

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6522812号
(P6522812)

(45) 発行日 令和1年5月29日(2019.5.29)

(24) 登録日 令和1年5月10日(2019.5.10)

(51) Int.Cl.	F I
G06F 12/00 (2006.01)	G06F 12/00 531R
	G06F 12/00 545A
	G06F 12/00 531J

請求項の数 15 外国語出願 (全 51 頁)

(21) 出願番号	特願2018-1662 (P2018-1662)	(73) 特許権者	507303550
(22) 出願日	平成30年1月10日 (2018.1.10)		アマゾン・テクノロジーズ・インコーポレ ーテッド
(62) 分割の表示	特願2016-501614 (P2016-501614) の分割		アメリカ合衆国・98108-1226・ ワシントン州・シアトル・パイオーボク ス・81226
原出願日	平成26年3月12日 (2014.3.12)	(74) 代理人	100098394
(65) 公開番号	特開2018-77895 (P2018-77895A)		弁理士 山川 茂樹
(43) 公開日	平成30年5月17日 (2018.5.17)	(74) 代理人	100064621
審査請求日	平成30年1月10日 (2018.1.10)		弁理士 山川 政樹
(31) 優先権主張番号	61/799,609	(72) 発明者	グプタ, アヌラグ・ウィンドラス
(32) 優先日	平成25年3月15日 (2013.3.15)		アメリカ合衆国・98109-5210・ ワシントン州・シアトル・テリー アヴェ ニュー ノース・410
(33) 優先権主張国	米国 (US)		
(31) 優先権主張番号	14/201,505		
(32) 優先日	平成26年3月7日 (2014.3.7)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 分散型データベースシステム用高速クラッシュ回復

(57) 【特許請求の範囲】

【請求項1】

分散型ストレージシステムを実装する複数のストレージノードであって、前記分散型ストレージシステムは、データベースのためにログ構造化データストレージを実装するように構成され、前記複数のストレージノードは、ネットワークを通じてデータベースシステムに接続し、前記ログ構造化データストレージの複数のリドゥログレコードが、前記複数のストレージノードにおいて、前記ネットワークを通じて前記データベースシステムから既に受信されており、前記複数のリドゥログレコードのそれぞれは、前記複数のストレージノードの中で前記データベースのために記憶されたデータに対する変更を記述する、複数のストレージノードと、

前記データベースシステムを実装するデータベースヘッドノードであって、前記データベースヘッドノードは、

前記ネットワークを通じて、前記複数のストレージノードとの接続を確立することと、

前記リドゥログレコードをリプレイすることなく、一つ以上のアクセス要求のための前記データベースへのアクセスを提供することと、

を含む、データベースヘッドノード故障からの回復を実行する、データベースヘッドノードと、

を備えるシステム。

【請求項2】

前記データベースヘッドノードはさらに、
前記データベースに対するアクセス要求を受信するように、
前記受信したアクセス要求に基づいて、前記ストレージノードに記憶されたデータページの現在の状態に対する要求を、前記ネットワークを通じて、前記複数のストレージノードのうちの一つに送信するように、及び
前記ストレージノードにおいて、その現在の状態で前記データページを生成するために、前記複数のリドゥログレコードのうちの一つ以上が、前記データページの既に保存された状態に適用された前記要求されたデータページを、その現在の状態で前記ストレージノードから受信するように構成される、請求項 1 に記載のシステム。

【請求項 3】

前記複数のストレージノードに送信された前記複数のリドゥログレコードのうち少なくともいくつかはシステムトランザクションを備え、前記複数のストレージノードのうちの一つのストレージノードは、
前記システムトランザクションが不完全であると決定するように、及び
前記少なくともいくつかのリドゥログレコードによって変更されたデータページの現在の状態を生成するときに適用されないとして前記複数のリドゥログレコードのうちの前記少なくともいくつかを識別するように、
構成される、請求項 1 に記載のシステム。

【請求項 4】

前記複数のストレージノードに送信された前記複数のリドゥログレコードのうち少なくともいくつかはシステムトランザクションを備え、前記データベースヘッドノードはさらに、
前記システムトランザクションが不完全であると決定するように、
前記少なくともいくつかのリドゥログレコードによって変更されたデータページの現在の状態を生成するときに適用されないとして前記複数のリドゥログレコードのうち少なくともいくつかを識別するように、及び
前記複数のストレージノードのうちの一つ以上に、適用されない、前記複数のリドゥログレコードのうちの前記識別された少なくともいくつかを示す通知を送信するように、
構成される、請求項 1 に記載のシステム。

【請求項 5】

データベースヘッドノードを実装する一つ以上のコンピューティング装置によって、
データベースヘッドノード故障からの回復時に、
データベースのためのデータを記憶する分散型ストレージシステムを実装する複数のストレージノードの一つ以上のストレージノードとの接続を確立することであって、前記複数のストレージノードは、前記データベースヘッドノードを実装する前記一つ以上のコンピューティング装置以外の複数のコンピューティング装置に実装され、前記分散型ストレージシステムは、前記データベースのためにログ構造化データストレージを実装するように構成され、前記ログ構造化データストレージの複数のリドゥログレコードが、前記複数のストレージノードにおいて既に受信されており、前記複数のリドゥログレコードのそれぞれは、前記データベースのために記憶されたデータに対する変更を、それが受信された前記それぞれのストレージノードで記述する、接続を確立することと、
前記リドゥログレコードのリプレイなしで、一つ以上のアクセス要求のための前記データベースへのアクセスを提供することと、
を実行することを含む方法。

【請求項 6】

前記データベースに対するアクセス要求を受信することと、
前記アクセス要求を受信することに応じて、前記一つ以上のストレージノードから、前記データベースのための前記データの部分を記憶する一つ以上のデータページの現在の状態を要求することと、
前記アクセス要求にサービスを提供するために、前記データページのための前記データ

10

20

30

40

50

の部分記憶する前記一つ以上のストレージノードから前記一つ以上のデータページの現在の状態を受信することと、
をさらに含む、請求項 5 に記載の方法。

【請求項 7】

受信された前記データベースのための前記データの部分を記憶する前記一つ以上のデータページのうちの少なくとも一つの前記現在の状態は、前記一つ以上のストレージノードのうちの一つが、前記少なくとも一つの前記データページの既に記憶されたバージョンまで前記複数のリドゥログレコードのうちの一つ以上をリプレイすることによって生成される、請求項 6 に記載の方法。

【請求項 8】

受信された前記データベースのための前記データの部分を記憶する前記一つ以上のデータページのうちの前記少なくとも一つからの異なるデータページの前記現在の状態は、前記データページの既に記憶されたバージョンまで前記複数のリドゥログレコードのうちの一つ以上をリプレイすることなく、前記一つ以上のストレージノードのうちの一つによって送信される、請求項 7 に記載の方法。

【請求項 9】

前記データベースヘッドノードは、前記複数のストレージノードにリドゥログレコードとして送信された変更をアンドゥするために複数のアンドゥログレコードを維持し、前記方法は、

前記一つ以上のストレージノードから受信した前記一つ以上のデータページのうちのの一つが、不完全なユーザトランザクションによって影響を及ぼされると決定することと、前記ユーザトランザクションは、前記一つのデータページを含む前記一つ以上のストレージノードに記憶された前記データに対して変更を指示することと、

前記ユーザトランザクションによって前記データページに対して指示された変更をアンドゥするために、前記データページに一つ以上のアンドゥログレコードを適用することと、
をさらに含む、請求項 6 に記載の方法。

【請求項 10】

前記データベースヘッドノードは、前記ユーザトランザクションを含む複数の不完全なユーザトランザクションを示すトランザクションテーブルを維持し、前記方法は、

前記トランザクションテーブルに少なくとも部分的に基づいて、前記複数の不完全なユーザトランザクションのうちの一つによって影響を及ぼされた一つ以上の追加のデータページを決定することと、

前記一つ以上のストレージノードから、一つ以上の追加のデータページの現在の状態を要求することと、

前記一つ以上の追加のデータページを受信することに応じて、追加の一つ以上のアンドゥログレコードを前記一つ以上の追加のデータページに適用して、前記少なくとも一つ不完全なユーザトランザクションによって前記一つ以上の追加のデータページに対して指示された変更をアンドゥすることと、

をさらに含む、請求項 9 に記載の方法。

【請求項 11】

前記一つ以上の追加のデータページを決定することと、前記一つ以上の追加のデータページを要求することと、及び前記一つ以上の追加のデータページに前記追加の一つ以上のアンドゥログレコードを適用することは、前記データベースヘッドノードでバックグラウンドプロセスの一部として実行され、また、前記アクセス要求を受信することと、前記一つ以上のデータページの前記現在の状態を要求することと、及び前記一つ以上のデータページの現在の状態を受信することは、フォアグラウンドプロセスの一部として実行される、請求項 10 に記載の方法。

【請求項 12】

前記データベースヘッドノード故障からの回復時に、

10

20

30

40

50

アクセスのために前記データベースを利用可能にする前に、前記データベースのための前記複数のストレージノードに記憶された前記データの、既に記録されたスナップショットに対応する状態への復元に対する要求を前記複数のストレージノードに送信することであって、前記復元が、前記複数のリドゥログのうちの一つ以上を、前記データの以前のバージョンに適用することを含むことと、
をさらに含む、請求項 5 に記載の方法。

【請求項 1 3】

一つ以上のコンピューティング装置による実行時に、
データベースヘッドノード故障からの回復時に、
データベースのためのデータを記憶する分散型ストレージシステムを実装する複数のストレージノードの一つ以上のストレージノードとの接続を確立することであって、前記複数のストレージノードは、前記データベースヘッドノードを実装する前記一つ以上のコンピューティング装置以外の複数のコンピューティング装置に実装され、前記分散型ストレージシステムは、前記データベースのためにログ構造化データストレージを実装するように構成され、前記ログ構造化データストレージの複数のリドゥログレコードが、前記複数のストレージノードにおいて既に受信されており、前記複数のリドゥログレコードのそれぞれは、前記データベースのために記憶されたデータに対する変更を、それが受信された前記それぞれのストレージノードで記述する、接続を確立することと、
前記リドゥログレコードのリプレイなしで、一つ以上のアクセス要求のための前記データベースへのアクセスを提供することと、
を実施するデータベースシステムのデータベースヘッドノードを実装するプログラム命令を記憶する、非一時的なコンピュータ可読記憶媒体。

【請求項 1 4】

前記複数のストレージノードで既に受信された前記複数のリドゥログレコードは、前記データベースヘッドノードとは異なるデータベースヘッドノードから受信された、請求項 1 3 に記載の非一時的なコンピュータ可読記憶媒体。

【請求項 1 5】

前記データベースヘッドノードは、
前記データベースのためのアクセス要求を受信することと、
前記アクセス要求を受信することに応じて、前記一つ以上のストレージノードから前記データベースのための前記データの部分を記憶する一つ以上のデータページの現在の状態を要求することと、
前記アクセス要求にサービスを提供するために前記データベースのための前記データの部分を記憶する前記一つ以上のデータページの現在の状態を受信することであって、前記一つ以上の受信したデータページのうちの前記少なくとも一つの前記現在の状態は、前記一つ以上のストレージノードのうちの一つが、前記少なくとも一つのデータページの既に記憶されたバージョンまで前記複数のリドゥログレコードのうちの一つ以上をリプレイすることによって生成されることと、
をさらに実施する、請求項 1 3 に記載の非一時的なコンピュータ可読記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

ソフトウェアスタックの多様な構成要素の分散は、いくつかの場合、（例えば複製によって）フォルトトレランス、より高い耐久性、及び（例えば、より少ない大型の高価な構成要素よりむしろ、多くのより小型でより安価な構成要素を使用することにより）より安価な解決策を提供する（または支援する）ことができる。ただし、データベースは、従来、分散の影響を最も受けにくいソフトウェアスタックの構成要素の中にある。例えば、データベースが提供すると期待されているいわゆる ACID プロパティ（例えば、原子性、一貫性、独立性、及び永続性）を保証しつつもデータベースを分散することは困難であることがある。

10

20

30

40

50

【背景技術】

【0002】

大部分の既存のリレーショナルデータベースは分散化されていないが、いくつかの既存のデータベースは、2つの共通モデル、つまり「シェアードナッシング」モデル及び「シェアードディスク」モデルの内の1つを使用して（より大型のモノリシックシステムを単に利用することによって「スケールアップ」されることと対照的に）「スケールアウト」される。一般的に、「シェアードナッシング」モデルでは、受信されたクエリーは（それぞれがクエリーの構成要素を含む）データベースシャードに分解され、これらのシャードはクエリー処理のために異なる計算ノードに送られ、結果は、結果が返される前に収集され、統合される。一般的に「シェアードディスク」モデルでは、クラスタのあらゆる計算ノードは同じ基礎的データにアクセスできる。このモデルを利用するシステムでは、キャッシュコヒーレンシーを管理するために細心の注意を払う必要がある。これらのモデルの両方において、大型のモノリシックデータベースは（スタンドアロンデータベースインスタンスの機能性のすべてを含んだ）複数のノードで複製され、それらを縫い合わせるために「グルー」ロジックが追加される。例えば、「シェアードナッシング」モデルでは、グルーロジックは、クエリーを再分割し、クエリーを複数の計算ノードに送信し、次いで結果を結合するディスパッチャーの機能性を提供してよい。「シェアードディスク」モデルでは、グルーロジックが（例えば、キャッシング層でコヒーレンシーを管理するために）複数のノードのキャッシュをともに融合させるのに役立つ。これらの「シェアードナッシング」データベースシステム及び「シェアードディスク」データベースシステムは

10

20

【0003】

実施形態は、いくつかの実施形態及び例示的な図面について一例として本明細書に説明されているが、当業者は実施形態が説明されている実施形態または図面に制限されないことを認識する。図面及び図面に対する詳細な説明は、開示されている特定の形式に実施形態を制限することを目的とするのではなく、逆に、添付の特許請求の範囲によって定められる精神及び範囲に入るすべての修正形態、同等物、及び変更形態を対象とすることを目的とすることが理解されるべきである。本明細書に使用される見出しは編成のためだけであり、明細書または特許請求項の範囲を制限するために使用されることを意図していない。本願を通して使用されるように、単語「してよい」は、強制の意味（つまり、しなければならないを意味する）よりむしろ、許可の意味（つまり、する可能性を有することを意味する）で使用される。単語「含む」、「含んだ」、及び「含む」はオープンエンド関係を示し、したがって含むが、これに限定されるものではないことを意味する。同様に、単語「有する」、「有している」、及び「有する」もオープンエンド関係を示し、したがって有するが、これに限定されるものではないことを意味する。本明細書で使用される用語「第1の」、「第2の」、「第3の」等は、それらが前に来る名詞に対するラベルとして使用され、いかなるタイプの順序付け（例えば、空間的、時間的、論理的等）も、係る順序付けがはっきりと特記されない限り暗示しない。

30

【0004】

多様な構成要素は、1つまたは複数のタスクを実行する「ように構成される」として記述されてよい。係る文脈では、「ように構成される」は、動作中に1つまたは複数のタスクを実行する「構造を有する」を概して意味する大まかな記述である。したがって、構成要素は、構成要素が現在そのタスクを実行していなくてもタスクを実行するように構成できる（例えば、コンピュータシステムは、動作が現在実行されていなくても動作を実行するように構成されてよい）。いくつかの文脈では、「ように構成される」は、動作中に1つまたは複数のタスクを実行する「回路網を有する」を概して意味する構造の大まかな記述であってよい。したがって、構成要素は、構成要素が現在オンでなくてもタスクを実行するように構成できる。一般的に、「ように構成される」に対応する構造を形成する回路網はハードウェア回路を含んでよい。

40

50

【0005】

多様な構成要素は、説明での便宜上、1つまたは複数のタスクを実行すると記述されてよい。係る説明は、言い回し「ように構成される」を含んでいるとして解釈されるべきである。1つまたは複数のタスクを実行するように構成される構成要素を記述することは、その構成要素に対して特許法第112条、第6項の解釈を行使することを明白に目的としていない。

【0006】

「に基づいて」。本明細書に使用されるように、この用語は、決定に影響を及ぼす1つまたは複数の要因を説明するために使用される。この用語は、決定に影響を及ぼすことがある追加の要因を除外しない。すなわち、決定は、それらの要因だけに基づいてよい、または少なくとも部分的にそれらの要因に基づいてよい。言い回し「Bに基づいてAを決定する」を考える。BがAの決定に影響を及ぼす要因であることがある一方、係る言い回しは、Aの決定がCにも基づいていることを除外しない。他の例では、AはBだけに基づいて決定されてよい。

10

【0007】

本開示の範囲は、本明細書に（明示的または暗示的のどちらかで）開示される任意の特徴または特徴の組合せまたはその任意の一般論を、それが本明細書で扱われる課題のいずれかまたはすべてを軽減するか否かに関わらず含む。したがって、特徴の係る任意の組合せに対して、本願（または本願に対する優先権を主張する出願）の手続き処理中に新しい特許請求の範囲が策定されることがある。特に、添付特許請求の範囲に関して、従属請求項からの特徴は独立請求項の特徴と組み合されてよく、それぞれの独立請求項からの特徴は任意の適切な方法で、及び単に添付の特許請求の範囲に列挙される特定の組合せではなく、組み合されてよい。

20

【図面の簡単な説明】

【0008】

【図1】一実施形態に係るデータベースソフトウェアスタックの多様な構成要素を示すブロック図である。

【図2】いくつかの実施形態に従って、ウェブサービスベースのデータベースサービスを実装するように構成されてよいサービスシステムアーキテクチャを示すブロック図である。

30

【図3】一実施形態に係る、データベースエンジン、及び別個の分散型データベースストレージサービスを含むデータベースシステムの多様な構成要素を示すブロック図である。

【図4】一実施形態に係る、分散型データベース最適化ストレージシステムを示すブロック図である。

【図5】一実施形態に係る、データベースシステムでの別個の分散型データベース最適化ストレージシステムの使用を示すブロック図である。

【図6】一実施形態に係る、分散型データベース最適化ストレージシステムの所与のノードにデータ及びメタデータがどのように記憶されてよいのかを示すブロック図である。

【図7】一実施形態に係る、データベースボリュームの例の構成を示すブロック図である。

40

【図8】いくつかの実施形態に係る、分散型データベースシステムでのシステム全体のチェックポイント回避のための方法を示す流れ図である。

【図9A】いくつかの実施形態に係る、分散型データベースシステムのための高速クラッシュ回復を実行する方法を明示する一連の図である。

【図9B】いくつかの実施形態に係る、分散型データベースシステムのための高速クラッシュ回復を実行する方法を示す流れ図である。

【図9C】いくつかの実施形態に係る、回復されたデータベースでアクセス要求を処理するための方法を示す流れ図である。

【図10】多様な実施形態に従って、データベースエンジン、及び別個の分散型データベースストレージサービスを含むデータベースシステムの少なくとも一部を実装するように

50

構成されるコンピュータシステムを示すブロック図である。

【発明を実施するための形態】

【0009】

分散型データベースシステムのためのシステム全体のチェックポイント回避の多様な実施形態が開示される。分散型ストレージシステムのストレージノードは、いくつかの実施形態では、ストレージノードに記憶される特定のデータページにリンクされる1つまたは複数のリドウログレコードをデータベースシステムから受信してよい。データページは、データベースのためのデータを記憶する複数のデータページの内の1つであってよい。合体イベントは、特定のデータページにリンクされる1つまたは複数のリドウログレコードに少なくとも部分的に基づいて特定のデータページについて検出されてよい。合体動作は、特定のデータページの以前に記憶されていたバージョンに1つまたは複数のログレコードを適用して、特定のデータページをその現在の状態で生成するために実行されてよい。

10

【0010】

分散型データベースシステムのための高速クラッシュ回復の多様な実施形態が開示される。データベースシステムヘッドノードは、いくつかの実施形態では、故障回復動作を実行してよい。システム故障から回復すると、データベースのためのデータを記憶する分散型ストレージシステムのストレージノードとの接続が確立されてよい。いくつかの実施形態では、ストレージノードとの接続の確立時、データベースヘッドノードは、データベースをアクセスのために利用できるようにしてよい。少なくともいくつかの実施形態では、1つまたは複数のアクセス要求が受信されてよく、1つまたは複数のデータページの現在の状態が要求され、ストレージノードから受信されてよい。

20

【0011】

明細書は、まず、システム全体のチェックポイント回避（例えば、作成、削除、使用、操作等）及び高速クラッシュ回復の技法を実装するように構成される例のウェブサービスベースのデータベースサービスを説明する。例のウェブサービスベースのデータベースサービスの説明に含まれているのは、データベースエンジン及び別個の分散型データベースストレージサービス等の、例のウェブサービスベースのデータベースサービスの多様な態様である。明細書は、次いでシステム全体のチェックポイント回避及び高速クラッシュ回復のための方法の多様な実施形態のフローチャートを説明する。次に、明細書は、開示されている技法を実装してよい例のシステムを説明する。明細書を通して多様な例が提供される。

30

【0012】

本明細書に説明されるシステムは、いくつかの実施形態では、クライアント（例えば、加入者）がクラウドコンピューティング環境でデータストレージシステムを操作できるようにするウェブサービスを実装してよい。いくつかの実施形態では、データストレージシステムは、高度にスケーラブル且つ拡張可能である企業クラスのデータベースシステムであってよい。いくつかの実施形態では、クエリーは複数の物理リソース全体で分散されるデータベースストレージに向けられてよく、データベースシステムは必要に応じてスケールアップ、またはスケールダウンされてよい。データベースシステムは、異なる実施形態で、多様なタイプ及び/または編成のデータベーススキーマと効果的に機能してよい。いくつかの実施形態では、クライアント/加入者は、例えばSQLインタフェースを介してデータベースシステムに対話的に等、いくつかの方法でクエリーを提出してよい。他の実施形態では、外部アプリケーション及びプログラムは、データベースシステムにオープンデータベースコネクティビティ（ODBC）ドライバインタフェース及び/またはJavaデータベースコネクティビティ（JDBC）ドライバインタフェースを使用してクエリーを提出してよい。

40

【0013】

すなわち、本明細書に説明されるシステムは、いくつかの実施形態では、単一のデータベースシステムの多様な機能構成要素が本質的に分散されるサービス指向型データベースアーキテクチャを実装してよい。これらのシステムは、例えば、（それぞれが、アプリケ

