。流れ図2100は、データベースにコミットしようとする次の完了されたトランザクションを検出するために、ブロック2102へとループする。

【0062】

他方で、ブロック2104が、現行状態が、トランザクションが開始された時点で存在した状態と同じでないと決定した場合、ブロック2104は、現行状態が変更されたものと決定する。ブロック2110は、トランザクションの書き込みバッファを解放する。ブロック2112は、トランザクションの打ち切りが完了したトランザクションポインタを削除する。前述のように、打ち切られたトランザクションは再開することが可能であり、デルタページは、トランザクションが完了し、コミットしようとする場合、トランザクションがデータベースにコミットされるべきであるかどうかを判別することになる。

【0063】

トランザクションが不必要に打ち切られるのを防ぐための読み取りバッファ

トランザクションが現行状態と一致しない状態をポイントする場合に必ず、デルタページが現行トランザクションを打ち切る代わりに、デルタページは、現行トランザクションが読み取ったデータを介入トランザクションが書き込んだ場合に、現行トランザクションを打ち切ることができる。現行状態が変更された場合に必ず現行トランザクションを打ち切ることによって、競合する変更がデータベースに適用されないことが保証される。しかしながら、同時に多くのトランザクションが実行している可能性がある場合、および、それらのトランザクションの一部が長いまたは複雑な計算を含む可能性がある場合、トランザクションの打ち切りは計算リソースを無駄にする。したがって、介入トランザクションが、現行トランザクションが読み取ったかまたは依拠したいかなるデータにもアクセスしなかったかまたはこれを変更しなかった場合、デルタページが現行トランザクションを打ち切る必要はない。

【0064】

デルタページのあるモードは、各トランザクションに対して読み取りバッファを作成する。読み取りバッファは、トランザクションがアクセスする任意のデータを追跡する。完了したトランザクションがデータベースへのコミットを試行するが、介入トランザクションがデータベースの状態を変更したことがわかった場合、デルタページは、読み取りバッファとデータベースの現行状態とを比較して、現行トランザクションの打ち切りが必要であるかどうかを判別することができる。他のモードでは、デルタページは、現行状態が変更され、現行トランザクションが読み取ったデータを介入トランザクションが変更または上書きした場合に、現行トランザクションを打ち切る。このモードでは、デルタページは、介入トランザクションが現行トランザクションによって読み取られたアドレスにデータを書き込んだ場合に、トランザクションを打ち切ることができるか、または、デルタページは、データの値が実際に変更したかどうか判別するために、現行トランザクションによって読み取られた実際のデータと、介入トランザクションによって書き込まれたデータとを、比較することができる。

【0065】

図22～24は、現行トランザクションによって読み取られたデータを介入トランザクションが書き込んだ場合に、読み取りバッファを使用してトランザクションを打ち切る、デルタページのモードを示す。図22は、データベースの現行状態に対して実行される、トランザクション1 2200およびトランザクション2 2250という2つのトランザクションを示す。現行状態ポインタ120はオリジナル状態110をポイントし、これによって現行状態またはデータベースの現行状態を示す。トランザクション1 2200の場合、デルタページは、読み取りバッファ2210、書き込みバッファ2220、

10

20

30

40

50

およびトランザクションポインタ 2 2 3 0 を作成する。前述のように、デルタページは、各トランザクションについて 1 つまたは複数の書き込みバッファを作成することができる。デルタページは、読み取りバッファがサイズまたはその他の理由で制限されている場合、複数の読み取りバッファを作成することもできる。デルタページは、読み取りバッファ 2 2 1 0 を現行状態にポイントさせ、書き込みバッファ 2 2 2 0 を読み取りバッファ 2 2 1 0 にポイントさせ、トランザクションポインタ 2 2 3 0 を書き込みバッファ 2 2 2 0 にポイントさせる。これに対応して、トランザクション 2 2 2 5 0 の場合、デルタページは読み取りバッファ 2 2 6 0、書き込みバッファ 2 2 7 0、およびトランザクションポインタ 2 2 8 0 を作成する。デルタページは、読み取りバッファ 2 2 6 0 を現行状態にポイントさせ、書き込みバッファ 2 2 7 0 を読み取りバッファ 2 2 6 0 にポイントさせ、トランザクションポインタ 2 2 8 0 を書き込みバッファ 2 2 7 0 にポイントさせる。

10

【 0 0 6 6 】

トランザクション 1 2 2 0 0 についての読み取りバッファ 2 2 1 0 は、トランザクション 1 2 2 0 0 がアドレス 3 でデータを読み取ったことを示す。トランザクション 2 2 2 5 0 についての読み取りバッファ 2 2 6 0 も、トランザクション 2 2 2 5 0 がアドレス 3 でデータを読み取ったことを示す。しかしながら、たとえそれぞれのトランザクションが同じデータを読み取ったとしても、どちらのトランザクションもそのアドレスのデータを変更しようとししない。したがって、デルタページはどちらのトランザクションも打ち切らなくてよい。

20

【 0 0 6 7 】

空の書き込みバッファと同様に、読み取りバッファは、データベースの状態を変更しない。したがって、デルタページの諸実施形態は、データが格納されたオブジェクトであるかのように、データベースの現行状態にない読み取りバッファを採用することができる。その代わりに、デルタページは、トランザクションが適用しようとする変更を格納するためにデルタページが作成できる部分マッピングの外部に格納することができる。読み取りバッファは、別の場所に格納することが可能であり、トランザクションポインタは、読み取りバッファの場所を示すために別のポインタを維持することができる。

【 0 0 6 8 】

介入トランザクションが現行トランザクションによって読み取られたアドレスのデータを変更する場合にのみ、デルタページが現行トランザクションを打ち切る場合、図 2 5 ~ 図 2 7 を参照しながら以下で説明するように、読み取りバッファは、読み取ったアドレスおよび値を格納するはずである。他方で、介入トランザクションが現行トランザクションによって読み取られたデータを上書きした場合に、デルタページがトランザクションを打ち切る場合、たとえ介入トランザクションが同じ値を書き込んだとしても、読み取りバッファは、図 2 2 ~ 図 2 4 に示されるように、トランザクションが読み取る 1 つまたは複数のアドレスを含むだけでよい。

30

【 0 0 6 9 】

図 2 3 は、デルタページがトランザクション 1 2 2 0 0 をデータベースにコミットした後のデータベースの状態を示す。デルタページは、トランザクション 1 2 2 0 0 について作成された最新の書き込みバッファ 2 2 2 0 にポイントするように、現行状態ポインタ 1 2 0 を変更し、トランザクションポインタ 2 2 3 0 を削除する。トランザクション 2 2 2 5 0 が完了すると、デルタページはトランザクション 2 2 2 5 0 をコミットすべきかどうかを判別する。ここでも、デルタページがトランザクション 2 2 2 5 0 をコミットする前にデータベースの現行状態が変更されているため、図 1 5 ~ 図 2 1 を参照しながら前述したトランザクションをコミットするかどうかを判別するモードで、デルタページはトランザクション 2 2 2 5 0 を打ち切ることになる。これに対して、現行のモードは、より現実的な評価をする。

40

【 0 0 7 0 】

このモードでは、たとえデータベースの現行状態が変更されていても、介入トランザク

50

ションが現行トランザクションによってアクセスされるアドレスにデータを書き込んでいない限り、デルタページは現行トランザクションを打ち切らない。デルタページは、トランザクション2 2250がデータを読み取ったアドレスに、介入トランザクションであるトランザクション1 2200が書き込まれたかどうかを判別するために、トランザクション2 2250の読み取りバッファ2260を比較し、それによって、トランザクション2 2250の処理が依拠した状態を変更する。トランザクション2 2250は、読み取りバッファ2260に従って、アドレス3からデータを読み取るだけである。しかしながら、現行状態ポインタ120によって示された現行状態をチェックすると、トランザクション1 2200は、書き込みバッファ2220に従ってアドレス4にデータを書き込んだのみである。したがってデルタページは、データベースにトランザクション2 2250をコミットする。

10

【0071】

図24は、データベースの現行状態にトランザクション2 2250をコミットしているデルタページを示す。デルタページは、トランザクション2 2250の読み取りバッファ2260を、現行状態ポインタによって示されたオブジェクトである、トランザクション1 2200の書き込みバッファ2220にポイントさせる。次にデルタページは、トランザクション2 2250の書き込みバッファ2270を指示するように、現行状態ポインタ120を変更し、トランザクション2 2250をデータベースにコミットする。

【0072】

20

図25～図27は、介入トランザクションがデータを書き込んだアドレスに介入トランザクションが書き込むだけでなく、現行トランザクションが読み取ったデータ値を介入トランザクションが実際に変更する場合に、デルタページが現行トランザクションを打ち切るだけの、それほどリスクを嫌わないモードを示す。図25は、データベースの現行状態に対して実行される、トランザクション1 2500およびトランザクション2 2550の、2つのトランザクションを示す。現行状態は、データベースのオリジナル状態110をポイントする、現行状態ポインタ120によって示される。トランザクション1 2500の場合、デルタページは、読み取りバッファ2510、書き込みバッファ2420、およびトランザクションポインタ2530を作成する。デルタページは、読み取りバッファ2510をその時点での現行状態にポイントさせ、書き込みバッファ2520を読み取りバッファ2510にポイントさせ、トランザクションポインタ2230を書き込みバッファ2520にポイントさせる。これに対応して、トランザクション2 2550の場合、デルタページは読み取りバッファ2560、書き込みバッファ2570、およびトランザクションポインタ2580を作成する。デルタページは、読み取りバッファ2560をその時点での現行状態にポイントさせ、書き込みバッファ2570を読み取りバッファ2560にポイントさせ、トランザクションポインタ2580を書き込みバッファ2570にポイントさせる。

30

【0073】

図22～図24の例とは対照的に、トランザクション1 2500についての読み取りバッファ2510は、トランザクション1 2500がアドレス3の値を読み取ったこと、およびアドレス3から読み取られた値がヌルであったことを示す。トランザクション2 2550についての読み取りバッファ2560は、トランザクション2 2500もアドレス3の値を読み取ったこと、および、アドレス3から読み取られた値もヌルであったことを見つけたことを示す。トランザクション2 2550についての書き込みバッファ2570は、アドレス3の値への変更を格納し、格納された値を文字列「CAT」に変更する。

40

【0074】

図26は、デルタページがトランザクション2 2550をデータベースにコミットした後のデータベースの現行状態を示す。したがってデルタページは、トランザクション2 2550がデータベースに適用するいずれかの変更を格納している書き込みバッ

50

ァ 2 5 7 0 をポイントするように、現行状態ポインタ 1 2 0 を変更し、トランザクションポインタ 2 5 8 0 を破棄する。このモードでは、トランザクション 1 2 5 0 0 が完了すると、デルタページは、トランザクション 2 2 5 5 0 が変更した値をトランザクション 1 2 5 0 0 が読み取ったかどうかを判別する。

【 0 0 7 5 】

デルタページは、トランザクション 1 2 5 1 0 の読み取りバッファ 2 5 1 0 と、現行状態ポインタ 1 2 0 によって示されるデータベースの現行状態とを比較する。読み取りバッファ 2 5 1 0 によれば、トランザクション 1 はアドレス 3 を読み取り、値がヌルであることがわかった。さらにデルタページは現行状態ポインタ 1 2 0 によって示されるデータベースの現行状態を追い、トランザクション 2 2 5 5 0 がアドレス 3 で格納された値を「C A T」に変更したことを見つける。したがって、トランザクション 1 2 5 0 0 がデータベースに適用しようとするいかなる変更も、今では古くなったデータに基づいている可能性がある。

10

【 0 0 7 6 】

したがって、デルタページは、図 2 7 に示されるようにトランザクション 1 2 5 0 0 を打ち切る。デルタページは、読み取りバッファ 2 5 1 0、書き込みバッファ 2 5 2 0、およびトランザクションポインタ 2 5 3 0 を解放または削除し、デルタページは、現行状態ポインタ 1 2 0 を、トランザクション 2 2 5 5 0 がデータベースに適用した変更を格納している書き込みバッファ 2 5 7 0 をポイントしたままにする。所望であれば、トランザクション 1 2 2 0 0 を前述のように再開することができる。

