

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2016-519379

(P2016-519379A)

(43) 公表日 平成28年6月30日 (2016. 6. 30)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 12/00 (2006.01)	G06F 12/00 535F	
	G06F 12/00 518A	
	G06F 12/00 545A	

審査請求 有 予備審査請求 未請求 (全 35 頁)

(21) 出願番号 特願2016-514043 (P2016-514043)
 (86) (22) 出願日 平成26年5月13日 (2014. 5. 13)
 (85) 翻訳文提出日 平成28年1月7日 (2016. 1. 7)
 (86) 国際出願番号 PCT/US2014/037901
 (87) 国際公開番号 W02014/186396
 (87) 国際公開日 平成26年11月20日 (2014. 11. 20)
 (31) 優先権主張番号 13/893, 004
 (32) 優先日 平成25年5月13日 (2013. 5. 13)
 (33) 優先権主張国 米国 (US)

(71) 出願人 507303550
 アマゾン・テクノロジーズ・インコーポレ
 ーテッド
 アメリカ合衆国・89507・ネバダ州・
 レノ・ピーオーボックス 8102
 (74) 代理人 100064621
 弁理士 山川 政樹
 (74) 代理人 100098394
 弁理士 山川 茂樹
 (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 ノース・410

最終頁に続く

(54) 【発明の名称】 トランザクションの順序付け

(57) 【要約】

データベースサービスのノードは、データベースサービスによって記憶されたレコードの読み取りを実行する読み取り要求、及びそのレコードに対してトランザクションを実行するトランザクション要求を受け取ることができる。第1の時刻指示と第2の時刻指示とを、読み取りとトランザクションとにそれぞれ関連付けることができる。潜在的な読み取り異常（例えば、ファジーリード、リードスキューなど）は、第1の時刻指示が、第2の時刻指示の閾値内にあると判定したことに少なくとも部分的に基づいて検出され得る。潜在的な読み取り異常を検出したことに応答して、第1の時刻指示が第2の時刻指示よりも早い時点を示しているかどうかに関わらず、トランザクション要求によって指定されたトランザクションの後に上記読み取りが実行され得る。

【選択図】 図3

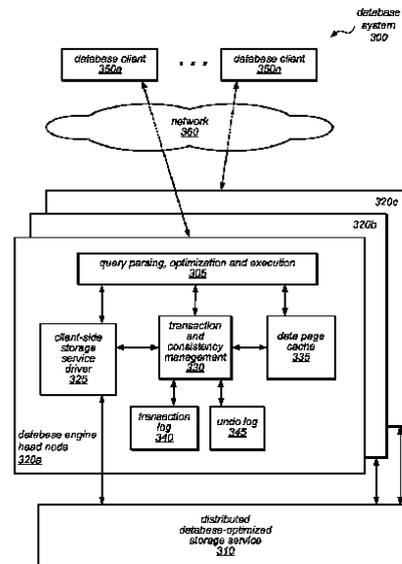


FIG. 3

【特許請求の範囲】**【請求項 1】**

システムであって、

それぞれが少なくとも1つのプロセッサ、及びメモリを含む複数のコンピューティングノードであって、データベースサービスを集合的に実施するように構成されている複数のコンピューティングノード、を含み、

前記複数のコンピューティングノードの第1のノードは、

データベーステーブル内の特定のデータページの特定のデータレコードを対象とする書き込み要求であって、前記特定のデータレコードに対して行うべき変更を指定する書き込み要求を受け取り、

前記変更をコミットする時刻を示すコミット時刻を前記書き込み要求に割り当てるように構成されており、

前記複数のコンピューティングノードの第2のノードは、

前記データベーステーブル内の前記特定のデータページを読み取る読み取り要求を受け取り、

読み取りの一貫性ポイントを示す読み取り一貫性ポイント時刻を前記読み取り要求に割り当て、

前記読み取り一貫性ポイント時刻が前記コミット時刻の精度ウィンドウ内にあると判定し、

前記読み取り一貫性ポイント時刻が前記コミット時刻よりも早い時点を示しているかどうかに関わらず、前記書き込み要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を行わせるように構成されている、システム。

【請求項 2】

前記複数のコンピューティングノードの前記第2のノードは、

前記書き込み要求に関連付けられた書き込みアンラッチ時刻が、前記読み取り要求に関連付けられた読み取りラッチ時刻の精度ウィンドウ内にあると判定するようにさらに構成されており、

前記書き込み要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を前記行わせることは、

前記読み取りラッチ時刻を後の時点に移動すること、及び

前記読み取り一貫性の時点以降に前記読み取り要求を再試行すること、を含む請求項1に記載のシステム。

【請求項 3】

前記書き込み要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を前記行わせることは、前記移動された読み取りラッチ時刻が前記書き込みアンラッチ時刻の精度ウィンドウ内にはないと判定すること、をさらに含む、請求項2に記載のシステム。

【請求項 4】

前記第1のノード及び前記第2のノードのそれぞれは、互いの前記精度ウィンドウ内に維持されている個別のクロックを保持するように構成されており、前記コミット時刻が前記第1のノードによって決定され、前記読み取り一貫性の時点が前記第2のノードによって決定される、請求項1に記載のシステム。

【請求項 5】

方法であって、

複数のコンピューティングノードによって、

1つ以上のクライアントから、記憶されたレコードの読み取りを実行する読み取り要求、及び前記レコードの更新を実行する更新要求を受け取ること、

第1の時刻指示と第2の時刻指示とを、前記読み取りと前記更新とにそれぞれ関連付けること、並びに

前記第1の時刻指示が前記第2の時刻指示の閾値内にあると判定したことに少なくとも

10

20

30

40

50

部分的に基づき、潜在的な読み取り異常を検出すること、
を実行すること、を含む方法。

【請求項 6】

前記潜在的な読み取り異常を前記検出したことに応答して、前記第 1 の時刻指示が前記第 2 の時刻指示よりも早い時点を示しているかどうかに関わらず、前記更新要求によって指定された前記更新の後に前記読み取り要求によって指定された前記読み取りを実行させること、をさらに含む、請求項 5 に記載の方法。

【請求項 7】

前記読み取りの再試行に関連付けられた再試行の時刻指示が前記第 1 の時刻指示よりも後の時刻を示すように、前記第 1 の時刻指示以降に前記読み取りを再試行すること、及び前記再試行された読み取りに対して読み取り異常が発生していないと判定すること、をさらに含む請求項 5 に記載の方法。

10

【請求項 8】

前記第 2 の時刻指示は前記更新をコミットする時刻を示す、請求項 5 に記載の方法。

【請求項 9】

前記潜在的な読み取り異常を前記検出することは、第 3 の時刻指示にさらに基づいており、前記第 3 の時刻指示も前記更新要求に関連付けられている、請求項 5 に記載の方法。

【請求項 10】

前記潜在的な読み取り異常を前記検出することは、前記第 3 の時刻指示が第 4 の時刻指示の閾値内にあると判定したことにさらに基づいており、前記第 4 の時刻指示も前記読み取り要求に関連付けられている、請求項 9 に記載の方法。

20

【請求項 11】

前記読み取り要求は前記複数のノードの第 1 のノードによって受け取られ、前記更新要求は前記複数のノードの別の第 2 のノードによって受け取られ、前記第 1 のノード及び前記第 2 のノードは個別のクロックをそれぞれ保持しており、前記読み取りに関連付けられた前記第 1 の時刻指示は前記第 1 のノードの前記個別のクロックによって決定され、前記更新に関連付けられた前記第 2 の時刻指示は前記第 2 のノードの前記個別のクロックによって決定される、請求項 5 に記載の方法。

【請求項 12】

前記潜在的な読み取り異常は潜在的なファジーリードである、請求項 5 に記載の方法。

30

【請求項 13】

再試行の頻度に少なくとも部分的に基づいて前記閾値を変更すること、をさらに含む、請求項 5 に記載の方法。

【請求項 14】

前記検出することは、前記レコードの前記読み取り、及び別のレコードの読み取りまたは前記レコードの第 2 の読み取りのいずれかを含む読み取り要求に対して実行される、請求項 5 に記載の方法。

【請求項 15】

システムであって、

1 つ以上のプロセッサと、

40

プログラム命令が記憶された 1 つ以上のメモリであって、前記プログラム命令が、データベースサービスのデータベースノードを実装するために前記 1 つ以上のプロセッサによってコンピュータで実行可能である 1 つ以上のメモリと、を含み、前記データベースノードは、

前記データベースサービスによって記憶されたレコードの読み取りを指定する読み取り要求によって指定された前記読み取りに第 1 の時刻指示に関連付け、

前記第 1 の時刻指示が、別のデータベースノードによって受け取られたトランザクション要求であって、前記レコードを変更するトランザクションを指定するトランザクション要求によって指定された前記トランザクションに関連付けられた第 2 の時刻指示の閾値内にあると判定し、

50

前記第1の時刻指示が前記第2の時刻指示よりも早い時点を示しているかどうかに関わらず、前記トランザクション要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を行わせるように構成されている、システム。

【発明の詳細な説明】

【背景技術】

【0001】

ソフトウェアスタックの様々な構成要素を分散させることにより、ある場合には、（例えば、複製による）耐障害性、高度な持続性、及び（例えば、大型で高価な構成要素を少数使用するのではなく、それよりも小型で安価な構成要素を多数使用することによって）費用がより少なくて済む解決策が得られる（またはそれらの助けになる）。しかしながら、データベースは、歴史的に見ると、分散に最も適用し難いソフトウェアスタックの構成要素間に存在してきた。例えば、データベースを分散させる一方で、データベースによる提供が期待されるいわゆるACID特性（例えば、原子性（Atomicity）、一貫性（Consistency）、独立性（Isolation）及び持続性（Durability））を依然として保証することは困難な場合がある。特に、一貫性及び独立性の各特性に関しては、分散データベースシステムのノード間を整合してノード全体の因果関係を保つことは、従来システムにとって非常に難しいことが証明されている。

10

【図面の簡単な説明】

【0002】

【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様々な構成要素を示すブロック図である。

20

【図2】いくつかの実施形態に係る、トランザクションの順序付けを行うように構成されたWebサービススペ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得る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。

【図3】一実施形態に係る、トランザクションの順序付けを行うように構成されたデータベースシステムの様々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、トランザクションの順序付けを行うように構成された分散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、トランザクションの順序付けを行うように構成されたデータベースシステム内の別個の分散データベース最適化ストレージシステムの利用を示すブロック図である。

30

【図6】トランザクションの順序付け方法の一実施形態を示すフロー図である。

【図7】図7A～7Cは、様々な実施形態に係る様々なトランザクションの順序付けを行うシナリオを示すタイミング図である。

【図8】様々な実施形態に係る、トランザクションの順序付けを実装するように構成されたコンピュータシステムを示すブロック図である。

【0003】

本明細書では、いくつかの実施形態及び説明図の実施例として実施形態を説明するが、当業者であれば、記載の実施形態または図面にこれらの実施形態が限定されないことを認識するであろう。図面及びそれに対する詳細な説明は、開示された特定の形態に実施形態を限定することを意図するものではなく、むしろ、本発明は、添付された特許請求の範囲によって定められるような概念及び趣旨に該当する全ての修正形態、均等形態及び代替形態を包含するものであることを理解すべきである。本明細書で用いられる見出しは、体系化を目的とするにすぎず、本明細書または特許請求の範囲を限定するために用いられることを意味するものではない。本出願全体を通して用いられるように、用語「し得る（may）」は、義務的な意味（すなわち、～しなければならないを意味する）ではなく、許容的な意味（すなわち、～する可能性を有することを意味する）で用いられる。用語「含む（include）」、「含んでいる（including）」及び「含む（includes）」は、非制限的な関係を示しており、従って、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、用語「有する（have）」、「有している（having）」

40

50

g)」及び「有する(has)」も非制限的な関係を示すものであり、従って、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。本明細書で用いられる用語「第1の(first)」、「第2の(second)」及び「第3の(third)」などは、それらに先行する名詞のラベルとして用いられるものであり、順序付けなどが明示的に示されない限り、(例えば、空間的、時間的、論理的などの)いかなる順序付けも意味しない。

【0004】

種々の構成要素は、1つ以上のタスクを実行「するように構成された(configured to)」と記載されることがある。このような文脈では、「するように構成された(configured to)」は、動作中に1つ以上のタスクを実行する「構造を備えている(having structure that)」ことを広く意味する広義の記述である。従って、構成要素がタスクを現在実行していないときでも、そのタスクを実行するようにその構成要素を構成することができる(例えば、コンピュータシステムは、現在動作が実行されていないときでも、その動作を実行するように構成され得る)。ある文脈では、「するように構成された(configured to)」は、動作中に1つ以上のタスクを実行する「回路を備えている(having circuitry that)」ことを広く意味する構造についての広義の記述であり得る。従って、構成要素が現在オンになっていないときでも、タスクを実行するようにその構成要素を構成することができる。通常、「するように構成された(configured to)」に相当する構造をなす回路は、ハードウェア回路を含み得る。

【0005】

種々の構成要素は、本明細書において便宜上、1つ以上のタスクを実行するものとして説明され得る。こうした説明は、「するように構成された(configured to)」という語句を含むものとして解釈されるべきである。1つ以上のタスクを実行するように構成されている構成要素について述べることは、その構成要素に対して米国特許法第112条第6段落の解釈を援用することを明示的に意図していない。