50

ーションサーバ、サーチ機能性、またはデータベースのコア機能を提供するために必要とされる機能性を超える他の機能性等の外來の機能性を含んでよい)複数の完全にモノリシックなデータベースインスタンスを束ねるよりむしろ、データベースの基本的な動作(例えば、クエリー処理、トランザクション管理、キャッシング、及び記憶)を、個々に且つ無関係にスケラブルであってよい階層に編成してよい。例えば、いくつかの実施形態では、本明細書に説明されるシステムの各データベースインスタンスは、(単一のデータベースエンジンヘッドノード及びクライアント側ストレージシステムドライバを含んでよい)データベース階層、及び(既存のシステムのデータベース階層で従来実行される動作のいくつかを集的に実行する複数のストレージノードを含んでよい)別個の分散されたストレージシステムを含んでよい。

10

【0014】

本明細書により詳細に説明されるように、いくつかの実施形態では、データベースの最低レベルの動作(例えば、バックアップ動作、復元動作、スナップショット動作、回復動作、ログレコード操作動作、及び/または多様なスペース管理動作)のいくつかは、データベースエンジンからストレージ層にオフロードされ、複数のノード及びストレージデバイス全体で分散されてよい。例えば、いくつかの実施形態では、データベースエンジンがデータベース(またはデータベースのデータページ)に変更を適用し、次いで修正されたデータページをストレージ層に送信するよりむしろ、記憶されているデータベース(及びデータベースのデータページ)に対する変更の適用は、ストレージ層自体の責任であってよい。係る実施形態では、修正されたデータページよりむしろ、リドゥログレコードがストレージ層に送信されてよく、その後リドゥ処理(例えば、リドゥログレコードの適用)はいくぶんゆったりと且つ(例えば、バックグラウンドプロセスによって等)分散された方法で実行されてよい。いくつかの実施形態では、クラッシュ回復(例えば、記憶されているリドゥログレコードからのデータページの再構築)は、ストレージ層によって実行されてもよく、分散された(及び、いくつかの場合、ゆったりとした)バックグラウンドプロセスによって実行されてもよい。

20

【0015】

いくつかの実施形態では、リドゥログだけ(及び修正されたデータページではない)がストレージ層に送信されるため、データベース階層とストレージ層との間にあるネットワークトラフィックは、既存のデータベースシステムにおいてよりもはるかに少なくてもよい。いくつかの実施形態では、各リドゥログは、各リドゥログが変更を指定する対応するデータページのサイズのほぼ10分の1であってよい。データベース階層及び分散型ストレージシステムから送信される要求が非同期であってよいこと、及び複数の係る要求が一度に送信中であってよいことに留意されたい。

30

【0016】

一般的に、1個のデータを与えられた後、データベースの主要な要件は、最終的のその1個のデータを返すことができることである。これを行うために、データベースはそれぞれが異なる機能を実行するいくつかの異なる構成要素(または階層)含んでよい。例えば、従来のデータベースは3つの階層、つまりクエリーパーシング、最適化、及び実行を実行するための第1の階層、トランザクション性(transactionality)、回復、及び耐久性を提供するための第2の階層、及びローカルでアタッチされたディスクでまたはネットワークでアタッチされたストレージのどちらかでストレージを提供する第3の階層を有すると見なされてよい。上述されたように、従来のデータベースをスケリングしようとする以前の試みは、通常、データベースの3つすべての階層を複製し、それらの複製されたデータベースインスタンスを複数のマシン全体で分散することを伴っていた。

40

【0017】

いくつかの実施形態では、本明細書に説明されるシステムは、従来のデータベースにおいてとは異なってデータベースシステムの機能性を仕切ってよく、スケリングを実装するために複数のマシン全体で(完全なデータベースインスタンスよりもむしろ)機能構成

50

要素のサブセットだけを分散してよい。例えば、いくつかの実施形態では、クライアントが面する階層は、どのデータが記憶されるべきなのか、または取り出されるべきなのかを指定するが、どのようにしてデータを記憶するのか、または取り出すのかは指定しない要求を受信するように構成されてよい。この階層は、要求のパーシング及び/または最適化（例えば、SQLのパーシング及び要求）を実行してよい。一方、別の階層が、クエリーの実行に責任を負ってよい。いくつかの実施形態では、第3の階層が結果のトランザクション性及び一貫性を提供することに責任を負ってよい。例えば、この階層は、いわゆるACIDプロパティのいくらか、特にデータベースをターゲットとするトランザクションの原子性を強化するよう構成されてよく、データベースの中で一貫性を維持し、データベースをターゲットとするトランザクション間で独立性を保証する。いくつかの実施形態では、第4の階層が次いで多様な種類の障害が存在する場合に記憶されているデータの耐久性を提供することに責任を負ってよい。例えば、この階層は、ロギングの変更、データベースクラッシュからの回復、基礎的な記憶ボリュームに対するアクセスの管理、及び/または基礎的な記憶ボリュームにおけるスペース管理に責任を負ってよい。

10

【0018】

ここで図を参照すると、図1は、一実施形態に係る、データベースソフトウェアスタックの多様な構成要素を示すブロック図である。この例に示されるように、データベースインスタンスは、それぞれがデータベースインスタンスの機能性の一部を提供する、複数の機能構成要素（または層）を含んでよい。この例では、データベースインスタンス100は、（110として示される）クエリーパーシング及びクエリー最適化層、（120として示される）クエリー実行層、（130として示される）トランザクション性及び一貫性管理層、並びに（140として示される）耐久性及びスペース管理層を含む。上述されたように、いくつかの既存のデータベースシステムでは、データベースインスタンスのスケールアップは、（図1に示される層のすべてを含んだ）データベースインスタンス全体を1回または複数回複製して、次いで層を互いに縫い合わせるためにグルーロジックを追加することを含んでよい。いくつかの実施形態では、本明細書に説明されるシステムは、代わりにデータベース階層から別個のストレージ層に耐久性及びスペース管理層140の機能性をオフロードしてよく、その機能性をストレージ層の複数のストレージノード全体で分散してよい。

20

【0019】

いくつかの実施形態では、本明細書に説明されるデータベースシステムは、図1に示されるデータベースインスタンスの上半分の構造の多くを保持してよいが、バックアップ動作、復元動作、スナップショット動作、回復動作、及び/または多様なスペース管理動作の少なくとも部分に対する責任を記憶階層に再配分してよい。このようにして機能性を再配分し、データベース階層と記憶階層との間でログ処理をしっかりと結合することは、スケラブルデータベースを提供する以前の手法と比較されるときに性能を改善し、可用性を高め、コストを削減してよい。例えば、（実際のデータページよりもサイズがはるかに小さい）リドゥログレコードだけがノード全体で送り出される、または書込み動作のレーテンシパスの中で持続してよいので、ネットワーク及び入出力帯域幅の要件が削減されてよい。さらに、データページの生成は、入信書込み動作を遮ることなく、（フォアグラウンド処理が許すので）各ストレージノードでバックグラウンドで独立して実行できる。いくつかの実施形態では、ログ構造化された非上書きストレージの使用が、例えばデータページの移動またはコピーよりむしろメタデータ操作を使用することによって、バックアップ動作、復元動作、スナップショット動作、ポイントインタイムリカバリ動作、及びボリューム増大動作をより効率的に実行できるようにしてよい。いくつかの実施形態では、ストレージ層は、複数のストレージノード全体でクライアントの代わりに記憶されたデータの複製（及び/またはリドゥログレコード等の、そのデータと関連付けられたメタデータ）に対する責任を負ってもよい。例えば、データ（及び/またはメタデータ）は、（例えば、ストレージノードの集合体が独自の物理的に別個の独立したインフラストラクチャで実行する単一の「可用性ゾーン」の中で等）ローカルに、及び/または単一の領域のもし

30

40

50

くは異なる領域の可用性ゾーン全体で複製されてよい。

【0020】

多様な実施形態では、本明細書に説明されるデータベースシステムは、さまざまなデータベース動作のために標準的なまたはカスタムのアプリケーションプログラミングインタフェース（API）をサポートしてよい。例えば、APIは、データベースの作成、テーブルの作成、テーブルの変更、ユーザーの作成、ユーザーの削除、テーブルでの1行または複数行の挿入、値のコピー、テーブルの中からのデータの選択（例えば、テーブルの問合せ）、クエリーの取消しまたはアボート、スナップショットの作成のための動作、及び/または他の動作をサポートしてよい。

【0021】

いくつかの実施形態では、データベースインスタンスのデータベース階層は、多様なクライアントプログラム（例えば、アプリケーション）及び/または加入者（ユーザー）からの読取り要求及び/または書込み要求を受信し、次いで要求をパースし、関連付けられたデータベース動作（複数の場合がある）実施するための実行計画を作成するデータベースエンジンヘッドノードサーバを含んでよい。例えば、データベースエンジンヘッドノードは、複雑なクエリー及び接合の結果を得るために必要な一連のステップを作成してよい。いくつかの実施形態では、データベースエンジンヘッドノードは、データベース階層と別個の分散型データベース最適化ストレージシステムとの間の通信だけではなく、データベースシステムのデータベース階層とクライアント/加入者との間の通信も管理してよい。

【0022】

いくつかの実施形態では、データベースエンジンヘッドノードは、JDBCインタフェースまたはODBCインタフェースを通してエンドクライアントからSQL要求を受信すること、及びローカルでSQL処理及び（ロッキングを含んでよい）トランザクション管理を実行することに責任を負ってよい。ただし、データベースエンジンヘッドノード（またはデータベースエンジンヘッドノードの多様な構成要素）は、データページをローカルで生成するよりむしろ、リドゥログレコードを生成してよく、リドゥログレコードを別個の分散型ストレージシステムの適切なノードに送り出してよい。いくつかの実施形態では、分散型ストレージシステムのためのクライアント側ドライバは、データベースエンジンヘッドノードでホストされてよく、それらのリドゥログレコードが向けられるセグメント（またはセグメントのデータページ）を記憶する1つのストレージシステムノード（または複数のストレージシステムノード）にリドゥログレコードを送ることに責任を負ってよい。例えば、いくつかの実施形態では、各セグメントは保護グループを形成する複数のストレージシステムノードでミラーリングされてよい（またはそれ以外の場合、耐久的にされてよい）。係る実施形態では、クライアント側ドライバは、各セグメントが記憶されるノードを追跡調査してよく、クライアント要求が受信されるときに（例えば非同期で、及び実質的にほぼ同時に並列で）セグメントが記憶されるノードのすべてにリドゥログを送ってよい。クライアント側ドライバが（リドゥログレコードがストレージノードに書き込まれていることを示すことがある）保護グループのストレージノードの書込み選抜グループ（quorum）から肯定応答を受信するとすぐに、クライアント側ドライバはデータベース階層に（例えば、データベースエンジンヘッドノードに）要求された変更の肯定応答を送信してよい。例えば、データが保護グループを使用することによって耐久的にされる実施形態では、データベースエンジンヘッドノードは、クライアント側ドライバが書込み選抜グループを構成するために十分なストレージノードインスタンスから回答を受信するまで及び受信しない限り、トランザクションをコミットできないことがある。同様に、特定のセグメントに向けられる読取り要求の場合、クライアント側ドライバは、（例えば非同期で、及び実質的に同時に並列で）セグメントが記憶されるノードのすべてに読取り要求を送ってよい。クライアント側ドライバは保護グループのストレージノードの読取り選抜グループから要求されたデータを受信するとすぐに、クライアント側ドライバはデータベース階層に（例えば、データベースエンジンヘッドノードに）要求されたデータを返

10

20

30

40

50

してよい。

【0023】

いくつかの実施形態では、データベース階層（またはより詳細には、データベースエンジンヘッドノード）は、最近アクセスされたデータページが一時的に保持されるキャッシュを含んでよい。係る実施形態では、係るキャッシュに保持されるデータページをターゲットとする書込み要求が受信されると、対応するリドゥログレコードをストレージ層に送り出すことに加えて、データベースエンジンはそのキャッシュに保持されているデータページのコピーに変更を適用してよい。ただし、他のデータベースシステムにおいてとは異なり、このキャッシュに保持されるデータページはストレージ層にフラッシュされることはなく、該データページはいつでも（例えば、キャッシュに入れられたコピーに最も最近に適用された書込み要求のリドゥログレコードがストレージ層に送信され、肯定応答された後のいつでも）廃棄されてよい。キャッシュは、異なる実施形態で、一度に多くても一人の書込み者（または複数の読取り者）によるキャッシュへのアクセスを制御するための多様なロッキング機構のいずれかを実装してよい。ただし、係るキャッシュを含む実施形態では、キャッシュは複数のノード全体で分散れるのではなく、所与のデータベースインスタンスのためにデータベースエンジンヘッドノードだけに存在してよいことに留意されたい。したがって、管理するキャッシュコピーレンシーまたは一貫性問題がないことがある。

10

【0024】

いくつかの実施形態では、データベース階層は、例えば、読取り要求を送ることができるデータベース階層の異なるノードでのデータの読取り専用コピー等、システムでの同期または非同期の読取りレプリカの使用をサポートしてよい。係る実施形態では、所与のデータベースのデータベースエンジンヘッドノードが特定のデータページに向けられる読取り要求を受信すると、データベースエンジンヘッドノードはこれらの読取り専用コピーの内のいずれか1つ（または特定の1つ）に要求を送ってよい。いくつかの実施形態では、データベースエンジンヘッドノードのクライアント側ドライバは、（例えば、これらの他のノードにそのキャッシュを無効にするように促すために）キャッシュに入れられたデータページに対する更新及び/または失効についてこれらの他のノードに通知するように構成されてよい（その後これらの他のノードはストレージ層から更新されたデータページの更新済みのコピーを要求してよい）。

20

30

【0025】

いくつかの実施形態では、データベースエンジンヘッドノードで実行中のクライアント側ドライバは、記憶階層にプライベートインタフェースを曝露してよい。いくつかの実施形態では、クライアント側ドライバは従来のiSCSIインタフェースを1つまたは複数の他の構成要素（例えば、他のデータベースエンジンまたは仮想コンピューティングサービス構成要素）に曝露してもよい。いくつかの実施形態では、記憶階層でのデータベースインスタンスのためのストレージは、制限なくサイズを増大することがあり、それと関連付けられた、制限されない数のIOPSを有することがある単一のボリュームとしてモデル化されてよい。ボリュームが作成されるとき、ボリュームは特定のサイズで、（例えば、ボリュームがどのように複製されるのかを指定する）特定の可用性/耐久性特徴で、及び/またはボリュームと関連付けられたIOPSレートで（例えば、ピークと持続の両方）作成されてよい。例えば、いくつかの実施形態では、さまざまな異なる耐久性モデルがサポートされてよく、ユーザー/加入者は自らのデータベースのために、複製コピー、ゾーン、もしくは領域の数、及び/またはその耐久性、性能、及びコストの目的に基づいて複製が同期であるのか、それとも非同期であるのかを指定できてよい。

40

【0026】

いくつかの実施形態では、クライアント側ドライバはボリュームについてのメタデータを維持してよく、ストレージノード間で追加のホップを必要とすることなく、読取り要求及び書込み要求を実行するために必要なストレージノードのそれぞれに非同期要求を直接的に送信してよい。例えば、いくつかの実施形態で、データベースに対する変更を行う要

50

求に応じて、クライアント側ドライバは、ターゲットとされたデータページのストレージを実装している1つまたは複数のノードを決定し、それらのストレージノードに対してその変更を指定するリドゥログレコード（複数の場合がある）を送るように構成されてよい。ストレージノードは、次いで、リドゥログレコードに指定される変更を将来のある時点でターゲットとされたデータページに適用することに責任を負ってよい。書込みはクライアント側ドライバに肯定応答されるので、クライアント側ドライバは、ボリュームが耐久的となる点を先に進めてよく、データベース階層に対してコミットを肯定応答してよい。上述されたように、いくつかの実施形態では、クライアント側ドライバはストレージノードサーバにデータページを絶対に送信しないことがある。これは、ネットワークトラフィックを削減するだけでなく、チェックポイントまたは以前のデータベースシステムでのフォアグラウンド処理スループットを制約するバックグラウンド書込み者スレッドの必要性を削除してもよい。

10

【0027】

いくつかの実施形態では、多くの読取り要求がデータベースエンジンヘッドノードキャッシュによって提供されてよい。ただし、大規模故障イベントは一般的すぎて、メモリ内複製だけを許可できないので、書込み要求は耐久性を必要としてよい。したがって、本明細書に説明されるシステムは、記憶階層内のデータストレージを2つの領域、つまりリドゥログレコードがデータベース階層から受信されるときにリドゥログレコードが書き込まれる小さなアペンド専用ログ構造化領域、及びバックグラウンドでデータページの新しいバージョンを作成するために、ログレコードがともに合体するより大きな領域として実装することによって、フォアグラウンドレーテンシパス内にあるリドゥログレコード書込み動作のコストを最小限に抑えるように構成されてよい。いくつかの実施形態では、メモリ内構造は、インスタンス化されたデータブロックが参照されるまで連鎖ログレコード後方へ、データページの前のリドゥログレコードを指すデータページごとに維持される。この手法は、読取りがおもにキャッシュに入れられるアプリケーション内を含んで、混合した読取り 書込みワークロードに優れた性能を提供してよい。

20

【0028】

いくつかの実施形態では、リドゥログレコードのためのログ構造化データストレージへのアクセスは、（ランダム入出力動作よりむしろ）一連の順次入出力動作から構成されてよい。そのため、行われている変更は互いに密接にパックされてよい。データページに変更するたびに、永続データストレージに対する2つの入出力動作（リドゥログのための動作及び修正されたデータページ自体のための動作）が生じる既存のシステムとは対照的に、いくつかの実施形態では、本明細書に説明されるシステムはリドゥログレコードの受信に基づいて分散型ストレージシステムのストレージノードでデータページを合体させることによってこの「書込み増幅」を回避してよい。

30

【0029】

上述されたように、いくつかの実施形態では、データベースシステムの記憶階層はデータベーススナップショットを撮ることに責任を負ってよい。ただし、記憶階層はログ構造化ストレージを実装するため、データページ（例えば、データブロック）のスナップショットを撮ることはデータページ/ブロックに最も最近適用されたリドゥログレコードと関連付けられたタイムスタンプ（またはデータページ/ブロックの新しいバージョンを作成するために複数のリドゥログレコードを合体させるための最も最近の動作と関連付けられたタイムスタンプ）を記録すること、及びページ/ブロックの以前のバージョン及び時間内に記録された点までのあらゆる以後のログエントリのガベージコレクションを妨げることを含んでよい。係る実施形態では、データベーススナップショットを撮ることは、オフボリュームバックアップ戦略を利用するときに必要なとされるだろう、データブロックの読取り、コピー、または書込みを必要としないことがある。いくつかの実施形態では、ユーザー/加入者はアクティブデータセットに加えてオンボリュームスナップショットのためにどれほど多くの追加スペースを保つことを希望するのかが選ぶことができる。修正されたデータだけが追加のスペースを必要とするので、スナップショットのスペ

40

50

ース要件は最小であってよい。異なる実施形態では、スナップショットは、不連続（例えば、各スナップショットは時間の特定の時点でのデータページ内のデータのすべてに対するアクセスを提供してよい）または連続（例えば、各スナップショットは2つの時点の間のデータページに存在するデータのすべてのバージョンに対するアクセスを提供してよい）であってよい。いくつかの実施形態では、以前のスナップショットに戻ることは、そのスナップショット以降のすべてのリドゥログレコード及びデータページが無効であり、ガベージコレクション可能であることを示すためにログレコードを記録すること、及びスナップショット点後のすべてのデータベースキャッシュエントリを廃棄することを含んでよい。係る実施形態では、ストレージシステムは、ストレージシステムが通常の順方向読取り/書き込み処理で行うのと同様に、要求されるように、及びすべてのノード全体でバックグラウンドで、ブロック単位でリドゥログレコードをデータブロックに適用するので、前進復帰は必要とされないことがある。クラッシュ回復は、それによってノード全体で並列且つ分散型にされてよい。

10

【0030】

ウェブサービスベースのデータベースサービスを実装するように構成されてよいサービスシステムアーキテクチャの一実施形態が図2に示される。示されている実施形態では、（データベースクライアント250aから250nとして示される）多くのクライアントがネットワーク260を介してウェブサービスプラットフォーム200と対話するように構成されてよい。ウェブサービスプラットフォーム200は、データベースサービス210、分散型データベース最適化ストレージサービス220、及び/または1つまたは複数の他の仮想コンピューティングサービス230の1つまたは複数のインスタンスとインタフェースをとるように構成されてよい。所与の構成要素の1つまたは複数が存在してよい場合、本明細書でのその構成要素に対する参照は単数形または複数形のどちらかで行われてよいことが留意される。ただしどちらの形の使用も他方を排除することを目的としていない。

20

【0031】

多様な実施形態では、図2に示される構成要素は、コンピュータハードウェア（例えば、マイクロプロセッサもしくはコンピュータシステム）によって直接的にまたは間接的に実行可能な命令として、またはこれらの技法の組合せを使用してコンピュータハードウェアの中で直接的に実装されてよい。例えば、図2の構成要素はそれぞれが図10に示され、以下に説明されるコンピュータシステム実施形態に類似してよい、いくつかのコンピューティングノード（つまり、単にノード）を含むシステムによって実装されてよい。多様な実装形態では、所与のサービスシステム構成要素（例えば、データベースサービスの構成要素またはストレージサービスの構成要素）の機能性は、特定のノードによって実装されてよい、またはいくつかのノード全体で分散されてよい。いくつかの実施形態では、所与のノードは複数のサービスシステム構成要素（例えば、複数のデータベースサービスシステム構成要素）の機能性を実装してよい。

30

【0032】

一般的に言えば、クライアント250は、データベースサービスに対する要求（例えば、スナップショットを生成する要求等）を含むウェブサービス要求を、ネットワーク260を介してウェブサービスプラットフォーム200に提出するように構成可能な任意のタイプのクライアントを包含してよい。例えば、所与のクライアント250は、ウェブブラウザの適切なバージョンを含んでよい、またはウェブブラウザによって提供される実行環境に対する拡張部として、またはウェブブラウザによって提供される実行環境の中で実行するように構成されるプラグインモジュールまたは他のタイプのコードモジュールを含んでよい。代わりに、クライアント250（例えば、データベースサービスクライアント）は、データベースアプリケーション（もしくはデータベースアプリケーションのユーザーインタフェース）、メディアアプリケーション、オフィスアプリケーション、または1つまたは複数のデータベースを記憶する、及び/または1つまたは複数のデータベースにアクセスするために永続記憶装置リソースを利用してよい任意の他のアプリケーション等のアプリケー