20

【 0 0 7 7 】

たとえトランザクション 2 2 5 5 0 がトランザクション 1 2 2 0 0 によって読み取られたアドレスに新しい値を書き込んだとしても、トランザクション 2 2 2 5 0 が書き込んだデータがトランザクション 1 2 5 0 0 の結果を変更している可能性はないことに留意されたい。しかしながら、トランザクション 1 2 5 0 0 の結果が、トランザクション 2 2 5 5 0 によって書き込まれたデータの結果として変更されるかどうかを判別するために、トランザクション 1 2 5 0 0 を再実行しなければならない。デルタページのこのモードは、現行状態が変更されているというだけで、トランザクション 1 2 5 0 0 を打ち切らず、トランザクション 1 2 5 5 0 の再実行にかかる時間が足りない現行トランザクションを打ち切るべきであるかどうかのいくつかの実質的な分析を提供する。

30

【 0 0 7 8 】

図 2 8 は、デルタページがトランザクションをコミットすべきであるかどうかを判別するより実質的なモードを容易にするために、読み取りバッファおよび単一書き込みバッファを作成することによって、デルタページがトランザクションを処理するために使用するモードを示す。デルタページは、読み取りバッファの使用と共に単一書き込みバッファの使用に限定されておらず、単一書き込みバッファの選択は、単なる 1 つの可能な代替である。

【 0 0 7 9 】

流れ図 2 8 0 0 は、ブロック 2 8 0 2 がトランザクションを受け取ると開始される。ブロック 2 8 0 4 は、トランザクションがデータベースに適用しようとする変更の部分マッピングを示すことになる、切り替えポインタを作成する。ブロック 2 8 0 6 は読み取りバッファを作成し、これを、トランザクションが開始された時点で存在した状態にポイントさせる。ブロック 2 8 0 8 は単一書き込みバッファを作成し、この単一書き込みバッファを読み取りバッファにポイントさせる。ブロック 2 8 1 0 は、トランザクションが適用しようとする変更の部分状態を示すために、切り替えポインタを単一書き込みバッファにポイントさせる。

40

【 0 0 8 0 】

ブロック 2 8 1 2 は、トランザクションによって読み取られたすべてのデータを読み取りバッファに格納する。前述のように、介入トランザクションが現行トランザクションと同じデータを読み取った場合、デルタページがトランザクションを打ち切ることができ

50

るモードでは、読み取りバッファは単に読み取ったデータのアドレスを格納すればよい。他方で、現行トランザクションによって読み取られたデータを介入トランザクションが書き込んだかどうかに基づいて、トランザクションをコミットするかどうかをデルタページャが判別する場合、読み取りバッファはアドレスとそのアドレスから読み取ったデータとの両方を格納すべきである。

【0081】

ブロック2814は、トランザクションがデータベースに変更を適用しようとするかどうかを判別する。適用しようとする場合、ブロック2816は、切り替えポインタを単一書き込みバッファを避けてポイントさせる。ブロック2818は、単一書き込みバッファへの変更を格納する。ブロック2820は、切り替えポインタを再度単一書き込みバッファにポイントさせる。

10

【0082】

ブロック2822は、トランザクションが完了するかまたは実行を継続するかを判別する。トランザクションが完了しない場合、流れ図2800はブロック2812にループして、トランザクションによって読み取られたデータがあれば読み取りバッファに格納する。他方で、ブロック2822がトランザクションを完了すると判別した場合、ブロック2824はトランザクションによって実行された変更のコミットを試行することになる。ブロック2826は他のトランザクションの受け取りを待つ。

【0083】

図29は、デルタページャが読み取りバッファを使用して、トランザクションをコミットするかまたは打ち切るかを判別するモードを示す。流れ図2900はブロック2902で開始される。ブロック2902は、完了されたトランザクションがデータベースをコミットしようとする場合を決定し、流れ図2900はトランザクションがデータベースにコミットしようとするまでブロック2902へとループする。ブロック2902が、完了されたトランザクションがコミットしようとするのを検出すると、ブロック2904は、トランザクションが開始されて以来、現行状態が変更されたかどうかを判別する。

20

【0084】

現行状態が変更されていない場合、ブロック2912は現行状態ポインタを、単一書き込みバッファ、または、トランザクションポインタも指示しているバッファである、トランザクションについて作成された複数の書き込みバッファのうちの最新のバッファにポイントさせる。ブロック2914は、トランザクションポインタを削除し、データベースの更新を完了する。データベースの現行状態が変更されていない場合、デルタページャは、読み取りバッファの検査またはデータベースへの適用の検索なしに、トランザクションをコミットすることができる。現行状態が変更されていない場合、いかなる介入トランザクションも現行トランザクションが依拠するデータを変更することはない。流れ図2900は、データベースにコミットしようとしている次の完了されたトランザクションを検出するために、ブロック2902へとループする。

30

【0085】

他方で、ブロック2904が、トランザクションが開始されて以来、現行状態が変更されていると決定した場合、ブロック2906は、読み取りバッファ内に列挙されたデータが現行状態で変更されているかどうかを判別する。ここでも、デルタページャのこのモードは、デルタページャが介入トランザクションをコミットする結果として、現行状態が変更されているかどうか、読み取りバッファを検査するのみである。読み取りバッファ内に列挙されたデータが現行状態で変更されていない場合、デルタページャはトランザクションをコミットすることになり、流れ図2900はブロック2912へと進む。

40

【0086】

他方で、ブロック2906が現行状態で、読み取りバッファ内のアドレスにデータが書き込まれたことを見つけた場合、または、そのアドレスでのデータの値が読み取りバッファ内に格納されたものと異なる場合、デルタページャはトランザクションを打ち切る。ブロック2908は、トランザクションについての読み取りおよび書き込みバッファを解放

50

する。次に、ブロック 2 9 1 0 は、現行トランザクションについてのトランザクションポインタを解放または削除し、トランザクションの打ち切りを完了する。次に、流れ図 2 9 0 0 は、ブロック 2 9 0 2 へとループして、コミットしようとする次のトランザクションを待つ。

【 0 0 8 7 】

トランザクションの合体および状態の保持

デルタページの書き込みバッファは、多くの利点を提供する。一例を挙げると、トランザクションがコミットできず、打ち切らなければならない場合、デルタページは、そのトランザクションについて作成した書き込みバッファを無視または削除し、デルタページは、そうでなければデータベースに適用された可能性のある誤った変更または競合する変更を再書き込みするかまたは取り消す必要はない。しかしながら、デルタページが多くの書き込みバッファを現行状態に追加した後では、現行状態におけるすべてのデータへのアクセスは非効率的になる可能性がある。たとえば比較的単純な図 2 0 の例であっても、トランザクションは、トランザクションが読み取ろうとしてデータの値を決定するために、書き込みバッファ 1 9 3 0、1 9 1 0、5 1 0、4 1 0、3 1 0、およびオリジナル状態 1 1 0 を通って再度読み取らなければならない可能性がある。デルタページの 1 つのモードでは、データベースの現行状態への効率的なアクセスを保証するために、書き込みバッファを合体する。

【 0 0 8 8 】

図 3 0 は、デルタページが図 2 0 の書き込みバッファを合体することができる 1 つのモードを示す。図 3 0 で、デルタページは、書き込みバッファに格納された変更をデータベースのオリジナル状態 1 1 0 に適用することによって、書き込みバッファを合体する。第 1 の段階 3 0 0 0 で、デルタページは、書き込みバッファが合体されることになるポイントを示す合体ポインタ 3 0 1 0 を認識する。合体ポインタ 3 0 1 0 の利点について、以下で詳細に説明する。デルタページは、次に、オリジナル状態 1 1 0 をポイントする第 1 の書き込みバッファ内に格納された変更を適用する。図 3 0 に示されるように、第 1 の書き込みバッファは空の書き込みバッファ 3 1 0 であった。デルタページは、この空の書き込みバッファを現行状態から除外することによって、書き込みバッファ 3 1 0 などの空のバッファを合体する。

【 0 0 8 9 】

空の書き込みバッファ 3 1 0 は、トランザクションがデータベースに適用しようとする変更を表す部分マッピングを変更しないことに留意されたい。したがって、デルタページが、トランザクションが適用しようとする変更を格納する書き込みバッファを作成するとすぐに、空の書き込みバッファが破棄される可能性があるか、またはデルタページが空の書き込みバッファ 3 1 0 を作成さえしない可能性もある。たとえば、アドレス 4 の値を変更するために書き込みバッファ 5 1 0 が追加されると、デルタページの実施形態は、書き込みバッファ 5 1 0 を、トランザクションが開始された時点で存在する状態にポイントさせ、空の書き込みバッファ 3 1 0 を破棄する可能性がある。

【 0 0 9 0 】

第 2 の段階 3 0 2 0 で、デルタページは、次の書き込みバッファ 4 1 0 に格納された変更をオリジナル状態 1 1 0 に適用し、オリジナル状態を更新済み状態 3 0 3 0 に置き換える。書き込みバッファ 4 1 0 に格納された、アドレス 2 の値を 1 に設定する変更は、更新済みバッキングストア 3 0 3 0 を作成するために、アドレス 2 3 0 4 0 に適用される。

【 0 0 9 1 】

最終結果 3 0 5 0 は、データベースの新しい状態 3 0 6 0 を示す。新しい状態 3 0 6 0 では、書き込みバッファ 5 1 0 に格納された変更が、アドレス 4 3 0 7 0 で格納された値を文字列「D O G」に変更するために適用され、書き込みバッファ 1 9 3 0 に格納された変更が、アドレス 3 3 0 8 0 で格納された値を文字列「C A T」に変更するために適用された。デルタページは、ここで、新しい状態 3 0 6 0 をポイントするように現行状

10

20

30

40

50

態ポインタ 1 2 0 を切り替える。その後デルタページは、合体ポインタ 3 0 1 0 を除外する。

【 0 0 9 2 】

デルタページのバッファの合体に関して注目する 4 つのポイントがある。第 1 に、たとえ合体プロセス中であっても、各ポインタはデータベースの同じ状態を不変的にポイントし続ける。たとえば第 2 の段階 3 0 2 0 で、デルタページが、更新済み状態 3 0 3 0 を作成するために、書き込みバッファ 4 1 0 内の変更をオリジナル状態 1 1 0 に適用した場合、他のポインタによって示される状態も同様に維持された。現行状態ポインタ 1 2 0、書き込みバッファ 5 1 0、1 9 1 0 および 1 9 3 0、ならびに合体ポインタ 3 0 1 0 は、すべて、すべて同じ値を提示するデータベースの状態を依然としてポイントしていた。

10

【 0 0 9 3 】

第 2 に、現行のデータベースを変更するために 1 つまたは複数の新しいトランザクションが開始されコミットされた場合、新しいトランザクションはデータベースの不変状態をポイントし続けることになった。新しいトランザクションを開始すると、デルタページは、現行状態ポインタ 1 2 0 によって示された現行状態をポイントする書き込みバッファを作成する。デルタページがトランザクションをコミットする場合、デルタページは、それぞれの後続トランザクションについて作成された最新の書き込みバッファをポイントするように、現行状態ポインタ 1 2 0 を切り替える。したがってデルタページは、たとえ書き込みバッファの合体中であっても、新しいトランザクションに関する状態の不変性を保持する。

20

【 0 0 9 4 】

第 3 に、合体ポインタ 3 0 1 0 は合体プロセスが停止するポイントを示す。デルタページが追加の書き込みバッファを現行状態にコミットする場合、デルタページは、デルタページが他の合体に従事するまで追加のバッファを合体しない。合体は無限に続けることができるが、数個の可能性のある新しい書き込みバッファを継続的に合体するために計算リソースを充当することは、計算リソースの効率的な使用ではない可能性がある。

【 0 0 9 5 】

第 4 に、図 3 1 に示されるように、デルタページは、データベースのオリジナル状態 1 1 0 を上書きすることなしにバッファを合体することができる。たとえばオリジナル状態 1 1 0 は、アクセスが非効率的な可能性があるディスクストレージに格納することができるか、または後続のトランザクションによって変更が適用されることなく、オリジナル状態 1 1 0 を維持することが望ましい可能性がある。したがって、デルタページは、1 つまたは複数の中間オブジェクトに書き込みバッファを合体することができる。

30

【 0 0 9 6 】

図 3 1 で、書き込みバッファを中間オブジェクトに合体する第 1 の段階 3 1 0 0 では、デルタページは、合体される第 1 の書き込みバッファである書き込みバッファ 3 1 0 と、オリジナル状態 1 1 0 との間に、中間の合体済みオブジェクト 3 1 1 0 を作成する。合体済みオブジェクト 3 1 1 0 の挿入は、データベース内のいかなる現行状態の不変性も変更せず、現行状態ポインタ 1 2 0 および書き込みバッファ 3 1 0、4 1 0、5 1 0、1 9 1 0、および 1 9 3 0 は、依然として同じデータベースの状態をポイントする。合体済みオブジェクト 3 1 1 0 ならびにオリジナル 1 1 0 は、キャッシュメモリ、メモリ、ディスク、または他の形のストレージを含む、任意の所望のストレージ内に維持することができる。