【0006】

「~に基づく(Base on)」。本明細書で使用する場合、本用語は、決定に影響を及ぼす1つ以上の要因を説明するために用いられる。本用語は、決定に影響を及ぼし得る追加の要因を除外しない。すなわち、決定は、これらの要因だけにに基づくものであってもよく、またはこれらの要因に少なくとも部分的に基づくものであってもよい。「Bに基づいてAを決定する(determine A based on B)」という語句について考察する。Bは、Aの決定に影響を及ぼす要因であり得るが、このような語句は、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のみに基づいて決定されてもよい。

【0007】

本開示の範囲は、本明細書において取り込まれる問題のいずれかまたは全てを軽減するものであろうとなかろうと、(明示的または暗示的に)本明細書に開示された任意の特徴もしくは特徴の任意の組み合わせ、またはそれらの任意の一般化を包含する。従って、本出願(または本出願の優先権を主張する出願)の手続き中に、任意のこうした特徴の組み合わせに対して新たな特許請求の範囲を作成してもよい。特に、添付された特許請求の範囲を参照すると、従属請求項からの特徴は、独立請求項の特徴と組み合わせてもよく、各独立請求項からの特徴は、添付された特許請求の範囲に列挙された特定の組み合わせのみならず、任意の適切な方法で組み合わせてもよい。

【発明を実施するための形態】

【0008】

トランザクションの順序付けを行う様々な実施形態を開示する。本実施形態の様々な実施形態は、記憶されたレコードの読み取りを実行する読み取り要求、及びそのレコードに対して(例えば、書き込みなどの)トランザクションを実行するトランザクション要求を受け取る(例えば、データベースサービスの)ノードを含み得る。本実施形態の様々な実施形態は、第1の時刻指示と第2の時刻指示とを、読み取りとトランザクションとにそれ

10

20

30

40

50

ぞれ関連付けるノードも含み得る。本実施形態の様々な実施形態は、第1の時刻指示が第2の時刻指示の閾値内にあると判定したことに少なくとも部分的に基づき、潜在的な読み取り異常（例えば、ファジーリード、リードスキューなど）を検出することをさらに含み得る。なお、いくつかの実施形態では、検出は、第1の時刻指示及び第2の時刻指示以外の時刻指示に基づくものであってもよい。潜在的な読み取り異常を検出したことに応答して、第1の時刻指示が第2の時刻指示よりも早い時点を示しているかどうかに関わらず、トランザクション要求によって指定されたトランザクションの後に上記読み取りが実行され得る。いくつかの事例では、読み取りを再試行することにより、その再試行については潜在的な読み取り異常が発生しないようにしてもよい。

【0009】

まず、本明細書では、開示されたトランザクションの順序付け技術を実装するように構成された例示的なWebサービスベースのデータベースサービスについて説明する。この例示的なWebサービスベースのデータベースサービスに関する説明には、データベースエンジン及び別個の分散データベースストレージサービスなどの、例示的なWebサービスベースのデータベースサービスの様々な態様が含まれる（なお、いくつかの実施形態では、データベースエンジンから上記ストレージサービスを分けなくてもよい）。次いで、本明細書では、トランザクションの順序付け方法の様々な実施形態のフローチャートについて説明する。次に、本明細書では、開示された技術を実装可能な例示的なシステムについて説明する。本明細書全体を通じて、様々な実施例が提示される。なお、開示されたトランザクションの順序付け技術は、図1～5の例示的なデータベースサービス以外のシステム（データの読み取り、書き込み及び記憶に使用可能な他のシステムなど）において用いられてもよい。例えば、開示された技術は、データの読み取りとそのデータに対する一連の更新とが、上記読み取りからそれらの更新が見えるようになった時点で行われ得る任意のシステムにおいて用いられてもよい。

【0010】

本明細書で説明されたシステムは、いくつかの実施形態では、クラウドコンピューティング環境においてクライアント（例えば、加入者）がデータストレージシステムを運用可能にするWebサービスを実装し得る。いくつかの実施形態では、データストレージシステムは、高度にスケラブルかつ拡張可能な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、複数の物理リソースにわたって分散されているデータベースストレージにクエリを送ってもよく、必要に応じてデータベースシステムをスケールアップまたはスケールダウンしてもよい。データベースシステムは、種々の実施形態では、様々な形式及び/または編成のデータベーススキーマを用いて有効に機能し得る。いくつかの実施形態では、クライアント/加入者は、いくつかの方法を用いて、例えば、SQLインターフェースを介して対話的に、データベースシステムにクエリを送出してもよい。他の実施形態では、外部のアプリケーション及びプログラムが、ODBC（Open Database Connectivity）ドライバインターフェース及び/またはJDBC（Java Database Connectivity）ドライバインターフェースを用いてデータベースシステムにクエリを送出してもよい。

【0011】

より具体的には、本明細書で説明されたシステムは、いくつかの実施形態では、1つのデータベースシステムの様々な機能要素が本質的に分散されているサービス指向データベースアーキテクチャを実装してもよい。例えば、複数の包括的かつ一体的なデータベースインスタンス（これらのインスタンスのそれぞれは、アプリケーションサーバ、検索機能、またはデータベースの基幹機能を提供するのに必要とされるもの以外の他の機能などの、外部機能を備え得る）を1つに結び付けるのではなく、これらのシステムは、データベースの基本動作（例えば、クエリ処理、トランザクション管理、キャッシング及びストレージ）を、個別のかつ独立的にスケラブルな各層に編成してもよい。例えば、いくつかの実施形態では、本明細書で説明されたシステム内の各データベースインスタンスは、データベース階層（この階層は、1つのデータベースエンジン・ヘッドノード、及びクライ

10

20

30

40

50

アント側ストレージシステムドライバを含み得る)、並びに別個の分散ストレージシステム(このシステムは、既存システムのデータベース階層において従来行われた動作の一部を集合的に実行する複数のストレージノードを含み得る)を含んでもよい。本明細書で既に述べたように、説明されたトランザクションの順序付け技術は、他のシステムにおいても同様に適用され得る。

【0012】

本明細書でより詳細に説明するように、いくつかの実施形態では、データベースの最下位レベルの動作の一部(例えば、バックアップ、復元、リカバリ、ログレコード操作及び/または様々な空き領域管理動作)をデータベースエンジンからストレージ層に肩代わりさせ、これらを、複数のノード及びストレージ装置にわたって分散させてもよい。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、データベーステーブル(またはそのデータページ)に変更を適用し、次いで変更されたデータページをストレージ層に送るのではなく、記憶されたデータベーステーブル(及びそのデータページ)に変更を適用することをストレージ層自体が担当してもよい。このような実施形態では、変更されたデータページではなくREDOログレコードをストレージ層に送ってもよく、その後、REDO処理(例えば、REDOログレコードの適用)を(例えば、バックグラウンドプロセスによって)ある程度緩慢に、かつ分散的な方法で実行してもよい。いくつかの実施形態では、クラッシュリカバリ(例えば、記憶されたREDOログレコードからデータページを再構築すること)は、ストレージ層によって行われてもよく、分散された(かつ、場合によっては、緩慢な)バックグラウンドプロセスによって行われてもよい。

10

20

【0013】

いくつかの実施形態では、REDOログ(及び未変更のデータページ)のみがストレージ層に送られるため、データベース階層とストレージ層の間のネットワークトラフィックが、既存のデータベースシステムにおけるネットワークトラフィックと比べて大幅に減少し得る。いくつかの実施形態では、各REDOログは、そのログが変更を指定するデータページに相当するサイズの1/10程度となり得る。なお、データベース階層及び分散ストレージシステムから送られる要求は非同期であってよく、こうした複数のかかる要求を一度に移送してもよい。

【0014】

一般に、1つのデータが与えられた上で、データベースにとって求められる最も重要なことは、データベースが最終的に同一のデータを戻すことができることである。これを行うために、データベースがいくつかの異なる構成要素(または階層)を含み、これらの構成要素のそれぞれが、異なる機能を実行するようにしてもよい。例えば、従来のデータベースは、3つの階層、すなわち、クエリを解析し、最適化し、実行する第1階層、トランザクション性、リカバリ及び持続性を提供する第2階層、並びにローカル接続のディスク上、またはネットワーク接続のストレージ上のいずれかにストレージを提供する第3階層を持っているものとみなすことができる。上述したように、従来のデータベースをスケールリングするための従来の試みでは、通常、上記データベースの3つの階層全てを複製すること、及びこ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にわたって分散させることが求められている。

30

40

【0015】

いくつかの実施形態では、本明細書で説明されたシステムは、スケールリングを実装するために、従来のデータベースを用いたシステムとは違ってデータベースシステムの機能を分割してもよく、(完全なデータベースインスタンスではなく)それらの機能要素のサブセットのみを複数のマシンにわたって分散させてもよい。例えば、いくつかの実施形態では、クライアントに直接対応する階層は、データの記憶手順または検索手順ではなく、どのデータを記憶または検索すべきかを指定する要求を受け取るように構成され得る。この階層が要求の解析及び/または最適化(例えば、SQLの解析及び最適化)を実行してもよく、それと共に別の階層がクエリを実行する役割を果たしてもよい。いくつかの実施形態では、第3階層は、トランザクション性及び結果の一貫性を提供する役割を果たしても

50

よい。例えば、この階層は、いわゆる A C I D 特性の一部、特に、データベースを対象とするトランザクションの原子性、データベース内の一貫性の維持、及びデータベースを対象とするトランザクション間の独立性を保証することを実現するように構成され得る。いくつかの実施形態では、この第 3 階層は、開示されたトランザクションの順序付け技術を実装してもよい。いくつかの実施形態では、次いで、第 4 階層が、様々な種類の障害の存在下で記憶データの持続性を提供する役割を果たしてもよい。例えば、この階層は、変更ロギング、データベースクラッシュからのリカバリ、下位のストレージ容量へのアクセスの管理、及び/または下位のストレージ容量における空き領域の管理を行う役割を果たしてもよい。

【 0 0 1 6 】

次に図面に目を向けると、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様々な構成要素を示すブロック図である。本実施例に示すように、データベースインスタンスが複数の機能要素（または機能層）を含み、これらの機能要素のそれぞれが、データベースインスタンスの機能の一部を提供するようにしてもよい。本実施例では、データベースインスタンス 1 0 0 は、(1 1 0 に示す) クエリ解析・クエリ最適化層、(1 2 0 に示す) クエリ実行層、(1 3 0 に示す) トランザクション性・一貫性管理層、及び (1 4 0 に示す) 持続性・空き領域管理層を含む。上述したように、いくつかの既存のデータベースシステムにおいてデータベースインスタンスをスケールアップするには、データベースインスタンス全体 (図 1 に示した層の全てを含む) を 1 回以上複製し、次いで、グループロジックを加えてそれらを統合することが求められる場合がある。いくつかの実施形態では、本明細書で説明されたシステムは、その代わりに持続性・空き領域管理層 1 4 0 の機能をデータベース階層から別個のストレージ層に肩代わりさせてもよく、ストレージ層内の複数のストレージノードにわたってその機能を分散させてもよい。なお、開示されたトランザクションの順序付け技術は、持続性・空き領域管理層 1 4 0 がデータベース階層の一部である実施形態において適用されてもよい。

【 0 0 1 7 】

様々な実施形態では、本明細書で説明されたデータベースシステムは、様々なデータベースの動作/トランザクションのために、標準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ーフェース (application programming interface : A P I) に対応し得る。例えば、この A P I は、データベースを作成すること、テーブルを作成すること、テーブルを変更すること、ユーザーを作成すること、ユーザーを削除すること、テーブルに 1 つ以上の行を挿入すること、値をコピーすること、テーブル内からデータを選択すること (例えば、テーブルにクエリを行うこと)、クエリをキャンセルすること、またはクエリを中止することを行う操作、及び/または他の操作に対応してもよい。

【 0 0 1 8 】

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、読み取り要求及び/または書き込み要求 (並びに/または他のトランザクション要求) を様々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) から受け取り、次いで、それらを解析し、(単数または複数の) 関連するデータベース動作を実行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 (プライマリノードと呼ぶこともある) を含んでもよい。例えば、このデータベースエンジン・ヘッドノードは、複合的なクエリ及びジョインに対する結果を得るのに必要な一連のステップを作成してもよい。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データベース最適化ストレージシステムとの間の通信に加えて、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を管理してもよい。いくつかの実施形態では、以下でより詳細に説明するように、データベースエンジン・ヘッドノードは、トランザクションの順序付けを行うように構成されてもよい。これにより、特定の分離レベル (例えば、読み取り一貫性など) を保持しやすくすることができる。

10

20

30

40

50

【0019】

いくつかの実施形態では、データベース階層（またはより具体的には、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を一時的に保持するキャッシュを含んでもよい。このような実施形態では、かかるキャッシュ内のデータページを対象とする書き込み要求を受け取った場合、対応するREDOログレコードをストレージ層に送り出すのに加えて、データベースエンジンは、自らのキャッシュ内のコピーにその変更を適用してもよい。しかしながら、他のデータベースシステムとは異なり、このキャッシュに保持されているデータページを、ストレージ層に絶えずフラッシュさせなくてもよく、このデータページは、任意の時点で（例えば、キャッシュ済みコピーに対して直近に適用された、書き込み要求のREDOログレコードがストレージ層に送られて承認された後の任意の時点で）破棄され得る。このキャッシュは、種々の実施形態では、多くとも1つのライター（または複数のリーダー）によるキャッシュへの同時アクセスを制御するための様々なロッキング機構のいずれを実装してもよい。しかしながら、かかるキャッシュを含む実施形態では、複数のノードにわたってキャッシュを分散させなくてもよく、所与のデータベースインスタンスに対するデータベースエンジン・ヘッドノードにのみキャッシュが存在し得ることに留意されたい。従って、キャッシュのコヒーレンシまたは一貫性の問題を扱わなくてもよい。また、そうは言っても、それぞれがデータベースエンジン・ヘッドノードを備えた複数のデータベースインスタンスが存在してもよいことに留意されたい。