40

50

ションを包含してよい。いくつかの実施形態では、係るアプリケーションは、必ずしもすべてのタイプのウェブベースのデータに対する完全なブラウザサポートを実装しなくてもウェブサービス要求を生成し、処理するための（例えば、ハイパテキスト転送プロトコル（HTTP）の適切なバージョンのための）十分なプロトコルサポートを含んでよい。すなわち、クライアント250は、ウェブサービスプラットフォーム200と直接的に対話するように構成されるアプリケーションであってよい。いくつかの実施形態では、クライアント250は、表象状態転送（Representational State Transfer）（REST）様式ウェブサービスアーキテクチャ、ドキュメントベースもしくはメッセージベースのウェブサービスアーキテクチャ、または別の適切なウェブサービスアーキテクチャに従ってウェブサービス要求を生成するよう構成されてよい。

10

【0033】

いくつかの実施形態では、クライアント250（例えば、データベースサービスクライアント）は、データベースのウェブサービスベースのストレージへのアクセスを、他のアプリケーションに、それらのアプリケーションにはトランスペアレントな方法で提供するように構成されてよい。例えば、クライアント250は、オペレーティングシステムまたはファイルシステムと統合して、本明細書に説明されるストレージモデルの適切な変形に従ってストレージを提供するように構成されてよい。ただし、オペレーティングシステムまたはファイルシステムは、ファイル、ディレクトリ、及び/またはフォルダの従来のファイルシステム階層等の、アプリケーションに異なるストレージインタフェースを提示してよい。係る実施形態では、アプリケーションは図1のストレージシステムサービスモデルを利用するために修正される必要はないことがある。代わりに、ウェブサービスプラットフォーム200へのインタフェースをとることの詳細は、オペレーティングシステム環境の中で実行するアプリケーションの代わりに、クライアント250及びオペレーティングシステムまたはファイルシステムによって調整されてよい。

20

【0034】

クライアント250は、ネットワーク260を介してウェブサービスプラットフォーム200にウェブサービス要求（例えば、スナップショット要求、スナップショット要求のパラメータ、読取り要求、スナップショットの復元等）を伝達し、ウェブサービスプラットフォーム200から応答を受信してよい。多様な実施形態では、ネットワーク260は、クライアント250とプラットフォーム200との間でウェブベースの通信を確立するために必要なネットワーキングハードウェア及びプロトコルの任意の適切な組合せを包含してよい。例えば、ネットワーク260は、集合的にインターネットを実装する多様な電気通信ネットワーク及びサービスプロバイダを概して包含してよい。また、ネットワーク260は、公衆無線ネットワークまたは構内無線ネットワークだけではなく、ローカルエリアネットワーク（LAN）または広域ネットワーク（WAN）等の構内ネットワークも含んでよい。例えば、所与のクライアント250とウェブサービスプラットフォーム200の両方も、独自の内部ネットワークを有する企業の中でそれぞれプロビジョニングされてよい。係る実施形態では、ネットワーク260は、インターネットとウェブサービスプラットフォーム200との間だけではなく、所与のクライアント250とインターネットとの間にネットワーキングリンクを確立するために必要なハードウェア（例えば、モデム、ルータ、開閉器、ロードバランサ、プロキシサーバ等）及びソフトウェア（例えば、プロトコルスタック、財務会計ソフト、ファイアウォール/セキュリティソフトウェア等）を含んでよい。いくつかの実施形態では、クライアント250は、公衆インターネットよりむしろ構内ネットワークを使用してウェブサービスプラットフォーム200と通信してよい。例えば、クライアント250は、データベースサービスシステム（例えば、データベースサービス210及び/または分散型データベース最適化ストレージサービス220を実装するシステム）と同じ企業の中でプロビジョニングされてよい。係る場合、クライアント250は、構内ネットワーク260（例えば、インターネットベースの通信プロトコルを使用してよいが、公にアクセス可能ではないLANまたはWAN）を通して完全にプラットフォーム200と通信してよい。

30

40

50

【 0 0 3 5 】

一般的に言えば、ウェブサービスプラットフォーム 2 0 0 は、データページ（またはデータページのレコード）にアクセスする要求等のウェブサービス要求を受信し、処理するように構成される 1 つまたは複数のサービスエンドポイントを実装するように構成されてよい。例えば、ウェブサービスプラットフォーム 2 0 0 は、特定のエンドポイントを実装するように構成されるハードウェア及び/またはソフトウェアを含んでよく、したがってそのエンドポイントに向けられた HTTP ベースのウェブサービス要求は適切に受信され、処理される。一実施形態では、ウェブサービスプラットフォーム 2 0 0 は、クライアント 2 5 0 からウェブサービス要求を受信し、ウェブサービス要求を、処理のためにデータベースサービス 2 1 0、分散型データベース最適化ストレージサービス 2 2 0、及び/または別の仮想コンピューティングサービス 2 3 0 を実装するシステムの構成要素に転送するように構成されるサーバシステムとして実装されてよい。他の実施形態では、ウェブサービスプラットフォーム 2 0 0 は、大規模なウェブサービス要求処理ロードを動的に管理するように構成されるロードバランス機能及び他の要求管理機能を実装する（例えば、クラスタポロジの）いくつかの別個のシステムとして構成されてよい。多様な実施形態では、ウェブサービスプラットフォーム 2 0 0 は、REST 様式またはドキュメントベースの（例えば、SOAP ベースの）タイプのウェブサービス要求をサポートするように構成されてよい。

10

【 0 0 3 6 】

いくつかの実施形態では、ウェブサービスプラットフォーム 2 0 0 は、クライアントのウェブサービス要求に対するアドレス可能なエンドポイントとして機能することに加えて、多様なクライアント管理機能を実装してよい。例えば、プラットフォーム 2 0 0 は、例えば要求側クライアント 2 5 0 のアイデンティティ、クライアント要求の数及び/または頻度、クライアント 2 5 0 の代わりに記憶されているまたは取り出されるデータテーブル（またはデータテーブルのレコード）のサイズ、クライアント 2 5 0 によって使用される全体的な記憶帯域幅、クライアント 2 5 0 によって要求されるストレージのクラス、または任意の他の測定可能なクライアント使用パラメータを追跡調査することによって、ストレージリソースを含むウェブサービスのクライアント使用の計量及びアカウンティングを調整してよい。プラットフォーム 2 0 0 は、財務会計システム及び請求書作成システムを実装してもよい、またはクライアント使用活動の報告及び請求書作成のために外部システムによって照会され、処理されてよい使用データのデータベースを維持してもよい。特定の実施形態では、プラットフォーム 2 0 0 は、クライアント 2 5 0 から受け取られる要求の割合及びタイプ、係る要求によって活用される帯域幅、係る要求のためのシステム処理レーテンシ、システム構成要素活用（例えば、ストレージサービスシステムの中のネットワーク帯域幅及び/またはストレージ活用）、要求から生じるエラーの割合及びタイプ、記憶され、要求されるデータページもしくはそのレコードの特徴（例えば、サイズ、データタイプ等）を反映する測定基準、または任意の他の適切な測定基準等、さまざまなストレージサービスシステム操作測定基準を収集する、監視する、及び/または統合するよう構成されてよい。いくつかの実施形態では、係る測定基準はシステム構成要素を調整し、維持するためにシステム管理者によって使用されてよい。一方、他の実施形態では、係る測定基準（または係る測定基準の関連性のある部分）は、係るクライアントがデータベースサービス 2 1 0、分散型データベース最適化ストレージサービス 2 2 0、及び/または別の仮想コンピューティングサービス 2 3 0（またはそれらのサービスを実装する基礎的なシステム）の使用を監視できるようにするためにクライアント 2 5 0 に曝露されてよい。

20

30

40

【 0 0 3 7 】

いくつかの実施形態では、プラットフォーム 2 0 0 は、ユーザー認証手順及びアクセス制御手順も実装してよい。例えば、特定のデータベースにアクセスする所与のウェブサービス要求の場合、プラットフォーム 2 0 0 は、要求と関連付けられるクライアント 2 5 0 が特定のデータベースにアクセスする権限を与えられているかどうかを確かめるように構成されてよい。プラットフォーム 2 0 0 は、例えばアイデンティティ、パスワード、もしくは他

50

の信用証明書を特定のデータベースと関連付けられた信用証明書に対して評価する、または特定のデータベースに対する要求されたアクセスを、特定のデータベースに対するアクセス制御リストに対して評価することによって係る権限付与を決定してよい。例えば、クライアント250が特定のデータベースにアクセスするほど十分な信用証明書を有していない場合、プラットフォーム200は、例えばエラー状態を示す応答を要求側クライアント250に返すことによって対応するウェブサービス要求を拒絶してよい。多様なアクセス制御方針は、データベースサービス210、分散型データベース最適化ストレージサービス220、及び/または他の仮想コンピューティングサービス230によってアクセス制御情報のレコードまたはリストとして記憶されてよい。

【0038】

ウェブサービスプラットフォーム200が、クライアント250がデータベースサービス210を実装するデータベースシステムの特徴にそれを通してアクセスしてよい一次インタフェースを表してよいが、ウェブサービスプラットフォーム200が係る特徴に対する単独のインタフェースを表す必要がないことが留意される。例えば、ウェブサービスインタフェースとは別個であってよい代替のAPIは、データベースシステムを提供する企業にとって内部のクライアントがウェブサービスプラットフォーム200を迂回できるようにするために使用されてよい。本明細書に説明される例の多くで、分散型データベース最適化ストレージサービス220が、クライアント250にデータベースサービスを提供するコンピューティングシステムまたは企業システムにとって内部であってよく、外部クライアント（例えば、ユーザーまたはクライアントアプリケーション）に曝露されないことがあることに留意されたい。係る実施形態では、内部「クライアント」（例えば、データベースサービス210）は、（例えば、これらのサービスを実装するシステムの間で直接的にAPIを通して）分散型データベース最適化ストレージサービス220とデータベースサービス210との間の実線として示されるローカルネットワークまたは構内ネットワーク上で分散型データベース最適化ストレージサービス220にアクセスしてよい。係る実施形態では、クライアント250の代わりにデータベースを記憶する上での分散型データベース最適化ストレージサービス220の使用はそれらのクライアントにとってトランスペアレントであってよい。他の実施形態では、分散型データベース最適化ストレージサービス220は、データベース管理のためにデータベースサービス210に依存するアプリケーション以外のアプリケーションに、データベースまたは他の情報のストレージを提供するために、ウェブサービスプラットフォーム200を通してクライアント250に曝露されてよい。これは、ウェブサービスプラットフォーム200と分散型データベース最適化ストレージサービス220の間の破線によって図2に示される。係る実施形態では、分散型データベース最適化ストレージサービス220のクライアントは、ネットワーク260を介して（例えば、インターネット上で）分散型データベース最適化ストレージサービス220にアクセスしてよい。いくつかの実施形態では、仮想コンピューティングサービス230は、クライアント250の代わりにコンピューティングサービス230を実行する上で使用されるオブジェクトを記憶するために（例えば、仮想コンピューティングサービス230と分散型データベース最適化ストレージサービス220との間で直接的にAPIを通して）分散型データベース最適化ストレージサービス220からストレージサービスを受信するように構成されてよい。これは、仮想コンピューティングサービス230と分散型データベース最適化ストレージサービス220との間の破線によって図2に示される。いくつかのケースでは、プラットフォーム200のアカウントングサービス及び/または信用証明書発行（*credentialing*）サービスは、管理クライアント等の内部クライアントにとって、または同じ企業の中のサービス構成要素間では不必要となつてよい。

【0039】

多様な実施形態では、異なる記憶方針が、データベースサービス210及び/または分散型データベース最適化ストレージサービス220によって実装されてよいことに留意されたい。係る記憶方針の例は、耐久性方針（例えば、記憶されるデータベース（またはデ

10

20

30

40

50

ータベースのデータページ)のインスタンスの数、及びデータベースが記憶される異なるノードの数を示す方針)、及び/または(要求トラフィックを一様にしようとしてデータベースまたはデータベースのデータページを、異なるノード、ボリューム、及び/またはディスク全体で分散してよい)ロードバランシング方針を含んでよい。さらに、異なる記憶方針は、サービスの多様な1つによって異なるタイプの記憶された項目に適用されてよい。例えば、いくつかの実施形態では、分散型データベース最適化ストレージサービス220は、データページに対するよりもリドゥログレコードに対してより高い耐久性を実装してよい。

【0040】

図3は、一実施形態に従って、データベースエンジン、及び別個の分散型データベースストレージサービスを含むデータベースシステムの多様な構成要素を示すブロック図である。この例では、データベースシステム300は、いくつかのデータベースのそれぞれのためのそれぞれのデータベースエンジンヘッドノード320、及び(データベースクライアント350aから350nとして示されるデータベースシステムのクライアントにとって可視であってよい、または可視でないことがある)分散型データベース最適化ストレージサービス310を含む。この例で示されるように、データベースクライアント350aから350nの内の1つまたは複数は、データベースヘッドノード320(例えば、それぞれがそれぞれのデータベースインスタンスの構成要素である、ヘッドノード320a、ヘッドノード320b、またはヘッドノード320c)に、ネットワーク360を介してアクセスしてよい(例えば、これらの構成要素はネットワークアドレス指定可能且つデータベースクライアント350aから350nにアクセス可能であってよい)。ただし、データベースクライアント350aから350nの代わりに、1つまたは複数のデータベースのデータページ(及びリドゥレコード及び/またはそれと関連付けられた他のメタデータ)を記憶し、本明細書に説明されるようにデータベースシステムの他の機能を実行するためにデータベースシステムによって利用されてよい分散型データベース最適化ストレージサービス310は、異なる実施形態では、ネットワークアドレス指定可能且つストレージクライアント350aから350nにアクセス可能であってよい、またはアクセス可能でないことがある。例えば、いくつかの実施形態では、分散型データベース最適化ストレージサービス310は、ストレージクライアント350aから350nに非可視である方法で多様な記憶動作、アクセス動作、ロギング変更動作、回復動作、ログレコード操作動作、及び/またはスペース管理動作を実行してよい。

【0041】

上述されたように、各データベースインスタンスは、多様なクライアントプログラム(例えばアプリケーション)及び/または加入者(ユーザー)から要求(例えば、スナップショット要求等)を受信し、次いで要求をパースし、要求を最適化し、関連付けられたデータベース動作(複数の場合がある)を実施するための実行計画を作成する単一のデータベースエンジンヘッドノード320を含んでよい。図3に示される例では、データベースエンジンヘッドノード320aのクエリーパーシング、最適化、及び実行構成要素305は、データベースクライアント350aから受信され、データベースエンジンヘッドノード320aがその構成要素であるデータベースインスタンスをターゲットとするクエリーのためにこれらの機能を実行してよい。いくつかの実施形態では、クエリーパーシング、最適化、及び実行構成要素305はデータベースクライアント350aに、書込み肯定応答、要求されたデータページ(及びデータページの部分)、エラーメッセージ、及びまたは他の応答を適宜に含んでよいクエリー応答を返してよい。この例に示されるように、データベースエンジンヘッドノード320aは、分散型データベース最適化ストレージサービス310の中で多様なストレージノードに読取り要求及び/またはリドゥログレコードを送り、分散型データベース最適化ストレージサービス310から書込み肯定応答を受信し、分散型データベース最適化ストレージサービス310から要求されたデータページを受信し、及び/またはデータページ、エラーメッセージ、または他の応答を(同様にそれらをデータベースクライアント350aに返してよい)クエリーパーシング、最適化、及

10

20

30

40

50

び実行構成要素 305 に返してよい、クライアント側ストレージサービスドライバ 325 も含んでよい。

【0042】

この例では、データベースエンジンヘッドノード 320 a は、最近アクセスされたデータページが一時的に保持されてよいデータページキャッシュ 335 を含む。図 3 に示されるように、データベースエンジンヘッドノード 320 a は、データベースエンジンヘッドノード 320 a が構成要素であるデータベースインスタンスでトランザクション性及び一貫性を提供することに責任を負ってよいトランザクション及び一貫性管理構成要素 330 も含んでよい。例えば、この構成要素は、データベースインスタンス及び該データベースインスタンスに向けられるトランザクションの原子性、一貫性、及び独立性のプロパティを保証することに責任を負ってよい。図 3 に示されるように、データベースエンジンヘッドノード 320 a は、多様なトランザクションのステータスを追跡調査し、コミットしないトランザクションのあらゆるローカルでキャッシュに入れられた結果をロールバックするためにトランザクション及び一貫性管理構成要素 330 によって利用されてよいトランザクションログ 340 及びアンドゥログ 345 も含んでよい。

10

【0043】

図 3 に示される他のデータベースエンジンヘッドノード 320 (例えば、320 b 及び 320 c) のそれぞれが類似する構成要素を含んでよく、データベースクライアント 350 a から 350 n の内の 1 つまたは複数によって受信され、それが構成要素であるそれぞれのデータベースインスタンスに向けられるクエリーのために類似する機能を実行してよいことに留意されたい。

20

【0044】

いくつかの実施形態では、本明細書に説明される分散型データベース最適化ストレージシステムは、1 つまたは複数のストレージノードでの記憶のために多様な論理ボリューム、セグメント、及びページでデータを編成してよい。例えば、いくつかの実施形態では、各データベースは論理ボリュームによって表され、各論理ボリュームはストレージノードの集合体上でセグメント化される。ストレージノード内の特定のストレージノード上で生きる各セグメントは、隣接ブロックアドレスのセットを含む。いくつかの実施形態では、各データページはセグメントに記憶され、したがって各セグメントは 1 つまたは複数のデータページの集合体及びそれが記憶する各データページの (リドゥログとも呼ばれる) 変更ログ (例えば、リドゥログレコードのログ) を記憶する。本明細書に詳細に説明されるように、ストレージノードは (本明細書で ULR とも呼ばれてよい) リドゥログレコードを受信し、リドゥログレコードを合体させて、(例えば、ゆったりと及び/またはデータページもしくはデータベースクラッシュに対する要求に応じて) 対応するデータページ及び/または追加のもしくは代替のログレコードの新しいバージョンを作成するように構成されてよい。いくつかの実施形態では、データページ及び/または変更ログは (クライアントによって指定されてよく、クライアントの代わりにデータベースシステムでデータベースが維持されている) 可変構成に従って複数のストレージノード全体でミラーリングされてよい。例えば、異なる実施形態では、データログまたは変更ログの 1 つのコピー、2 つのコピー、または 3 つのコピーがデフォルト構成、アプリケーションに特有の耐久性優先度、またはクライアントによって指定される耐久性優先度に従って、1 つ、2 つ、または 3 つの異なる可用性ゾーンもしくは領域のそれぞれに記憶されてよい。

30

40

【0045】

本明細書に使用されるように、以下の用語は、多様な実施形態に従って分散型データベース最適化ストレージシステムによってデータの編成を説明するために使用されてよい。

【0046】

ボリューム：ボリュームは、ストレージシステムのユーザー/クライアント/アプリケーションが理解するストレージのきわめて耐久性のある単位を表す論理概念である。すなわち、ボリュームはデータベースの多様なユーザーページに対する書込み動作の単一の一貫性がある順序付けられたログとしてユーザー/クライアント/アプリケーションに見え

50

る分散型ストアである。各書込み動作は、ボリュームの中で単一のユーザーページのコンテンツに対する論理的な順序付けられた変形を表すユーザーログレコード（ULR）で符号化されてよい。上述されたように、ULRは、本明細書でリドゥログレコードと呼ばれてもよい。各ULRは、一意の識別子（例えば、論理シーケンス番号（LSN））を含んでよい。各ULRは、ULRに高い耐久性及び可用性を提供するために、保護グループ（PG）を形成する、分散型ストア内の1つまたは複数の同期セグメントに持続してよい。ボリュームは、バイトの可変サイズの連続範囲にLSN型の読取り/書込みインタフェースを提供してよい。

【0047】

いくつかの実施形態では、ボリュームはそれぞれが保護グループを通して耐久的にされた複数のエクステントから構成されてよい。係る実施形態では、ボリュームはボリュームエクステントの変わりやすい連続シーケンスから構成されるストレージの単位を表してよい。ボリュームに向けられる読取り及び書込みは、構成するボリュームエクステントに対する対応する読取り及び書込みにマッピングされてよい。いくつかの実施形態では、ボリュームのサイズは、ボリュームエクステントを追加することにより、又は、ボリュームの端部からボリュームエクステントを除去することにより変更されてもよい。

【0048】

セグメント：セグメントは、単一ストレージノードに割り当てられるストレージの制限される耐久性の単位である。すなわち、セグメントは、特有の固定サイズバイト範囲のデータに、限られたベストエフォート型の耐久性（例えば、ストレージノードである、故障の永続的であるが冗長ではない単一点）を提供する。多様な実施形態では、このデータは、いくつかの場合では、ユーザーアドレス指定可能なデータのミラーであってよい、またはこのデータはボリュームメタデータまたはイレイジャーコーディングされたビット等の他のデータであってよい。所与のセグメントは、正確に1つのストレージノード上で生きてよい。ストレージノードの中で、複数のセグメントが各SSD上で生きてよく、各セグメントは1つのSSDに制限されてよい（例えば、セグメントは複数のSSDに及ばないことがある）。いくつかの実施形態では、セグメントはSSD上で連続領域を占有するように要求されないことがある。むしろ、各SSDにセグメントのそれぞれによって所有される領域を記述する割当てマップがあってよい。上述されたように、保護グループは複数のストレージノードに渡って拡散される複数のセグメントから構成されてよい。いくつかの実施形態では、セグメントは、（サイズが作成時に定義される）バイトの固定サイズの隣接範囲に、LSN型読取り/書込みインタフェースを提供してよい。いくつかの実施形態では、各セグメントはセグメントUUIID（例えば、セグメントの汎用一意識別子）によって識別されてよい。

【0049】

記憶ページ：記憶ページは、概して固定サイズのメモリのブロックである。いくつかの実施形態では、各ページは、オペレーティングシステムによって定義されるサイズのメモリの（例えば、バーチャルメモリ、ディスク、または他の物理メモリの）ブロックであり、本明細書では用語「データブロック」によって参照されてもよい。すなわち、記憶ページは隣接セクタのセットであってよい。記憶ページは、ヘッダ及びメタデータがあるログページでの単位だけではなく、SSDでの割当ての単位としても役立ってよい。いくつかの実施形態では、及び本明細書に説明されるデータベースシステムの文脈では、用語「ページ」または「記憶ページ」は、通常、4096バイト、8192バイト、16384バイト、または32768バイト等の2の倍数であってよいデータベース構成によって定義されるサイズの類似したブロックを指してよい。

【0050】

ログページ：ログページは、ログレコード（例えば、リドゥログレコードまたはアンドゥログレコード）を記憶するために使用される記憶ページのタイプである。いくつかの実施形態では、ログページは、サイズが記憶ページと同一であってよい。各ログページは、例えばそれが属するセグメントを識別するメタデータ等、そのログページについてのメタ