40

【 0 0 9 7 】

最終結果 3 1 2 0 は、デルタページが、書き込みバッファ 4 1 0、5 1 0、および 1 9 3 0 に格納されたすべての変更を、合体済みオブジェクト 3 1 1 0 に合体したことを示す。したがって、書き込みバッファ 4 1 0 に格納された、アドレス 2 の値を 1 に変更するという変更は、オリジナル状態 1 1 0 に適用される代わりに、アドレス 2 3 1 3 0 で合体済みオブジェクト 3 1 1 0 に格納される。同様に、デルタページは、書き込みバッファ 5 1 0 および 1 9 3 0 からの変更を、それぞれアドレス 3 3 1 4 0 およびアドレス 4

50

3150で、合体済みオブジェクト3110に格納する。デルタページは、合体済みオブジェクト3110を指示するように現行状態ポインタ120を変更し、この合体済みオブジェクト3110はオリジナル状態110をポイントする。したがって、データベースの状態は決して変更されず、現行状態ポインタ120は、デルタページが書き込みバッファを合体済みオブジェクト3110に合体する前に存在したデータベースと同じ状態を引き続きポイントする。

【0098】

有利なことに、データを読み取ろうとしている新しいトランザクションは、ここで、合体済みオブジェクト3110およびオリジナル状態110という2つのデータストアをチェックするだけでよい。後者が望ましい場合、デルタページは、合体済みオブジェクト3110のコンテンツをオリジナル状態に合体し、更新済みバックリングストアを作成することができる。別の方法として、デルタページは、合体済みオブジェクトを維持し、同じ合体済みオブジェクト3110への書き込みバッファの後続の合体を実行することができる。さらに別の方法として、スナップショットに関して以下で詳細に説明するように、デルタページは追加の合体済みオブジェクトを作成することができる。

【0099】

前述のように、デルタページはポインタの不変性を維持する。データベース内の合体済みオブジェクトでは、各オブジェクトが追加された時点で存在した中間状態を表すバッファ間のポインタは、現行状態が保持される限り消去される。しかしながら、後続バッファからのポインタ以外のポインタがバッファをポイントする場合、デルタページによって認識される不変性は、ポインタによって示される状態がバッファの合体後も存続しなければならないことを保証する。デルタページのポインタの不変性を利用して、デルタページは、選択された状態をポイントするスナップショットポインタを遵守し、たとえばスナップショットポインタの先行または後続のバッファをデルタページが合体する場合であっても、その状態を保持することを遵守する。

【0100】

図32は、デルタページがスナップショットポインタ3210と共に保持することになるデータベースの現行状態を示す。現行状態ポインタ120は、書き込みバッファ3220、3230、および3240によってオリジナル状態110が修正された、現行状態をポイントする。第1の書き込みバッファ3220は、文字列「PEAR」をアドレス0に書き込む。第2の書き込みバッファ3230は、値1をアドレス2に書き込む。第3の書き込みバッファ3240は、文字列「DOG」をアドレス4に書き込む。

【0101】

デルタページが、現行状態を保持するように命じられた場合、デルタページは、現行状態をポイントするスナップショットポインタ3210を作成する。デルタページは、たとえば定期的なバックアップまたはユーザ要求を自動的に保持するプログラムが、現行状態などの状態のスナップショットを要求した場合、状態を保持するためにスナップショットポインタ3210を挿入する。スナップショットポインタ3210は選択された状態をポイントすることから、デルタページによって遵守されるポインタの不変性の結果として、デルタページはポインタおよび対応する状態を維持することになる。

【0102】

図33は、追加のトランザクションが現行状態を変更した後の、データベースのその後の状態を示す。第4の書き込みバッファ3310は、文字列「PEACH」をアドレス0に書き込む。第5の書き込みバッファ3320は、文字列「CAT」をアドレス4に書き込む。デルタページは、第5の書き込みバッファ3320をポイントするように現行状態ポインタ120を切り替える。図33で、第5の書き込みバッファ3320は、アドレス4に格納された書き込みバッファ3240の文字列「DOG」を、文字列「CAT」で上書きする。しかしながら、スナップショットポインタ3310は、図32の状態を表す書き込みバッファ3240を依然としてポイントする。

【0103】

10

20

30

40

50

図34は、デルタページャによって許可された合体を示す。デルタページャがオリジナル状態110を上書きしないケースを想定すると、デルタページャは、書き込みバッファ3220、3230、および3240によって適用された変更を集めることによって、スナップショットポインタ3210と、第1の合体済みオブジェクト3410内のオリジナル状態のバックングストアとの間で、データベースの現行状態におけるオブジェクトを合体する。次にデルタページャは、書き込みバッファ3310および3320によって適用された変更を第2の合体済みオブジェクト3320内に集める。現行状態ポインタ120は、この第2の合体済みオブジェクト3420をポイントする。第2の合体済みオブジェクト3420は、第1の合体済みオブジェクト3410をポイントし、これはオリジナル状態110をポイントする。

10

【0104】

デルタページャは、スナップショットオブジェクト3210のポインタの不変性を保持するために、別々の合体済みオブジェクト3410および3420を使用する。スナップショットオブジェクト3210が存在する限り、デルタページャは図32の選択された状態を維持することになる。したがって、たとえば図33および図34の現行状態では、デルタページャは、書き込みバッファ3240によってアドレス4に格納された文字列「DOG」を、書き込みバッファ3320によってアドレス4に格納された文字列「CAT」で置き換える。しかしながら、スナップショットポインタ3210によって保持された選択された状態では、文字列「DOG」は依然として第1の合体済みオブジェクト3410内のアドレス4に格納されたままである。

20

【0105】

スナップショットポインタ3210に関して注目する5つのポイントがある。第1にデルタページャは、スナップショットポインタを追加することによって、任意の選択された時点でデータベースの状態を保持することができる。第2にデルタページャは、複数の状態を保持するために複数のスナップショットポインタを含むことができる。第3に、合体に関して前述したように、デルタページャは、ポインタの不変性により、1つまたは複数の以前の状態を保持しながら、トランザクションの受け取りおよびデータベースへの変更のコミットを続行することができる。第4に、スナップショットポインタ3210によって保持された状態がもはや必要でない場合、スナップショットポインタ3210は除去され、後続の合体によって保持された状態が解放される。

30

【0106】

第5に、スナップショットポインタは、デルタページャが複数の異なる状態を単一のデータストア内で維持できるようにする。現行状態を維持するために従来のデータベースのバックアップイメージが保存される場合、後続のイメージが以前のイメージを上書きするため、各イメージは別々に格納される。複数の状態は、大量のストレージを消費する可能性があり、バックアップイメージは、高速デバイス内のストレージスペースを保持するために、しばしば低速ストレージデバイスへ追いやられる。しかしながら、デルタページャは以前の状態の上に後続の状態を構築することから、スナップショットポインタは単一のデータストア内に以前の状態を保持する。

40

【0107】

図35は、デルタページャが、バッファまたは、以前に作成されてここでさらに合体されることになる合体済みオブジェクトなどの、他のオブジェクトを合体する際に使用する1つのモードを示す。流れ図3500はブロック3502で開始される。ブロック3502は、データベースの現行状態でオブジェクトの合体を開始する。プログラムが周期的に合体を開始するか、またはユーザが合体を開始することができる。ブロック3504は、図34の例にあるようなオリジナル状態、最新の合体済みオブジェクト、またはデルタページャがデータベースに追加されたバッファを合体することになる他のポイントを含むことができる、合体開始ポイントを識別する。ブロック3506は、次の合体されていないオブジェクトに進む。図33の例で、ブロック3504は、オリジナル状態110から始まり、次の合体されていないオブジェクトである書き込みバッファ3220として、合体

50

のための開始ポイントを識別する。

【 0 1 0 8 】

ブロック 3 5 0 8 は、次に作成されるバッファから以外のバッファへのいずれかのポイントが存在するかどうかを判別する。デルタページは、データベースの状態が変更されない場合、一連のバッファまたは他のオブジェクトを合体する際にポイントを消去することができるが、デルタページは、スナップショットポイント 3 2 1 0 からのポイントなどの他のポイントは消去しない。ブロック 3 5 0 8 が、次のバッファから以外のポイントが存在しないと判別した場合、ブロック 3 5 1 0 は、このバッファと次に作成されるバッファとを合体済みオブジェクト内に組み入れる。ブロック 3 5 1 2 は、合体する可能性のある追加のオブジェクトが存在するかどうかを判別する。追加のオブジェクトは次に作成されるバッファを含む可能性があるか、または、データベース内のオブジェクトが事前に合体されている場合は、さらに合体されることになる合体済みオブジェクトが存在する可能性がある。合体する可能性のある追加のオブジェクトが存在する場合、流れ図 3 5 0 0 はブロック 3 5 0 6 へとループして、次の合体されていないオブジェクトに進む。しかしながらブロック 3 5 1 2 は、合体する他のオブジェクトがないと決定する場合がある。たとえば、デルタページがデータベースを完全に合体した場合、デルタページは、合体されていない現行状態ポイント 1 2 0 に到達することになる。ブロック 3 5 1 2 が、合体する追加のオブジェクトが存在しないと決定した場合、ブロック 3 5 1 4 は合体を完了する。合体の完了により、たとえば動作の完了を確認するメッセージを送信またはログ記録する場合がある。

10

20

【 0 1 0 9 】

他方で、ブロック 3 5 0 8 は、状態内の次のオブジェクト以外からのポイントが存在すると決定する場合がある。たとえば図 3 2 ~ 図 3 4 に示されるように、スナップショットポイント 3 2 1 0 は、オブジェクトが表す状態を保持するためのオブジェクトをポイントすることができる。ブロック 3 5 0 8 が、他のポイントが存在すると決定した場合、流れ図 3 5 0 0 はブロック 3 5 1 2 へと進み、合体する可能性のある追加のオブジェクトが存在するかどうかを判別する。

【 0 1 1 0 】

状態の永続化

デルタページは、データベースの選択された部分を永続的にすることができる。デルタページは、データベースの選択された部分を保持するために、不揮発性ストレージにコミットすることになる。

30

【 0 1 1 1 】

図 3 6 は、永続媒体であるディスク 3 6 1 0 上に部分的に、および揮発性媒体であるメモリ 3 6 2 0 内に部分的に格納された、オブジェクトを含む、データベースの現行状態を示す。現行状態ポイント 1 2 0 は、ディスク 3 6 1 0 上に格納されたオリジナル状態 1 1 0 を含み、デルタページが書き込みバッファ 3 6 3 0 および 3 6 4 0 内の変更をコミットした、現行状態をポイントする。

【 0 1 1 2 】

1 つのモードで、デルタページは、永続ポイント 3 6 5 0 を挿入することによって、プログラムまたはユーザがデータベースの選択した部分を永続的にできるようにする。永続ポイント 3 6 5 0 は、永続的にされるオブジェクトを含むデータベースの状態をポイントする。図 3 6 で、永続ポイント 3 6 5 0 は、第 1 の書き込みバッファ 3 6 3 0 を含む状態をポイントする。

40

【 0 1 1 3 】

図 3 7 が示すように、デルタページは、永続ポイント 3 2 5 0 によって示される状態に含まれたオブジェクトがあれば、ディスクストレージ 3 7 1 0 などの永続格納可能ストレージに移動する。この例では、第 1 の書き込みバッファ 3 6 3 0 はディスクストレージ 3 7 1 0 に移動される。デルタページは、データベースの状態を、任意数のストレージ媒体にわたって常駐させることができる。したがって、デルタページが第 1 の書き込み

50

バッファ 3 7 3 0 を書き込むディスクストレージ 3 7 1 0 は、オリジナル状態 1 1 0 が常駐するディスクストレージ 3 6 1 0 と同じストレージデバイスとするか、または別のストレージデバイスとすることができる。永続ポインタ 3 6 5 0 がポイントする状態に含まれていない第 2 の書き込みバッファ 3 6 4 0 は、メモリ 3 6 2 0 内に残される。

【 0 1 1 4 】

デルタページャが、永続ポインタ 3 6 5 0 によって示されるデータベースの現行状態の一部を永続ストレージに格納すると、デルタページャは永続ポインタを解放することができる。デルタページャは、データベースの選択された部分が永続的にされていることを確認するメッセージを生成することができる。