10

【0020】

いくつかの実施形態では、データベース階層は、システム内の読み取り用レプリカ（例えば、読み取り要求がルーティングされる可能性のある、データベース階層の種々のノード上のデータの読み取り専用コピー）を同期または非同期で利用することに対応してもよい。このような実施形態では、所与のデータベーステーブルに対するデータベースエンジン・ヘッドノードは、特定のデータページを対象とする読み取り要求を受け取った場合、その要求を、これらの読み取り専用コピーのうちの任意の1つ（または特定の1つ）にルーティングしてもよい。いくつかの実施形態では、データベースエンジン・ヘッドノード内のクライアント側ドライバは、これらの他のノードに、（例えば、これらのノードに内部キャッシュの無効化を促し、その後、これらのノードが、更新済みデータページの更新済みコピーをストレージ層から要求できるようにするために）キャッシュ済みのデータページに対する更新及び/または無効化について通知するように構成されてもよい。

20

30

【0021】

いくつかの実施形態では、クライアント側ドライバは、ボリュームに関するメタデータを保持してもよく、ストレージノード間にホップを追加せずとも、読み取り要求と書き込み要求を実行するの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送ってもよい。例えば、いくつかの実施形態では、データベーステーブルを変更する要求に回答して、クライアント側ドライバは、対象データページ用のストレージを実装している1つ以上のノードを決定し、上記変更を指定する（単数または複数の）REDOログレコードをこれらのストレージノードにルーティングするように構成されてもよい。次いで、ストレージノードは、REDOログレコード内で指定された変更を、将来のある時点で対象データページに適用する役割を果たしてもよい。書き込みが承認されてクライアント側ドライバに返されると、クライアント側ドライバは、ボリュームが持続的である時点を進めて、データベース階層に返すコミットを承認してもよい。先に述べたように、いくつかの実施形態では、クライアント側ドライバは、データページをストレージノードサーバに絶えず送らなくてもよい。これにより、ネットワークトラフィックが減少し得るのみならず、前述のデータベースシステムでのフォアグラウンド処理能力を制約するチェックポイント・スレッドまたはバックグラウンド・ライタースレッドが不要ともなり得る。

40

【0022】

いくつかの実施形態では、クライアント側ドライバは、本明細書で説明されるように、読み取り要求を受け取って複数のレコードを検索するデータベースエンジン・ヘッドノード

50

ドのために、開示されたトランザクションの順序付けを実行してもよい。例えば、データベースサービスのデータベースエンジン・ヘッドノードは、読み取り要求を受け取って、データベースサービスによって記憶されたレコードの読み取りを実行してもよい。別のデータベースエンジン・ヘッドノードは、トランザクション要求を受け取って、そのレコードに対してトランザクション（例えば、書き込みなど）を実行してもよい。読み取り要求を受け取ったデータベースエンジン・ヘッドノードは、その読み取りに関連付けられた時刻指示が、上記トランザクションに関連付けられた第2の時刻指示の閾値内にあると判定したことに基づき、潜在的な読み取り異常（例えば、ファジーリード、リードスキューなど）を検出し得る。潜在的な読み取り異常を検出したことに応答して、第1の時刻指示が第2の時刻指示よりも早い時点を示しているかどうかに関わらず、トランザクション要求によって指定されたトランザクションの実行後に、上記読み取りが実行され得る。ある事例では、読み取りを再試行することにより、その再試行については潜在的な読み取り異常が発生しないようにしてもよい。なお、データベースエンジン・ヘッドノードは、ある時刻に読み取り要求を受け取ってデータテーブルにクエリを行い、別の時刻にトランザクション要求を受け取ってデータテーブルを変更してもよい。以下で説明するように、通常リード、ファジーリード及びリードスキューに関する様々な例示的なタイミング図を図7A～Cに示す。

10

【0023】

いくつかの実施形態では、データベースエンジン・ヘッドノードのキャッシュによって多くの読み取り要求が提供されてもよい。しかしながら、書き込み要求には持続性が必要となる場合がある。これは、大規模な障害イベントがあまりにもよく起こるため、インメモリレプリケーションのみでは対処することができない場合があるためである。従って、本明細書で説明されたシステムは、ストレージ階層内のデータストレージを2つの領域として実装することにより、フォアグラウンド・レイテンシーパスに含まれるREDOログレコードの書き込み動作のコストを最小化するように構成され得る。この2つの領域とは、すなわち、データベース階層からREDOログレコードを受け取ったときにそれらが書き込まれる小さな追記専用のログ構造化領域、及びバックグラウンドにおいてデータページの新たなバージョンを作成するためにログレコードを1つに併合させたより大きな領域である。いくつかの実施形態では、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、当該ページの後ろ向き連鎖のログレコードにおける最後のREDOログレコードを指している各データページを対象としてインメモリ構造を維持してもよい。この手法により、読み取りが主にキャッシュされたアプリケーションにおけるものを含む、読み取り・書き込みが混在した作業負荷に対して良好な性能を得ることができる。

20

30

【0024】

Web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可能なサービスシステムアーキテクチャの一実施形態を図2に示す。例示した実施形態では、（データベースクライアント250a～250nと示した）いくつかのクライアントは、ネットワーク260を介してWebサービスプラットフォーム200と相互作用するように構成され得る。Web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の1つ以上のインスタンス、分散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つ以上の他の仮想コンピューティングサービス230と相互作用するように構成され得る。なお、所与の構成要素の1つ以上のインスタンスが存在し得る場合、本明細書においてその構成要素への言及は、単数形または複数形のいずれでもなされることがある。しかしながら、いずれかの形式を用いることは、他の形式を排除することを意図するものではない。

40

【0025】

様々な実施形態では、図2に示した構成要素は、コンピュータハードウェア内に直接実装され、コンピュータハードウェア（例えば、マイクロプロセッサまたはコンピュータシステム）によって直接的もしくは間接的に実行可能な命令として実装され、またはこれらの技術の組み合わせを用いて実装され得る。例えば、図2の構成要素は、いくつかのコンピューティングノード（または単に、ノード）を含むシステムによって実装され得る。こ

50

これらのノードのそれぞれは、図8に示され、以下で説明されるコンピュータシステムの実施形態と同様のものであってよい。様々な実施形態では、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）の機能は、特定のノードによって実装されてもよく、またはいくつかのノードにわたって分散されてもよい。いくつかの実施形態では、所与のノードは、2つ以上のサービスシステム構成要素（例えば、2つ以上のデータベースサービスシステム構成要素）の機能を実装してもよい。

【0026】

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を求める要求（例えば、トランザクション要求、読み取り要求など）を含むWebサービス要求を、Webサービスプラットフォーム200にネットワーク260を介して送信するように構成可能な任意の種類のクライアントを含んでもよい。例えば、所与のクライアント250は、適切なバージョンのWebブラウザを搭載してもよく、または、Webブラウザによって提供される実行環境に対する拡張機能、もしくはその環境内の拡張機能として実行するように構成されたプラグインモジュールもしくは他の種類のコードモジュールを搭載してもよい。あるいは、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、永続的記憶装置のリソースを利用して1つ以上のデータベーステーブルを記憶し、かつ/またはそれらにアクセスし得る、データベースアプリケーションなどのアプリケーション（またはそのユーザーインターフェース）、メディアアプリケーション、オフィスアプリケーション、またはその他のアプリケーションを含んでもよい。いくつかの実施形態では、こうしたアプリケーションは、（例えば、適切なバージョンのHTTP（Hypertext Transfer Protocol）に対して）十分なプロトコルサポートを含むことにより、全ての種類のWebベースのデータ用に完全なブラウザサポートを必ずしも実装せずに、Webサービス要求を生成及び処理できるようにしてもよい。すなわち、クライアント250は、Webサービスプラットフォーム200と直接相互作用するように構成された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、REST（Representational State Transfer）スタイルのWebサービスアーキテクチャ、ドキュメントベースのもしくはメッセージベースのWebサービスアーキテクチャ、または他の適切なWebサービスアーキテクチャに従って、Webサービス要求を生成するように構成されてもよい。

【0027】

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、上記アプリケーションに透過的な方法で、データベーステーブルの、Webサービスベースのストレージに対するアクセスを他のアプリケーションに提供するように構成されてもよい。例えば、クライアント250は、本明細書で説明されたストレージモデルの適切な変形例に従ってストレージを提供するオペレーティングシステムまたはファイルシステムと一体化されるように構成されてもよい。しかしながら、この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/またはフォルダで構成される従来のファイルシステム階層構造などの、別のストレージインターフェースをアプリケーションに提示してもよい。こうした実施形態では、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アプリケーションを修正しなくてもよい。むしろ、Webサービスプラットフォーム200へのインターフェースの詳細については、オペレーティングシステム環境内で実行するアプリケーションに代わり、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得る。

【0028】

クライアント250は、ネットワーク260を介し、Webサービス要求（例えば、トランザクション要求、読み取り要求など）をWebサービスプラットフォーム200に伝達してもよく、そこから応答を受け取ってもよい。様々な実施形態では、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の間のWebベース通信を確立するのに必要なネットワーク接続ハードウェアとネットワーク接続プロトコルの任意の適切

10

20

30

40

50

な組み合わせを含んでもよい。例えば、ネットワーク260は、インターネットを集合的に実装する様々な遠距離通信ネットワーク及びサービスプロバイダを通常含んでもよい。ネットワーク260は、パブリックまたはプライベートな無線ネットワークに加えて、ローカルエリアネットワーク(local area network: LAN)またはワイドエリアネットワーク(wide area network: WAN)などのプライベートネットワークを含んでもよい。例えば、所与のクライアント250とWebサービスプラットフォーム200の両方を、内部ネットワークを独自に有する企業内にそれぞれ提供してもよい。こうした実施形態では、ネットワーク260は、インターネットとWebサービスプラットフォーム200の間に加えて、所与のクライアント250とインターネットの間でネットワーク接続リンクを確立するのに必要な、ハードウェア(例えば、モデム、ルーター、スイッチ、負荷分散装置、プロキシサーバなど)、及びソフトウェア(例えば、プロトコルスタック、会計ソフトウェア、ファイアウォール/セキュリティソフトウェアなど)を含んでもよい。なお、いくつかの実施形態では、クライアント250は、パブリックインターネットではなくプライベートネットワークを用いてWebサービスプラットフォーム200と通信してもよい。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)として同一企業内に提供されてもよい。その場合、クライアント250は、プライベートネットワーク260(例えば、インターネットベースの通信プロトコルを利用し得るが、公にアクセス可能ではないLANまたはWAN)全体を通じてプラットフォーム200と通信してもよい。

10

20

【0029】

一般的に言えば、Webサービスプラットフォーム200は、データページ(またはそのレコード)にアクセスする要求などのWebサービス要求を受け取って処理するように構成された1つ以上のサービスエンドポイントを実装するように構成され得る。例えば、Webサービスプラットフォーム200は、特定のエンドポイントを実装するように構成されたハードウェア及び/またはソフトウェアを含むことにより、そのエンドポイントを対象とするHTTPベースのWebサービス要求を正しく受け取って処理できるようにしてもよい。一実施形態では、Webサービスプラットフォーム200は、クライアント250からWebサービス要求を受け取り、それらを、データベースサービス210、分散データベース最適化ストレージサービス220及び/または別の仮想コンピューティングサービス230を実装するシステムの構成要素に転送して処理できるように構成されたサーバシステムとして実装されてもよい。他の実施形態では、Webサービスプラットフォーム200は、大量のWebサービス要求処理負荷を動的に管理するように構成された負荷分散機能及び他の要求管理機能を実装するいくつかの別個のシステムとして(例えば、クラスタポロジを用いて)構成されてもよい。様々な実施形態では、Webサービスプラットフォーム200は、RESTスタイルまたはドキュメントベースの(例えば、SOAPベースの)種類のWebサービス要求に対応するように構成されてもよい。

30

【0030】

クライアントのWebサービス要求に対してアドレス可能エンドポイントとして機能することに加えて、いくつかの実施形態では、Webサービスプラットフォーム200は、様々なクライアント管理機能を実装してもよい。例えば、プラットフォーム200は、要求中のクライアント250の識別情報、クライアント要求の数及び/または頻度、クライアント250に代わって記憶もしくは検索されたデータテーブル(もしくはそのレコード)のサイズ、クライアント250によって使用されるストレージ帯域全体、クライアント250によって要求されたストレージのクラス、またはその他の測定可能なクライアントの使用状況パラメータを追跡することなどにより、ストレージリソースを含めた、クライアントによるWebサービスの使用状況の測定と計算とを連携して行うようにしてもよい。プラットフォーム200は、いくつかの実施形態では、クライアントのWebサービス要求を、このプラットフォームの各データベースインスタンスの特定のデータベースエンジン・ヘッドノードに分配するように構成されてもよい。簡単な例として、時刻1にて、