10

20

30

40

50

データを含むヘッダを含んでよい。ログページが編成の単位であり、必ずしも書込み動作に含まれるデータの単位ではないことがあることに留意されたい。例えば、いくつかの実施形態では、標準的な転送処理の間、書込み動作は、一度の1つのセクタをログの末尾に書き込んでよい。

【0051】

ログレコード：ログレコード（例えば、ログページの個々の要素）はいくつかの異なるクラスであってよい。例えば、ストレージシステムのユーザー/クライアント/アプリケーションによって作成され、理解されるユーザーログレコード（ULR）は、ボリューム内のユーザーデータに対する変更を示すために使用されてよい。ストレージシステムによって生成される制御ログレコード（CLR）は、現在の無条件ボリューム耐久性（unc
o
n
d
i
t
i
o
n
a
l
v
o
l
u
m
e
d
u
r
a
b
l
e）LSN（VDL）等のメタデータを追跡調査するために使用される制御情報を含んでよい。ヌルログレコード（NLR）は、いくつかの実施形態では、ログセクタまたはログページの未使用のスペースを充填するためのパディングとして使用されてよい。いくつかの実施形態では、これらのクラスのそれぞれの中に多様なタイプのログレコードがあつてよく、ログレコードのタイプはログレコードを解釈するために呼び出される必要がある関数に対応してよい。例えば、1つのタイプは特定の圧縮フォーマットを使用する圧縮フォーマットのユーザーページのすべてのデータを表してよく、第2のタイプは、ユーザーページの中のバイト範囲の新しい値を表してよく、第3のタイプは、整数として解釈されるバイトのシーケンスに対する増分動作を表してよく、第4のタイプはページの中の別の場所に1バイト範囲をコピーすること
10
20

【0052】

ペイロード：ログレコードのペイロードは、ログレコードに、または特定のタイプのログレコードに特有であるデータまたはパラメータ値である。例えば、いくつかの実施形態では、大部分（またはすべての）ログレコードが含み、ストレージシステム自体が理解するパラメータまたは属性のセットがあつてよい。これらの属性は、セクタサイズに比較して相対的に小さくてよい共通のログレコードヘッダ/構造の部分であつてよい。さらに、大部分のログレコードは、そのログレコードタイプに特有の追加のパラメータまたはデータ
30

【0053】

セグメントログでログレコードを記憶する際に、いくつかの実施形態では、ペイロードはログヘッダとともに記憶されてよいことに留意されたい。他の実施形態では、ペイロードは別の場所に記憶されてよく、そのペイロードが記憶される場所に対するポインタはログヘッダとともに記憶されてよい。さらに他の実施形態では、ペイロードの一部はヘッダに記憶されてよく、ペイロードの残りは別個の場所に記憶されてよい。ペイロード全体がログヘッダとともに記憶される場合、これは帯域内ストレージと呼ばれてよい。それ以外の場合、ストレージは帯域外であると呼ばれてよい。いくつかの実施形態では、大部分の大きなAULRのペイロードは（以下に説明される）ログのコールドゾーンで帯域外で記憶されてよい。

【0054】

ユーザーページ：ユーザーページは、（固定サイズの）バイト範囲、及びストレージシステムのユーザー/クライアントに可視である特定のボリュームのためのそのアラインメントである。ユーザーページは論理概念であり、特定のユーザーページのバイトは任意の
40
50

記憶ページにそのまま記憶されてよい、または記憶されないことがある。特定のボリュームのユーザーページのサイズは、そのボリュームの記憶ページサイズとは無関係であってよい。いくつかの実施形態では、ユーザーページサイズはボリュームごとに設定可能であってよく、ストレージノード上の異なるセグメントは異なるユーザーページサイズを有してよい。いくつかの実施形態では、ユーザーページサイズは、セクタサイズ（例えば、4 K B）の倍数となるように制約されてよく、上限（例えば、6 4 K B）を有してよい。他方、記憶ページサイズは、ストレージノード全体にとって固定であってよく、基礎的なハードウェアに対する変更がない限り変化しないことがある。

【 0 0 5 5 】

データページ：データページは、圧縮された形式でユーザーページデータを記憶するために使用される記憶ページのタイプである。いくつかの実施形態では、データページに記憶されるあらゆる1個のデータがログレコードと関連付けられ、各ログレコードは（データセクタとも呼ばれる）データページの中のセクタに対するポインタを含んでよい。いくつかの実施形態では、データページは各セクタによって提供されるメタデータ以外の任意の埋込みメタデータを含まないことがある。データページ内のセクタ間には関係性がなくてよい。代わりに、ページへの編成は、セグメントへのデータの割当ての粒度の表現としてのみ存在してよい。

【 0 0 5 6 】

ストレージノード：ストレージノードは、ストレージノードサーバコードが配備される単一のバーチャルマシンである。各ストレージノードは、複数のローカルにアタッチされたSSDを含んでよく、1つまたは複数のセグメントへのアクセスにネットワークAPIを提供してよい。いくつかの実施形態では、多様なノードはアクティブリスト上または（例えば、ノードが応答するには低速である、またはそれ以外の場合、正常に機能しないが、完全に使用不可ではない場合等）劣化したリスト上にあってよい。いくつかの実施形態では、クライアント側ドライバは、ノードが交換されるべきかどうか、及びいつノードが交換されるべきかを判断するため、及び/または観察された性能に基づいて、いつ及びどのようにして多様なノードの間でデータを再配分するのかを決定するために、ノードをアクティブまたは劣化として分類するのを支援してよい（または、分類するのに責任を負ってよい）

【 0 0 5 7 】

SSD：本明細書において参照されるように、用語「SSD」は、例えばディスク、ソリッドステートドライブ、電池によって支援されるRAM、不揮発性RAMデバイス（例えば、1つまたは複数のNV-DIMM）、または別のタイプの永続ストレージデバイス等の、その記憶ボリュームによって利用されるストレージのタイプに関わりなく、ストレージノードによって見られるローカルブロック記憶ボリュームを指してよい。SSDは、必ずしも直接的にハードウェアにマッピングされない。例えば、異なる実施形態では、単一のソリッドステートストレージデバイスは、各ボリュームが複数のセグメントに分割され、複数のセグメントに渡ってストライピングされる複数のローカルボリュームに分けられる可能性がある、及び/または単一ドライブは単に管理の容易さのために複数のボリュームに分割されてよい。いくつかの実施形態では、各SSDは単一の固定場所で割当てマップを記憶してよい。このマップは、特定のセグメントによってどの記憶ページが所有されているのか、及び（データページと対照的に）これらのページの内のどれがログページであるのかを示してよい。いくつかの実施形態では、記憶ページは、転送処理が割当てを待機する必要がなくてよいように各セグメントに事前に割り当てられてよい。割当てマップに対するあらゆる変更は、新規に割り当てられた記憶ページがセグメントによって使用される前に耐久的にされる必要があることがある。

【 0 0 5 8 】

分散型データベース最適化ストレージシステムの一実施形態は、図4のブロック図によって示される。この例では、データベースシステム400は、相互接続460上でデータベースエンジンヘッドノード420と通信する分散型データベース最適化ストレージシ

10

20

30

40

50

テム 4 1 0 を含む。図 3 に示される例でのように、データベースエンジンヘッドノード 4 2 0 は、クライアント側ストレージサービスドライバ 4 2 5 を含んでよい。この例では、分散型データベース最適化ストレージシステム 4 1 0 は (4 3 0、4 4 0、及び 4 5 0 として示されるストレージシステムサーバノードを含んだ) 複数のストレージシステムサーバノードを含み、複数のストレージシステムサーバノードのそれぞれは、それが記憶するセグメント (複数の場合がある) のためのデータページ及びリドゥログのストレージ、多様なセグメント管理機能を実行するように構成されるハードウェア及び/またはソフトウェアを含む。例えば、各ストレージシステムサーバノードは以下の動作、つまり、複製 (例えば、ストレージノードの中で等ローカルに)、データページを生成するためのリドゥログの合体、スナップショット (例えば、作成、復元、削除等)、ログ管理 (例えば、ログレコードの操作)、クラッシュ回復、及び/または (例えば、セグメントの) スペース管理の内のいずれかまたはすべての少なくとも一部を実行するように構成されるハードウェア及び/またはソフトウェアを含んでよい。各ストレージシステムサーバノードは、データブロックがクライアント (例えば、ユーザー、クライアントアプリケーション、及び/またはデータベースサービス加入者) の代わりに記憶されてよい (例えば、SSD 等) 複数のアタッチされたストレージデバイスも有してよい。

10

【 0 0 5 9 】

図 4 に示される例では、ストレージシステムサーバノード 4 3 0 は、データページ (複数の場合がある) 4 3 3、セグメントリドゥログ (複数の場合がある) 4 3 5、セグメント管理機能 4 3 7、及びアタッチされた SSD 4 7 1 から 4 7 8 を含む。再び、ラベル「SSD」はソリッドステートドライブを指してよい、または指さないこともあるが、基礎的なハードウェアに関わりなく、より概してローカルブロック記憶ボリュームを指してよいことに留意されたい。同様に、ストレージシステムサーバノード 4 4 0 は、データページ (複数の場合がある) 4 4 3、セグメントリドゥログ (複数の場合がある) 4 4 5、セグメント管理機能 4 4 7、及びアタッチされた SSD 4 8 1 から 4 8 8 を含み、ストレージシステムサーバノード 4 5 0 は、データページ (複数の場合がある) 4 5 3、セグメントリドゥログ (複数の場合がある) 4 5 5、セグメント管理機能 4 5 7、及びアタッチされた SSD 4 9 1 から 4 9 8 を含む。

20

【 0 0 6 0 】

上述されたように、いくつかの実施形態では、セクタは、SSDでのアラインメントの単位であり、書込みが部分的だけに完了されるリスクなしに書き込むことができる SSDでの最大サイズであってよい。例えば、多様なソリッドステートドライブ及びスピニングメディアのセクタサイズは 4 K B であってよい。本明細書に説明される分散型データベース最適化ストレージシステムのいくつかの実施形態では、ありとあらゆるセクタは、セクタがその一部であるより高レベルのエンティティに関わりなく、セクタの始まりに 6 4 ビット (8 バイト) の CRC を含んで有してよい。係る実施形態では、(セクタが SSD から読み取られるたびに確認されてよい) この CRC は破損を検出する際に使用されてよい。いくつかの実施形態では、ありとあらゆるセクタは、その値がセクタをログセクタ、データセクタ、または初期化されていないセクタとして該セクタを識別する「セクタタイプ」バイトを含んでもよい。例えば、いくつかの実施形態では、0 のセクタタイプバイト値は、セクタが初期化されていないことを示してよい。

30

40

【 0 0 6 1 】

いくつかの実施形態では、分散型データベース最適化ストレージシステムのストレージシステムサーバノードのそれぞれは、例えばリドゥログを受信し、データページ等を送り返すために、データベースエンジンヘッドノードとの通信を管理するノードサーバのオペレーティングシステムで実行中のプロセスのセットを実装してよい。いくつかの実施形態では、分散型データベース最適化ストレージシステムに書き込まれるすべてのデータブロックは、(例えば、リモートキー値耐久性バックアップストレージシステムで) 長期の及び/またはアーカイブのストレージにバックアップされてよい。

【 0 0 6 2 】

50

図5は、一実施形態に係る、データベースシステムでの別個の分散型データベース最適化ストレージシステムの使用を示すブロック図である。この例では、1つまたは複数のクライアントプロセス510が、データベースエンジン520及び分散型データベース最適化ストレージシステム530を含むデータベースシステムによって維持される1つまたは複数のデータベースにデータを記憶してよい。図5に示される例では、データベースエンジン520がデータベース階層構成要素560、及び(分散型データベース最適化ストレージシステム530とデータベース階層構成要素560との間のインタフェースとして働く)クライアント側ドライバ540を含む。いくつかの実施形態では、データベース階層構成要素560は、図3のクエリーパーシング、最適化、及び実行構成要素305、並びにトランザクション及び一貫性管理構成要素330によって実行される機能等の機能を実行してよい、及び/またはデータページ、トランザクションログ、及び/またはアンドゥログ(例えば、図3のデータページキャッシュ335、トランザクションログ340、及びアンドゥログ345によって記憶されるもの)を記憶してよい。

10

【0063】

この例では、1つまたは複数のクライアントプロセス510は、データベース階層構成要素560に(ストレージノード535aから535nの内の1つまたは複数に記憶されるデータをターゲットとする読取り要求及び/または書込み要求を含んでよい)データベースクエリー要求515を送信してよく、データベース階層構成要素560からデータベースクエリー応答517(例えば、書込み肯定応答及び/または要求されたデータを含む応答)を受信してよい。データページに書き込む要求を含む各データベースクエリー要求515は、分散型データベース最適化ストレージシステム530への以後のルーティングのためにクライアント側ドライバ540に送信されてよい、1つまたは複数のレコード書込み要求541を生成するためにパースされ、最適化されてよい。この例では、クライアント側ドライバ540は、それぞれのレコード書込み要求541に対応する1つまたは複数のリドゥログレコード531を生成してよく、リドゥログレコード531を分散型データベース最適化ストレージシステム530のストレージノード535の特定のストレージノードに送信してよい。分散型データベース最適化ストレージシステム530は、データベースエンジン520に(具体的には、クライアント側ドライバ540に)各リドゥログレコード531の対応する書込み肯定応答532を返してよい。クライアント側ドライバ540は、これらの書込み肯定応答をデータベース階層構成要素560に(書込み応答542として)渡してよく、データベース階層構成要素560は次いでデータベースクエリー応答517の内の1つとして1つまたは複数のクライアントプロセス510に対応する応答(例えば、書込み肯定応答)を送信してよい。

20

30

【0064】

この例では、データページを読み込む要求を含む各データベースクエリー要求515は、1つまたは複数のレコード読取り要求543を生成するためにパースされ、最適化されてよく、レコード読取り要求543は分散型データベース最適化ストレージシステム530への以後のルーティングのためにクライアント側ドライバ540に送信されてよい。この例では、クライアント側ドライバ540は、分散型データベース最適化ストレージシステム530のストレージノード535の特定のストレージノードにこれらの要求を送信してよく、分散型データベース最適化ストレージシステム530はデータベースエンジン520に(具体的には、クライアント側ドライバ540に)要求されたデータページ533を返してよい。クライアント側ドライバ540は、戻りデータレコード544としてデータベース階層構成要素560に返されたデータページを送信してよく、データベース階層構成要素560は次いでデータベースクエリー応答517として1つまたは複数のクライアントプロセス510にデータページを送信してよい。

40

【0065】

いくつかの実施形態では、多様なエラーメッセージ及び/またはデータ損失メッセージ534が、分散型データベース最適化ストレージシステム530からデータベースエンジン520に(具体的には、クライアント側ドライバ540に)送信されてよい。これらの

50

メッセージは、クライアント側ドライバ540から、エラー報告メッセージ及び/または損失報告メッセージ545として、データベース階層構成要素560に、及び次いで1つまたは複数のクライアントプロセス510に、データベースクエリー応答517とともに（または代わりに）渡されてよい。

【0066】

いくつかの実施形態では、分散型データベース最適化ストレージシステム530のAPI531から534、及びクライアント側ドライバ540のAPI541から545は、データベースエンジン520が分散型データベース最適化ストレージシステム530のクライアントであるかのように、分散型データベース最適化ストレージシステム530の機能性をデータベースエンジン520に曝露してよい。例えば、データベースエンジン520は、データベースエンジン520及び分散型データベース最適化ストレージシステム530の組合せによって実装されるデータベースシステムの多様な動作（例えば、記憶動作、アクセス動作、ロギング変更動作、回復動作、及び/またはスペース管理動作）を実行するために（またはそれらの実行を容易にするために）（クライアント側ドライバ540を通して）リドゥログレコードまたは要求データページをこれらのAPIを通して書き込んでよい。図5に示されるように、分散型データベース最適化ストレージシステム530は、それぞれが複数のアタッチされたSSDを有してよいストレージノード535aから535nにデータブロックを記憶してよい。いくつかの実施形態では、分散型データベース最適化ストレージシステム530は、多様なタイプの冗長性方式の適用によって、記憶されているデータブロックに高い耐久性を提供してよい。

【0067】

多様な実施形態では、図5のデータベースエンジン520と分散型データベース最適化ストレージシステム530との間のAPI呼出し及び応答（例えば、API531から534）、及び/またはクライアント側ドライバ540とデータベース階層構成要素560との間のAPI呼出し及び応答（例えば、API541から545）は、（例えば、ゲートウェイ制御プレーンによって管理される）安全なプロキシ接続上で実行されてよい、または公衆ネットワーク上でもしくは代わりにバーチャルプライベートネットワーク（VPN）接続等のプライベートチャネル上で実行されてよいことに留意されたい。本明細書に説明されるデータベースシステムの構成要素への、及びデータベースシステムの構成要素の間のこれらの及び他のAPIは、シンプルオブジェクトアクセスプロトコル（SOAP）技術及び表象状態転送（REST）技術を含むが、これに限定されるものではない異なる技術に従って実装されてよい。例えば、これらのAPIは、SOAP APIまたはRESTful APIとして実装されてよいが、必ずしも実装されない。SOAPは、ウェブベースのサービスとの関連で情報を交換するためのプロトコルである。RESTは分散型ハイパーメディアシステム用のアーキテクチャスタイルである。（RESTfulウェブサービスとも呼ばれてよい）RESTful APIは、HTTP及びREST技術を使用して実装されるウェブサービスAPIである。本明細書に説明されるAPIは、いくつかの実施形態では、データベースエンジン520及び/または分散型データベース最適化ストレージシステム530との統合をサポートするために、C、C++、Java、C#、及びPerlを含むが、これに限定されるものではない多様な言語でクライアントライブラリでラップされてよい。

【0068】

上述されたように、いくつかの実施形態では、データベースシステムの機能構成要素は、データベースエンジンによって実行される構成要素と、別個の分散されたデータベース最適化ストレージシステムで実行される構成要素との間で仕切られてよい。1つの特定の例では、（例えば、単一のデータブロックを、そのデータブロックにレコードを追加することによって更新するために）何かをデータベースに挿入する要求をクライアントプロセス（またはクライアントプロセスのスレッド）から受信することに応じて、データベースエンジンヘッドノードの1つまたは複数の構成要素は、クエリーパーシング、最適化、及び実行を実行してよく、クエリーの各部分をトランザクション及び一貫性管理構成要素に

10

20

30

40

50

送信してよい。トランザクション及び一貫性管理構成要素は、他のクライアントプロセス（またはクライアントプロセスのスレッド）が同時に同じ行を修正しようとしていないことを保証してよい。例えば、トランザクション及び一貫性管理構成要素は、この変更がデータベースにおいて原子的に、一貫して、耐久的に、及び独立して実行されることを保証することに責任を負ってよい。例えば、トランザクション及び一貫性管理構成要素は、分散型データベース最適化ストレージサービスのノードの1つに送信されるリドゥログレコードを生成し、ACIDプロパティがこのトランザクションについて満たされていることを保証する順序で及び/またはタイミングでリドゥログレコードを（他のクライアント要求に応じて生成される他のリドゥログとともに）分散型データベース最適化ストレージサービスに送信するために、データベースエンジンヘッドノードのクライアント側ストレージサービスドライバとともに機能してよい。対応するストレージノードは、（ストレージサービスによって「更新レコード」と見なされてよい）リドゥログレコードを受信すると、データブロックを更新し、データブロックのリドゥログを更新してよい（例えば、データブロックに向けられるすべての変更のレコード）。いくつかの実施形態では、データベースエンジンは、この変更のためにアンドゥログレコードを生成することに責任を負ってよく、アンドゥログのためのリドゥログレコードを生成することにも責任を負ってよく、この両方ともトランザクション性を保証するために（データベース階層で）ローカルに使用されてよい。ただし、従来のデータベースシステムにおいてとは異なり、本明細書に説明されるシステムは、（変更をデータベース階層で適用し、修正されたデータブロックをストレージシステムに送るよりむしろ）データブロックに変更を適用するための責任をストレージシステムに移してよい。さらに、図8から図9Bで本明細書に説明されるように、多様な実施形態では、システム全体のチェックポイントは、ストレージシステムによっても実行されてよい多様なログレコード演算に起因するデータベースシステムクラッシュからの高速回復とともに、データベースシステムで回避されてよい。

【0069】

異なる実施形態で、さまざまな割当てモデルがSSDのために実装されてよい。例えば、いくつかの実施形態では、ログエントリページ及び物理アプリケーションページが、SSDデバイスと関連付けられたページの単一のヒープから割り当てられてよい。この手法は、未指定のままとなるために、及び自動的に使用に適合するためにログページ及びデータページによって消費される相対的な記憶量を残すという優位点を有してよい。また、手法は、ページが使用され、準備なしに随意に転用されるまでページを準備されないままにできるという優位点も有してよい。他の実施形態では、割当てモデルはストレージデバイスをログエントリ及びデータページのための別々のスペースに仕切ってよい。一度係る割当てモデルが図6のブロック図に示され、以下に説明される。

【0070】

図6は、一実施形態に係る、分散型データベース最適化ストレージシステムの所与のストレージノード（または永続ストレージデバイス）にデータ及びメタデータがどのように記憶されてよいのかを示すブロック図である。この例では、SSDストレージスペース600は、610と名前が付けられたスペースの部分にSSDヘッダ及び他の固定メタデータを記憶する。SSDストレージスペース600は、620と名前が付けられたスペースの部分にログページを記憶し、追加のログページのために初期化され、確保される、630と名前が付けられたスペースを含む。（640として示される）SSDストレージスペース600の一部分は初期化されているが、割り当てられておらず、（650として示される）スペースの別の部分は初期化されておらず、割り当てられていない。最後に、660と名前が付けられたSSDストレージスペース600の部分はデータページを記憶する。

【0071】

この例では、最初の使用可能なログページスロットは615として示され、最後の使用されたログページスロット（一時的）は625として示される。最後の確保されたログページスロットは635として示され、最後の使用可能なログページスロットは645とし

10

20

30

40

50

て示される。この例では、最初の使用されたデータページスロット（一時的）は665として示される。いくつかの実施形態では、SSDストレージスペース600の中でのこれらの要素（615、625、635、645、及び665）のそれぞれの位置は、それぞれのポインタによって識別されてよい。