【 0 1 1 5 】

図 3 8 は、デルタページャがバッファを永続的にする際の 1 つのモードを示す。流れ図 3 5 0 0 はブロック 3 8 0 2 で開始される。ブロック 3 8 0 2 は永続要求を受け取る。ブロック 3 8 0 4 は、永続的にするデータベースの部分を決定するために永続ポインタを配置する。ブロック 3 8 0 6 は、永続ポインタと、永続的にされるデータベース内の最新の永続ストレージ部分との間の、すべてのオブジェクトを識別する。ブロック 3 8 0 8 は、永続ポインタと最新の永続ストアとの間のすべてのオブジェクトを、永続ストレージオブジェクトにコピーする。ブロック 3 8 1 0 は、永続的にされる現行状態の部分をポイントするいかなるポインタも、永続ストレージオブジェクトに変更する。図 3 7 の例では、デルタページャは、書き込みバッファ 3 6 4 0 のポインタを永続オブジェクトに変更する。別の方法として、たとえば、現行状態ポインタ 1 2 0 またはスナップショットポインタ 3 2 1 0 などの、永続ポインタ 3 6 5 0 がポイントする状態への他のポインタが存在する場合、また、これらのポインタは、ブロック 3 8 0 8 によって作成された永続オブジェクトへポイントされる。

【 0 1 1 6 】

永続ポインタ 3 6 5 0 がポイントする状態内のオブジェクトが永続オブジェクトにコピーされ、この状態へのポインタが変更されると、ブロック 3 8 1 2 は、メモリを他の用途に解放するために、永続ポインタがポイントするメモリ内のオブジェクトを解放する。ブロック 3 8 1 4 は、永続ポインタを削除する。ブロック 3 1 6 は、たとえば、永続要求を実行するメッセージをユーザに送信すること、または、メッセージをシステムログに追加することによって、永続要求の完了を確認する。

【 0 1 1 7 】

データのキャッシュ

デルタページャは、性能を向上させるためにデータのキャッシュも提供する。デルタページャは、図 7 ~ 図 9 を参照しながら説明した一連の累積バッファを使用するなどにより、データの重複ストレージを可能にする。状態内でのデータの重複ストアが状態を変更することはない。したがってデルタページャは、状態を変更することなく、ディスクストレージまたは他の大容量記憶装置などの低速ストア上に格納されたデータを、メモリなどの高速ストア内にキャッシュすることができる。

【 0 1 1 8 】

図 3 9 は、キャッシュオブジェクト 3 9 1 0 を含むデータベースの現行状態を示す。キャッシュオブジェクトポインタ 3 9 1 0 は、ディスクストレージ 3 9 5 0 に格納された永続オブジェクト 3 9 3 0 などの、キャッシュされるオブジェクトをポイントする。デルタページャは、このケースでは読み取りバッファ 3 9 6 2 である、永続オブジェクト 3 9 3 0 に格納されていない第 1 のオブジェクトを、キャッシュオブジェクト 3 9 1 0 にポイントさせる。

【 0 1 1 9 】

現行状態では、トランザクションがアドレス 0 の値を検索すると、トランザクションは、文字列「CAT」を現行状態に書き込む書き込みバッファ 3 9 6 8 内のアドレス 0 の現行値を見つけることになる。トランザクションは、書き込みバッファ 3 9 6 8 以外は探す必要がなく、トランザクションはメモリ取り出し速度でアドレス 0 の値を受け取ることに

10

20

30

40

50

なる。他方で、トランザクションがアドレス2に格納された値を求める場合、デルタページはこれを見つけるまで現行状態全体にわたってこの値を求めることになる。キャッシュオブジェクト3910がない場合、デルタページはディスクストレージ3950内の永続オブジェクト3930まで探しに行くことになる。ディスクストレージ3950からのデータへのアクセスは、メモリ3960へのアクセスに比べて遅い。キャッシュオブジェクト3910をメモリ3930内に作成することで、アクセス効率を向上させることができる。

【0120】

図39の例では、永続オブジェクト3930は、値1を格納するアドレス2 3932、文字列「CAT」を格納するアドレス3 3934、および文字列「DOG」を格納するアドレス4 3936などの、アドレスの範囲で値を格納する。永続オブジェクト3930は、文字列「ROCK」を格納するアドレス99 3938までの、データが格納されないヌル値の範囲3940も格納する。

【0121】

デルタページは、メモリ3960内にキャッシュオブジェクト3910を作成する。デルタページは、キャッシュオブジェクト3910を、このケースでは永続オブジェクト3930である、キャッシュされているオブジェクトにポイントさせる。次にデルタページは、第1の非永続オブジェクトである読み取りバッファ3962を、キャッシュオブジェクト3910にポイントさせる。

【0122】

デルタページがキャッシュオブジェクトを作成し、現行状態にキャッシュオブジェクトを挿入するようにポインタを変更すると、デルタページはキャッシュを読み込む（populate）。1つのモードでは、デルタページが永続オブジェクト3930からデータを取り出す場合、デルタページは同じデータをキャッシュオブジェクト3910に格納する。また、永続オブジェクトへのアクセスに時間を費やした後、デルタページは永続オブジェクト3930からデータのブロックを取り出し、これをキャッシュオブジェクト3910に格納することもできる。たとえば、キャッシュオブジェクト3910が読み込まれる前に、デルタページが読み取りバッファ3962を作成したトランザクションはアドレス4からのデータを求め、これが永続オブジェクト3930内のアドレス4

3936に格納される。永続オブジェクト3930にアクセスしている間、デルタページは、アドレス2 3932、アドレス3 3934、およびアドレス99 3938までの他のアドレスからのデータを含むブロックに関するデータを取り出し、そのデータをキャッシュオブジェクト3910内のそれらのアドレスに格納することもできる。

【0123】

デルタページが永続オブジェクト3930から取り出したデータをキャッシュオブジェクト3910に格納した後、キャッシュオブジェクト3910にコピーされたデータに関する後続の要求は、かなり速いメモリ速度で実行されることになる。したがって、その後、読み取りバッファ3962が作成されたトランザクションがアドレス99にデータを要求した場合、デルタページはキャッシュオブジェクト3910内のアドレス99 3918からデータを取り出すことができる。同様に、読み取りバッファ3964が作成されたトランザクションがアドレス3にデータを要求した場合、デルタページは、ディスクストレージ3950からデルタページがデータを取り出すのを待つことなく、キャッシュオブジェクト内のアドレス3 3914からデータを取り出すことができる。

【0124】

デルタページのキャッシュオブジェクト3910に関して注目する3つの特徴がある。第1に、キャッシュオブジェクト3910の挿入によって、デルタページがデータ状態を作成および使用する方法が変更されないことである。たとえば、たとえキャッシュオブジェクト3910内のアドレス4 3916が文字列「DOG」を格納している場合でも、書き込みバッファ3966は、書き込みバッファ3966内のアドレス4に文字列「BIRD」を書き込むことによって、データベースの状態を変更することができる。後続

10

20

30

40

50

のトランザクションでは、現行状態ポインタ120がポイントする現行状態は、キャッシュオブジェクト3910がアドレス4 3916に何を格納しているかにかかわらず、文字列「B I R D」がアドレス4に格納されていることを見つけることになる。

【0125】

第2に、キャッシュオブジェクト3910を使用する状態が永続的になっている場合、デルタページは、キャッシュオブジェクト3910を再開および再読み込みする必要がない。その代わりに、データが永続ストレージにコピーされるため、ディスクにコピーされている状態内のデータとキャッシュオブジェクト3910内に格納されたデータとが異なる場合、デルタページはキャッシュオブジェクト3910内のデータを更新することができる。別の方法としてデルタページは、永続ストレージにコピーされている状態と比した場合、もはや現行ではない、キャッシュオブジェクト3910内のエントリを無効にすることができる。どちらの技法も、古い状態のキャッシュを修正して、更新済み状態の有効キャッシュとする。このように実行することによって、デルタページはキャッシュオブジェクトを無期限に再使用する。キャッシュが更新全体にわたって保持できない場合、更新が永続的にされた場合に必ずキャッシュを再構築することで効率性が失われる。

【0126】

第3に、デルタページは、範囲の照会を容易にするためにキャッシュオブジェクト3910を使用することができる。キャッシュオブジェクトは、範囲の照会をバッキングストアに発行することによって、ヌルのラン(runs of null)を効率的に発見し、デルタページによって維持されるキャッシュオブジェクトまたは他のオブジェクトの一部でこうしたランをコンパクトに表すことができる。たとえば、ヌルのランの表現をキャッシュオブジェクトのフィールドに含めることができる。メモリに格納されたオブジェクト内でこの表現を維持することによって、範囲の照会を容易にするためのヌルのランの表現へのアクセス速度が改善される。たとえば、トランザクションが「B I R D」、「C A T」、または「D O G」の後に次のペットの文字列を見つけようとする、「D O G」の後の永続オブジェクト3930内にヌルデータの長い文字列3940が存在した場合、キャッシュオブジェクト3910は、範囲の照会をサポートするために、アドレス5 3920のように、ヌル値を格納しているラン内のアドレスの最後、または、非ヌル値を格納しているアドレスの最初を示す、第1のヌルアドレスフィールド内のエントリを含むことができる。

【0127】

図40は、デルタページがキャッシュオブジェクトを維持する際の1つのモードを示す。流れ図4000はブロック4002で開始される。ブロック4002は、メモリまたは他の高速アクセスストレージ内にキャッシュオブジェクトを作成する。ブロック4004はこのキャッシュオブジェクトを、キャッシュされているオブジェクトにポイントさせる。ブロック4006は、キャッシュされているオブジェクトをポイントするオブジェクトを識別し、このオブジェクトをキャッシュオブジェクトにポイントさせる。ブロック4008は、キャッシュオブジェクト内にキャッシュされているオブジェクトからトランザクションによって読み取られたアドレスおよびデータを格納する。前述のように、デルタページが、キャッシュされているオブジェクトからデータを取り出す場合、好ましくは、デルタページは、後続のディスクアクセス動作の必要性を潜在的に回避するために、読み取られたアドレスに隣接するアドレスのブロックからデータを取り出すことになる。

【0128】

ブロック4010は、キャッシュオブジェクトをポイントする状態が永続的とされるかどうかを判別する。永続的とされない場合、流れ図4000は、トランザクションによって読み取られたアドレスおよびデータを引き続きキャッシュオブジェクトに格納するために、ブロック4008へとループする。しかしながら、ブロック4010が、キャッシュをポイントする状態が永続的とされるものと決定した場合、ブロック4012は、永続的とされるオブジェクト内の変更を使用してキャッシュエントリを更新する。別の方法として、デルタページは、古くなったキャッシュエントリを更新する代わりに無効化するこ

とができる。ブロック 4 0 1 2 がキャッシュエントリを更新すると、ブロック 4 0 1 4 は、永続的とされない第 1 のオブジェクトをキャッシュオブジェクトにポイントさせる。

【 0 1 2 9 】

デルタページャによって使用されるオブジェクト

要約すると、図 4 1 は、前述の例示的動作を容易にするためにデルタページャが使用するオブジェクトを示す。この要約に表された汎用オブジェクトを参照するために、新しい参照番号が使用されている。図 4 1 は、永続ディスクストレージ 4 1 2 0 およびメモリ 4 1 3 0 に格納されたようなオブジェクトの範囲を含む、例示的状況 4 1 0 0 を示す。

【 0 1 3 0 】

デルタページャは、ディスクストレージ 4 1 2 0 内の 1 つまたは複数の永続オブジェクト 4 1 2 2 および 4 1 2 4 を使用することができる。これらのオブジェクトは、オリジナル状態および前述のように永続的とされた任意の他の状態を格納することができる。たとえば永続ポインタ 4 1 2 6 は永続的とされており、それがポイントする状態が永続オブジェクト 4 1 2 4 内に顕在化 (manifest) された後に除去されることになるため、破線で示されている。永続オブジェクトは、1 つまたは複数のデバイスに格納することができる。

10

【 0 1 3 1 】

デルタページャは、メモリ 4 1 3 0 に格納された複数のオブジェクトを含むことができる。デルタページャ 4 1 4 0 は、キャッシュ可能な永続オブジェクト 4 1 2 2 および 4 1 2 4 と後続のオブジェクトとの間に挿入された、キャッシュオブジェクト 4 1 4 0 を含むことができる。スナップショットポインタ 4 1 5 2 によって保持された状態を含むスナップショットオブジェクト 4 1 5 0 は、キャッシュオブジェクト 4 1 4 0 をポイントする。スナップショットオブジェクト 4 1 5 0 は、合体済みオブジェクトまたは一連の合体されていないオブジェクトを含む、選択された状態を含む。前述のように、スナップショットポインタ 4 1 5 2 は、たとえば保持された状態に含まれるかまたは保持された状態に後で追加されるオブジェクトが合体される場合であっても、ある状態を保持することができる。