40

50

プラットフォーム 200 は、読み取り要求をデータベースエンジン・ヘッドノード 1 に分配してもよく、時刻 3 にて、プラットフォームは、書き込み要求をデータベースエンジン・ヘッドノード 2 に分配してもよい。プラットフォーム 200 は、財務会計システム及び課金システムを実装してもよく、または使用状況データのデータベースを保持してもよい。この使用状況データは、クライアントの利用活動に対する報告及び課金を行う外部システムによって照会及び処理が行われる。ある実施形態では、プラットフォーム 200 は、様々なストレージサービスシステムの運用上の測定基準を収集、監視、及び/または集約するように構成されてもよい。このような測定基準には、クライアント 250 から受け取った要求の割合及び種類を反映している測定基準、かかる要求によって利用される帯域、かかる要求に対するシステムの処理レイテンシー、システム構成要素の利用率（例えば、ストレージサービスシステム内のネットワーク帯域及び/もしくはストレージの利用率）、要求に起因するエラーの割合及び種類、記憶され、要求されたデータページもしくはそのレコードの特性（例えば、サイズ、データ種類など）、またはその他の適切な測定基準などがある。いくつかの実施形態では、こうした測定基準は、システム構成要素の調整及び保守を行うためにシステム管理者によって用いられてもよい。一方、他の実施形態では、こうした測定基準（または、当該測定基準の関連部分）をクライアント 250 に公開することにより、当該クライアントが、データベースサービス 210、分散データベース最適化ストレージサービス 220 及び/または別の仮想コンピューティングサービス 230（またはそれらのサービスを実装する下位システム）における自らの使用状況を監視できるようにしてもよい。

10

20

【0031】

いくつかの実施形態では、プラットフォーム 200 は、ユーザー認証手順及びユーザーアクセス制御手順を実装してもよい。例えば、特定のデータベーステーブルにアクセスする所与の Web サービス要求に対して、プラットフォーム 200 は、その要求に関連付けられたクライアント 250 が当該特定のデータベーステーブルにアクセスする権限を与られているかどうかを確認するように構成され得る。プラットフォーム 200 は、こうした認証を、例えば、当該特定のデータベーステーブルに関連付けられた証明書と照合して識別情報、パスワードもしくは他の証明書を評価すること、または当該特定のデータベーステーブル用のアクセス制御リストと照合して、要求された当該特定のデータベーステーブルへのアクセスを評価することによって判定してもよい。例えば、当該特定のデータベーステーブルにアクセスするための正当な証明書をクライアント 250 が持っていない場合、プラットフォーム 200 は、例えば、要求中のクライアント 250 に対してエラー状態を示す応答を返すことにより、対応する Web サービス要求を拒絶してもよい。様々なアクセス制御ポリシーが、アクセス制御情報のレコードまたはリストとして、データベースサービス 210、分散データベース最適化ストレージサービス 220 及び/または仮想コンピューティングサービス 230 によって記憶されてもよい。

30

【0032】

なお、Web 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 を実装するデータベースシステムの機能にクライアント 250 がアクセスできるようにするための主要なインターフェースに相当し得るが、かかる機能への唯一のインターフェースに相当する必要はない。例えば、Web サービスインターフェースとは異なる代替 API を用いることにより、データベースシステムを提供する企業内部のクライアントが、Web サービスプラットフォーム 200 を回避できるようにしてもよい。なお、本明細書で説明された実施例の多くでは、分散データベース最適化ストレージサービス 220 は、データベースサービスをクライアント 250 に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の内部にあってよく、外部のクライアント（例えば、ユーザーまたはクライアントアプリケーション）に公開されなくてもよい。このような実施形態では、内部の「クライアント」（例えば、データベースサービス 210）は、分散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 の間の実線で示すように、ローカルネットワークまたはプライベートネットワークを介して分散データベース最適化ストレージサ

40

50

ービス 220 に（例えば、API を通じて、これらのサービスを実装するシステム間で直接的に）アクセスしてもよい。このような実施形態では、データベーステーブルに記憶する際にクライアント 250 に代わって分散データベース最適化ストレージサービス 220 を利用することは、それらのクライアントにとって透過的に行われてもよい。他の実施形態では、Web サービスプラットフォーム 200 を通じてクライアント 250 に分散データベース最適化ストレージサービス 220 を公開することにより、データベースを管理するための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対し、データベーステーブルまたは他の情報のストレージを提供してもよい。これは、図 2 において Web サービスプラットフォーム 200 と分散データベース最適化ストレージサービス 220 の間の点線によって示されている。このような実施形態では、分散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、分散データベース最適化ストレージサービス 220 にネットワーク 260 経由で（例えば、インターネットを介して）アクセスしてもよい。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、コンピューティングサービス 230 を実行するのに用いられるオブジェクトをクライアント 250 に代わって記憶する分散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（例えば、API を通じて、仮想コンピューティングサービス 230 と分散データベース最適化ストレージサービス 220 の間で直接的に）受け取るように構成され得る。これは、図 2 において仮想コンピューティングサービス 230 と分散データベース最適化ストレージサービス 220 の間の点線によって示されている。ある場合には、プラットフォーム 200 の会計サービス及び/または証明書発行サービスは、管理用クライアントなどの内部のクライアントにとっては、または同一企業内のサービス構成要素間では不要となる場合がある。

10

20

30

40

50

【0033】

なお、様々な実施形態では、種々のストレージポリシーが、データベースサービス 210 及び/または分散データベース最適化ストレージサービス 220 によって実装され得る。このようなストレージポリシーの実施例としては、持続性ポリシー（例えば、記憶されるデータベーステーブル（もしくはそのデータページ）のインスタンスの数、及びそれらが記憶される種々のノードの数を示すポリシー）、並びに/または負荷分散ポリシー（このポリシーにより、データベーステーブルもしくはそのデータページが、要求のトラフィックを一様にするために種々のノード、ボリューム及び/またはディスクにわたって分散され得る）を挙げることができる。加えて、上記サービスの中の多様なサービスにより、各種の記憶済み項目に異なるストレージポリシーを適用してもよい。例えば、いくつかの実施形態では、分散データベース最適化ストレージサービス 220 は、データページの持続性よりも REDO ログレコードの持続性が高くなるように実装されてもよい。

【0034】

図 3 は、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様々な構成要素を示すブロック図である。本実施例では、データベースシステム 300 は、いくつかのデータベーステーブルのそれぞれに対する各データベースエンジン・ヘッドノード 320、及び分散データベース最適化ストレージサービス 310（これは、データベースクライアント 350 a ~ 350 n 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見えても見えなくてもよい）を含む。本実施例に示すように、データベースクライアント 350 a ~ 350 n の 1 つ以上が、データベースヘッドノード 320（例えば、ヘッドノード 320 a、ヘッドノード 320 b またはヘッドノード 320 c であって、これらのノードのそれぞれは、各データベースインスタンスの構成要素である）にネットワーク 360 経由でアクセスしてもよい（例えば、これらの構成要素は、ネットワークアドレス可能かつデータベースクライアント 350 a ~ 350 n にとって利用可能であってよい）。しかしながら、異なる実施形態では、分散データベース最適化ストレージサービス 310 は、データベースシステムによって利用されて、1 つ以上のデータベーステーブルのデータページ（並びに REDO ログレコード及び/またはそれに関連付けられた他のメタデータ）をデータベースクライアン

ト 3 5 0 a ~ 3 5 0 n に代わって記憶し、本明細書で説明されるデータベースシステムの他の機能を実行し得るが、ネットワークアドレス可能かつストレージクライアント 3 5 0 a ~ 3 5 0 n にとって利用可能であってもよく、そうでなくてもよい。例えば、いくつかの実施形態では、分散データベース最適化ストレージサービス 3 1 0 は、ストレージクライアント 3 5 0 a ~ 3 5 0 n にとって見えないようにして、様々な記憶、アクセス、変更ロギング、リカバリ、ログレコード操作、及び/または空き領域管理動作を行ってもよい。

【 0 0 3 5 】

先に述べたように、各データベースインスタンスは、1つのデータベースエンジン・ヘッドノード 3 2 0 を含み得る。このヘッドノードは、要求（例えば、トランザクション要求など）を様々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）から受け取り、次いでそれらを解析し、それらを最適化し、実行計画に展開することにより、関連するデータベース動作（単数または複数）を実行する。図 3 に示した例では、データベースエンジン・ヘッドノード 3 2 0 a の、クエリの解析・最適化・実行用構成要素 3 0 5 は、データベースクライアント 3 5 0 a から受け取ったクエリであって、データベースエンジン・ヘッドノード 3 2 0 a が構成要素になっているデータベースインスタンスを対象とするクエリに対してこれらの機能を実行し得る。いくつかの実施形態では、クエリの解析・最適化・実行用構成要素 3 0 5 は、データベースクライアント 3 5 0 a にクエリ応答を返してもよい。この応答には、必要に応じて、書き込み承認、要求されたデータページ（もしくはその各部）、エラーメッセージ及びまたは他の応答が含まれ得る。本実施例に示すように、データベースエンジン・ヘッドノード 3 2 0 a は、クライアント側ストレージサービスドライバ 3 2 5 を含んでもよい。このストレージサービスドライバは、読み取り要求及び/または（例えば、書き込みの）REDO ログレコードを分散データベース最適化ストレージサービス 3 1 0 内の様々なストレージノードにルーティングし、書き込み承認を分散データベース最適化ストレージサービス 3 1 0 から受け取り、要求されたデータページを分散データベース最適化ストレージサービス 3 1 0 から受け取り、かつ/またはデータページ、エラーメッセージもしくは他の応答をクエリの解析・最適化・実行用構成要素 3 0 5 に返してもよい（さらに、この構成要素は、このデータページ、エラーメッセージまたは他の応答をデータベースクライアント 3 5 0 a に返してもよい）。

【 0 0 3 6 】

本実施例では、データベースエンジン・ヘッドノード 3 2 0 a はデータページキャッシュ 3 3 5 を含んでおり、その内部には、最近アクセスされたデータページが一時的に保持され得る。図 3 に示すように、データベースエンジン・ヘッドノード 3 2 0 a は、トランザクション・一貫性管理構成要素 3 3 0 を含んでもよい。この構成要素は、データベースエンジン・ヘッドノード 3 2 0 a が構成要素になっているデータベースインスタンスにおいてトランザクション性及び一貫性を提供する役割を果たしてもよい。例えば、この構成要素は、データベースインスタンスの原子性、一貫性及び独立性の諸特性、並びにそのデータベースインスタンス対象とするトランザクションを確保する役割を果たしてもよい。例えば、データベースサービスのデータベースエンジン・ヘッドノードは、読み取り要求を受け取って、データベースサービスによって記憶されたレコードの読み取りを実行してもよい。別のデータベースエンジン・ヘッドノードは、トランザクション要求を受け取って、そのレコードに対してトランザクション（例えば、書き込みなど）を実行してもよい。次いで、読み取り要求を受け取ったデータベースエンジン・ヘッドノードのトランザクション・一貫性管理構成要素 3 3 0 は、その読み取りに関連付けられた時刻指示が、上記トランザクションに関連付けられた第 2 の時刻指示の閾値内にあると判定したことに基づき、潜在的な読み取り異常（例えば、ファジーリード、リードスキューなど）を検出し得る。潜在的な読み取り異常を検出したことに応答して、第 1 の時刻指示が第 2 の時刻指示よりも早い時点を示しているかどうかに関わらず、トランザクション要求によって指定されたトランザクションの後に上記読み取りが実行され得る。場合によっては、読み取りを

10

20

30

40

50

再試行することにより、その再試行については潜在的な読み取り異常が発生しないようにしてもよい。

【0037】

図3に示すように、データベースエンジン・ヘッドノード320aは、トランザクションログ340及びUNDOログ345を含んでもよい。これらのログをトランザクション・一貫性管理構成要素330によって利用することにより、様々なトランザクションのステータスを追跡し、コミットされていないトランザクションに関する、ローカルにキャッシュされた任意の結果をロールバックしてもよい。

【0038】

なお、図3に示した他のデータベースエンジン・ヘッドノード320のそれぞれ（例えば、320b及び320c）は同様の構成要素を含んでもよく、当該データベースエンジン・ヘッドノードが構成要素になっている各データベースインスタンスを対象とし、かつデータベースクライアント350a～350nの1つ以上によって受け取られるクエリ及び/または他のトランザクションに対して同様の機能を実行してもよい。例えば、開示されたトランザクションの順序付け技術は、本明細書で説明するように、2つの異なるデータベースエンジン・ヘッドノードが、闕時間内に同一データにアクセスしている（例えば、一方が読み取り中であり、一方が書き込み中である）シナリオにおいて実施され得る。

【0039】

分散データベース最適化ストレージシステムの一実施形態をブロック図によって図4に示す。本実施例では、データベースシステム400は、分散データベース最適化ストレージシステム410を含んでおり、このストレージシステムは、インターコネクト460を介してデータベースエンジン・ヘッドノード420と通信する。図3に示した実施例と同様に、データベースエンジン・ヘッドノード420は、クライアント側ストレージサービスドライバ425を含んでもよい。本実施例では、分散データベース最適化ストレージシステム410は、（430、440及び450として示したものを含む）複数のストレージシステムサーバノードを含む。これらのノードのそれぞれは、各ノードが記憶する（単数または複数の）セグメントに対応したREDOログ及びデータページ用のストレージ、並びに様々なセグメント管理機能を実行するように構成されたハードウェア及び/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：（ローカルに、例えば、ストレージノード内で行われる）複製、データページを生成するためのREDOログの併合、ログ管理（例えば、ログレコードの操作）、クラッシュリカバリ、及び/または（例えば、セグメントに対する）空き領域管理、のいずれかまたは全ての少なくとも一部を実行するように構成されたハードウェア及び/またはソフトウェアを含んでもよい。各ストレージシステムサーバノードは、複数の付属ストレージ装置（例えば、SSD）を備えてもよい。このストレージ装置には、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の加入者）に代わってデータブロックが記憶されてもよい。