【0072】

図6に示される割当て手法では、有効なログページはフラットストレージスペースの始まりにバックされてよい。ログページが解放されるために開く穴は、アドレススペースのさらに先に入る追加のログページスロットが使用される前に再使用されてよい。例えば、最悪の場合、最初のn個のログページスロットが有効なログデータを含み、この場合、nは今まで同時に存在した有効なログページの最大数である。この例では、有効データページはフラットストレージスペースの最後にバックされてよい。データページが解放されることにより開く穴は、アドレススペースでより下方の追加のデータページスロットが使用される前に再使用されてよい。例えば、最悪の場合、最後のmのデータページが有効なデータを含み、この場合mは今まで同時に存在した有効なデータページの最大数である。

【0073】

いくつかの実施形態では、ログページスロットが有効なログページエントリの潜在的なセットの部分になることができる前に、ログページスロットは有効な将来のログエントリページのために混同できない値に初期化される必要がある。廃棄されたログページは新しい有効なログページについて絶対に混同されることがないほど十分なメタデータを有するので、これは、リサイクルされるログページスロットに暗黙に当てはまる。ただし、ストレージデバイスが最初に初期化される時、またはアプリケーションデータページを記憶するために潜在的に使用されたスペースが再利用される時、ログページスロットは、ログページスロットがログページスロットプールに加えられる前に初期化される必要がある。いくつかの実施形態では、ログスペースのバランスを取り戻す/再利用することは、バックグラウンドタスクとして実行されてよい。

【0074】

図6に示される例では、カレントログページスロットプールは（615で）最初の使用可能なログページスロットと最後の確保されたログページスロット（625）との間に領域を含む。いくつかの実施形態では、このプールは、（例えば、最後の確保されたログページスロット635を識別するポインタに対する更新を持続させることによって）新しいログページスロットの再初期化なしに最後の使用可能なログページスロット（625）まで安全に増大してよい。この例では、（ポインタ645によって識別される）最後の使用可能なログページスロットを超えて、プールは、初期化されたログページスロットを持続し、最後の使用可能なログページスロット（645）のためのポインタを持続的に更新することによって、（ポインタ665によって識別される）最初の使用されたデータページスロットまで成長してよい。この例では、650として示される、SSDストレージスペース600の以前に初期化されておらず、割り当てられていない部分は、ログページを記憶するためにとりあえず利用されてよい。いくつかの実施形態では、カレントログページスロットプールは、最後の確保されたログページスロット（635）のポインタに対する更新を持続することによって（ポインタによって識別される）最後の使用されたログページスロットの位置まで縮小されてよい。

【0075】

図6に示される例では、カレントデータページスロットプールは、（ポインタ645によって識別される）最後の使用可能なログページスロットと、SSDストレージスペース600の最後との間に領域を含む。いくつかの実施形態では、データページプールは、最後の使用可能なログページスロット（645）のポインタに対する更新を持続することによって、最後の確保されたログページスロット（635）に対するポインタによって識別される位置まで安全に成長してよい。この例では、640として示される、SSDストレージスペース600の以前に初期化されたが、割り当てられていない部分は、データページを記憶するためにとりあえず利用されてよい。これを超えて、プールは、最後の確保された

10

20

30

40

50

ログページスロット(635)及び最後の使用可能なログページスロット(645)のポインタに対する更新を継続し、ログページよりむしろデータページを記憶するために、630及び640として示されるSSDストレージスペース600の部分を効果的に割り当てし直すことによって、最後の使用されたログページスロット(625)のポインタによって識別される位置まで安全に成長してよい。いくつかの実施形態では、データページスロットプールは、追加のログページスロットを初期化し、最後の使用可能なログページスロット(645)のポインタに対する更新を継続することによって、最初の使用されたデータページスロット(665)のポインタによって識別される位置まで安全に縮小されてよい。

【0076】

図6に示される割当て手法を利用する実施形態では、ログページプール及びデータページプールのページサイズは、優れたパッキング挙動を容易にしつつも、独立して選択されてよい。係る実施形態では、有効なログページが、アプリケーションデータによって形成されるスプーフィングされたログページにリンクする可能性はないことがあり、壊れたログと依然として書き込まれていない次のページにリンクする有効なログテールとを区別することが可能なことがある。図6に示される割当て手法を利用する実施形態では、起動時、最後の確保されたログページスロット(635)に対するポインタによって識別される位置までのログページスロットのすべてが迅速に且つ連続して読み取られてよく、(推論されるリンクング/順序付けを含む)ログインデックス全体が再構築されてよい。係る実施形態では、すべてはLSN順序制御制約から推論できるので、ログページ間の明示的なリンクングの必要性がないことがある。

【0077】

いくつかの実施形態では、セグメントは3つの主要な部分(またはゾーン)、つまり、ホットログを含む部分、コールドログを含む部分、及びユーザーページデータを含む部分から構成されてよい。ゾーンは、必ずしもSSDの隣接領域ではない。むしろ、ゾーンは、記憶ページの粒度で点在することがある。さらに、セグメント及びそのプロパティについてのメタデータを記憶するセグメントごとにルートページがあってよい。例えば、セグメントのルートページはセグメントのためのユーザーページサイズ、セグメント内のユーザーページの数、(フラッシュ番号(flush number)の形で記録されてよい)ホットログゾーンの現在の始まり/ヘッド、ボリュームエポック、及び/またはアクセス制御メタデータを記憶してよい。

【0078】

いくつかの実施形態では、ホットログゾーンは、それらがストレージノードによって受信されるにつれ、クライアントからの新しい書込みを受け入れてよい。ページの以前のバージョンからのデルタの形をとるユーザーページ/データページに対する変更を指定するデルタユーザーログレコード(DULR)及び完全なユーザーページ/データページのコンテンツを指定する絶対ユーザーログレコード(AULR)の両方とも、ログに完全に書き込まれてよい。ログレコードは、ほぼ、ログレコードが受信される順序でこのゾーンに追加されてよく(例えば、ログレコードがLSNでソートされるのではない)、それらはログページに渡って広がることがある。例えばログレコードは独自のサイズの表示を含んでよい等、ログレコードは自己記述的である必要がある。いくつかの実施形態では、ガベージコレクションはこのゾーンで実行されない。代わりに、スペースは、すべての必要とされるログレコードがコールドログにコピーされた後にログの始まりから切り詰めることによって再利用されてよい。ホットゾーンのログセクタは、セクタが作成されるたびに最も最近の既知の無条件VDLで注釈されてよい。条件付きのVDL CLRは、それらが受信されるにつれホットゾーンに書き込まれてよいが、最も最近に書き込まれたVDL CLRだけが意味を持ってよい。

【0079】

いくつかの実施形態では、新しいログページが書き込まれるたびに、新しいログページにはフラッシュ番号が割り当てられる。フラッシュ番号は、各ログページの中のあらゆる

10

20

30

40

50

セクタの部分として書き込まれてよい。フラッシュ番号は、2つのログページを比較するときに、どのログページが後に書き込まれたのかを決定するために使用されてよい。フラッシュ番号は単調に増加し、SSD（またはストレージノード）に対して調べられて（scoped）よい。例えば、単調に増加するフラッシュ番号のセットは、SSD上のすべてのセグメント（またはストレージノード上のすべてのセグメント）の間で共有される。

【0080】

いくつかの実施形態では、コールドログゾーンで、ログレコードはそのLSNの昇順で記憶されてよい。このゾーンでは、AULRはそのサイズに応じて必ずしもデータをインラインで記憶しないことがある。例えば、AULRが大きなペイロードを有する場合、ペイロードのすべてまたは一部がデータゾーンに記憶されてよく、AULRはそのデータがデータゾーンのどこに記憶されているのかを指してよい。いくつかの実施形態では、コールドログゾーンのログページは、セクタ単位でよりむしろ、一度に1全ページ、書き込まれてよい。コールドゾーンのログページは一度に全ページ書き込まれるため、全セクタ内のフラッシュ番号が同一ではないコールドゾーンのどのようなログページも不完全に書き込まれたページと見なされてよく、無視されてよい。いくつかの実施形態では、コールドログゾーンでは、DULRは（最大2ログページまで）複数のログページに及ぶことができることがある。しかし、AULRは、例えば合体動作が単一の原子的な書込みでDULRをAULRで置き換えることができるように、複数のログセクタに及ぶことができないことがある。

【0081】

いくつかの実施形態では、コールドログゾーンは、ホットログゾーンからログレコードをコピーすることによってポピュレートされる。係る実施形態では、LSNが現在の無条件ボリューム耐久性LSN（VDL）以下であるログレコードだけがコールドログゾーンにコピーされる資格があつてよい。ホットログゾーンからコールドログゾーンにログレコードを移動するとき、（多くのCLR等の）いくつかのログレコードは、それらがもはや必要ではないため、コピーされる必要がないことがある。さらにユーザーページのなんらかの追加の合体がこの点で実行されてよく、このことが必要とされるコピーの量を削減してよい。いくつかの実施形態では、いったん所与のホットゾーンログページが完全に書き込まれ、もはや最新のホットゾーンログページではなく、ホットゾーンログページ上のすべてのULRがコールドログゾーンに無事にコピーされると、ホットゾーンログページは解放され、再使用されてよい。

【0082】

いくつかの実施形態では、例えば記憶階層のSSDにもはや記憶される必要のないログレコード等、もはやサポートされていないログレコードによって占められているスペースを再利用するために、ガベージコレクションがコールドログゾーンで行われてよい。例えば、ログレコードは同じユーザーページに対する以後のAULRがあるときにサポートされなくなつてよく、ログレコードによって表されるユーザーページのバージョンはSSDでの保持に必要とされない。いくつかの実施形態では、ガベージコレクションプロセスは、2つ以上の隣接するログページをマージし、2つ以上の隣接するログページをそれらのページが置き換えているログページからの旧式ではないログレコードのすべてを含むより少ない新しいログページで置き換えることによってスペースを再利用してよい。新しいログページには、それらが置き換えているログページのフラッシュ番号よりも大きい新しいフラッシュ番号が割り当てられてよい。これらの新しいログページの書込みが完了した後、置き換えられたログページが空きページプールに加えられてよい。いくつかの実施形態では、あらゆるポインタを使用するログページの明示的な連鎖がないことがあることに留意されたい。代わりに、ログページのシーケンスはそれらのページに対するフラッシュ番号によって暗黙に決定されてよい。ログレコードの複数のコピーが検出されるたびに、最高のフラッシュ番号のログページに存在するログレコードが有効であると見なされてよく、他はもはやサポートされないと見なされてよい。

【0083】

いくつかの実施形態では、例えば、データゾーン（セクタ）の中で管理されるスペースの粒度がデータゾーン（記憶ページ）の外の粒度とは異なってよいため、なんらかのフラグメンテーションがあってよい。いくつかの実施形態では、このフラグメンテーションを管理するために、システムは各データページによって使用されるセクタの数を追跡調査してよく、ほぼ全データページから優先的に割り当ててよく、（データを新しい場所に、それが依然として関連している場合に移動することを必要としてよい）ほぼ空のデータページのガベージコレクションを優先的に行ってよい。セグメントに割り当てられるページが、いくつかの実施形態では3つのゾーンの間で転用されてよいことに留意されたい。例えば、セグメントに割り当てられていたページが解放されると、ページはある期間そのセグメントと関連付けられたままとなってよく、後にそのセグメントの3つのゾーンのいずれかで使用されてよい。あらゆるセクタのセクタヘッダは、セクタが属するゾーンを示してよい。いったんページ内のすべてのセクタが空くと、ページは、ゾーンに渡って共有される共通の空き記憶ページプールに返されてよい。この空き記憶ページの共有は、いくつかの実施形態では、フラグメンテーションを削減（または回避）してよい。

10

【0084】

いくつかの実施形態では、本明細書に説明される分散型データベース最適化ストレージシステムは、メモリ内に多様なデータ構造を維持してよい。例えば、セグメントに存在するユーザーページごとに、ユーザーページテーブルが、このユーザーページが「クリアされる」かどうか（つまり、このユーザーページがすべてのゼロを含んでいるかどうか）、該ページのためのコールドログゾーンからの最新のログレコードのLSN、及びページのホットログゾーンからのすべてのログレコードの場所のレイ/リストを示すビットを記憶してよい。ログレコードごとに、ユーザーページテーブルはセクタ番号、そのセクタ中のログレコードのオフセット、そのログページの中で読み取るセクタの数、（ログレコードが複数のログページに及ぶ場合）第2のログページのセクタ番号、及びそのログページの中で読み取るセクタの数を記憶してよい。いくつかの実施形態では、ユーザーページテーブルは、コールドログゾーンからのあらゆるログレコードのLSN、及び/またはAULRがコールドログゾーンにある場合、最新のAULRのペイロードのセクタ番号のレイを記憶してもよい。

20

【0085】

本明細書に説明される分散型データベース最適化ストレージシステムのいくつかの実施形態では、LSNインデックスはメモリに記憶されてよい。LSNインデックスは、コールドログゾーンの中のログページにLSNをマッピングしてよい。コールドログゾーンのログレコードがソートされていることを考えれば、それはログページあたり1つのエントリを含むためであってよい。ただし、いくつかの実施形態では、あらゆる旧式ではないLSNがインデックスに記憶され、対応するセクタ番号、オフセット、及びログレコードごとのセクタの数にマッピングされてよい。

30

【0086】

本明細書に説明される分散型データベース最適化ストレージシステムのいくつかの実施形態では、ログページテーブルはメモリに記憶されてよく、ログページテーブルはコールドログゾーンのガベージコレクションの間に使用されてよい。例えば、ログページテーブルはどのログレコードがもはやサポートされていないのか（例えば、どのログレコードのガベージコレクションを行うことができるのか）、及び各ログページでどれほど多くの空きスペースが使用できるのかを識別してよい。

40

【0087】

本明細書に説明されるストレージシステムでは、エクステントは、ボリュームを表すために他のエクステントと結合できる（連結できる、またはストライピングできるのかのどちらか）ストレージの高度に耐久性の単位を表す論理概念であってよい。各エクステントは、単一の保護グループでのメンバーシップによって耐久的にされてよい。エクステントは、LSN型の読取り/書込みインタフェースを、作成時に定義される固定サイズを有する隣接バイトサブレンジに提供してよい。エクステントに対する読取り/書込み動作は、

50

含む側の保護グループによって1つまたは複数の適切なセグメント読取り/書込み動作にマッピングされてよい。本明細書に使用されるように、用語「ボリュームエクステント」は、ボリュームの中のバイトの特有のサブレンジを表すために使用されるエクステントを指してよい。

【0088】

上述されたように、ボリュームは、それぞれが1つまたは複数のセグメントから構成される保護グループによって表される複数のエクステントから構成されてよい。いくつかの実施形態では、異なるエクステントに向けられるログレコードはインタリーブされたLSNを有してよい。ボリュームに対する変更が特定のLSNまで耐久的となるためには、そのLSNまでのすべてのログレコードが、それらが属しているエクステントに関わりなく耐久的である必要があつてよい。いくつかの実施形態では、クライアントは、まだ耐久的にされていない未決ログレコードを追跡調査してよく、いったん特定のLSNまでのすべてのULRが耐久的にされると、クライアントはボリュームの保護グループの内の1つにボリューム耐久性LSN(VDL)メッセージを送信してよい。VDLは、保護グループのすべての同期ミラーセグメントに書き込まれてよい。これは「無条件VDL」と呼ばれることがあり、それはセグメントで起こる書込み活動とともに多様なセグメントに(またはより詳細には、多様な保護グループに)周期的に持続されてよい。いくつかの実施形態では、無条件VDLはログセクタヘッダに記憶されてよい。

【0089】

多様な実施形態では、セグメントで実行されてよい動作は、(ホットログゾーンの末尾にDULRまたはAULRを書き込み、次いでユーザーページテーブルを更新することを含んでよい)クライアントから受信されたDULRまたはAULRを書き込むこと、(ユーザーページのデータセクタの位置を突き止め、あらゆる追加のDULRを適用する必要なしにデータセクタを返すことを含んでよい)コールドユーザーページを読み取ること、(ユーザーページの最も最新のAULRのデータセクタの位置を突き止めることを含み、ユーザーページに、それを返す前にあらゆる以後のDULRを適用してよい)ホットユーザーページを読み取ること、(適用された最後のDULRを置き換えるAULRを作成するためにユーザーページのDULRを合体させることを含んでよい)DULRをAULRで置き換えること、ログレコードを操作すること等を含んでよい。本明細書に説明されるように、合体は、ユーザーページのより最近のバージョンを作成するためにユーザーページの初期のバージョンにDULRを適用するプロセスである。(別のDULRが書き込まれるまで)合体の前に書き込まれたすべてのDULRは要求に応じて読み取られ、適用される必要はないことがあるため、ユーザーページを合体させることは読取りレーテンシを削減するのに役立ってよい。また、合体は、(ログレコードが存在することを必要とするスナップショットがないならば)古いAULR及びDULRをもはやサポートされなくすることによってストレージスペースを再利用するのに役立ってよい。いくつかの実施形態では、合体動作は、最も最新のAULRを場所を見つけ、DULRのいずれも省略することなく、あらゆる以後のDULRを順番に適用することを含んでよい。上述されたように、いくつかの実施形態では、合体はホットログゾーンの中で実行されないことがある。代わりに、合体はコールドログゾーンの中で実行されてよい。いくつかの実施形態では、合体は、ログレコードがホットログゾーンからコールドログゾーンにコピーされるにつれて実行されてもよい。

【0090】

いくつかの実施形態では、ユーザーページを合体させる決定は、(例えば、DULRチェーンの長さが合体動作の所定の閾値を超える場合、システム全体での方針、アプリケーション特有の方針、またはクライアントによって指定される方針に従って)、またはクライアントに読み取られているユーザーページごとに、ページの未決のDULRチェーンのサイズによってトリガされてよい。

【0091】

図7は、一実施形態に係る、データベースボリューム710の例の構成を示すブロック

10

20

30

40

50

図である。この例では、(アドレス範囲715 aから715 eとして示される)多様なアドレス範囲715のそれぞれに対応するデータが(セグメント745 aから745 nとして示される)異なるセグメント745として記憶される。すなわち、多様なアドレス範囲715のそれぞれに対応するデータは(エクステント725 aから725 b、及びエクステント735 aから735 hとして示される)異なるエクステントに編成されてよく、これらのエクステントの多様なエクステントが、(ストライプセット720 a及びストライプセット720 bとして示されるもの等の)ストライピングを行って、または行わないで(730 aから730 fとして示される)異なる保護グループ730に含まれてよい。この例では、保護グループ1はイレイジャーコーディングの使用を示す。この例では、保護グループ2及び3、並びに保護グループ6及び7は互いのミラーリングされたデータセットを表す。一方、保護グループ4は単一インスタンス(非冗長)データセットを表す。この例では、保護グループ8は、他の保護グループを結合する複数階層保護グループを表す(例えば、これは複数領域保護グループを表してよい)。この例では、ストライプセット1(720 a)及びストライプセット2(720 b)は、いくつかの実施形態で、エクステント(例えば、エクステント725 a及び725 b)がどのようにしてボリュームの中にストライピングされてよいのかを示す。

10

【0092】

すなわち、この例では、保護グループ1(730 a)は、それぞれ範囲1から3(715 aから715 c)のデータを含むエクステントaからc(735 aから735 c)を含み、これらのエクステントはセグメント1から4(745 aから745 d)にマッピングされる。保護グループ2(730 b)は、範囲4(715 d)からストライピングされたデータを含むエクステントd(735 d)を含み、このエクステントはセグメント5から7(745 eから745 g)にマッピングされる。同様に、保護グループ3(730 c)は、範囲4(715 d)からストライピングされたデータを含むエクステントe(735 e)を含み、セグメント8から9(745 hから745 i)にマッピングされ、保護グループ4(730 d)は、範囲4(715 d)からストライピングされたデータを含むエクステントf(735 f)を含み、セグメント10(745 j)にマッピングされる。この例では、保護グループ6(730 e)は、範囲5(715 e)からストライピングされたデータを含むエクステントg(735 g)を含み、セグメント11から12(745 kから745 l)にマッピングされ、保護グループ7(730 f)は、やはり範囲5(715 e)からストライピングされたデータを含むエクステントh(735 h)を含み、セグメント13-14(745 mから745 n)にマッピングされる。

20

30

【0093】

ここで図8を参照すると、多様な実施形態では、上述されたように、データベースシステムは、ストレージノードのデータページの中に記憶されているデータに対する多様なアクセス要求(例えば、書込み要求)に応じてリドゥログレコードを生成し、リドゥログレコードが生成されたそれぞれのデータページを記憶するストレージノードにリドゥログレコードを送信するように構成されてよい。ストレージノードは、特定のデータページのための合体イベントを検出し、それに応じて特定のデータページのために合体動作を実行してよい。典型的なデータベースシステムは、一方、周期的な間隔で記憶されるデータに適用される生成されたリドゥログのすべてをフラッシュし、このようにしてデータベースによって実行されるアクセス要求及び他のタスクの処理を中断させるシステム全体のチェックポイントを適用してよい。

40

【0094】

図8の方法は、分散型データベース最適化ストレージシステム410(例えば、ストレージシステムサーバノード(複数の場合がある)430、440、450等)のログ構造化ストレージシステムの多様な構成要素によって実行されているとして説明されてよいが、方法はいくつかの場合、いずれの特定の構成要素によっても実行される必要はない。例えば、いくつかの場合、図8の方法は、いくつかの実施形態に従ってなんらかの他の構成要素またはコンピュータシステムによって実行されてよい。また、いくつかの場合、デー

50

データベースシステム400の構成要素は、図4の例に示されるのとは異なって組み合わされてよい、または存在してよい。多様な実施形態では、図8の方法は分散型データベース最適化ストレージシステムの1台または複数のコンピュータによって実行されてよく、その内の1つは図10のコンピュータシステムとして示される。図8の方法は、システム全体のチェックポイント回避のための方法の1つの例の実装として示される。他の実装では、図8の方法は追加のブロック、または図示されるよりも少ないブロックを含んでよい。

【0095】

810に示されるように、データベースのために記憶される特定のデータページにリンクされるリドゥログレコードが維持されてよい。これらのリドゥログレコード(上述されたように、URLと呼ばれることがある)は、ユーザーデータに対する変更を記述してよい。リドゥログレコードは、データページ等のユーザーデータの特定の部分にリンクされてよい。例えば、いくつかの実施形態では、リドゥログレコードは、特定のデータページに最終的にリンクされるリドゥログレコードの連鎖を形成し、各リドゥログレコードはデータページのための以前に受信されたリドゥログレコードを指す。この例を使用すると、3つのリドゥログレコードが特定のデータページにリンクされる場合には、最も最近に受信されたリドゥログレコードは次に最も最近に受信されたリドゥログレコードを指し、次に最も最近に受信されたリドゥログレコードは同様に3番目に最も最近に受信されたリドゥログレコードを指し、3番目に最も最近に受信されたリドゥログレコードはデータページの最も最近に保存された状態を指す。前のリドゥログレコードに対する各ポイントによって示されるリドゥログレコードの論理的な順序付けが、係るリドゥログレコードが係る順序で物理的に記憶されることを暗示しないことに留意されたい。図6に関して上述されたように、これらのリドゥログレコードは、いくつかの実施形態では、ユーザーデータの他の部分にリンクされた他のリドゥログレコードとインタリーブされてよい。したがって、前の例は制限的となることを目的としていない。