20

【 0 1 3 2 】

最新の合体動作で合体済みオブジェクト 4 1 6 0 が生成されて以来、読み取りバッファ 4 1 7 0 および書き込みバッファ 4 1 8 0 などの複数のバッファが現在の状態に追加されている。現行状態ポインタ 4 1 9 0 はこの現行状態を示す。

30

【 0 1 3 3 】

例示的諸実施形態を実施するための動作環境

図 4 2 は、デルタページャを実施するための例示的動作環境 4 2 0 0 を示す。動作環境 4 2 0 0 は好適な動作環境の一例に過ぎず、前述のようなデルタページャの諸実施形態例または他の諸実施形態の使用または機能の範囲に関して、いかなる制限をも示唆することを意図するものではない。さらに動作環境 4 2 0 0 は、例示的動作環境 4 2 0 0 に示された構成要素のうちのいずれか 1 つまたはいずれかの組み合わせに関する、いかなる依存性または要件を有するものとも解釈されるべきではない。

【 0 1 3 4 】

デルタページャの実施プロセスは、動作環境 4 2 0 0 内で実行される、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コンテキストで説明することができる。通常、プログラムモジュールは、特定のタスクを実行するかまたは特定の抽象データ型を実施する、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。さらに、当業者であれば、デルタページャを実施するプロセスが、ハンドヘルドデバイス、マルチプロセッサシステム、マイクロプロセッサベースかまたはプログラム可能な大衆消費電化製品、ミニコンピュータ、メインフレームコンピュータ、などを含む、様々なコンピュータシステム構成を使用して実施可能であることを理解されよう。デルタページャを実施するプロセスは、通信ネットワークを介してリンクされたりリモート処理デバイスによってタスクが実行される、分散コンピューティング環境でも実施可能である。分散コンピューティング環境では、メモリストレージデバイスを含むローカルとリモートの両方の

40

50

コンピュータストレージ媒体内に、プログラムモジュールを配置することができる。

【 0 1 3 5 】

図 4 2 を参照すると、デルタページのプロセスを実施するための例示的動作環境 4 2 0 0 は、処理ユニット 4 2 2 0 と、システムメモリ 4 2 3 0 と、システムメモリ 4 2 3 0 を含む様々なシステム構成要素を処理ユニット 4 2 2 0 に結合するシステムバス 4 2 2 1 とを含む、コンピュータ 4 2 1 0 を含む。

【 0 1 3 6 】

コンピュータ 4 2 1 0 は、通常、様々なコンピュータ読み取り可能媒体を含む。例を挙げると、コンピュータ読み取り可能媒体は、コンピュータストレージ媒体および通信媒体を備えることができるが、これらに限定されることはない。コンピュータストレージ媒体の例には、ランダムアクセスメモリ (R A M)、読み取り専用メモリ (R O M)、電気的消去可能プログラム可能読み取り専用メモリ (E E P R O M)、フラッシュメモリまたは他のメモリ技術、C D R O M、デジタル汎用ディスク (D V D) あるいは他の光またはホログラフィディスクストレージ、磁気カセット、磁気テープ、磁気ディスクストレージ、または他の磁気ストレージデバイス、あるいは、所望の情報を格納するために使用することおよびコンピュータ 4 2 1 0 によってアクセスすることが可能な任意の他の媒体、が含まれるが、これらに限定されることはない。システムメモリ 4 2 3 0 は、R O M 4 2 3 1 および R A M 4 2 3 2 などの揮発性および / または不揮発性メモリの形のコンピュータストレージ媒体を含む。コンピュータ 4 2 1 0 内の諸要素間で (起動時などの) 情報の転送を助ける基本ルーチンを含む、基本入出力システム 4 2 3 3 (B I O S) は、通常、R O M 4 2 3 1 に格納される。R A M 4 2 3 2 は、通常、処理ユニット 4 2 2 0 によって即時アクセス可能および / または現在その上で動作中の、データおよび / またはプログラムモジュールを含む。例を挙げると、図 4 2 は、オペレーティングシステム 4 2 3 4、アプリケーションプログラム 4 2 3 5、他のプログラムモジュール 4 2 3 6、およびプログラムデータ 4 2 3 7 を示すが、これらに限定されることはない。

【 0 1 3 7 】

コンピュータ 4 2 1 0 は、他の取り外し可能 / 取り外し不能で揮発性 / 不揮発性のコンピュータストレージ媒体も含むことができる。単なる例として挙げると、図 4 2 は、取り外し不能で不揮発性の磁気媒体から読み取るかまたはこれに書き込むハードディスクドライブ 4 2 4 1 と、取り外し可能で不揮発性の磁気ディスク 4 2 5 2 から読み取るかまたはこれに書き込む磁気ディスクドライブ 4 2 5 1 と、C D - R O M または他の光媒体などの取り外し可能で不揮発性の光ディスク 4 2 5 6 から読み取るかまたはこれに書き込む光ディスクドライブ 4 2 5 5 と、を示す。例示的動作環境で使用可能な他の取り外し可能 / 取り外し不能で揮発性 / 不揮発性のコンピュータストレージには、磁気テープカセット、フラッシュメモリユニット、デジタル汎用ディスク、デジタルビデオテープ、ソリッドステート R A M、ソリッドステート R O M、などが含まれるが、これらに限定されることはない。ハードディスクドライブ 4 2 4 1 は、通常、インターフェース 4 2 4 0 などの取り外し不能メモリインターフェースを介して、システムバス 4 2 2 1 に接続される。磁気ディスクドライブ 4 2 5 1 および光ディスクドライブ 4 2 5 5 は、通常、インターフェース 4 2 5 0 などの取り外し可能メモリインターフェースによってシステムバス 4 2 2 1 に接続される。

【 0 1 3 8 】

前述し、図 4 2 に示された、ドライブおよびそれらに関連付けられたコンピュータストレージ媒体は、コンピュータ読み取り可能命令、データ構造、プログラムモジュール、およびコンピュータ 4 2 1 0 向けの他のデータの、ストレージを提供する。たとえばハードディスクドライブ 4 2 4 1 は、オペレーティングシステム 4 2 4 4、アプリケーションプログラム 4 2 4 5、他のプログラムモジュール 4 2 4 6、およびプログラムデータ 4 2 4 7 を格納するように示される。これらの構成要素は、オペレーティングシステム 4 2 3 4、アプリケーションプログラム 4 2 3 5、他のプログラムモジュール 4 2 3 6、およびプログラムデータ 4 2 3 7 と同じであるか、または異なるものとすることができることに留

10

20

30

40

50

意されたい。通常、RAMに格納されたオペレーティングシステム、アプリケーションプログラムなどは、対応するシステム、プログラム、またはハードディスクドライブ4241から読み取られたデータの一部であり、この部分は、所望の機能に応じてサイズおよび範囲が異なる。オペレーティングシステム4244、アプリケーションプログラム4245、他のプログラムモジュール4246、およびプログラムデータ4247には、それらが少なくとも異なるコピーであることが可能なことを示すために、ここでは異なる番号が与えられている。ユーザは、キーボード4262、一般にはマウス、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デバイス4261、無線入力受信構成要素4263、あるいは、リモートコントロールなどの無線ソースなどの、入力デバイスを介して、コマンドおよび情報をコンピュータ4210に入力することができる。他の入力デバイス（図示せず）は、マイクロフォン、ジョイスティック、ゲームパッド、衛星放送用パラボラアンテナ、スキャナなどを含むことができる。これらおよび他の入力デバイスは、しばしば、システムバス4221に結合されるユーザ入力インターフェース4260を介して、処理ユニット4220に接続されるが、パラレルポート、ゲームポート、IEEE4294ポート、またはユニバーサルシリアルバス（USB）4298、または赤外線（IR）バス4299などの、他のインターフェースおよびバス構造による接続が可能である。前述のように、入力/出力機能は、通信ネットワークを介した分散様式で容易にすることができる。

【0139】

ディスプレイデバイス4291も、ビデオインターフェース4290などのインターフェースを介してシステムバス4221に接続される。ディスプレイデバイス4291は、モニタ、LCDスクリーン、TFTスクリーン、フラットパネルディスプレイ、従来型テレビジョン、またはスクリーンプロジェクタに限定されることのない、コンピュータ4210の出力を表示するための任意のデバイスとすることができる。コンピュータは、ディスプレイデバイス4291に加えて、出力周辺インターフェース4295を介して接続可能な、スピーカ4297およびプリンタ4296などの他の周辺出力デバイスを含むこともできる。

【0140】

コンピュータ4210は、リモートコンピュータ4280などの1つ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する、ネットワーク化環境で動作することになる。リモートコンピュータ4280は、パーソナルコンピュータとすることが可能であり、通常、コンピュータ4210に関して上記で説明した要素のうちの多くまたはすべてを含むが、図42ではメモリストレージデバイス4281のみが示されている。図42に示された論理接続は、ローカルエリアネットワーク（LAN）4271およびワイドエリアネットワーク（WAN）4273を含むが、メトロポリタンエリアネットワーク（MAN）、イントラネット、またはインターネットへの接続などの、他のネットワークも含むことができる。

【0141】

LANネットワーク環境で使用される場合、コンピュータ4210は、ネットワークインターフェースまたはアダプタ4270を介してLAN4271に接続される。WANネットワーク環境で使用される場合、コンピュータ4210は、通常、モデム4272、またはインターネットなどのWAN4273を介した通信を確立するための他の手段を含む。モデム4272は内蔵型または外付け型とすることが可能であり、ネットワークインターフェース4270または他の適切なメカニズムを介してシステムバス4221に接続される。モデム4272は、ケーブルモデム、DSLモデム、または他のブロードバンドデバイスとすることができる。ネットワーク化環境では、コンピュータ4210に関して示されたプログラムモジュール、またはその一部を、リモートメモリストレージデバイスに格納することができる。たとえば、図42はリモートアプリケーションプログラム4285をメモリデバイス4281上に常駐するものとして示しているが、これに限定されることはない。図示されたネットワーク接続は例示的なものであり、コンピュー

10

20

30

40

50

タ間に通信リンクを確立する他の手段も使用可能である。

【0142】

コンピュータ4210の多くの他の内部構成要素は示されていないが、当業者であれば、こうした構成要素および相互接続は周知であることを理解されよう。たとえば、テレビジョンチューナカードおよびネットワークインターフェースカードなどの様々な拡張カードをコンピュータ4210内に含めることが一般的である。したがって、デルタページャの実施プロセスの例示的諸実施形態を説明する上で、コンピュータ4210の内部構造に関する追加の細部を開示する必要はない。

【0143】

コンピュータ4210が電源投入またはリセットされた場合、ROM 4231に格納されたBIOS 4233は、オペレーティングシステムまたはその必要部分を、ハードディスクドライブ4241からRAM 4232へとロードするように、処理ユニット4220に命じる。オペレーティングシステム4244として指定されたオペレーティングシステムのコピーされた部分がRAM 4232にロードされると、処理ユニット4220はオペレーティングシステムコードを実行し、オペレーティングシステム4234のユーザインターフェースに関連付けられた視覚要素をディスプレイデバイス4291上に表示させる。通常、アプリケーションプログラム4245がユーザによって開かれた場合、プログラムコードおよび関連データがハードディスクドライブ4241から読み取られ、必要な部分がRAM 4232にコピーされるが、このコピーされた部分は、本明細書では参照番号4235によって表される。

【0144】

結論

以上、構造的特徴および/または方法論的動作に特有の言葉で例示的諸実施形態について説明してきたが、添付の特許請求の範囲は、必ずしも前述の特定の特徴または動作に限定されるものでないことを理解されよう。むしろ、この特定の機能および動作は、例示的諸実施形態として開示されている。