【0040】

図4に示した実施例では、ストレージシステムサーバノード430は、（単数または複数の）データページ433、（単数または複数の）セグメントREDOログ435、セグメント管理機能437及び付属SSD471～478を含む。さらになお、「SSD」という表示は、その下位ハードウェアに関係なくソリッドステートドライブを表しても表さなくてもよいが、より一般的にはローカルのブロックストレージボリュームを表し得る。同様に、ストレージシステムサーバノード440は、データページ（単数または複数）443、（単数または複数の）セグメントREDOログ445、セグメント管理機能447及び付属SSD481～488を含み、ストレージシステムサーバノード450は、（単数または複数の）データページ453、（単数または複数の）セグメントREDOログ455、セグメント管理機能457及び付属SSD491～498を含む。

【0041】

いくつかの実施形態では、分散データベース最適化ストレージシステム内のストレージ

10

20

30

40

50

システムサーバノードのそれぞれは、一式のプロセスを実装し得る。これらのプロセスは、ノードサーバのオペレーティングシステム上で動作するものであり、このオペレーティングシステムは、データベースエンジン・ヘッドノードとの通信を管理して、例えば、REDOログを受け取り、データページを送り返すことなどを行う。いくつかの実施形態では、分散データベース最適化ストレージシステムに書き込まれた全データブロックを、（例えば、キー値を利用した永続的なりモートバックアップストレージシステムを用いて）長期ストレージ及び/またはアーカイブストレージにバックアップしてもよい。

【0042】

図5は、一実施形態に係る、データベースシステム内の別個の分散データベース最適化ストレージシステムの利用を示すブロック図である。本実施例では、1つ以上のクライアントプロセス510は、データベースエンジン520及び分散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保持される1つ以上のデータベーステーブルにデータを記憶し得る。図5に示した実施例では、データベースエンジン520は、データベース階層要素560、及びクライアント側ドライバ540（このドライバは、分散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要素560の間のインターフェースとして機能する）を含む。いくつかの実施形態では、データベース階層要素560は、図3のクエリの解析・最適化・実行用構成要素305及びトランザクション・一貫性管理構成要素330（例えば、トランザクションの順序付け）によって実行されるものなどの機能を実行してもよく、かつ/またはデータページ、トランザクションログ及び/もしくはUNDOログ（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UNDOログ345によって記憶されるものなど）を記憶してもよい。

【0043】

本実施例では、1つ以上のクライアントプロセス510は、データベース・クエリ要求515（この要求は、ストレージノード535a～535nの1つ以上に記憶されたデータを対象とする、読み取り要求、及び/または書き込み要求、及び/または他のトランザクション要求、を含み得る）をデータベース階層要素560に送ってもよく、（例えば、書き込み承認及び/または要求されたデータを含む）データベース・クエリ応答517をデータベース階層要素560から受け取ってもよい。データページに書き込む要求を含む各データベース・クエリ要求515を解析及び最適化して、1つ以上の書き込みレコード要求541が生成され得る。これらの要求は、引き続いて分散データベース最適化ストレージシステム530にルーティングできるように、クライアント側ドライバ540に送られてもよい。本実施例では、クライアント側ドライバ540は、各書き込みレコード要求541に対応する1つ以上のREDOログレコード531を生成してもよく、それらを、分散データベース最適化ストレージシステム530を構成するストレージノード535の特定のストレージノードに送ってもよい。分散データベース最適化ストレージシステム530は、各REDOログレコード531について、対応する書き込み承認532を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）返してもよい。クライアント側ドライバ540は、これらの書き込み承認をデータベース階層要素560に（書き込み応答542として）渡してもよく、次いで、このデータベース階層要素は、対応する応答（例えば、書き込み承認）を1つ以上のクライアントプロセス510にデータベース・クエリ応答517の1つとして送ってもよい。

【0044】

本実施例では、データページを読み取る要求を含む各データベース・クエリ要求515を解析及び最適化して、1つ以上の読み取りレコード要求543が生成され得る。これらの要求は、引き続いて分散データベース最適化ストレージシステム530にルーティングできるように、クライアント側ドライバ540に送られてもよい。本実施例では、クライアント側ドライバ540は、これらの要求を、分散データベース最適化ストレージシステム530を構成するストレージノード535の特定のストレージノードに送ってもよく、分散データベース最適化ストレージシステム530は、要求されたデータページ533を

10

20

30

40

50

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）返してもよい。クライアント側ドライバ 540 は、返されたデータページをデータベース階層要素 560 に返却データレコード 544 として送ってもよく、次いで、データベース階層要素 560 は、そのデータページを 1 つ以上のクライアントプロセス 510 にデータベース・クエリ応答 517 として送ってもよい。

【0045】

いくつかの実施形態では、様々なエラー及び/またはデータ損失のメッセージ 534 を分散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）送ってもよい。これらのメッセージは、クライアント側ドライバ 540 からデータベース階層要素 560 にエラー及び/または損失の報告メッセージ 545 として渡されてもよく、次いで、データベース・クエリ応答 517 と共に（またはその代わりに）1 つ以上のクライアントプロセス 510 に渡されてもよい。

10

【0046】

いくつかの実施形態では、分散データベース最適化ストレージシステム 530 の API 531 ~ 534、及びクライアント側ドライバ 540 の API 541 ~ 545 は、データベースエンジン 520 が分散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ように、分散データベース最適化ストレージシステム 530 の機能をデータベースエンジン 520 に公開してもよい。例えば、（クライアント側ドライバ 540 を介した）データベースエンジン 520 は、これらの API を通じて REDO ログレコードを書き込むこと、またはデータページを要求することにより、データベースエンジン 520 と分散データベース最適化ストレージシステム 530 の組み合わせによって実装されるデータベースシステムの様々な動作（例えば、記憶、アクセス、変更ロギング、リカバリ、及び/または空き領域管理動作）を実行する（またはその性能を促進する）ようにしてもよい。図 5 に示すように、分散データベース最適化ストレージシステム 530 は、ストレージノード 535 a ~ 535 n にデータブロックを記憶してもよく、これらのノードのそれぞれは、複数の付属 SSD を備えてもよい。いくつかの実施形態では、分散データベース最適化ストレージシステム 530 は、各種の冗長方式を適用することにより、記憶されたデータブロックに対して高度な持続性を提供してもよい。

20

【0047】

なお、様々な実施形態では、図 5 における、データベースエンジン 520 と分散データベース最適化ストレージシステム 530 の間の API コール及び API 応答（例えば、API 531 ~ 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要素 560 の間の API コール及び API 応答（例えば、API 541 ~ 545）は、安全なプロキシ接続（例えば、ゲートウェイ制御プレーンによって管理される接続）を介して実行されてもよく、または、パブリックネットワークを介して、もしくは代替的に、バーチャルプライベートネットワーク（virtual private network: VPN）接続などのプライベートチャネルを介して実行されてもよい。本明細書で説明されたデータベースシステムの構成要素に対する、かつ/またはこれらの構成要素間の上記の及び他の API は、SOAP (Simple Object Access Protocol) 技術及び REST (Representational state transfer) 技術を含む種々の技術に従って実装されてもよいが、これらに限定されることはない。例えば、これらの API は、SOAP API または RESTful API として実装されてもよいが、必ずしもそのように実装されるとは限らない。SOAP は、Web ベースのサービスに関連して情報を交換するためのプロトコルである。REST は、分散ハイパーメディアシステム用の構築様式である。RESTful API (RESTful Web サービスと呼ぶこともある) は、HTTP 及び REST 技術を用いて実装される Web サービス API である。本明細書で説明された API は、いくつかの実施形態では、C、C++、Java、C# 及び Perl を含むがこれらに限定されない様々な言語を用いたクライアントライブラリによってラップされて、データベースエンジン 5

30

40

50

20及び/または分散データベース最適化ストレージシステム530との一体化に対応し得る。

【0048】

上述したように、いくつかの実施形態では、データベースシステムの機能要素を、データベースエンジンによって実行される機能要素と、別個の分散データベース最適化ストレージシステムにおいて実行される機能要素との間で分割してもよい。ある具体的な例では、あるものをデータベーステーブルに挿入するために（例えば、1つのデータブロックにレコードを追加することによってそのデータブロックを更新するために）クライアントプロセス（またはそのスレッド）から要求を受け取ったことに応答して、データベースエンジン・ヘッドノードの1つ以上の構成要素は、クエリの解析、最適化及び実行を行ってもよく、そのクエリの各部を、トランザクション・一貫性管理構成要素に送ってもよい。

10

【0049】

トランザクション・一貫性管理構成要素は、同時に同一行を修正しようとしているクライアントプロセス（またはそのスレッド）が他にないことを保証し得る。例えば、トランザクション・一貫性管理構成要素は、この変更が原子的に、一貫して、持続的に、かつそのデータベース内で独立して実行されることを保証する役割を果たし得る。例えば、トランザクション・一貫性管理構成要素は、データベースエンジン・ヘッドノードのクライアント側ストレージサービスと連携して動作することにより、分散データベース最適化ストレージサービス内のノードの1つに送るべきREDOログレコードを生成し、それを、このトランザクションに対してACID特性が確実に満たされる順序及び/またはタイミングで、（他のクライアント要求に応答して生成された他のREDOログと共に）分散データベース最適化ストレージサービスに送ってもよい。そのREDOログレコードを受け取ると、対応するストレージノードは、該当するデータブロックを更新してもよく、そのデータブロックに対するREDOログ（例えば、そのデータブロックを対象とする全ての変更が記録されたレコード）を更新してもよい。いくつかの実施形態では、データベースエンジンは、この変更に対するUNDOログレコードを生成する役割を果たしてもよく、そのUNDOログに対するREDOログレコードを生成する役割もまた果たしてもよい。その上で、それらのレコードの両方をローカル的に（データベース階層内で）利用して、トランザクション性を保証するようにしてもよい。さらに、様々な実施形態では、トランザクション・一貫性管理構成要素は、トランザクションの順序付けを行うように構成され得る。例えば、トランザクション・一貫性管理構成要素は、ほぼ同時に（例えば、互いの閾値内の一貫性ポイントで）複数のデータベースエンジン・ヘッドノードがトランザクションの実行（例えば、読み取り及びコミット）を試みる状況において、潜在的な読み取り異常（例えば、ファジーリード、リードスキューなど）を検出するように構成され得る。潜在的な読み取り異常を検出したことに応答して、トランザクション・一貫性管理構成要素はさらに、その読み取りがより早い時刻に関連付けられている場合であっても、他のトランザクションの後に上記読み取りを発生させるように構成され得る。

20

30

【0050】

次に、図6に目を向けると、様々な実施形態では、データベースシステム300（または、データの読み取り、書き込み及び記憶に使用可能なデータベースサービス以外のあるシステム）は、トランザクションの順序付けを行うように構成されてもよい。図6の方法は、データベースエンジン・ヘッドノード320a、320b、320cなどのトランザクション・一貫性管理構成要素330及び/またはクライアント側ドライバなどの、分散データベースシステムの様々な構成要素（例えばノード）によって行われるものとして説明され得るが、この方法は、ある場合には特定の構成要素によって行われる必要はない。例えば、ある場合には、図6の方法は、いくつかの実施形態に係る、他の構成要素またはコンピュータシステムによって行われてもよい。あるいは、場合によっては、データベースシステム300の構成要素は、図3の実施例に示したものと異なるように、組み合わせられてもよく、または存在してもよい。様々な実施形態では、図6の方法は、分散データベースシステムの1つ以上のノードによって行われてもよく、これらのノードの1つが図

40

50

8のコンピュータシステムとして示されている。図6の方法は、トランザクションの順序付けを行う方法の例示的な実施態様として示されている。他の実施態様では、図6の方法は、図示したブロック以外の追加ブロック、または図示したブロックよりも少ないブロックを含み得る。

【0051】

610では、(例えば、データベースサービスまたは他のサービスによって記憶された)レコードの読み取りを実行する読み取り要求、及びそのレコードに対するトランザクションを実行するトランザクション要求を、例えば、(例えば、データベースサービスまたは他のサービスの)1つ以上のクライアントから受け取ってもよい。一実施形態では、この読み取り要求を、SELECTステートメントまたは他の要求として受け取ってもよい。この読み取り要求は、レコードのスナップショットを見ることが出来るスナップショットポイントの時点を求める要求と呼ばれることもある。上記トランザクション要求は、UPDATE、INSERT、または(例えば、データベースの)レコードを変更するのに使用可能な他のトランザクション(例えば、書き込みトランザクション)であってもよく、そのトランザクション要求をコミットしてもよい。様々な実施形態では、トランザクション要求及び読み取り要求は、記憶されているレコードに対して同時にアクセスし得る種々のノード(例えば、読み取り/書き込みを行うことができる複数のプライマリノード及び/または1つのプライマリノード、並びにその時点でレコードの読み取りのみが可能な読み取り用レプリカ)によって受け取られてもよい。例えば、一実施形態では、図2のWebサービスプラットフォーム200が、読み取り要求及びトランザクション要求を受け取り、種々のデータベースインスタンスの種々のノードにそれらをルーティングしてもよい。なお、図1~5における上記の実施例では、データベース階層とストレージ階層とが分かれている場合について説明したが、他の例では、ストレージ階層は、データベースインスタンスから分かれていなくてもよい。さらに、他の実施例では、本システムはデータベースサービスでなくともよく、その代わりに、記憶データの読み取り及び書き込みを実行可能な別のシステムであってもよい。さらになお、ライターである2つのノードは、両者の間でロックしていてもよい。ただし、同時処理を許可するために、リーダーであるノードとライターであるノードとは、両者の間でロックしていなくてもよい。