【0096】

多様な実施形態では、リドゥログレコードは、ストレージノード430、440、450等のストレージノードにデータが記憶されてよい、1つまたは複数のデータベースを管理してよいデータベースエンジンヘッドノード420等のデータベースシステムから受信されてよい。しかしながら、少なくともいくつかの実施形態では、ストレージノードは、ストレージノードがデータを記憶するための1つまたは複数の追加のデータベースシステムまたはノードからリドゥログレコードを受信してよい。これらの他のデータベースシステムまたはノードは、ストレージノードにそのそれぞれのデータベースのために記憶されているデータ特定の部分にリンクされたリドゥログレコードを送信してもよい。

【0097】

いくつかの実施形態では、受信されたリドゥログレコードが次いで記憶されてよい。図6は、係るリドゥログレコードがどのようにして受信され、処理され、ストレージノードに記憶されてよいのかの多様な実施形態を説明する。多様な形式のメタデータが、データページ等の特定の部分データにリンクされるリドゥログレコードの数つまりカウント等の記憶されているリドゥログレコードのために維持されてよい。例えば、上記に示された例のように、3つのリドゥログレコードが特定のデータページにリンクされる場合、次いで特定のデータページのリドゥログレコードカウントは3で維持されてよい。多様な他のログレコードに対するポイントまたはデータページの最も最近に保存された状態に対するポイント等、サイズまたは物理的な場所、及びリドゥログレコードがリンクされるデータの部分等のリドゥログレコードに関する他のメタデータが維持されてよい。

【0098】

記憶されているリドゥログレコードのために維持されるメタデータに対する更新は、リドゥログレコード自体に対する変更、それらがリンクされる特定のデータページに対する変更、またはリドゥログレコードを活用することによって、またはリドゥログレコードに関して実行される動作もしくは他の方法もしくは技法に応じて行われてよい。例えば、830で示されるように、合体動作が実行され、データページの現在の状態を生成するため

10

20

30

40

50

に特定のデータページにリンクされる1つまたは複数のリドゥログレコードを適用する場合、次いでリドゥログレコードカウントは特定のデータページに対するリドゥログレコードカウントからそれらの適用されたリドゥログレコードを削除するために更新されてよい。

【0099】

多様な実施形態では、特定のデータページのための合体イベントは、特定のデータページにリンクされる1つまたは複数のリドゥログレコードに少なくとも部分的に基づいて、820で示されるように検出されてよい。検出された合体イベントは、合体動作が特定のデータページに対して実行されてよいことを示してよい。少なくともいくつかの実施形態では、特定のデータページのための合体イベントを検出することは、他のデータページについて検出された合体イベントとは関係なく、または他のデータページについて検出された合体イベントを考慮せずに発生してよい。特定のデータページが、多くのリドゥログレコードが受信される「ホット」データページであってよいシナリオを考える。リドゥログレコードはめったに他のデータページのために受信されることはない。合体イベントを検出することは、合体閾値を超えるそれぞれのデータページにリンクされるリドゥログレコードの数に基づいてよく、したがって、このシナリオでは、合体イベントは他のデータページについてより、特定の「ホット」データページについてより頻繁に検出されてよい。

【0100】

合体イベントを検出することは、バックグラウンドプロセスとして実行してよいストレージノード監視構成要素またはプロセスの一部として実行されてよく、読取り要求、書込み要求、及び他のアクセス要求を処理するフォアグラウンドプロセスは、合体イベントの検出の前に（または合体イベントの検出を遅延させて）実行されてよい。合体イベントの検出は、ストレージノードの作業負荷が作業負荷閾値未満であるとき等、周期的な間隔または非周期的な間隔で発生してよい。

【0101】

特定のデータページにリンクされたリドゥログレコードに少なくとも部分的に基づいて合体イベントを検出するための多様な方法及び技法が実装されてよい。例えば、少なくともいくつかの実施形態では、合体閾値は合体イベントを検出するために活用されてよい。合体閾値は、合体イベントが検出される前に特定のデータページにリンクされてよいリドゥログレコードの数を定義してよい。例えば、特定のデータページが、10リドゥログレコードの合体閾値を超える11リドゥログレコードを有する場合、次いで合体イベントが検出されてよい。異なる合体閾値は異なるデータページに活用されてよい。例えば、データページにリンクされた頻繁なリドゥログレコードを受信する「ホット」データページのシナリオを再度考える。リドゥログレコードをあまり頻繁に受信しないデータページよりも高い合体閾値は、「ホット」データページに活用され、このようにして「ホット」データページに対して実行される合体動作の数を削減してよい。代わりに、いくつかの実施形態では、同じ合体閾値または類似する合体閾値が活用されてよい。合体閾値は、多様な他の技法または構成部品と結合されてもよい。例えば、他の構成部品を使用していつ合体閾値が超えられる可能性があるのかを計算し、タイマまたは他の構成要素を設定して、合体イベント検出を実行するバックグラウンドモジュールまたは他のプロセスに対し、特定のデータページのリドゥログレコードカウントが調べられるべきであることを示すこと。

【0102】

少なくともいくつかの実施形態では、特定のデータページに対する（またはデータページの特定のセットに対する）合体閾値が決定されてよい。例えば、いくつかの実施形態では、合体閾値はユーザー定義の合体閾値に従って決定されてよい。ユーザー定義の合体閾値は、要求され、決定され、もしくはデータベースエンジンヘッドノード420等のデータベースシステムからストレージノードに対して示される合体閾値であってよい、またはデータベースシステムのクライアントは合体イベントを検出するために使用される合体閾値を与えてよい。いくつかの実施形態では、合体閾値はストレージノードの作業負荷または性能に基づいて決定されてよい。例えば、いくつかの実施形態では、作業負荷測度また

10

20

30

40

50

は性能測度が、合体動作を実行するための能力が低いことを示す場合、次いで合体閾値は、検出される合体イベントの数がストレージノードによってその現在の作業負荷で処理され得るように増加されてよい。いくつかの実施形態では、リドゥログレコードが特定のデータページについて受信されるレートつまり頻度が計算され、合体閾値を決定するために使用されてよい。少なくともいくつかの実施形態では、リドゥログレコードのサイズ、物理記憶でのリドゥログレコードの場所、リドゥログレコードを記憶するために利用可能なスペース、及び/または合体動作がデータページの以前に記憶されたバージョンにリドゥログレコードを適用するために実行されてよい時刻等の多様な他の特徴が合体閾値を決定するために使用されてよい。

【0103】

特定のデータページに対する合体イベントを検出することに応じて、特定のデータページにリンクされる1つまたは複数のリドゥログレコードが、830で示されるように特定のデータページをその現在の状態で生成するために特定のデータの以前に記憶されたバージョンに適用されてよい。少なくともいくつかの実施形態では、特定のデータページにリンクされるリドゥログレコードを適用することは合体動作の一部として実行される。上述されたような合体動作つまり合体は、ユーザーページのより最近のバージョンを作成するためにDULR等のリドゥログレコードをユーザーページの初期のバージョンに適用してよい。いくつかの実施形態では、合体動作は、最も最近のAULR（例えば、データページの以前に記憶されたバージョン）の位置を突き止め、DULRのいずれも省略することなくあらゆる以後のDULRを順に適用することを含んでよい。例えば、3つのDULRが受信され、AULRにリンクされている場合、最初に受信されたDULRがAULRに適用される（このようにして、以前に記憶されたデータページを基準にして最初に受信された変更を適用する）。次いで、次に受信されたDULRが適用され、最後に最も最近のDULRが適用され、記憶ノードでのDULRの受信に基づいて決定される順にDULRを適用する。いくつかの実施形態では、新しいAULRは特定のデータページの現在の状態として生成される。リドゥログレコードカウント等の上述されたメタデータは、リドゥログレコードの適用を反映し、リドゥログレコードカウントに関して、その数をカウントから削除するために更新されてよい。

【0104】

少なくともいくつかの実施形態では、遅延は、820で示される合体イベントの検出と830で示されるリドゥログレコードの適用との間で発生してよい、または実行されてよい。例えば、該検出及び該適用を実行するストレージノードの作業負荷が、リドゥログレコードを適用することの実行と、合体イベントの検出との間の遅延を決定してよい。同様に、合体イベントの検出に応えるリドゥログレコードの適用はバックグラウンドプロセスの一部として実行されてよい、すなわち削減される、つまり多様なアクセス要求（例えば、読取り要求または書込み要求）の処理等、フォアグラウンドプロセスを実行しないときにだけ実行される。遅延した合体動作またはデータページのためのリドゥログレコードの適用は、データページがいつリドゥログレコードを適用させるべきであるのかの順序、シーケンス、またはタイミングを決定する、先入先出し（FIFO）待ち行列または優先順位待ち行列等のデータ構造に入れられてよい。例えば、上述されたシナリオでのように、「ホット」データページが検出された合体イベントを有する場合、別のデータページの代わりに「ホット」データページに対するリドゥログレコードの適用を実行する方がより効率的であることがある。バックグラウンドプロセスとしてリドゥログレコードの適用を遅延するまたは実行する結果として、合体イベントが検出されたデータページにリンクされる1つまたは複数の追加のリドゥログレコードが受信されてよい。少なくともいくつかの実施形態では、これらの追加のリドゥログレコードは、他のリドゥログレコードがデータページの以前に記憶されたバージョンに適用されるときに適用されてよい。

【0105】

図4に示されるように、複数のストレージノード430、440、450他は、分散型ストレージサービスの一部として実装されてよい。図8に関して上述された多様な方法及

10

20

30

40

50

び技法は、これらの複数のストレージノードによって互いと無関係に実行されてよい。各ストレージノードは、合体イベントを検出すること、及びそれに応じて同時にまたは互いと異なるときに1つまたは複数のリドゥログレコードを適用することを実行するだけでなく、異なる合体閾値または同じ合体閾値を決定してもよい。

【0106】

ここで、いくつかの実施形態に従って、分散型データベースシステムのための高速クラッシュ回復を実行するための方法を明示する一連の図を示す図9Aを参照する。典型的なデータベースシステムにおけるクラッシュ回復は達成が困難なプロセスである。これらの典型的なシステムでは、データベースシステム故障からの回復時、データベースのクリーンなバージョンが得られ、次いでディスクに記憶されていないトランザクションからのリドゥログレコードのすべてが、データベースをデータベースシステム故障の前のその現在の状態に復元するためにリプレイされなければならない、データベースにアクセスできるようになる前に多大な復元時間を生じさせる。図9Aは、一方、クラッシュ回復を実行するためのより高速且つより効率的な技法を提供してよい分散型データベースシステム用の高速クラッシュ回復の説明を提供する。

10

【0107】

シーン992で、図2に関して上述されたデータベースクライアント250等のデータベースクライアント906は、図2に上述されたネットワーク260上で、データベースを実装する、図4に関して上述されたデータベースヘッドノード430等のデータベースヘッドノード902と通信する。ストレージノード908は、データベースヘッドノード902によって実装されるデータベースのためのログ構造化データストレージを実装する1つまたは複数のストレージノードであってよい。多様なアクセス要求が受信され、その後ストレージノード908からアクセスされたデータを取り出すと、データベースヘッドノード902によってサービスを提供されてよい。図8に関して上述されたもの等のリドゥログレコードが生成され、ユーザーデータを送信する代わりにストレージノード908に送信されてよい。リドゥログレコードはストレージノード908で維持されてよい。少なくともいくつかの実施形態では、合体動作は、図8に関して上述されたように等、合体イベントの検出に応じて実行されてよい。

20

【0108】

シーン994は、データベースヘッドノード902の故障を示す。データベースヘッドノード故障は、電源喪失、利用可能なメモリなし、システム障害等の、データベースヘッドノードが機能を続行できないようにさせる任意のタイプのシステム故障であることがある。データベースクライアント906とデータベースヘッドノード902との間の通信は、図に示されるように送信または受信されないことがある。したがって、データベースに対するアクセスは提供され得ない。同様に、ストレージノード908とデータベースヘッドノード902との間の通信が送信または受信されないことがあり、したがってデータベースのために記憶されているデータに対する要求が処理されないことがある。

30

【0109】

シーン996では、回復動作が示されてよい。同じシステムハードウェアで再起動されたヘッドノードアプリケーションプログラムのバージョン、または異なるハードウェアで起動されたヘッドノードの別のインスタンスであってよい新しいデータベースヘッドノード904がオンラインにされてよい。ストレージノード908との接続は、示されるように、データベースヘッドノード904によって確立されてよい。シーン998は、ストレージノード908との接続の確立時、データベースヘッドノード902で実装されたのと同じデータベースが、新しいデータベースヘッドノード904でのアクセスのために利用可能にされてよいことを示す。読取り要求または書込み要求等のアクセス要求は、ネットワーク260を介してデータベースクライアント906から新しいデータベースヘッドノード904に送信されてよい。リドゥログレコードはすでに、アクセス要求にサービスを提供するために新しいデータベースヘッドノード908にデータベースのために記憶されているデータのカルレントバージョンを提供してよいストレージノード908に送信されて

40

50

いたので、新しいデータベースヘッドノード904は、データベースヘッドノード故障の前にデータの現在の状態を入手するためにこれらのリドゥログレコードをリプレイする必要がないことがある。ストレージノード908は、特定のデータに対する要求が受信されるとき特定のデータの以前に記憶されていたバージョンにリドゥログレコードを適用してよい。代わりに、特定のデータの現在の状態は、図8に関して上述されたように合体イベントが検出されるとき等、あらゆるリドゥログレコードがすでに適用されている特定のデータに向けられた状態でストレージノードにすでに記憶されていてよい。

【0110】

図9Bは、いくつかの実施形態に係る、分散型データベースシステムのための高速クラッシュ回復を実行する方法を示す流れ図である。多様な実施形態では、データベースヘッドノード故障が発生することがある。このヘッドノード故障はあらゆる通信、修正、または故障したデータベースヘッドノードによって実装され、管理されるデータベースへの他の形のアクセスを妨げることがある。例えば、図2に説明されるデータベースクライアント250等のデータベースシステムクライアントは、故障したデータベースヘッドノードに読取り要求または書込み要求を送信できないことがある。データベースヘッドノードの故障は、例えば図2に上述されたウェブサービスプラットフォーム200、または何らかの他のシステムもしくは構成要素によって検出されてよい。ヘッドノードの故障に応じて、再起動されたデータベースヘッドノードまたは新しいデータベースヘッドノード（例えば、以前に故障したヘッドノードと同じまたは異なるハードウェア上でホストされる新しいデータベースヘッドノード仮想インスタンス）が、回復動作を実行するように命令されてよい。いくつかの実施形態では、この回復動作はこれらの要素に制限されていないが、回復動作は図9Bに示される多様な要素を含んでよい。

【0111】

データベースヘッドノード故障からの回復は、910に示されるように発生してよい。回復は実行され、さまざまな方法で完了していると決定されてよい。例えば、データベースヘッドノードアプリケーションは、多様なテストを実行すること、多様な装置を有効にすること等、実行するために準備するとき多様な状態を有することがある。このプロセスの一部として、ノード故障からの回復の完了を示してよいデータベースヘッドノードについて準備完了した状態が決定されてよい。910に示されるように、データベースノード故障からの回復時、920に示されるように、データベースのためにデータを記憶する1台または複数のストレージノードとの接続が確立されてよい。

【0112】

図9A及び上記の多様な他の図に関して上述されたように、データベースは、図3及び図4に説明されるデータベースヘッドノード320または440等のデータベースヘッドノードによって実装され、管理されてよい。上述された読取り要求または書込み要求等のデータベースアクセス要求を実装することの一部として、データベースヘッドノードで処理されてよい。少なくともいくつかの実施形態では、データベースに対する変更を反映するリドゥログレコードは、ストレージノードに記憶されるデータに対する変更を反映する、図4で上述されたストレージノード450等の1つまたは複数のストレージノードに送信される。特定のデータページまたはデータの他の部分等の、変更されるデータを記憶するストレージノードは、変更される、データページ等のデータの部分にリンクされるリドゥログレコードを受信してよい。これらのリドゥログレコードは、データページのカレントバージョンに対する要求に応じて、または合体イベントの検出に応じて等、なんらかの他の時に、データページ等のデータの部分の以前に記憶されていたバージョンに適用されてよい（例えば、合体動作）。データベースのためのリドゥログレコードは、上述された多様な方法で、データベースヘッドノードで実装されるデータベースのために維持されるので、いくつかの実施形態では、ストレージノードはデータベースヘッドノードに、データベースヘッドノード故障の時刻まで最新であると保証されるデータの現在の状態を送信してよい。

【0113】

10

20

30

40

50

接続の確立先のストレージノードが識別されてよい。例えば、図4で上述されたクライアント側ストレージサービスドライバ425は、どのストレージノードがデータベースのためにデータを記憶するのか、及びデータベースのどの部分がストレージノードに記憶されるのかを示す情報を維持してよい。接続要求、または何らかの他の通信メッセージは、図4に関して上述された多様な通信方法の1つを使用して送信されてよい。同様に、肯定応答、及びストレージノード及び/またはデータベースヘッドノードのステータスについての他の情報が交換されてよい。

【0114】

920に示されるように、1つまたは複数のストレージノードとの接続の確立時、データベースは、930に示されるように、アクセスのために利用可能にされてよい。いくつかの実施形態では、アクセスは1つまたは複数のアクセス要求（例えば、読取り要求、書込み要求）に提供されてよい。データベースの可用性の表示が生成され、クライアントに送信されてよい。例えば、データベースがアクセスに利用可能である旨のメッセージがデータベースクライアントに送信されてよい。係るメッセージは、図2に説明されるウェブサービスプラットフォーム200、またはなんらかの他の通信プラットフォームもしくは装置を介して送信されてよい。上述されたように、典型的なデータベースシステムでは、リドゥログレコードのリプレイは、データベースを利用可能にする前に実行されなければならない。しかし、少なくともいくつかの実施形態では、データベースはリドゥログレコードをリプレイせずに利用可能にされてよい。リドゥログレコードとともに使用されるとき用語「リプレイ」が概してデータの以前に記憶されていたバージョンに対して1つまたは複数のリドゥログレコードを適用することを意味することに留意されたい。

【0115】

少なくともいくつかの実施形態では、ストレージノードは、データベースヘッドノード故障を検出できてよい、またはそれ以外の場合データベースヘッドノード故障を認識させられてよい。データベースヘッドノード故障の検出に応じて、ストレージノードは、ストレージノードで受信されたリドゥログレコードに対する切り詰め演算を実行してよい。切り詰め演算は、データベースヘッドノードの故障の前に完了しなかったシステムトランザクションの一部であるリドゥログレコードを決定してよい、または識別してよい。これらの識別されたリドゥログレコードは、それらがリンクされているデータページにそれらが適用され得ないように、削除されてよい、またはそれ以外の場合、マークされてよい、移動されてよい、もしくは識別されてよい。例えば、記憶ページが特定のデータページのために5リドゥログレコードを維持し、最も最近の3リドゥログレコードが、データベースヘッドノード故障の前に完了しなかったシステムトランザクションの一部である場合、次いでストレージノードは、2つの最も古いリドゥログレコードだけを適用することによってデータページの現在の状態を生成するとき最も最近の3リドゥログレコードを無視してよい。少なくともいくつかの実施形態では、切り詰め演算は、回復されたデータベースヘッドノードと接続を確立できるようになる前に、影響を受けたリドゥログレコードがあるストレージノードに対して実行されてよい。データベースエンジンヘッドノードは、いくつかの実施形態では、データベースヘッドノードの故障前に完了しなかったシステムトランザクションの一部であるリドゥログレコードを同様に決定し、または識別し、これらの識別されたリドゥログレコードが、それらがリンクされているデータページにそれらが適用され得ないように削除されてよい、またはそれ以外の場合マークされてよい、移動されてよい、または識別されてよい旨の通知をストレージノードに送信するように構成されてよい。例えば、図3に関して上述されたクライアント側ストレージサービスドライバ325等のクライアント側ストレージサービスドライバは、上述された技法を実行してよい。切り詰め演算を説明するこれらの技法は、いくつかの実施形態では、バックグラウンドプロセスの一部として実行されてよい。

【0116】

少なくともいくつかの実施形態では、システムトランザクションは、ユーザートランザクションを実行する、または実装するための動作または他の形の1つもしくは複数のタス

10

20

30

40

50

クであってよい。ユーザートランザクションは、受信されたアクセス要求から多様なタスクまたは動作を実行するために複数のシステムトランザクションを含んでよい。例えば、データベースに対する挿入命令が受信されてよい。ユーザートランザクションとして、この挿入命令は、挿入を実行するために、例えばb-ツリー等のデータベースデータ構造に作用する等、挿入を実行するための複数のシステムトランザクションを含んでよい。少なくともいくつかの実施形態では、不完全なユーザートランザクションは、ユーザートランザクションであり、該ユーザートランザクションに含まれるシステムトランザクションのすべてが完了していない（または耐久的にされていない）可能性がある。同様に、システムトランザクションは不完全なことがある。ユーザートランザクション及びシステムトランザクションの一部としてデータベースのために記憶されたデータに対して行われた変更を反映するリドゥログレコードは、いくつかの実施形態では、特定のユーザートランザクション及び/またはシステムトランザクションで識別されてよい。

10

【0117】

図9Cは、いくつかの実施形態に係る、回復されたデータベースでアクセス要求を処理するための方法を示す流れ図である。上述されたように、少なくともいくつかの実施形態では、アクセスのためにデータベースを利用できるようにしたデータベースヘッドノードで、アクセス要求が受信されてよい。アクセス要求は、読取り要求、書込み要求、またはデータベースのために記憶されているデータ入手するもしくは修正するための任意の他の要求であってよい。図9Cが示すように、アクセス要求は、940で示されるようにデータベースに対して受信されてよい。それに応じて、950示されるように1つまたは複数のストレージノードからの1つまたは複数のデータページに対する要求が行われてよい（クライアントからのアクセス要求及びデータベースヘッドノードからのデータ要求の両方とも、上記図5に関してより詳細に扱われている）。要求された1つまたは複数のデータページの現在の状態は、960に示されるように、ストレージノードから受信されてよい。上述されたように、この現在の状態は、データページの以前に記憶されたバージョンまで以前に受信されたリドゥログレコードをリプレイする、もしくはデータページの以前に記憶されたバージョンに以前に受信されたリドゥログレコードを適用することによって、または現在の状態であるデータページの以前に記憶されたバージョンを返すことによって生成されてよい。多様な実施形態では、各データページまたは要求されたデータの一部は、（例えばゆったりと）データに対する要求を受信することに応じて、その現在の状態