【図面の簡単な説明】

【0145】

【図1】データベースの初期の、現行状態を示す図である。

【図2】データベースを対象とする第1のトランザクションを示す図である。

【図3】バッファを使用してデータベースに対して実行およびコミットされているトランザクションを示す図である。

【図4】バッファを使用してデータベースに対して実行およびコミットされているトランザクションを示す図である。

【図5】バッファを使用してデータベースに対して実行およびコミットされているトランザクションを示す図である。

【図6】バッファを使用してデータベースに対して実行およびコミットされているトランザクションを示す図である。

【図7】累積書き込みバッファを使用して実行されているトランザクションを示す図である。

【図8】累積書き込みバッファを使用して実行されているトランザクションを示す図である。

【図9】累積書き込みバッファを使用して実行されているトランザクションを示す図である。

【図10】累積書き込みバッファを使用して実行されているトランザクションを示す図である。

【図11】累積書き込みバッファを使用して実行されているトランザクションを示す図である。

【図12】書き込みバッファを使用するトランザクションの実行モードを示す図である。

【図13】単一書き込みバッファを使用して実行されているトランザクションを示す図である。

10

20

30

40

50

ある。

【図 1 4】単一書き込みバッファを使用するトランザクションの実行モードを示す図である。

【図 1 5】第 1 のトランザクションが介入トランザクションとして実行される間に、データベースに対して実行されている第 2 のトランザクションを示す図である。

【図 1 6】介入トランザクションを受けての、現行トランザクションのコミットメントおよび打ち切りを示す図である。

【図 1 7】介入トランザクションを受けての、現行トランザクションのコミットメントおよび打ち切りを示す図である。

【図 1 8】介入トランザクションを受けての、現行トランザクションのコミットメントおよび打ち切りを示す図である。

10

【図 1 9】介入トランザクションを受けての、現行トランザクションのコミットメントおよび打ち切りを示す図である。

【図 2 0】介入トランザクションを受けての、現行トランザクションのコミットメントおよび打ち切りを示す図である。

【図 2 1】介入トランザクションを受けて、トランザクションをコミットするかまたは打ち切るタイミングを決定するモードを示す図である。

【図 2 2】読み取りおよび書き込みのバッファを使用してデータベースに対して実行およびコミットされている、競合する可能性のあるトランザクションを示す図である。

【図 2 3】読み取りおよび書き込みのバッファを使用してデータベースに対して実行およびコミットされている、競合する可能性のあるトランザクションを示す図である。

20

【図 2 4】読み取りおよび書き込みのバッファを使用してデータベースに対して実行およびコミットされている、競合する可能性のあるトランザクションを示す図である。

【図 2 5】読み取りおよび書き込みのバッファを使用してデータベースに対して実行およびコミットされている、競合する可能性のあるトランザクションを示す図である。

【図 2 6】読み取りおよび書き込みのバッファを使用してデータベースに対して実行およびコミットされている、競合する可能性のあるトランザクションを示す図である。

【図 2 7】読み取りおよび書き込みのバッファを使用してデータベースに対して実行およびコミットされている、競合する可能性のあるトランザクションを示す図である。

【図 2 8】読み取りおよび書き込みのバッファを使用するトランザクションの実行モードを示す図である。

30

【図 2 9】読み取りおよび書き込みのバッファを使用する介入トランザクションを受けて、トランザクションをコミットするかまたは打ち切るタイミングを決定するモードを示す図である。