【0052】

一実施形態では、様々なノードは、トランザクションの順序付けに用いられ得るクロックをノードごとに保持してもよい。これらのクロックは、ノード全体で同期されていてもよく、互いの範囲内で±の精度を有し得る。これらのクロックの精度は、マルチノードシステムにおいてゼロではない場合がある。そのため、イベントがほぼ同時に発生する可能性があるとして、これらのイベントにおいて正確に順序付けを行うこと(因果関係)及び特定の分離レベルを維持することが難しくなる場合がある。例えば、読み取り一貫性の分離レベルには、以下の特性が含まれ得る。すなわち、ステートメントの開始時において、別のノードからコミットされたものは何でもその時点において見える時点がある、コミットされていない別のノードからの変更はその時点において見えない、及びコミットされようとなかろうと、そのノード自体からの変更を見ることが出来る、という特性である。2つの時刻、すなわちAとBを所与として、次の3つのシナリオが発生し得る。A<B(Bの前にAが起こる)、A>B(Bの後にAが起こる)、及びA=B(AとBがほぼ同時に起こり、AとBが精度ウィンドウ内にあるようになっている)。

【0053】

いくつかの実施形態では、これらのクロックによって保持される時刻は、タイムスタンプ(例えば、2012年15日、20:00.35 GMT)であってもよく、他の実施形態では、これらの時刻は、ログシーケンス番号(log sequence number: LSN)などの、時刻を示す単調増加的な値であってもよい。この値は、システムのノードにわたる通信が発生すると増加し得る。LSNの実施例では、値が単調増加しているため、LSN100はLSN105よりも早い時点を示し得る。なお、LSN時刻の間において各数値を用いる必要はない。従って、LSN100及び105は、一実施例で

10

20

30

40

50

は、最も密接して割り当てられた2つのLSNであり得る。別の実施例では、LSN100~105のそれぞれを用いてもよい。

【0054】

620に示すように、第1の時刻指示と第2の時刻指示とを、読み取りとトランザクションとにそれぞれ関連付けてもよい。例えば、いくつかの実施形態では、要求(例えば、トランザクション、読み取りなど)を受け取ると、その要求を受け取ったノードは、そのヘッドノードの個別のクロックに基づき、そのトランザクションに1つ以上の時刻を割り当ててもよい。例えば、ヘッドノードは、読み取り要求をLSN100にて受け取ってもよく、一貫性の時点(スナップショット時刻)Tsとして100を割り当ててもよい。Tsは、スナップショット時刻が作成された時を表し得る。別の実施例として、ヘッドノード(例えば、別のヘッドノード)は、テーブルを更新する要求をLSN101にて受け取ってもよい。上記ノードは、LSN102にてテーブルを更新し、次いで、その更新をLSN103にてコミットしてもよい。このような実施例では、上記ノードは、トランザクションがコミットされた時を表すコミット時刻Tcとして103を割り当ててもよい。なお、トランザクションに関連付けられた他の時刻を割り当ててもよい。例えば、Tcに加えて、書き込みに関連付けられたTwが別の時刻にあってもよい。Twは、変更すべき最後のページがアンラッチ(例えば、ページの開放/アンロック)された直後の時刻を表し得るが、そのトランザクション要求を受け取ったヘッドノードによって割り当てられてもよい。別の実施例として、Tsに加えて、読み取りに関連付けられたTrが別の時刻にあってもよい。Trは、最初のデータページが読み取りラッチされた直後の時刻を表し得るが、その読み取り要求を受け取ったヘッドノードによって割り当てられてもよい。

10

20

【0055】

630に示すように、潜在的な読み取り異常が検出され得る。一実施形態では、読み取り要求(及びその読み取りの実行)を受け取ったヘッドノードは、ブロック630にて検出を実行し得る。いくつかの実施形態では、このような検出は、第1の時刻指示(例えば、Ts及び/またはTr)が、第2の時刻指示(例えば、Tc及び/またはTw)の閾値(例えば、精度ウィンドウ)内にあるとの判定に基づき得る。各種の潜在的な読み取り異常が発生する可能性がある。例えば、ファジーリードは、更新及び(再)読み取りがほぼ同時に起こるときに発生する可能性があり、結果として同一レコードについて異なる値を読み取る可能性がある。として5を用いた表1の実施例では、読み取りの一貫性ポイントよりも前に起きたときにコミットを処理することにより、適切な値を読み取れるようにその読み取りが正しく動作する。表1の実施例を図7Aに図示する。

30

【0056】

【表1】

表1

ノード1		ノード2	
101	トランザクション開始	100	スナップショットの作成 (Ts=100)
102	更新 X=10	101	
103	コミット (Tc=103)	102	
104		103	読み取り X(10) [100≠103]

40

【0057】

ファジーリードを表2の実施例に示し、図7Bに図示する。これも同様に、として5を用いている。ファジーリードでは、リードトランザクションにより、LSN101にて値(1)が読み取られ、次いで、同一レコードをLSN104にて再び読み取ると異なる値(10)が見える。本明細書で説明されるように、開示された技術は、表2の潜在的なファジーリードを検出することができ、読み取りを調節することによってそのファジーリードが実際には発生しないようにすることができる。

50

【 0 0 5 8 】

【表 2】

表 2

ノード 1		ノード 2	
1 0 1	トランザクション開始	1 0 0	スナップショットの作成 ($T_s = 1 0 0$)
1 0 2		1 0 1	読み取り X (1)
1 0 3	更新 X = 1 0	1 0 2	
1 0 4	コミット ($T_c = 1 0 4$)	1 0 3	
1 0 5		1 0 4	読み取り X (1 0) [$1 0 0 \neq 1 0 4$]

10

【 0 0 5 9 】

なお、表 2 のファジーリードは、(例えば、(単数または複数の)クライアント上で)動作がシリアル化されている場合には発生しない。これは、更新を行ってから初めて行を読み取ることになるためである。

【 0 0 6 0 】

別の異常読み取りはリードスキューである。これは、(例えば、複数の異なるレコードの)一貫性のないデータを読み取った状況である。(図 7 C に示した)表 3 の実施例について考察する。これも同様に を 5 としている。表 3 及び図 7 C の実施例では、データテーブルは、初期値を $X = 1$ 及び $Y = 2$ として、 $X = 2 Y$ である不変関係を有する。図示の通り、この初期値 1 は、LSN 1 0 1 にて X について読み取られるが、Y の更新値は、LSN 1 0 4 にて 2 0 と読み取られる。これは、 $X = 2 Y$ と矛盾している。

20

【 0 0 6 1 】

【表 3】

表 3

ノード 1		ノード 2	
1 0 1	トランザクション開始	1 0 0	スナップショットの作成 ($T_s = 1 0 0$)
1 0 2	更新 Y = 2 0	1 0 1	読み取り X (1)
1 0 3	更新 X = 1 0	1 0 2	
1 0 4	コミット ($T_c = 1 0 4$)	1 0 3	
1 0 5		1 0 4	読み取り Y (2 0) [$1 0 0 \neq 1 0 4$]

30

【 0 0 6 2 】

いくつかの実施形態では、潜在的な読み取り異常は、読み取りの一貫性ポイント T_s の精度ウィンドウ () 内にコミット時刻 T_c がある場合に検出され得る。潜在的な読み取り異常は、本明細書では、読み取り異常が発生し得る可能性が存在することを示すのに用いられるが、潜在的な読み取り異常は、読み取り異常が確実に発生することを必ずしも意味するものではないことに留意されたい。従って、読み取り異常の可能性が存在する場合、ブロック 6 4 0 にて以下で説明するように、システムは、それを検出し、潜在的な読み取り異常の回避を試みることができる。

40

【 0 0 6 3 】

上述したように、いくつかの実施形態では、 T_c 及び T_s 以外の時刻も同様に、読み取り及び/または他のトランザクションに関連付けてもよい。例えば、 T_w を用いた実施形態では、 $T_c > T_s$ である場合、読み取りを行ってもトランザクションによってなされた変更が見えず、異常読み取りは発生しない。 $T_c < T_s$ である場合、読み取りを行うとトランザクションによってなされた変更が見える。 $T_c \leq T_s$ かつ $T_w < T_s$ である場合、トランザクションによってなされた変更が読み取りの前に行われていたため、潜在的な読

50

み取り異常は存在しない。トランザクションによってなされた変更は、読み取りを行うことによって見える。さもなければ、 $T_c > T_s$ かつ $T_w > T_s$ の場合には、潜在的な読み取り異常が存在する。

【0064】

いくつかの実施形態では、 T_r を用いて潜在的な読み取り異常を検出するようにしてもよい。このような実施形態では、 $T_c > T_s$ である場合、トランザクションによってなされた変更は読み取りを行っても見えず、異常読み取りは発生しない。 $T_c < T_s$ である場合、読み取りを行うとトランザクションによってなされた変更が見える。 $T_c > T_s$ かつ $T_w < T_r$ である場合、トランザクションによってなされた変更が読み取りの前に行われていたため、潜在的な読み取り異常は存在しない。さもなければ、 $T_c > T_s$ かつ $T_w > T_r$ である場合には、潜在的な読み取り異常が存在し、ブロック630にて検出され得る。

10

【0065】

なお、いくつかの実施形態では、リードスキュー及びファジーリードは、最初のページをラッチした後にのみ起こる場合がある。そのため、1つのレコードを検索するステートメントでは、こうした異常が発生しない。従って、いくつかの実施形態では、ブロック630の検出ロジックは、(同一レコードを複数回検索しようと複数の異なるレコードを検索しようと)複数のレコードを検索する場合にのみ実行されてもよい。

【0066】

640に示すように、潜在的な読み取り異常を上記検出したことに応答して、その読み取り要求を受け取ったノードは、第1の時刻指示が第2の時刻指示よりも早い時点を示しているかどうかに関わらず、その読み取り要求によって指定された読み取りを、トランザクション要求によって指定されたトランザクションの後に実行させることができる。

20

【0067】

いくつかの実施形態では、トランザクションの後に読み取りを実行させることは、その読み取りの再試行に関連付けられた再試行の時刻指示が第1の時刻指示よりも後の時刻を示すようにその読み取りを再試行する読み取り要求を受け取ったノードを含んでもよい。例えば、潜在的なファジーリードまたはリードスキューが検出された場合、読み取りステートメントは、 T_s を維持する一方で T_r をリセットする(例えば、 T_r の時刻を進める)ことによって抽出され得る。再試行時に T_s を維持することにより、フォワード進行を達成することができる。なぜなら、 T_r が進められるため、トランザクションと読み取りが、比較ロジックにとって既知の事例の1つ(例えば、潜在的な読み取り異常が発生しない状況)に結果として収まり得るためである。

30

【0068】

なお、一実施形態では、再試行の場合、ブロック630の検出ロジックは、以前の T_r を置き換えた再試行時刻(更新された T_r)を用いて再度適用され得る。従って、その再試行が成功した場合、検出ロジックは、再試行された読み取りに対しては読み取り異常が発生しなかったと判定する。例えば、更新された T_r を用いることにより、異常読み取りが起こり得ない区分の1つ(例えば、 $T_c > T_s$ かつ $T_w < T_r$)の範囲内に再試行が収まっていると検出ロジックに判定させてもよく、異常読み取りを伴わずに再試行が行われる。一方、別の潜在的な異常読み取りが発生する場合があるため、別の再試行を行う際には、別の再試行の時刻を時間的にさらに進めるようにする。検出ロジックを適用し、潜在的な異常読み取りが存在すると判定し、そのステートメントを再試行することは、再試行が成功するまで何回でも発生する場合がある。あるいは、いくつかの実施形態では、ノードは、その読み取りが正確でない可能性があることを示すエラーメッセージを実際の読み取り値と共に返す前に、当該ステートメントをある最大回数(例えば、2回、3回、10回など)再試行してもよい。

40

【0069】

一実施形態では、この閾値を変更してもよい。例えば、システムにおいて既に発生したトランザクションに対する再試行の頻度に基づいて閾値を変更してもよい。ある実施例と

50

して、再試行が頻繁に発生している場合には、閾値を減少させて控えめの水準まで下げてもよい。同様に、閾値を増加させてもよい。

【0070】

いくつかの実施形態では、ブロック620、630及び640は、同一レコードを複数回読み取ろうと複数の異なるレコードを読み取ろうと、複数のレコードの検索を含む読み取り要求に対して行われてもよい。従って、このような実施形態では、読み取り異常の有無の確認は、1つのレコードの検索に対してではなく、複数のレコードの検索に対して行われてもよい。そのため、第1の読み取りは制約のない読み取りとなり得るが、これは、1回の読み取りのみを実行中である場合、(上記のような読み取り異常と解釈するには、第1の読み取りと矛盾する第2の読み取りが必要となるため)リードスキューまたはファジーリードが起こる可能性がないためである。従って、いくつかの実施形態では、本システムは、互いの精度ウィンドウ(閾値)内にある、1つのレコードの検索を求める読み取り要求と、トランザクション要求とを処理する際に図6の方法を常に適用しなくてもよい。