20

30

【0118】

少なくともいくつかの実施形態では、アンドゥログレコードは、データベースヘッドノードで維持されてよい。上述されたようなアンドゥログレコードは、不完全なユーザートランザクションが発生した場合に等、データに対して行われた変更をアンドゥするためにデータベースのために記憶されるデータに適用される変更を記録してよい。ユーザートランザクションは、（複数のシステムトランザクション等の）データベースのために記憶されるデータに対する複数の変更を含み、1つまたは複数のリドゥログレコード及び1つまたは複数のアンドゥログレコードを生成してよい。ユーザートランザクションは、ユーザートランザクションの変更のすべてがコミットされなかった（例えば、耐久的にされなかった）ときに不完全であることがある。図3に関して上述されたトランザクションログ340等のトランザクションテーブルは、どのユーザートランザクション、及びストレージノードに記憶されているデータのその関連付けられた部分が、データベースヘッドノード故障前にコミットされず、したがって不完全であるのかを示すために実装されてよい。970で示されるように、受信されたデータページがトランザクションテーブルによって示される等、不完全なユーザートランザクションによって影響を及ぼされるかどうかに関して決定が下されてよい。はいである場合、肯定の出口が示すように、次いでアンドゥログレコードの1つまたは複数が、不完全なトランザクションによって行われた変更をアンドゥして、972に示すように、データページの新しい現在の状態を生成するためにデータページに適用されてよい。アンドゥログレコードが適用された、つまり不完全なユーザ

40

50

トランザクションによってデータページが影響を及ぼされなかった後、次いでデータページの現在の状態が、980で示されるようにアクセス要求にサービスを提供するために提供されてよい。

【0119】

少なくともいくつかの実施形態で、トランザクションテーブルに基づいて、不完全なユーザートランザクションによって影響を受けたデータの部分を決定する、または識別するバックグラウンドプロセスが実行されてよい。不完全なユーザートランザクションによって影響を受けた、データページ等のデータの現在の状態に対する要求が送受されてよい。アンドウログレコードは、次いで、不完全なユーザートランザクションによってこれらのデータページに向けられた変更をアンドウするために適用されてよい。多様な実施形態では、データベースキャッシュが、アンドウログレコードが適用された後にこれらのデータページで更新されてよい。

10

【0120】

少なくともいくつかの実施形態では、以前に記録されたスナップショットが、データベースの状態を初期の状態に復元するために使用されてよい。例えば、アクセスのためにデータベースを利用可能にする前に、要求は、データベースのためのデータを以前に記録されたスナップショットに対応する状態に復元するためにストレージノードに送信されてよい。スナップショットは、以前に受信されたリドゥログレコードを、記録されたスナップショット点（例えば、タイムスタンプまたはマーカ）までリプレイできるようにする、ストレージノードに記憶されるリドゥログのためのタイムスタンプまたは他のマーカまたはインジケータを識別することによって記録されてよく、該復元は複数のリドゥログの1つまたは複数データを以前のバージョンに適用することを含む。ストレージノードにスナップショットを実装する追加説明が上記に示される。

20

【0121】

図9Bから図9Cの方法及び技法は、データベースエンジンヘッドノード420等のデータベースシステムの多様な構成要素によって実行されるとして説明されてよいが、方法は、いくつかの場合、いずれの特定の構成要素によっても実行される必要はない。例えば、いくつかの場合、図9Bから図9Cは、いくつかの実施形態に従って、なんらかの他の構成要素またはコンピュータシステムによって実行されてよい。また、いくつかの場合、データベースシステム400の構成部品は、データベースシステム400の構成要素は、図4の例に示されるのとは異なって組み合されてよい、または存在してよい。多様な実施形態では、図9Bから図9Cの方法は分散型データベースシステムの1台または複数のコンピュータによって実行されてよく、その内の1つは図10のコンピュータシステムとして示される。図9Bから図9Cの方法は、分散型データベースシステムの高速度クラッシュ回復のための方法の例の実装として示される。他の実装では、図9Bから図9Cの方法は追加のブロック、または図示されるよりも少ないブロックを含んでよい。

30

【0122】

本明細書に説明される方法は、多様な実施形態では、ハードウェア及びソフトウェアの任意の組合せによって実装されてよい。例えば、一実施形態では、方法は、プロセッサに結合されたコンピュータ可読記憶媒体に記憶されるプログラム命令を実行する1台または複数のプロセッサを含むコンピュータシステム（例えば、図10のコンピュータシステム）によって実装されてよい。プログラム命令は、本明細書に説明される機能性（例えば、本明細書に説明されるデータベースサービス/システム及び/またはストレージサービス/システムを実装する多様なサーバ及び他の構成要素の機能性）を実装するように構成されてよい。

40

【0123】

図10は、多様な実施形態に従って、本明細書に説明されるデータベースシステムの少なくとも一部を実装するように構成されるコンピュータシステムを示すブロック図である。例えば、コンピュータシステム1000は、異なる実施形態で、データベース階層のデータベースエンジンヘッドノード、またはデータベース階層のクライアントの代わりにデ

50

ータベース及び関連付けられたメタデータを記憶する別個の分散型データベース最適化ストレージシステムの複数のストレージノードの内の1つを実装するように構成されてよい。コンピュータシステム1000は、パーソナルコンピュータシステム、デスクトップコンピュータ、ラップトップコンピュータまたはノートパソコン、メインフレームコンピュータシステム、ハンドヘルドコンピュータ、ワークステーション、ネットワークコンピュータ、消費者装置、アプリケーションサーバ、ストレージデバイス、電話、携帯電話、または一般的に任意のタイプのコンピューティング装置を含むが、これに限定されることがない多様なタイプの装置のいずれかであってよい。

【0124】

コンピュータシステム1000は、入出力(I/O)インタフェース1030を介してシステムメモリ1020に結合される(いずれかが、単一スレッドまたはマルチスレッドであってよい複数のコアを含んでよい)1台または複数のプロセッサ1010を含む。コンピュータシステム1000は、I/Oインタフェース1030に結合されるネットワークインタフェース1040をさらに含む。多様な実施形態では、コンピュータシステム1000は、1台のプロセッサ1010を含んだユニプロセッサシステム、または数台のプロセッサ1010(例えば、2、4、8、または別の適切な数)を含んだマルチプロセッサシステムであってよい。プロセッサ1010は、命令を実行できる任意の適切なプロセッサであってよい。例えば、多様な実施形態では、プロセッサ1010は、x86、PowerPC、SPARC、もしくはMIPS ISA等のさまざまな命令セットアーキテクチャ(ISA)または任意の他の適切なISAのいずれかを実装する汎用プロセッサまたは組み込みプロセッサであってよい。マルチプロセッサシステムでは、プロセッサ1010のそれぞれが、一般に同じISAを実装してよいが、必ずしも同じISAを実装しないこともある。コンピュータシステム1000は、通信ネットワーク(例えば、インターネット、LAN等)上で他のシステム及び/または構成要素と通信するための1台または複数のネットワーク通信装置(例えば、ネットワークインタフェース1040)も含む。例えば、システム1000で実行中のクライアントアプリケーションは、単一のサーバ上、または本明細書で説明されるデータベースシステムの構成要素の内の1つまたは複数の実装するサーバのクラスタ上で実行中のサーバアプリケーションと通信するためにネットワークインタフェース1040を使用してよい。別の例では、コンピュータシステム1000上で実行中のサーバアプリケーションのインスタンスは、他のコンピュータシステム(例えば、コンピュータシステム1090)の上で実装されてよいサーバアプリケーション(または別のサーバアプリケーション)の他のインスタンスと通信するために、ネットワークインタフェース1040を使用してよい。

【0125】

示されている実施形態では、コンピュータシステム1000は、1台または複数の永続ストレージデバイス1060及び/または1台または複数のI/Oデバイス1080も含む。多様な実施形態では、永続ストレージデバイス1060は、ディスクドライブ、テープドライブ、ソリッドステートメモリ、他の大容量記憶装置、または任意の他の永続ストレージデバイスに相当してよい。コンピュータシステム1000(または、コンピュータシステム1000上で動作する分散アプリケーションもしくはオペレーティングシステム)は、所望されるように、命令及び/またはデータを永続ストレージデバイス1060に記憶してよく、必要に応じて記憶されている命令及び/またはデータを取出してよい。例えば、いくつかの実施形態では、コンピュータシステム1000は、ストレージシステムサーバノードをホストしてよく、永続記憶装置1060はそのサーバノードにアタッチされるSSDを含んでよい。

【0126】

コンピュータシステム1000は、プロセッサ(複数の場合がある)1010によってアクセス可能な命令及びデータを記憶するように構成される1つまたは複数のシステムメモリ1020を含む。多様な実施形態では、システムメモリ1020は、任意の適切なメモリ技術(例えば、キャッシュ、スタティックランダムアクセスメモリ(SRAM)、D

10

20

30

40

50

RAM、RDRAM、EDO RAM、DDR 10 RAM、同期ダイナミックRAM (SDRAM)、Rambus RAM、EEPROM、不揮発性/フラッシュタイプメモリ、または任意の他のタイプのメモリの内の1つまたは複数)を使用して実装されてよい。システムメモリ1020は、本明細書に説明される方法及び技法を実装するためにプロセッサ(複数の場合がある)1010によって実行可能であるプログラム命令1025を含んでよい。多様な実施形態では、プログラム命令1025は、プラットフォームネイティブバイナリ、Java(商標)バイトコード等の任意のインタープリター型言語で、またはC/C++、Java(商標)等の任意の他の言語で、またはその任意の組合せで符号化されてよい。例えば、示されている実施形態では、プログラム命令1025は、データベース階層のデータベースエンジンヘッドノードの、または異なる実施形態で、データ階層のクライアントの代わりにデータベース及び関連付けられたメタデータを記憶する別個の分散型データベース最適化ストレージシステムの複数のストレージノードの内の1つの機能性を実装するために実行可能なプログラム命令を含む。いくつかの実施形態では、プログラム命令1025は、複数の別個のクライアント、サーバノード、及び/または他の構成要素を実装してよい。

【0127】

いくつかの実施形態では、プログラム命令1025が、UNIX(登録商標)、Linux、Solaris(商標)、MacOS(商標)、Windows(商標)等の多様なオペレーティングシステムの内いずれかであってよいオペレーティングシステム(不図示)を実装するために実行可能な命令を含んでよい。プログラム命令1025のいずれかまたはすべては、多様な実施形態に従ってプロセスを実行するためにコンピュータシステム(または他の電子機器)をプログラミングするために使用されてよい、その上に記憶されている命令を有する非一過性のコンピュータ可読記憶媒体を含んでよいコンピュータプログラム製品、つまりソフトウェアとして提供されてよい。非一過性のコンピュータ可読記憶媒体は、マシン(例えば、コンピュータ)によって読取り可能な形(例えば、ソフトウェア、処理アプリケーション)をとる情報を記憶するための任意の機構を含んでよい。一般的に言えば、非一過性のコンピュータアクセス可能記憶媒体は、例えばI/Oインタフェース1030を介してコンピュータシステム1000に結合される、ディスクまたはDVD/CD-ROM等の磁気媒体または光学媒体等の、コンピュータ可読記憶媒体または記憶媒体を含んでよい。また、非一過性のコンピュータ可読記憶媒体は、コンピュータシステム1000のいくつかの実施形態では、システムメモリ1020または別のタイプのメモリとして含まれてよい、RAM(例えば、SDRAM、DDR SDRAM、RDRAM、SRAM等)、ROM等の任意の揮発性媒体または不揮発性媒体を含んでもよい。他の実施形態では、プログラム命令は、ネットワークインタフェース1040を介して実装されてよい等、ネットワークリンク及び/または無線リンク等の通信媒体を介して伝達される、光信号、音響信号、または他の形の伝搬信号(例えば、搬送波、赤外線信号、デジタル信号等)を使用して通信されてよい。

【0128】

いくつかの実施形態では、システムメモリ1020は、本明細書に説明されるように構成されてよいデータストア1045を含んでよい。例えば、本明細書に説明されるデータベース階層の機能を実行する際に使用されるトランザクションログ、アンドゥログ、キャッシュに入れられたページデータ、または他の情報等の、データベース階層によって(例えば、データベースエンジンヘッドノード上に)記憶されるとして本明細書に説明される情報は、データストア1045にもしくは1つまたは複数のノード上のシステムメモリ1020の別の部分に、永続記憶装置1060に、及び/または1つまたは複数のリモートストレージデバイス1070に異なるときに及び多様な実施形態で記憶されてよい。同様に、記憶階層によって記憶されているとして本明細書に説明される情報(例えば、本明細書に説明される分散型ストレージシステムの機能を実行する上で使用されるリドゥログレコード、合体データページ、及び/または他の情報)は、データストア1045にもしくは1つまたは複数のノード上のシステムメモリ1020の別の部分に、永続記憶装置10

10

20

30

40

50

60に、及び/または1つまたは複数のリモートストレージデバイス1070に異なるときに及び多様な実施形態で記憶されてよい。一般に、システムメモリ1020(例えば、システムメモリ1020の中のデータストア1045)、永続記憶装置1060、及び/またはリモートストレージ1070は、データブロック、データブロックのレプリカ、データブロックと関連付けられたメタデータ、及び/またはその状態、データベース構成情報、及び/または本明細書に説明される方法及び技法を実装する上で使用できる任意の他の情報を記憶してよい。

【0129】

一実施形態では、I/Oインタフェース1030は、プロセッサ1010と、システムメモリ1020と、ネットワークインタフェース1040または他の周辺インタフェースを通してを含んだシステムのあらゆる周辺装置との間のI/Oトラフィックを調整するように構成されてよい。いくつかの実施形態では、I/Oインタフェース1030は、1つの構成要素(例えば、システムメモリ1020)から別の構成要素(例えば、プロセッサ1010)による使用に適したフォーマットにデータ信号を変換するために任意の必要なプロトコル、タイミング、または他のデータ変形を実行してよい。いくつかの実施形態では、I/Oインタフェース1030は、例えばペリフェラルコンポーネントインターコネクト(PCI)バス規格、またはユニバーサルシリアルバス(USB)規格の変形等の多様なタイプの周辺バスを通してアタッチされるデバイスに対するサポートを含んでよい。いくつかの実施形態では、I/Oインタフェース1030の機能は、例えばノースブリッジ及びサウスブリッジ等、2つ以上の別々の構成要素に分割されてよい。また、いくつかの実施形態では、システムメモリ1020へのインタフェース等、I/Oインタフェース1030の機能性のいくつかまたはすべては、プロセッサ1010の中に直接的に組み込まれてよい。

【0130】

ネットワークインタフェース1040は、例えば、コンピュータシステム1000と、(本明細書に説明される1つまたは複数のストレージシステムサーバノード、データベースエンジンヘッドノード、及び/またはデータベースシステムのクライアントを実装してよい)他のコンピュータシステム1090等の、ネットワークにアタッチされる他のデバイスとの間でデータを交換できるように構成されてよい。さらに、ネットワークインタフェース1040は、コンピュータシステム1000と多様なI/O装置1050及び/またはリモートストレージ1070との間の通信を可能にするように構成されてよい。入出力装置1050は、いくつかの実施形態では、1つまたは複数のディスプレイ端末、キーボード、キーパッド、タッチパッド、スキャン装置、音声認識装置もしくは光学認識装置、または1つまたは複数のコンピュータシステム1000によってデータを入力するまたは取り出すために適した任意の他の装置を含んでよい。複数の入出力装置1050は、コンピュータシステム1000に存在してよい、またはコンピュータシステム1000を含む分散型システムの多様なノードで分散されてよい。いくつかの実施形態では、類似する入出力装置はコンピュータシステム1000とは別個であってよく、ネットワークインタフェース1040上で等、有線接続または無線接続を通してコンピュータシステム1000を含む分散型システムの1つまたは複数のノードと対話してよい。ネットワークインタフェース1040は、一般に1つまたは複数の無線ネットワークプロトコル(例えば、Wi-Fi/IEEE 802.11、または別の無線ネットワーク規格)をサポートしてよい。ただし、多様な実施形態では、ネットワークインタフェース1040は、例えば他のタイプのイーサネット(登録商標)ネットワーク等、任意の適切な有線汎用データネットワークまたは無線汎用データネットワークを介する通信をサポートしてよい。さらに、ネットワークインタフェース1040は、Fibre Channel SAN等のストレージエリアネットワークを介して、または任意の他の適切なタイプのネットワーク及び/またはプロトコルを介して、アナログ音声ネットワークまたはデジタルファイバ通信ネットワーク等の電気通信ネットワーク/電話網を介する通信をサポートしてよい。多様な実施形態では、コンピュータシステム1000は、図10に示される構成要素より多

10

20

30

40

50

い、少ない、または異なる構成要素（例えば、ディスプレイ、ビデオカード、オーディオカード、周辺装置、ATMインタフェース、イーサネットインタフェース、フレームリレーインタフェース等の他のネットワークインタフェース等）を含んでよい。

【0131】

本明細書に説明される分散型システムの実施形態のいずれも、またはその構成要素のいずれも1つまたは複数のウェブサービスとして実装されてよいことに留意されたい。例えば、データベースシステムのデータベース階層の中のデータベースエンジンヘッドノードは、データベースサービス、及び/または本明細書に説明される分散型ストレージシステムを利用する他のタイプのデータストレージサービスをウェブサービスとしてのクライアントに提示してよい。いくつかの実施形態では、ウェブサービスは、ネットワーク上で相互運用可能なマシン対マシンの対話をサポートするように設計されたソフトウェアシステム及び/またはハードウェアシステムによって実装されてよい。ウェブサービスは、ウェブサービス記述言語（WSDL）等のマシン処理可能なフォーマットで記述されるインタフェースを有してよい。他のシステムは、ウェブサービスのインタフェースの記述によって規定される方法でウェブサービスと対話してよい。例えば、ウェブサービスは、他のシステムが呼び出してよい多様な動作を定義してよく、多様な動作を要求するとき他のシステムが準拠することを期待されてよい特定のアプリケーションプログラミングインタフェース（API）を定義してよい。

10

【0132】

多様な実施形態では、ウェブサービスは、ウェブサービス要求と関連付けられるパラメータ及び/またはデータを含むメッセージを使用することによって要求されてよい、または呼び出されてよい。係るメッセージは、拡張マークアップ言語（XML）等の特定のマークアップ言語に従ってフォーマットされてよい、及び/またはシンプルオブジェクトアクセスプロトコル（SOAP）等のプロトコルを使用してカプセル化されてよい。ウェブサービス要求を実行するために、ウェブサービスクライアントは、要求を含むメッセージをアSEMBルし、ハイパテキスト転送プロトコル（HTTP）等のインターネットベースのアプリケーション層転送プロトコルを使用して、メッセージをウェブサービスに対応するアドレス可能なエンドポイント（例えば、ユニフォームリソースロケータ（URL））に伝達してよい。

20

【0133】

いくつかの実施形態では、ウェブサービスは、メッセージベースの技法よりむしろ、表象状態転送（「RESTful」）技法を使用して実装されてよい。例えば、RESTful技法に従って実装されるウェブサービスは、SOAPメッセージの中でカプセル化されるよりむしろ、PUT、GET、またはDELETE等のHTTP方法の中に含まれるパラメータを通して呼び出されてよい。

30

【0134】

以下の実施形態は以下の節を鑑みてさらによく理解されてよい。

1. 分散型ストレージシステムを実装する複数のストレージノードであって、分散型ストレージシステムがデータベースのためにログ構造化データストレージを実装するように構成され、複数のリドゥログレコードが複数のストレージノードでデータベースシステムから以前に受信されたことがあり、リドゥログレコードのそれぞれが複数のストレージノードの中でデータベースのために記憶されるデータに対する変更を記述する、複数のストレージノードと、

40

データベースシステムを実装するデータベースヘッドノードであって、
 複数のストレージノードとの接続を確立する、及び
 複数のストレージノードとの接続の確立時に、1つまたは複数のアクセス要求のためのデータベースへのアクセスを提供する

ための故障回復動作を実行するように構成される、データベースヘッドノードと、
 を備えるシステム。

2. 複数のリドゥログレコードをリプレイすることなく、アクセスがデータベースに提

50

供される、節 1 に記載のシステム。

3 . データベースシステムヘッドノードが、
データベースに対するアクセス要求を受信する、
受信されたアクセス要求に基づいて、ストレージノードに記憶されるデータページの現在の状態に対する要求を複数のストレージノードの内の 1 つに送信する、及び
複数のリドゥログレコードの 1 つまたは複数が、ストレージノードでデータページをその現在の状態で生成するためにデータページの以前に保存された状態に適用された、要求されたデータページをその現在の状態でストレージノードから受信する、
ようにさらに構成される、節 2 に記載のシステム。

4 . 複数のストレージノードに送信される複数のリドゥログレコードの少なくともいくつかはシステムトランザクションを含み、複数のストレージノードの内の 1 つのストレージノードが、
システムトランザクションが不完全であると決定する、及び
少なくともいくつかのリドゥログレコードによって変更されるデータページの現在の状態を生成するときに適用されないとして複数のリドゥログレコードの少なくともいくつかを識別する
ように構成される、節 1 に記載のシステム。

5 . 複数のストレージノードに送信される複数のリドゥログレコードの少なくともいくつかはシステムトランザクションを含み、データベースシステムヘッドノードが、
システムトランザクションが不完全であると決定する、
少なくともいくつかのリドゥログレコードによって変更されるデータページの現在の状態を生成するときに適用されないとして複数のリドゥログレコードの少なくともいくつかを識別する、及び
複数のストレージノードの 1 つまたは複数に、適用されない、複数のリドゥログレコードの識別された少なくともいくつかを示す通知を送信する
ようにさらに構成される、節 1 に記載のシステム。

6 . データベースヘッドノードを実装する 1 台または複数のコンピューティング装置によって、
データベースヘッドノード故障からの回復時に、
データベースのためのデータを記憶する分散型ストレージシステムを実装する複数のストレージノードの 1 つまたは複数のストレージノードと接続を確立することであって、分散型ストレージシステムが、データベースのためにログ構造化データストレージを実装するように構成され、複数のリドゥログレコードが複数のストレージノードで以前に受信されたことがあり、リドゥログレコードのそれぞれが、データベースのために記憶されたデータに対する変更を、それが受信されたそれぞれのストレージノードで記述する、接続を確立することと、
複数のストレージノードの 1 つまたは複数のストレージノードとの接続の確立時に、アクセスのためにデータベースを利用可能にすることと、
を実行することを、
含む方法。