【図 3 0】バッファの合体を示す図である。

【図 3 1】バッファの合体を示す図である。

【図 3 2】データベースの選択された状態が後続の状態に合体されないように保持するための、スナップショットオブジェクトの使用を示す図である。

【図 3 3】データベースの選択された状態が後続の状態に合体されないように保持するための、スナップショットオブジェクトの使用を示す図である。

40

【図 3 4】データベースの選択された状態が後続の状態に合体されないように保持するための、スナップショットオブジェクトの使用を示す図である。

【図 3 5】1 つまたは複数の選択された状態を保持しながらバッファを合体するモードを示す図である。

【図 3 6】永続的にされているバッファを示す図である。

【図 3 7】永続的にされているバッファを示す図である。

【図 3 8】オブジェクトを永続的にするモードを示す図である。

【図 3 9】キャッシュオブジェクトの使用を示す図である。

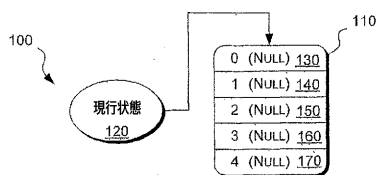
【図 4 0】キャッシュオブジェクトを作成および維持するモードを示す図である。

【図 4 1】前述のオブジェクトの集まりを示す図である。

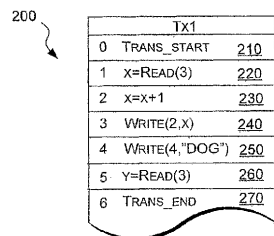
50

【図４２】前述のデータベースのトランザクションおよび動作を実行するために好適な、例示的動作環境を示す図である。

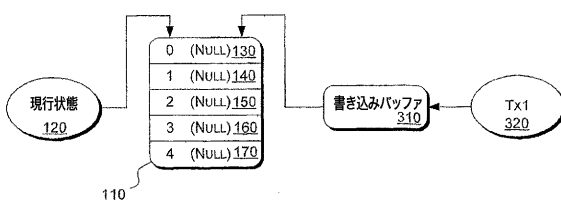
【図１】



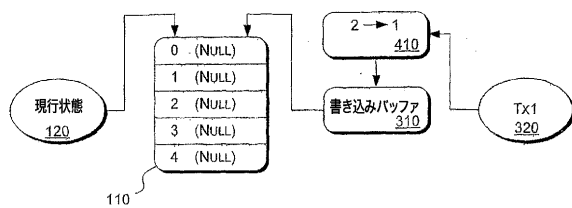
【図２】



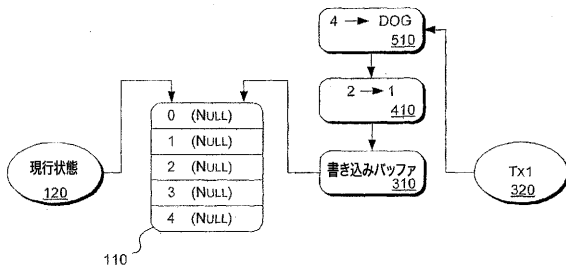
【図３】



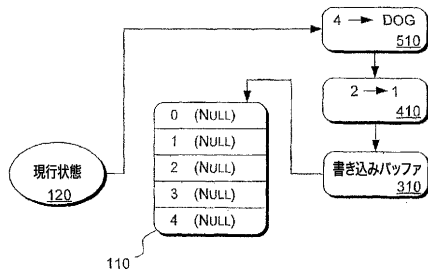
【図４】



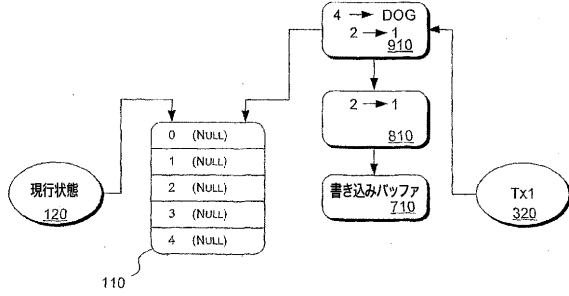
【図５】



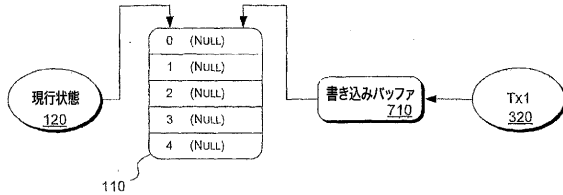
【図 6】



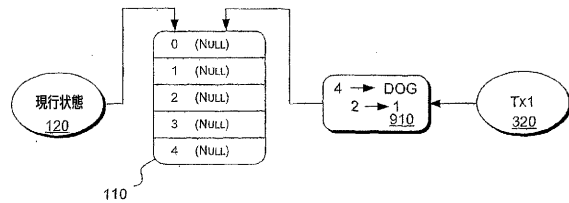
【図 9】



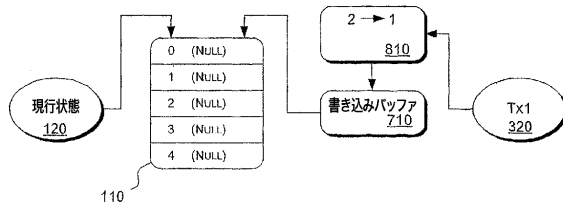
【図 7】



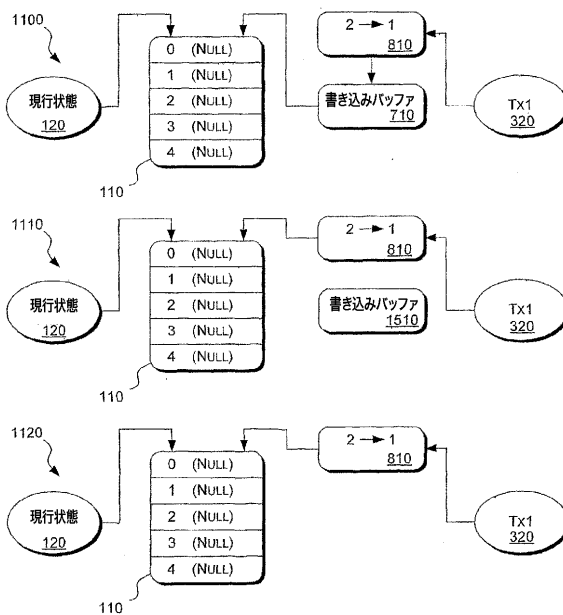
【図 10】



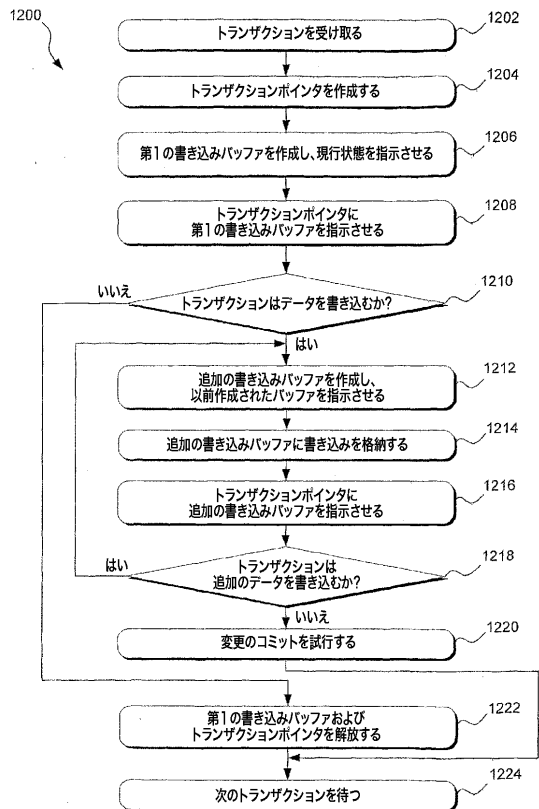
【図 8】



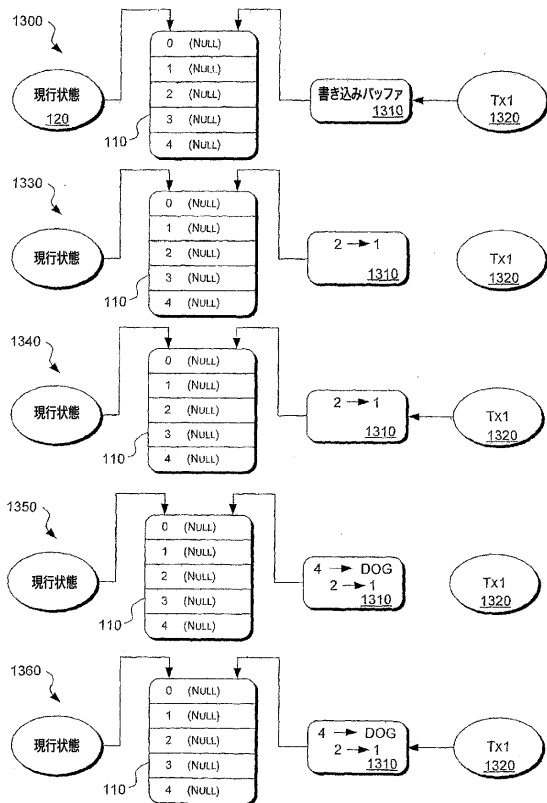
【図 11】



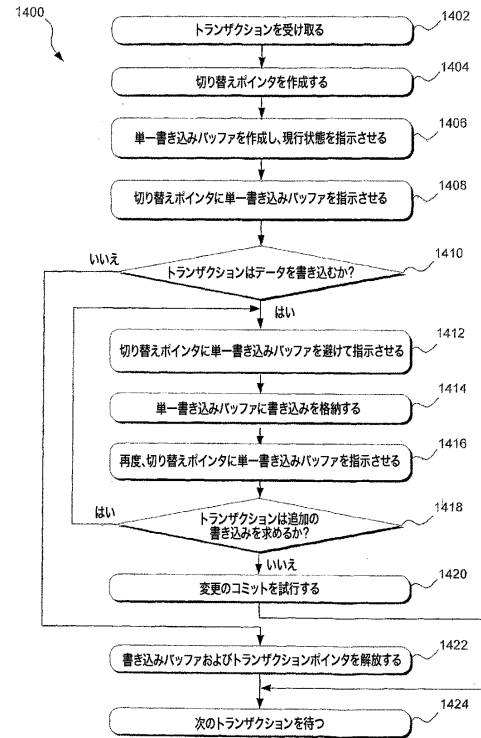
【図 12】



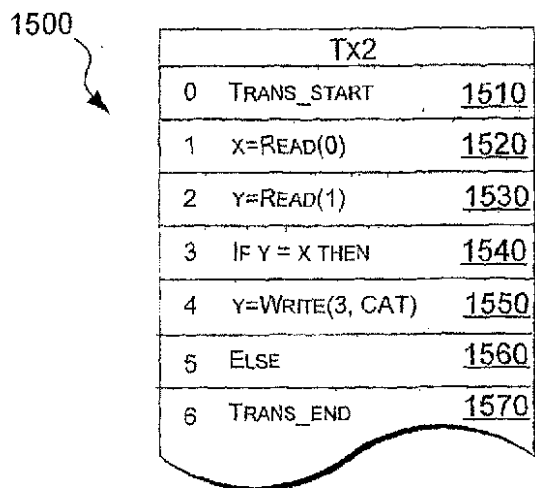
【図 13】



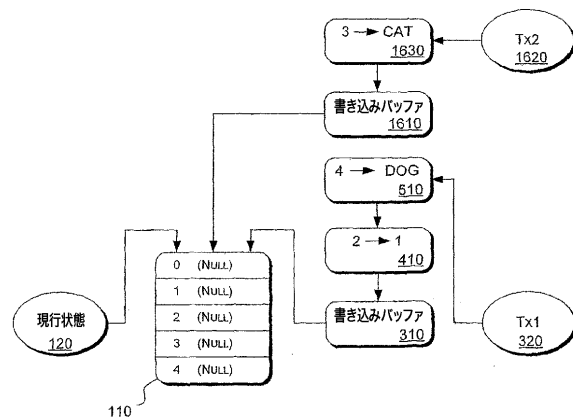
【図 14】



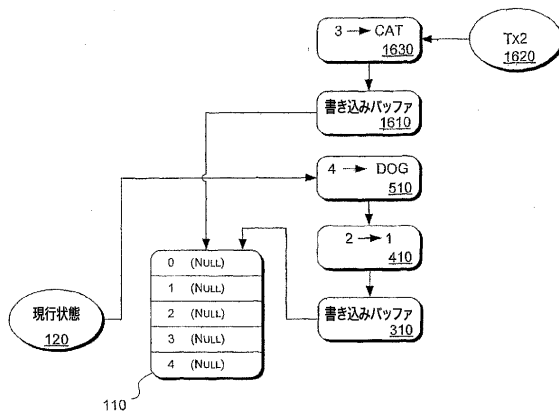
【図 15】



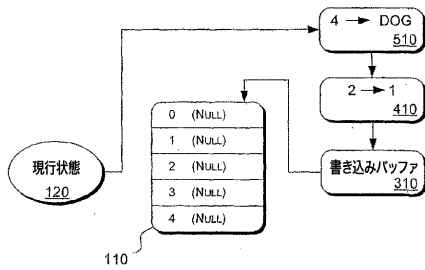
【図 16】



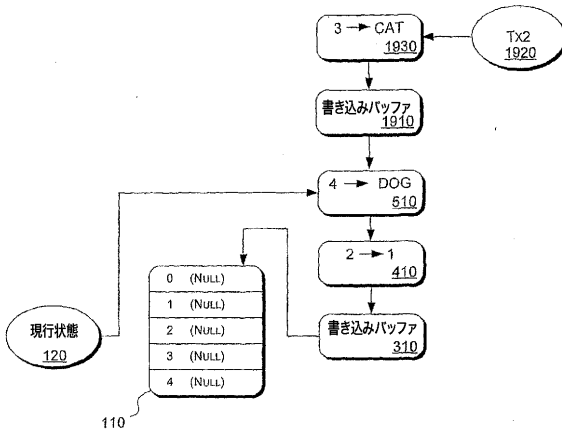
【図 17】



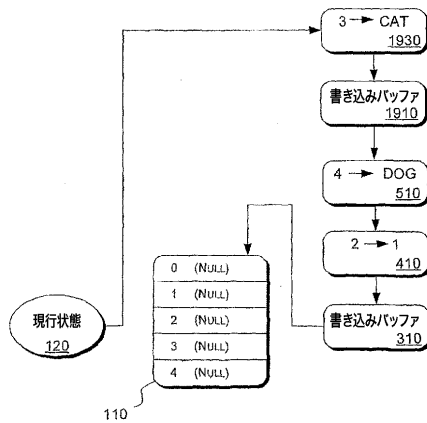