10

【0071】

いくつかの実施形態では、読み取り要求を受け取ったノードは、別のノードが書き込み要求を受け取って当該レコードを更新中であると認識する場合がある。例えば、本システムは、当該データへの変更を他のノードが認識することを保証する下位副構造を含んでもよい。例示的な下位副構造には、(例えば、各データページキャッシュ335または他のキャッシュに対するWebサービスプラットフォーム200のレベルでの)キャッシュフュージョン、または共有ディスクが含まれる。一実施形態では、共通ストレージ(例えば、ストレージ階層)の上にコヒーレントキャッシュを置いてもよい。コヒーレントキャッシュを用いることにより、あるものをノードが書き込んだ場合、別のノードがそれを見ることが保証され得る。一実施形態では、様々な時刻指示がトランザクションテーブルに記憶されてもよい。例えば、Twは、コミットの時刻以降にトランザクションのコミット時刻と共に記憶されてもよく、各時刻の値を書き込み、トランザクション識別子を生成してもよい。トランザクション識別子は、所与のトランザクションがアクティブか、それともコミットされたかを示し得る。所与のレコードに対してトランザクションがアクティブである場合、そのレコードに対して読み取りを実行しているノードは、そのレコードに対して以前の値を生成し得る(例えば、ロールバック、UNDOなど)。

20

30

【0072】

本明細書で説明された方法は、様々な実施形態では、ハードウェアとソフトウェアの任意の組み合わせによって実装され得る。例えば、一実施形態では、本方法は、プログラム命令を実行する1つ以上のプロセッサを含むコンピュータシステム(例えば、図8と同様のコンピュータシステム)によって実行され得る。これらのプログラム命令は、上記プロセッサに接続されたコンピュータ可読記憶媒体に記憶されている。このプログラム命令は、本明細書で説明された機能(例えば、本明細書で説明されたサービス/システム及び/またはストレージサービス/ストレージシステムを実装する様々なサーバ及び他の構成要素の機能)を実装するように構成され得る。

【0073】

開示されたトランザクションの順序付け技術を用いることにより、強力な、かつ理にかなった分離レベルを顧客に提供することができる。この技術により、システムのスケラビリティを向上させることができる。なぜなら、スナップショット(読み取り)を迅速に作成することが可能であり、ネットワーク通信を必要としなくてもよいためである。さらに、トランザクションが重複しない作業負荷も適切にスケールアウトすることができる。シングルton行の検索は迅速に行うことができる。というのも、それらの読み取りは、読み取り異常の検出ロジックを実行させずに制約なく行われ得るためである。(例えば、Tc及びTsに加えて、Tr及び/またはTwを計算に入れることによって)正確な精度ウィンドウを用いた実施形態では、ステートメントが再試行される可能性を減少させることができる。

40

50

【 0 0 7 4 】

図 8 は、種々の実施形態に係る、本明細書で説明されたシステムの少なくとも一部を実装するように構成された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 8 0 0 は、種々の実施形態では、（例えば、データベース階層もしくは同程度のシステムの）ノード、またはクライアントに代わってレコード及び関連メタデータを記憶する複数のストレージノードの 1 つを実装するように構成され得る。コンピュータシステム 8 0 0 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もしくはノートブックコンピュータ、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、コンシューマ機器、アプリケーションサーバ、ストレージ装置、電話、携帯電話、または一般的な任意の種類のコピューティング装置を含む各種の装置のいずれでもあってよいが、これらに限定されることはない。

10

【 0 0 7 5 】

コンピュータシステム 8 0 0 は、入出力 (input/output : I/O) インターフェース 8 3 0 を介してシステムメモリ 8 2 0 に接続された 1 つ以上のプロセッサ 8 1 0 を含む（これらのプロセッサはいずれも複数のコアを含んでよく、それらのコアはシングルスレッド方式でもマルチスレッド方式でもよい）。コンピュータシステム 8 0 0 は、I/O インターフェース 8 3 0 に接続されたネットワークインターフェース 8 4 0 をさらに含む。様々な実施形態では、コンピュータシステム 8 0 0 は、1 つのプロセッサ 8 1 0 を含むユニプロセッサシステム、またはいくつかのプロセッサ 8 1 0（例えば、2 個、4 個、8 個もしくは他の適切な数）を含むマルチプロセッサシステムであってよい。プロセッサ 8 1 0 は、命令を実行可能な任意の適切なプロセッサであってよい。例えば、様々な実施形態では、プロセッサ 8 1 0 は、様々な命令セットアーキテクチャ (instruction set architectures : 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 (x86、PowerPC、SPARC もしくは MIPS の ISA またはその他の適切な ISA など) 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 8 1 0 のそれぞれは、通常は同じ ISA を実装し得るが、必ずしもそのように実装されるとは限らない。コンピュータシステム 8 0 0 は、1 つ以上のネットワーク通信装置（例えば、ネットワークインターフェース 8 4 0）も含み、他のシステム及び/または構成要素と通信ネットワーク（例えば、インターネット、LAN など）を介して通信できるようになっている。例えば、システム 8 0 0 上で実行されるクライアントアプリケーションは、ネットワークインターフェース 8 4 0 を用いて、本明細書で説明されたデータベースシステムの構成要素の 1 つ以上を実装する 1 つのサーバまたはサーバのクラスター上で実行されるサーバアプリケーションと通信してもよい。別の例では、コンピュータシステム 8 0 0 上で実行されるサーバアプリケーションのインスタンスは、ネットワークインターフェース 8 4 0 を用いて、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 8 9 0）に実装され得る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してもよい。

20

30

【 0 0 7 6 】

例示した実施形態では、コンピュータシステム 8 0 0 は、1 つ以上の永続的記憶装置 8 6 0 及び/または 1 つ以上の I/O 装置 8 8 0 も含む。様々な実施形態では、永続的記憶装置 8 6 0 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、他の大容量ストレージ装置、またはその他の永続的記憶装置に対応し得る。コンピュータシステム 8 0 0（または分散アプリケーションもしくはそのシステム上で動作するオペレーティングシステム）は、所望により、命令及び/またはデータを永続的記憶装置 8 6 0 に記憶してもよく、記憶された命令及び/またはデータを必要に応じて検索してもよい。例えば、ある実施形態では、コンピュータシステム 8 0 0 は、ストレージシステムサーバノードをホスティングしてもよく、永続的記憶装置 8 6 0 は、そのサーバノードに付属する SSD を含んでもよい。

40

【 0 0 7 7 】

50

コンピュータシステム 800 は、1つ以上のシステムメモリ 820 を含んでおり、これらのメモリは、(単数または複数の)プロセッサ 810 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ている。様々な実施形態では、システムメモリ 820 は、任意の適切なメモリ技術(例えば、1つ以上のキャッシュ、SRAM (static random access memory)、DRAM、RDRAM、EDORAM、DDR 10 RAM、SDRAM (synchronous dynamic RAM)、Rambus RAM、EEPROM、不揮発性/フラッシュ型メモリ、またはその他の種類のメモリ)を用いて実装され得る。システムメモリ 820 は、プログラム命令 825 を含み得る。これらの命令は、本明細書で説明された方法及び技術を実装するために、(単数または複数の)プロセッサ 810 によって実行可能である。様々な実施形態では、プログラム命令 825 は、プラットフォームのネイティブバイナリ、Java (登録商標) バイトコードなどの任意のインタープリター型言語、もしくは C/C++、Java (登録商標) などのその他の言語、またはそれらの任意の組み合わせを用いてコード化され得る。例えば、プログラム命令 825 は、例示した実施形態では、データベース階層のデータベースエンジン・ヘッドノードの機能を実装するために実行可能なプログラム命令を含み、または、異なる実施形態では、データベース階層のクライアントに代わってデータベーステーブル及び関連メタデータを記憶する別個の分散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うちの1つの機能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 825 は、複数の独立したクライアント、サーバノード及び/または他の構成要素を実装してもよい。

10

20

【0078】

いくつかの実施形態では、プログラム命令 825 は、オペレーティングシステム(図示せず)を実装するために実行可能な命令を含んでもよく、このオペレーティングシステムは、UNIX (登録商標)、LINUX (登録商標)、Solaris (登録商標)、MacOS (登録商標)、Windows (登録商標) などの様々なオペレーティングシステムのいずれでもよい。プログラム命令 825 のいずれかまたは全ては、コンピュータプログラム製品、またはソフトウェアとして提供され得る。このコンピュータプログラム製品またはソフトウェアは、命令が記憶された非一時的コンピュータ可読記憶媒体を含んでもよく、この媒体を用いてコンピュータシステム(または他の電子装置)をプログラムすることにより、様々な実施形態に応じたプロセスを実行してもよい。非一時的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)によって読み取り可能な形態(例えば、ソフトウェア、処理アプリケーション)で情報を記憶するための任意の機構を含んでもよい。一般的に言えば、非一時的コンピュータアクセス可能媒体は、コンピュータ可読記憶媒体またはメモリ媒体を含み得る。このような媒体としては、磁気媒体または光学媒体(例えば、I/O インターフェース 830 を介してコンピュータシステム 800 に接続されたディスクまたは DVD/CD-ROM)などがある。非一時的コンピュータ可読記憶媒体は、RAM (例えば、SDRAM、DDR SDRAM、RDRAM、SRAM など)、ROM などの任意の揮発性または不揮発性媒体も含み得る。これらの媒体は、コンピュータシステム 800 のいくつかの実施形態では、システムメモリ 820 または別の種類のメモリとして含まれ得る。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ーフェース 840 を介して実装され得るようなネットワーク及び/または無線リンクなどの通信媒体を経由して伝わる光学的、音響的または他の形態の伝搬信号(例えば、搬送波、赤外線信号、デジタル信号など)を用いて伝達され得る。

30

40

【0079】

いくつかの実施形態では、システムメモリ 820 は、データ記憶部 845 を含んでもよく、この記憶部は、本明細書で説明されたように構成され得る。例えば、本明細書で説明されたデータベース階層の機能を実行する際に用いられるトランザクションログ、UNDO ログ、キャッシュ済みのページデータまたは他の情報などの、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)記憶されているものとして本明細書で説明された情報は、別々の時刻にて、及び様々な実施形態において、1つ以上のノ

50

ード上のデータ記憶部 8 4 5 もしくはシステムメモリ 8 2 0 の別の部分、永続的記憶装置 8 6 0、及び/または1つ以上のリモートストレージ装置 8 7 0 に記憶され得る。同様に、データベース階層によって記憶されているものとして本明細書で説明された情報(例えば、本明細書で説明された分散ストレージシステムの機能を実行する際に用いられる R E D O ログレコード、併合されたデータページ、及び/または他の情報など)は、別々の時刻にて、及び様々な実施形態において、1つ以上のノード上のデータ記憶部 8 4 5 もしくはシステムメモリ 8 2 0 の別の部分、永続的記憶装置 8 6 0、及び/または1つ以上のリモートストレージ装置 8 7 0 に記憶され得る。概して、システムメモリ 8 2 0 (例えば、システムメモリ 8 2 0 内のデータ記憶部 8 4 5)、永続的記憶装置 8 6 0 及び/またはリモートストレージ 8 7 0 は、データブロック、データブロックのレプリカ、データブロック及び/もしくはそれらの状態に関連付けられたメタデータ、データベース構成情報、並びに/または本明細書で説明された方法及び技術を実装するのに使用可能なその他の情報を記憶し得る。

10

【0080】

一実施形態では、I/Oインターフェース 8 3 0 は、ネットワークインターフェース 8 4 0 または他の周辺インターフェースを通るものを含む、プロセッサ 8 1 0、システムメモリ 8 2 0 及びシステム内の任意の周辺装置の間のI/Oトラフィックを調整するように構成され得る。いくつかの実施形態では、I/Oインターフェース 8 3 0 は、任意の必要なプロトコル変換、タイミング変換または他のデータ変換を実行することにより、ある構成要素(例えば、システムメモリ 8 2 0)からのデータ信号を別の構成要素(例えば、プロセッサ 8 1 0)での使用に適した形式に変換し得る。いくつかの実施形態では、I/Oインターフェース 8 3 0 は、例えば、P C I (P e r i p h e r a l C o m p o n e n t I n t e r c o n n e c t) パス規格またはU S B (U n i v e r s a l S e r i a l B u s) 規格を変型したものなどの様々な種類の周辺バスを介して接続された装置にも対応し得る。いくつかの実施形態では、I/Oインターフェース 8 3 0 の機能は、例えば、ノースブリッジ及びサウスブリッジなどの、2つ以上の別々の構成要素に分割され得る。また、いくつかの実施形態では、システムメモリ 8 2 0 とのインターフェースなどの、I/Oインターフェース 8 3 0 の機能の一部または全てをプロセッサ 8 1 0 に直接組み込んでよい。