7 . データベースに対するアクセス要求を受信することと、
アクセス要求を受信することに応じて、1 つまたは複数のストレージノードからデータベースのためのデータの部分を記憶する 1 つまたは複数のデータページの現在の状態を要求することと、
アクセス要求にサービスを提供するためにデータベースのためのデータの部分を記憶する 1 つまたは複数のストレージノードから 1 つまたは複数のデータページの現在の状態を受信することと、
をさらに含む、節 6 に記載の方法。

8 . 受信されたデータベースのためのデータの部分を記憶する 1 つまたは複数のデータページの少なくとも 1 つの現在の状態が、1 つまたは複数のストレージノードの 1 つが、

10

20

30

40

50

少なくとも1つのデータページの以前に記憶されたバージョンまで複数のリドゥログレコードの1つまたは複数を読み替えることによって生成される、節7に記載の方法。

9. 受信されたデータベースのためのデータの部分を記憶する1つまたは複数のデータページの少なくとも1つからの異なるデータページの現在の状態が、データページの以前に記憶されたバージョンまで複数のリドゥログレコードの1つまたは複数を読み替えることなく、1つまたは複数のストレージノードの内の1つによって送信される、節8に記載の方法。

10. データベースヘッドノードが、複数のストレージノードにリドゥログレコードとして送信される変更をアンドゥするために複数のアンドゥログレコードを維持し、方法が、

1つまたは複数のストレージノードから受信される1つまたは複数のデータページの1つが不完全なユーザートランザクションによって影響を及ぼされると決定することであって、ユーザートランザクションが、1つのデータページを含んだ1つまたは複数のストレージノードに記憶されるデータに対して変更を向ける、決定することと、

ユーザートランザクションによってデータページに向けられた変更をアンドゥするためにデータページに1つまたは複数のアンドゥログレコードを適用することと、をさらに含む、節7に記載の方法。

11. データベースヘッドノードが、ユーザートランザクションを含んだ複数の不完全なユーザートランザクションを示すランザクションテーブルを維持し、方法が、

ランザクションテーブルに少なくとも部分的に基づいて、複数の不完全なユーザートランザクションの少なくとも1つによって影響を及ぼされる1つまたは複数の追加のデータページを決定することと、

1つまたは複数のストレージノードから1つまたは複数の追加のデータページの現在の状態を要求することと、

1つまたは複数の追加のデータページを受信することに応じて、少なくとも1つの不完全なユーザートランザクションによって1つまたは複数の追加のデータページに向かって向けられる変更をアンドゥするために、1つまたは複数の追加のデータページに追加の1つまたは複数のアンドゥログレコードを適用することと、

をさらに含む、節10に記載の方法。

12. 1つまたは複数の追加のデータページを該決定すること、1つまたは複数の追加のデータページを該要求すること、及び1つまたは複数の追加のデータページに追加の1つまたは複数のアンドゥログレコードを該適用することが、データベースヘッドノードでバックグラウンドプロセスの一部として実行され、アクセス要求を該受信すること、1つまたは複数のデータページの現在の状態を該要求すること、及び1つまたは複数のデータページの現在の状態を該受信することがフォアグラウンドプロセスの一部として実行される、節11に記載の方法。

13. データベースヘッドノード故障からの回復時に、

アクセスのためにデータベースを利用可能にする前に、データベースのために複数のストレージノードに記憶されたデータの、以前に記録されたスナップショットに対応する状態への復元に対する要求を複数のストレージノードに送信することであって、該復元が複数のリドゥログの1つまたは複数を読み替えるデータの以前のバージョンに適用することを含む、送信することと、

をさらに含む、節6に記載の方法。

14. データベースが複数のリドゥログレコードを読み替えることなくアクセスのために利用可能にされる、節6に記載の方法。

15. 1台または複数のコンピューティング装置による実行時に、

データベースヘッドノード故障からの回復時に、

データベースのためのデータを記憶する分散型ストレージシステムを実装する複数のストレージノードの1つまたは複数のストレージノードとの接続を確立することであって、分散型ストレージシステムがデータベースのためにログ構造化データストレージを実装す

10

20

30

40

50

るように構成され、複数のリドゥログレコードが複数のストレージノードで以前に受信されたことがあり、リドゥログレコードのそれぞれが、データベースのために記憶されたデータに対する変更を、それが受信されたそれぞれのストレージノードで記述する、接続を確立すること、及び

複数のストレージノードの1つまたは複数のストレージノードとの接続の確立時に、1つまたは複数のアクセス要求のためにデータベースへのアクセスを提供すること、を実装するデータベースシステムのデータベースヘッドノードを実装するプログラム命令を記憶する非一過性のコンピュータ可読記憶媒体。

16. 複数のストレージノードで以前に受信された複数のリドゥログレコードが、該データベースヘッドノードとは異なるデータベースヘッドノードから受信された、節15に記載の非一過性のコンピュータ可読記憶媒体。

10

17. 複数のリドゥログレコードをリプレイすることなく、アクセスがデータベースに提供される、節15に記載の非一過性のコンピュータ可読記憶媒体。

18. データベースシステムヘッドノードが、データベースに対するアクセス要求を受信することと、アクセス要求を受信することに応じて、1つまたは複数のストレージノードからデータベースのためのデータの部分を記憶する1つまたは複数のデータページの現在の状態を要求することと、

アクセス要求にサービスを提供するためにデータベースのためのデータの部分を記憶する1つまたは複数のデータページの現在の状態を受信することであって、1つまたは複数の受信されたデータページの少なくとも1つの現在の状態が、1つまたは複数のストレージノードの1つが、少なくとも1つのデータページの以前に記憶されたバージョンまで複数のリドゥログレコードの1つまたは複数を実行することによって生成される、受信することと、

20

をさらに実装する、節15に記載の非一過性のコンピュータ可読記憶媒体。

19. 受信されたアクセス要求が読取り要求または書込み要求である、節18に記載の非一過性のコンピュータ可読記憶媒体。

20. データベースヘッドノードが、リドゥログレコードとして複数のストレージノードに送信された変更をアンドゥするために、複数のアンドゥログレコードを維持し、データベースヘッドノードが、

30

1つまたは複数のストレージノードから受信される1つまたは複数のデータページの1つが不完全なユーザートランザクションによって影響を及ぼされると決定することであって、ユーザートランザクションが1つのデータページを含んだ1つまたは複数のストレージノードに記憶されるデータに対して変更を向ける、決定することと、

ユーザートランザクションによってデータページに向けられた変更をアンドゥするためにデータページに1つまたは複数のアンドゥログレコードを適用することと、をさらに実装する、節18に記載の非一過性のコンピュータ可読記憶媒体。

21. データベースヘッドノードが、ユーザートランザクションを含んだ複数の不完全なユーザートランザクションを示すトランザクションテーブルを維持し、データベースシステムヘッドノードが、

40

トランザクションテーブルに少なくとも部分的に基づいて、複数の不完全なユーザートランザクションの少なくとも1つによって影響を及ぼされる1つまたは複数の追加のデータページを決定することと、

1つまたは複数のストレージノードから1つまたは複数の追加のデータページの現在の状態を要求することと、

1つまたは複数の追加のデータページを受信することに応じて、少なくとも1つの不完全なユーザートランザクションによって1つまたは複数の追加のデータページに向かって向けられる変更をアンドゥするために、1つまたは複数の追加のデータページに追加の1つまたは複数のアンドゥログレコードを適用することと、

をバックグラウンドプロセスとして実行すること

50

をさらに実装する、節 20 に記載の非一過性のコンピュータ可読記憶媒体。

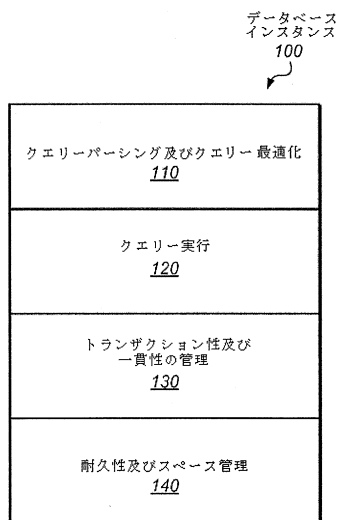
【0135】

図に示され、本明細書に説明される多様な方法は、方法の例の実施形態を表す。方法は、ソフトウェアで、ハードウェアで、またはソフトウェア及びハードウェアの組合せで手動で実装されてよい。任意の方法の順序は変更されてよく、多様な要素が追加、再順序付け、結合、省略、修正等、されてよい。

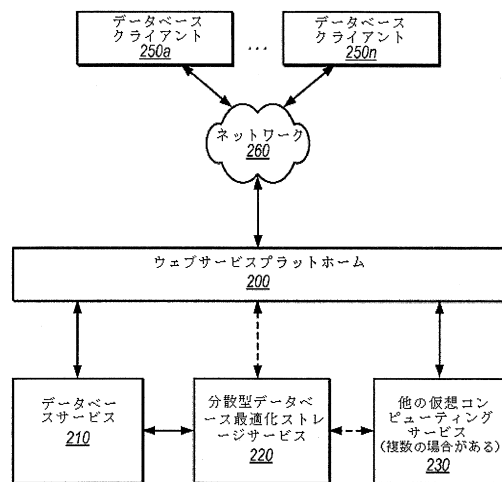
【0136】

上記実施形態はかなり詳細に説明されているが、いったん上記開示が完全に理解されると当業者に明らかになるように、多数の変形形態及び修正形態が加えられてよい。続く特許請求の範囲が、すべての係る修正形態及び変更を包含すると解釈され、したがって上記説明は制限的な意味よりむしろ例示的な意味で考えられることが意図される。

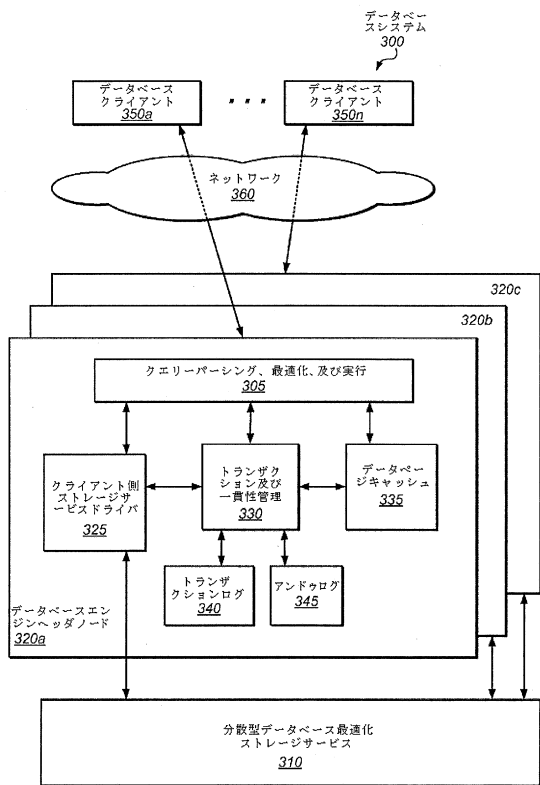
【図 1】



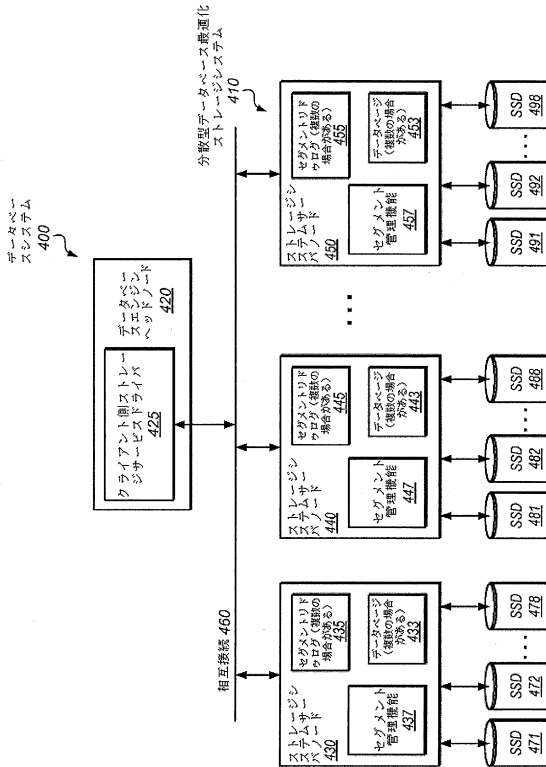
【図 2】



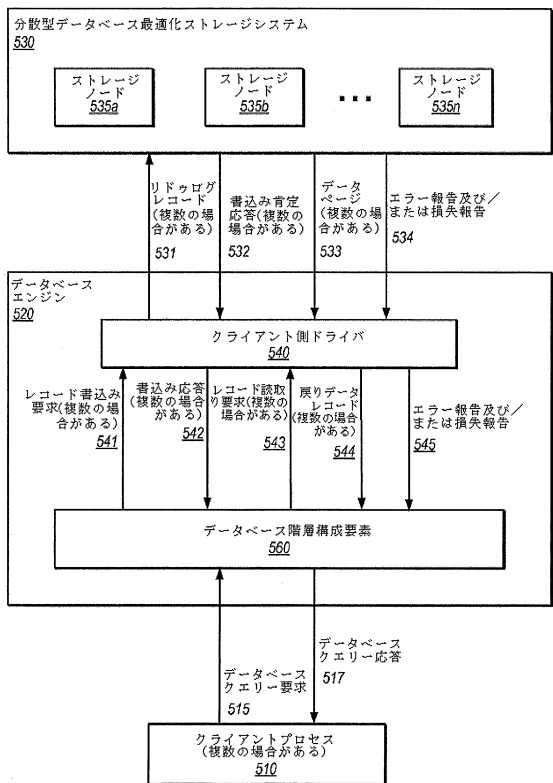
【図3】



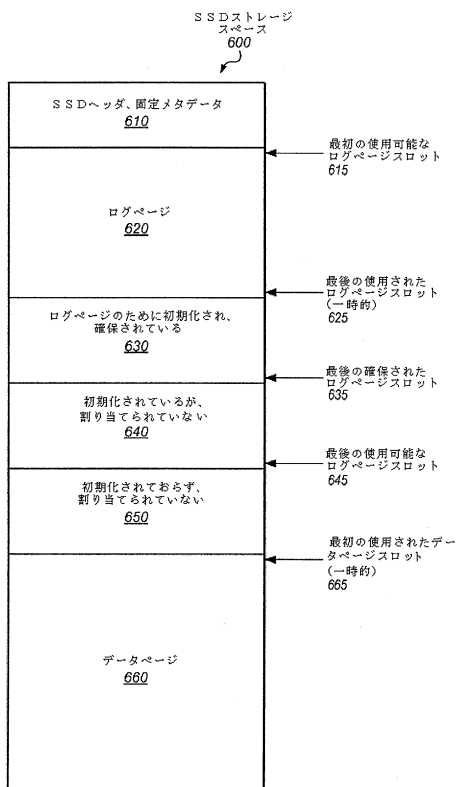
【図4】



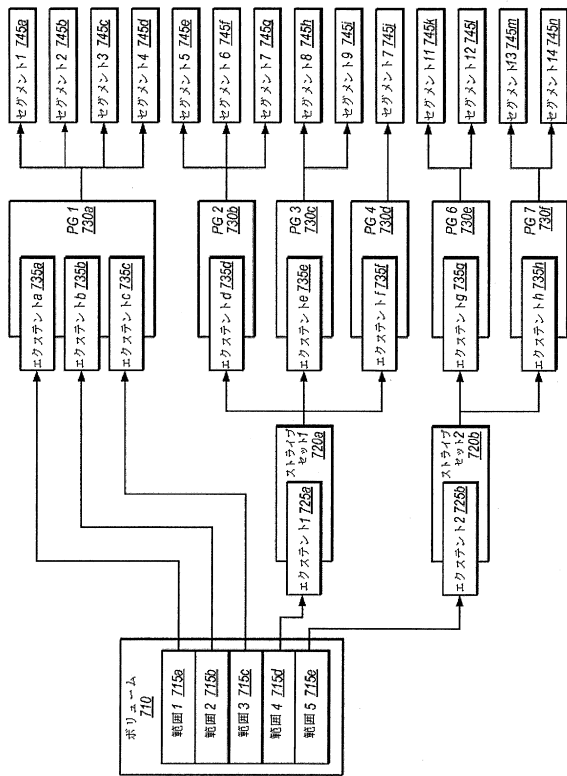
【図5】



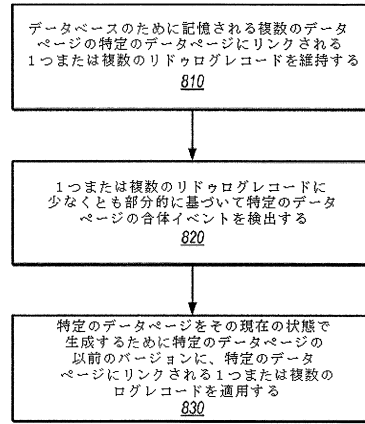
【図6】



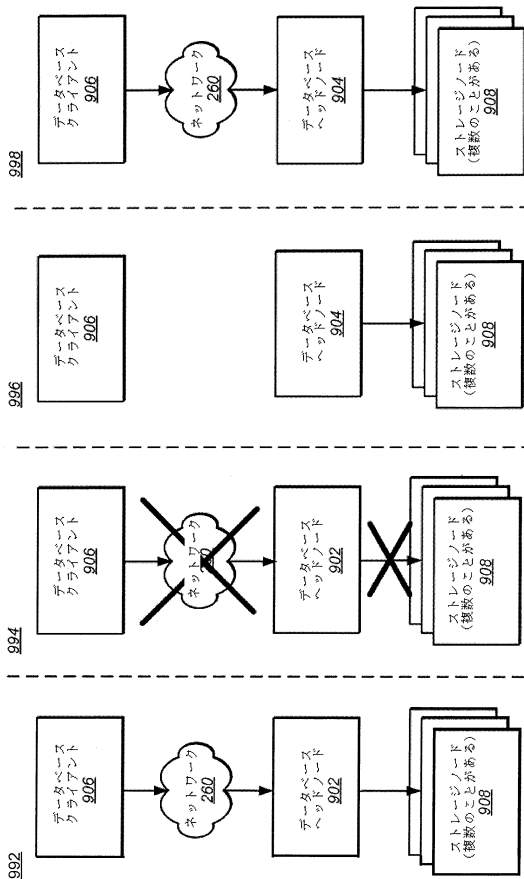
【 図 7 】



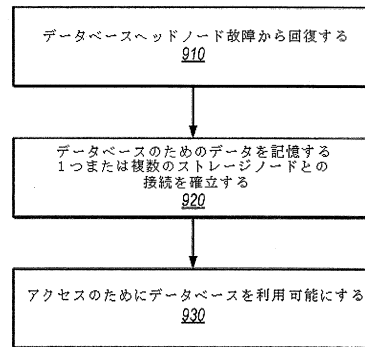
【 図 8 】



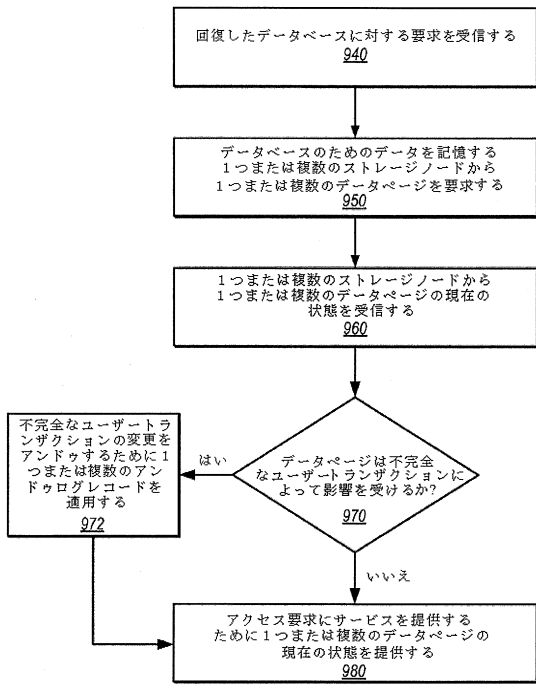
【 図 9 A 】



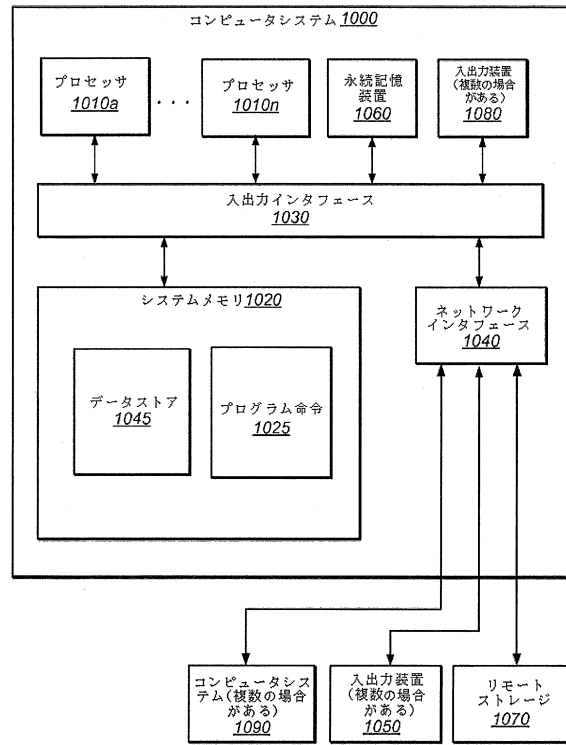
【 図 9 B 】



【図9C】



【図10】



フロントページの続き

- (72)発明者 バーチャル, ローリオン・ダレル
アメリカ合衆国・98109-5210・ワシントン州・シアトル・テリー アヴェニュー ノース
・410
- (72)発明者 マダヴァラブ, プラディーブ・ジュニャーナ
アメリカ合衆国・98109-5210・ワシントン州・シアトル・テリー アヴェニュー ノース
・410
- (72)発明者 ファハン, ニール
アメリカ合衆国・98109-5210・ワシントン州・シアトル・テリー アヴェニュー ノース
・410

審査官 漆原 孝治

- (56)参考文献 特表2012-507086(JP,A)
特開2004-227169(JP,A)
特開2004-348701(JP,A)
特開平06-195250(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G06F 12/00