【図 18】



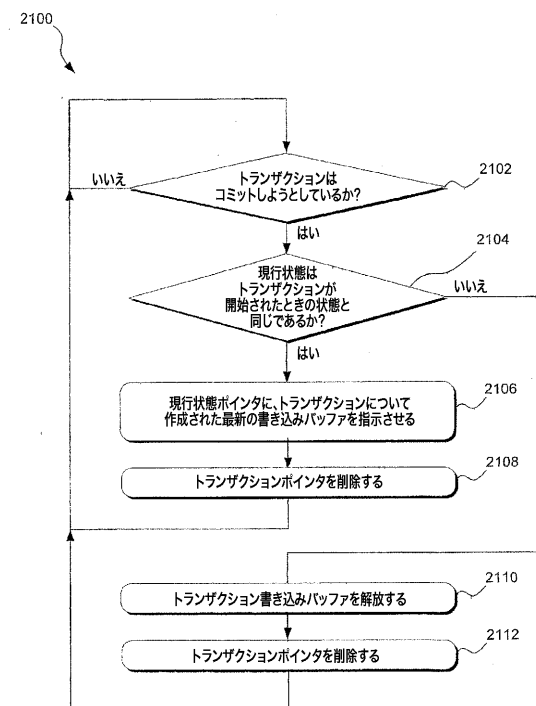
【図 19】



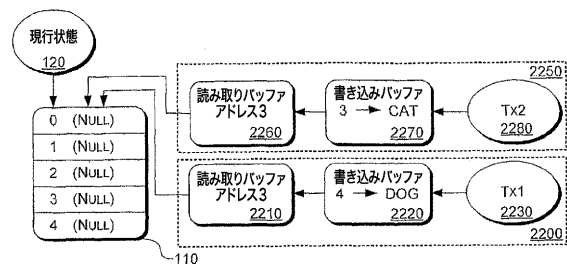
【図 20】



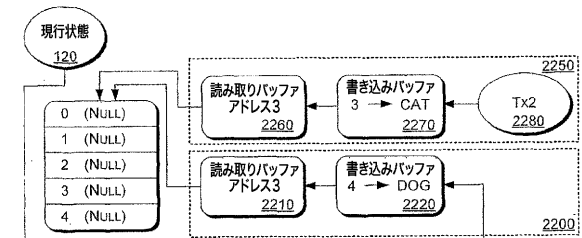
【図 21】



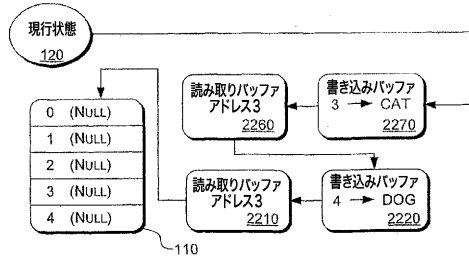
【図 22】



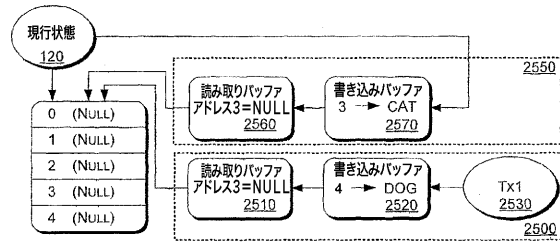
【図 23】



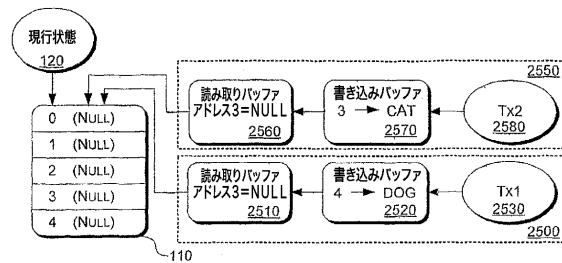
【図 24】



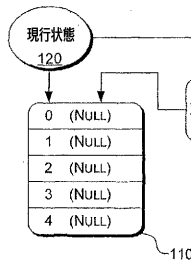
【図 26】



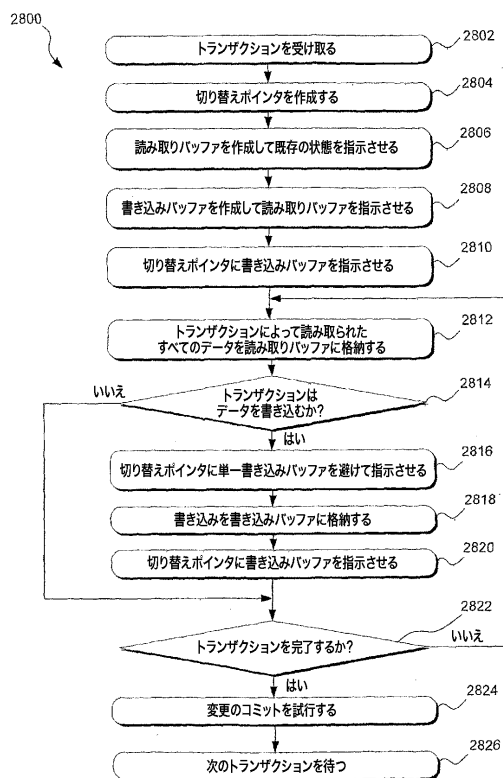
【図 25】



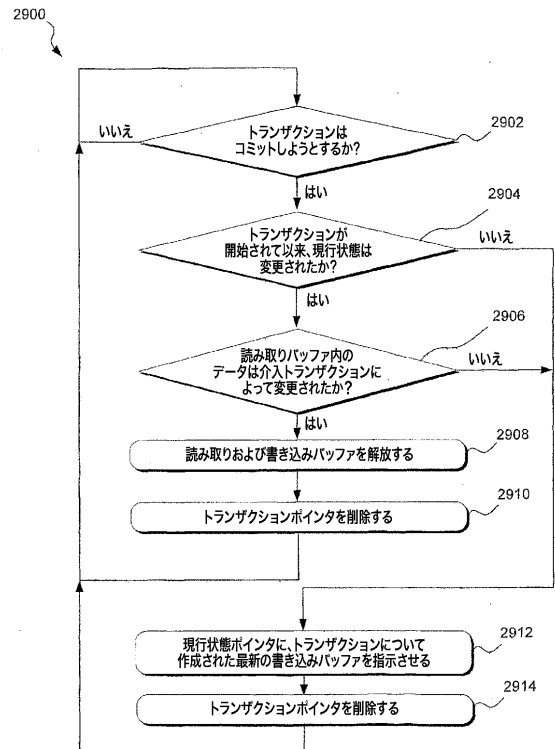
【図 27】



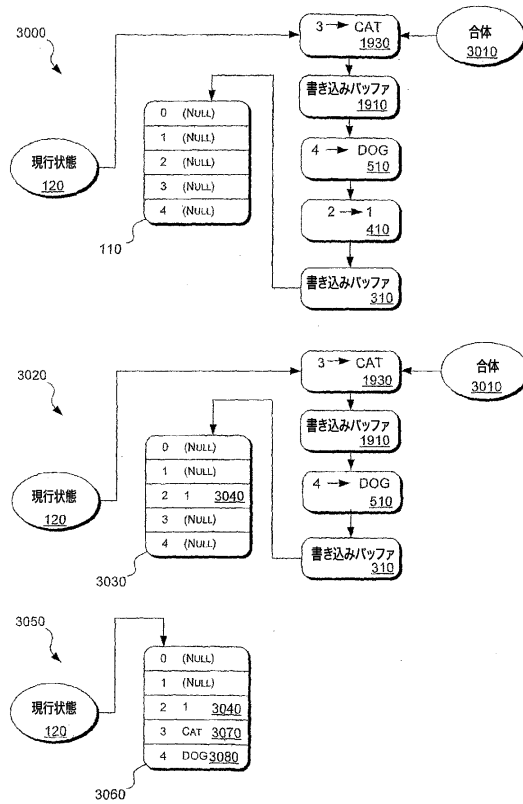
【図 28】



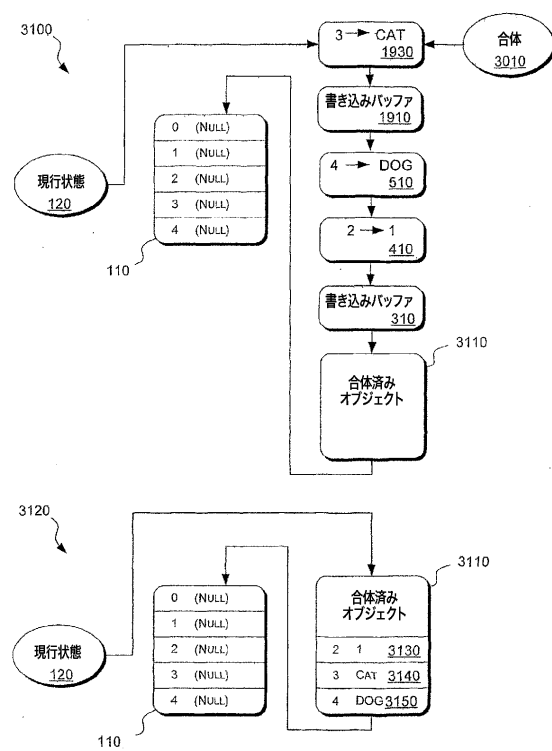
【図 29】



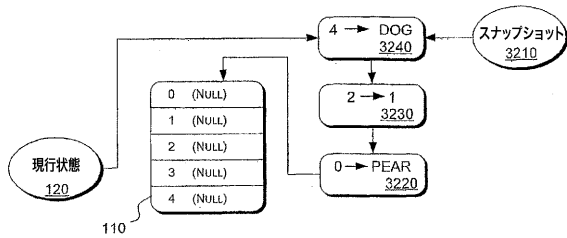
【図 30】



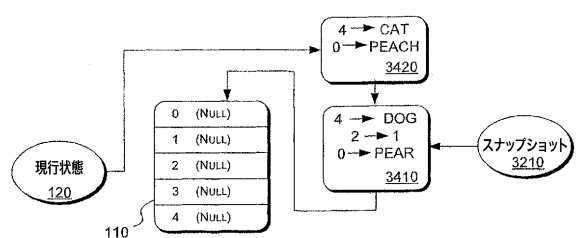
【図 31】



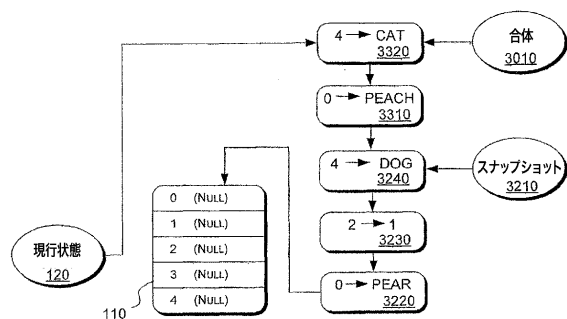
【図 32】



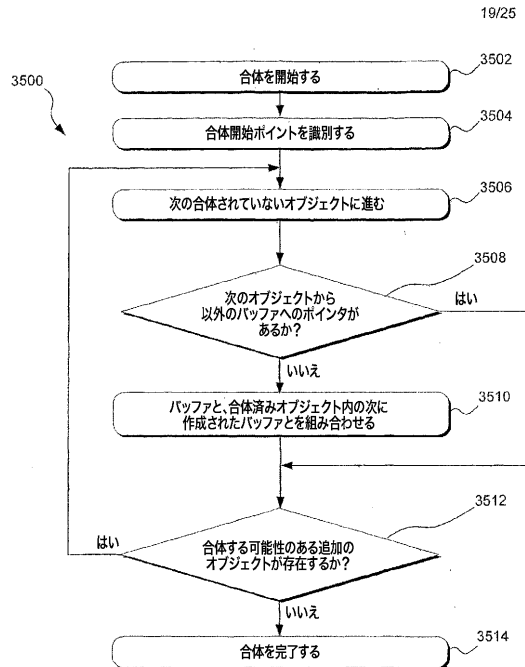
【図 34】



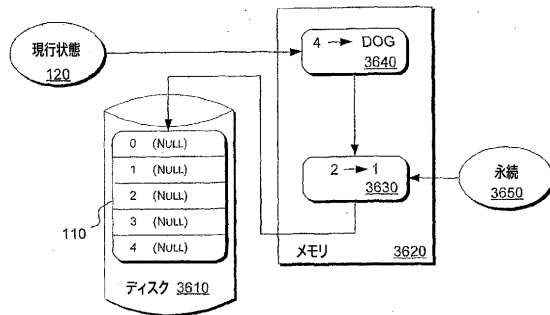
【図 33】



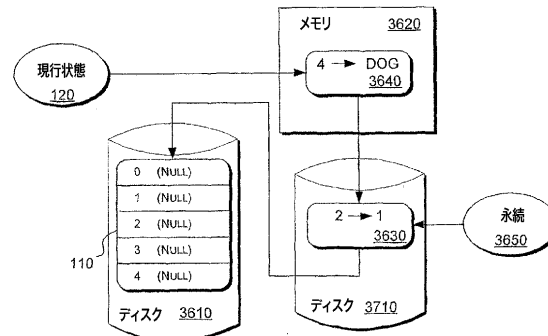
【図 35】



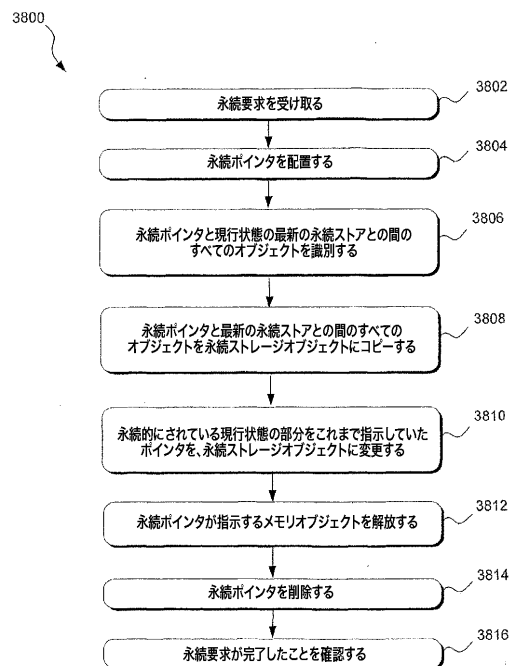
【図 36】



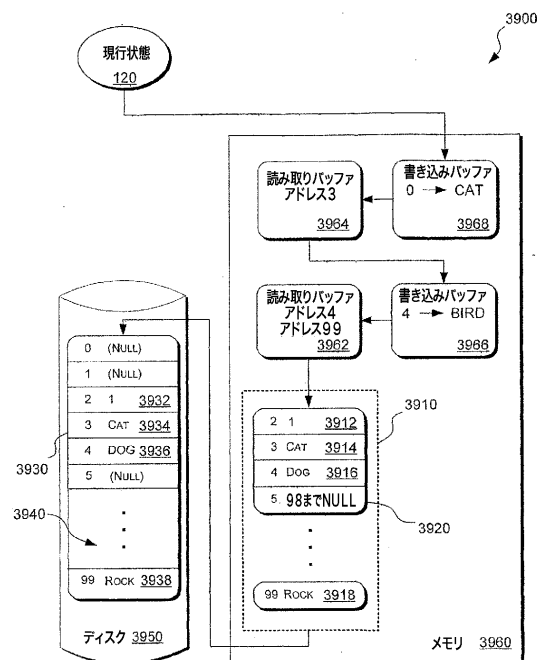
【図 37】



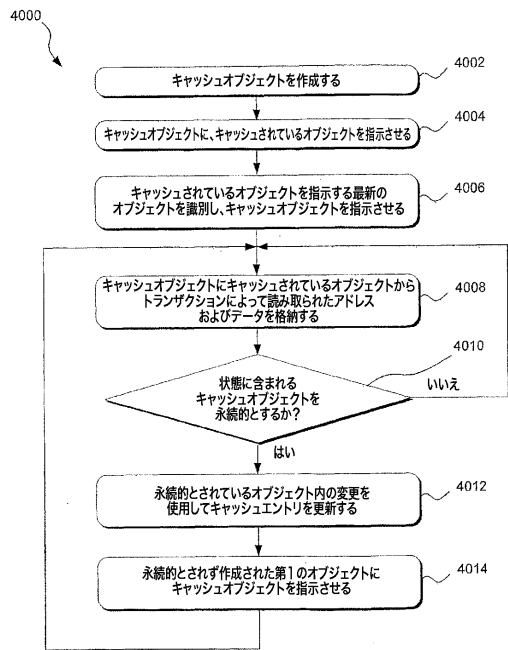
【図 38】



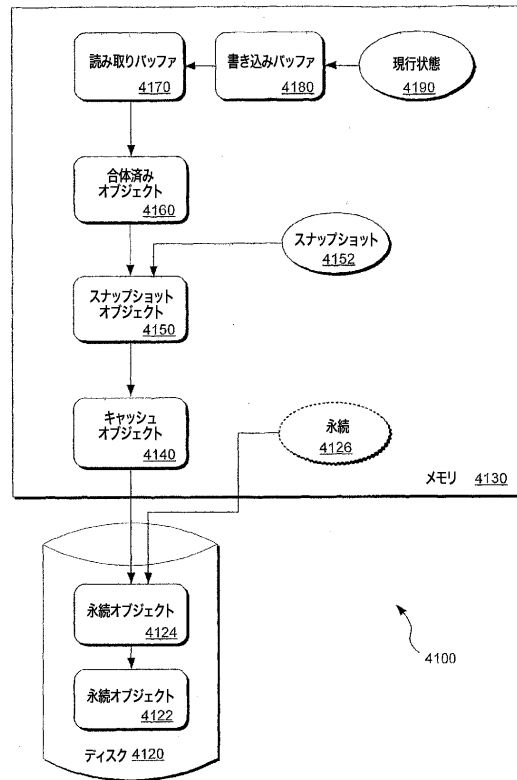
【図 39】



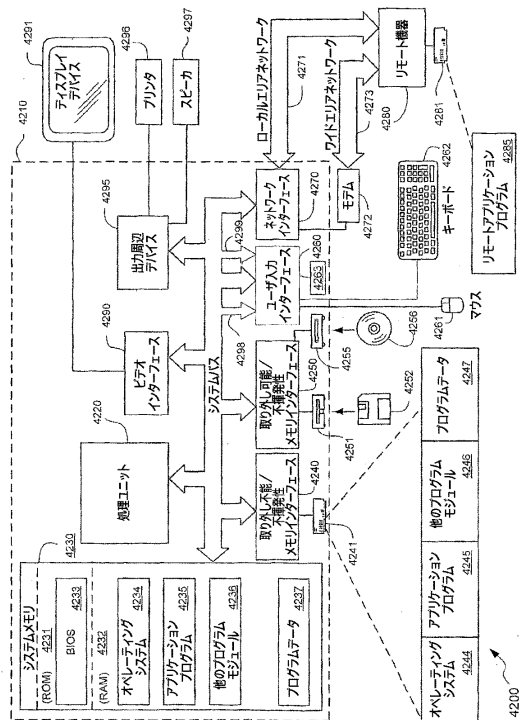
【図40】



【図41】



【図42】



フロントページの続き

(72)発明者 ジョン アール・ドゥサー

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ
マイクロソフト コーポレーション インターナショナル パテント内

審査官 桜井 茂行

(56)参考文献 特開平09-062554(JP,A)

国際公開第93/003436(WO,A1)

米国特許第6594751(US,B1)

国際公開第96/09730(WO,A1)

米国特許第6199141(US,B1)

米国特許第05983225(US,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00

G06F 17/30

G06F 17/40

JSTPlus(JDreamII)