20

【0081】

ネットワークインターフェース 8 4 0 は、例えば、コンピュータシステム 8 0 0 と、他のコンピュータシステム 8 9 0 (このシステムは、本明細書で説明されたデータベースシステムの1つ以上のストレージシステムサーバード、データベースエンジン・ヘッドノード及び/またはクライアントを実装し得る)などの、ネットワークに接続された他の装置との間でデータを交換できるように構成され得る。加えて、ネットワークインターフェース 8 4 0 は、コンピュータシステム 8 0 0 と、各種I/O装置 8 5 0 及び/またはリモートストレージ 8 7 0 との間の通信を行うことができるように構成され得る。入出力装置 8 5 0 は、いくつかの実施形態では、1つ以上の表示端末、キーボード、キーパッド、タッチパッド、スキャニング装置、音声認識装置もしくは光学的認識装置、または1つ以上のコンピュータシステム 8 0 0 によるデータの入力もしくは検索に適したその他の装置を含んでもよい。複数の入出力装置 8 5 0 は、コンピュータシステム 8 0 0 内に存在してもよく、またはコンピュータシステム 8 0 0 を含む分散システムの様々なノードに分散されてもよい。いくつかの実施形態では、同様の入出力装置は、コンピュータシステム 8 0 0 から分かれていてもよく、コンピュータシステム 8 0 0 を含む分散システムの1つ以上のノードと、ネットワークインターフェース 8 4 0 を介するような有線または無線接続を通じて相互作用してもよい。ネットワークインターフェース 8 4 0 は、通常1つ以上の無線ネットワーク接続プロトコル(例えば、W i - F i / I E E E 8 0 2 . 1 1、または別の無線ネットワーク接続規格)に対応し得る。しかしながら、様々な実施形態では、ネットワークインターフェース 8 4 0 は、例えば、他の種類のイーサネット(登録商標)ネットワークなどの、任意の適切な有線または無線の一般的なデータネットワークを経由した通

30

40

50

信に対応し得る。さらに、ネットワークインターフェース 840 は、アナログ音声ネットワークもしくはデジタルファイバー通信ネットワークなどの遠距離通信ネットワーク/電話網、ファイバーチャネル SAN などのストレージエリアネットワーク、またはその他の適切な種類のネットワーク及び/またはプロトコルを経由した通信に対応し得る。様々な実施形態では、コンピュータシステム 800 は、図 8 に示した構成要素よりも多い構成要素、これよりも少ない構成要素、またはこれとは異なる構成要素（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 インターフェース、イーサネットインターフェース、フレームリレーインターフェースなどの他のネットワークインターフェース）を含んでもよい。

【0082】

なお、本明細書で説明された分散システムの任意の実施形態、またはそれらの任意の構成要素は、1つ以上の Web サービスとして実装され得る。例えば、データベースシステムのデータベース階層内のデータベースエンジン・ヘッドノードは、本明細書で説明された分散ストレージシステムを利用するデータベースサービス及び/または他の種類のデータストレージサービスを、Web サービスとしてクライアントに提供し得る。いくつかの実施形態では、Web サービスは、ネットワークを介して相互運用可能なマシン対マシン相互作用に対応するように設計されたソフトウェア及び/またはハードウェアシステムによって実装され得る。Web サービスは、Web サービス記述言語 (Web Services Description Language: WSDL) などの、マシン処理可能な書式を用いて記述されたインターフェースを備えてもよい。他のシステムは、Web サービスのインターフェースを記述することによって定められる方法で Web サービスと相互作用してもよい。例えば、Web サービスは、他のシステムが呼び出すことができる様々な動作を規定してもよく、様々な動作を要求する際に他のシステムが適合していることを見込むことができる特定のアプリケーションプログラミングインターフェース (API) を規定してもよい。

【0083】

上記は、前述の条項に照らしてより理解され得る。

条項 1

システムであって、

それぞれが少なくとも 1 つのプロセッサ、及びメモリを含む複数のコンピューティングノードであって、データベースサービスを集合的に実施するように構成されている複数のコンピューティングノード、を含み、

前記複数のコンピューティングノードの第 1 のノードは、

データベーステーブル内の特定のデータページの特定のデータレコードを対象とする書き込み要求であって、前記特定のデータレコードに対して行うべき変更を指定する書き込み要求を受け取り、

前記変更をコミットする時刻を示すコミット時刻を前記書き込み要求に割り当てるように構成されており、

前記複数のコンピューティングノードの第 2 のノードは、

前記データベーステーブル内の前記特定のデータページを読み取る読み取り要求を受け取り、

読み取りの一貫性ポイントを示す読み取り一貫性ポイント時刻を前記読み取り要求に割り当て、

前記読み取り一貫性ポイント時刻が前記コミット時刻の精度ウィンドウ内にあると判定し、

前記読み取り一貫性ポイント時刻が前記コミット時刻よりも早い時点を示しているかどうかに関わらず、前記書き込み要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を行わせるように構成されている、システム。

条項 2

前記複数のコンピューティングノードの前記第 2 のノードは、

10

20

30

40

50

前記書き込み要求に関連付けられた書き込みアンラッチ時刻が、前記読み取り要求に関連付けられた読み取りラッチ時刻の精度ウィンドウ内にあると判定するようにさらに構成されており、

前記書き込み要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を前記行わせることは、

前記読み取りラッチ時刻を後の時点に移動すること、及び

前記読み取り一貫性の時点以降に前記読み取り要求を再試行すること、を含む条項 1 に記載のシステム。

条項 3

前記書き込み要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を前記行わせることは、前記移動された読み取りラッチ時刻が前記書き込みアンラッチ時刻の精度ウィンドウ内にはないと判定すること、をさらに含む、条項 2 に記載のシステム。

10

条項 4

前記第 1 のノード及び前記第 2 のノードのそれぞれは、互いの前記精度ウィンドウ内に維持されている個別のクロックを保持するように構成されており、前記コミット時刻が前記第 1 のノードによって決定され、前記読み取り一貫性の時点が前記第 2 のノードによって決定される、条項 1 に記載のシステム。

条項 5

方法であって、

20

複数のコンピューティングノードによって、

1 つ以上のクライアントから、記憶されたレコードの読み取りを実行する読み取り要求、及び前記レコードの更新を実行する更新要求を受け取ること、

第 1 の時刻指示と第 2 の時刻指示とを、前記読み取りと前記更新とにそれぞれ関連付けること、並びに

前記第 1 の時刻指示が前記第 2 の時刻指示の閾値内にあると判定したことに少なくとも部分的に基づき、潜在的な読み取り異常を検出すること、

を実行すること、を含む方法。

条項 6

前記潜在的な読み取り異常を前記検出したことに応答して、前記第 1 の時刻指示が前記第 2 の時刻指示よりも早い時点を示しているかどうかに関わらず、前記更新要求によって指定された前記更新の後に前記読み取り要求によって指定された前記読み取りを実行させること、をさらに含む条項 5 に記載の方法。

30

条項 7

前記読み取りの再試行に関連付けられた再試行の時刻指示が前記第 1 の時刻指示よりも後の時刻を示すように、前記第 1 の時刻指示以降に前記読み取りを再試行すること、及び前記再試行された読み取りに対して読み取り異常が発生していないと判定すること、をさらに含む条項 5 に記載の方法。

条項 8

前記第 2 の時刻指示は前記更新をコミットする時刻を示す、条項 5 に記載の方法。

40

条項 9

前記潜在的な読み取り異常を前記検出することは、第 3 の時刻指示にさらに基づいており、前記第 3 の時刻指示も前記更新要求に関連付けられている、条項 5 に記載の方法。

条項 10

前記潜在的な読み取り異常を前記検出することは、前記第 3 の時刻指示が第 4 の時刻指示の閾値内にあると判定したことにさらに基づいており、前記第 4 の時刻指示も前記読み取り要求に関連付けられている、条項 9 に記載の方法。

条項 11

前記読み取り要求は前記複数のノードの第 1 のノードによって受け取られ、前記更新要求は前記複数のノードの別の第 2 のノードによって受け取られ、前記第 1 のノード及び前

50

記第 2 のノードは個別のクロックをそれぞれ保持しており、前記読み取りに関連付けられた前記第 1 の時刻指示は前記第 1 のノードの前記個別のクロックによって決定され、前記更新に関連付けられた前記第 2 の時刻指示は前記第 2 のノードの前記個別のクロックによって決定される、条項 5 に記載の方法。

条項 1 2

前記潜在的な読み取り異常は潜在的なファジーリードである、条項 5 に記載の方法。

条項 1 3

前記潜在的な読み取り異常は潜在的なリードスキューである、条項 5 に記載の方法。

条項 1 4

再試行の頻度に少なくとも部分的に基づいて前記閾値を変更すること、をさらに含む条項 5 に記載の方法。

条項 1 5

前記検出することは、前記レコードの前記読み取り、及び別のレコードの読み取りまたは前記レコードの第 2 の読み取りのいずれかを含む読み取り要求に対して実行される、条項 5 に記載の方法。

条項 1 6

データベースサービスのデータベースノードを実装するためにコンピュータで実行可能なプログラム命令を記憶する非一時的コンピュータ可読記憶媒体であって、前記データベースノードは、

前記データベースサービスによって記憶されたレコードの読み取りを指定する読み取り要求によって指定された前記読み取りに第 1 の時刻指示に関連付け、

前記第 1 の時刻指示が、別のデータベースノードによって受け取られたトランザクション要求であって、前記レコードを変更するトランザクションを指定するトランザクション要求によって指定された前記トランザクションに関連付けられた第 2 の時刻指示の閾値内にあると判定し、

前記第 1 の時刻指示が前記第 2 の時刻指示よりも早い時点を示しているかどうかに関わらず、前記トランザクション要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を行わせるように構成されている、非一時的コンピュータ可読記憶媒体。

条項 1 7

前記データベースノードは、

第 3 の時刻指示を前記読み取りに関連付け、

前記第 3 の時刻指示が、前記トランザクションに関連付けられた第 4 の時刻指示の閾値内にあると判定するようにさらに構成されており、

前記トランザクション要求によって指定された前記変更がコミットされた後に前記読み取り要求の実行を前記行わせることは、

前記第 3 の時刻指示を後の時点に置き換えること、及び

前記第 1 の時刻指示以降に前記読み取り要求を再試行すること、を含む条項 1 6 に記載の非一時的コンピュータ可読記憶媒体。

条項 1 8

前記データベースノードは、前記後の時点が、前記第 4 の時刻指示の閾値内にはないと判定するようにさらに構成されている、条項 1 7 に記載の非一時的コンピュータ可読記憶媒体。

条項 1 9

前記第 1 の時刻指示及び前記第 2 の時刻指示は、単調増加する時刻指標である、条項 1 6 に記載の非一時的コンピュータ可読記憶媒体。

条項 2 0

前記読み取り要求は、前記レコードを含む複数のレコードを検索する要求である、条項 1 6 に記載の非一時的コンピュータ可読記憶媒体。

【 0 0 8 4 】

10

20

30

40

50

様々な実施形態では、Webサービスは、Webサービス要求に関連付けられたパラメータ及び/またはデータを含むメッセージの利用を通じて要求または起動され得る。こうしたメッセージは、XML (Extensible Markup Language) など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もよく、かつ/またはSOAP (Simple Object Access Protocol) などのプロトコルを用いてカプセル化されてもよい。Webサービス要求を実行するために、Webサービスクライアントは、その要求を含むメッセージを集めてもよく、そのメッセージを、HTTP (HyperText Transfer Protocol) などの、インターネットを利用したアプリケーション層の転送プロトコルを用いて、Webサービスに対応するアドレス可能エンドポイント (例えば、URL (Uniform Resource Locator)) に伝達してもよい。

10

【0085】

いくつかの実施形態では、Webサービスは、メッセージを利用した技術ではなく、「RESTful (Representational State Transfer)」技術を用いて実装されてもよい。例えば、RESTful技術に従って実装されたWebサービスは、SOAPメッセージ内にカプセル化されたものではなく、PUT、GETまたはDELETEなどのHTTPメソッド内に含まれるパラメータを通じて起動され得る。

【0086】

各図面に図示され、本明細書で説明された様々な方法は、方法の例示的な実施形態を表す。これらの方法は、ソフトウェア、ハードウェアまたはそれらの組み合わせを用いて手作業で実施されてもよい。あらゆる方法の順序は変更可能であり、様々な要素の追加、並べ替え、組み合わせ、省略、修正などを行ってもよい。

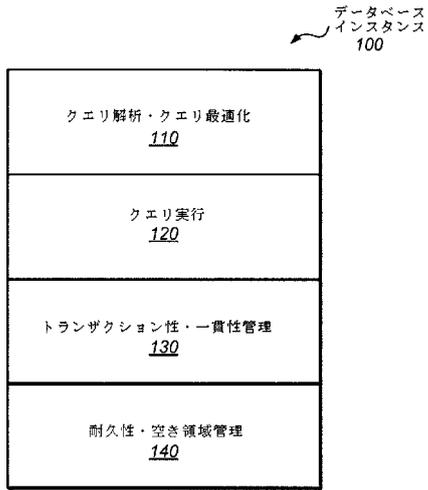
20

【0087】

上記のような実施形態についてかなり詳細に説明してきたが、上記の開示を一旦完全に理解すれば、当業者にとって明らかとなるような多くの変形及び修正をなすことが可能である。以下の特許請求の範囲は、このような変形及び修正の全て、並びに従って、限定的な意味ではなく例示とみなされるべき上述の説明を包含するように解釈されることが意図される。

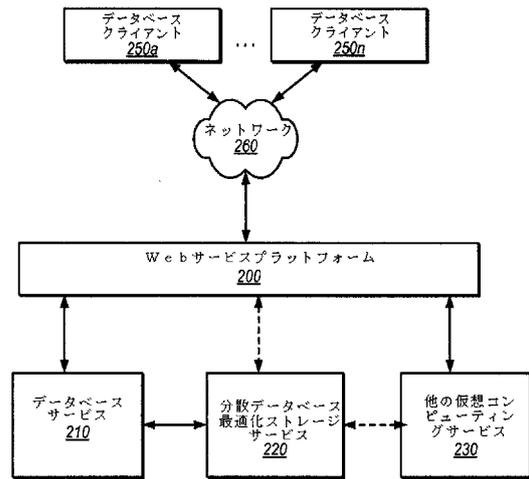
【 図 1 】

図 1



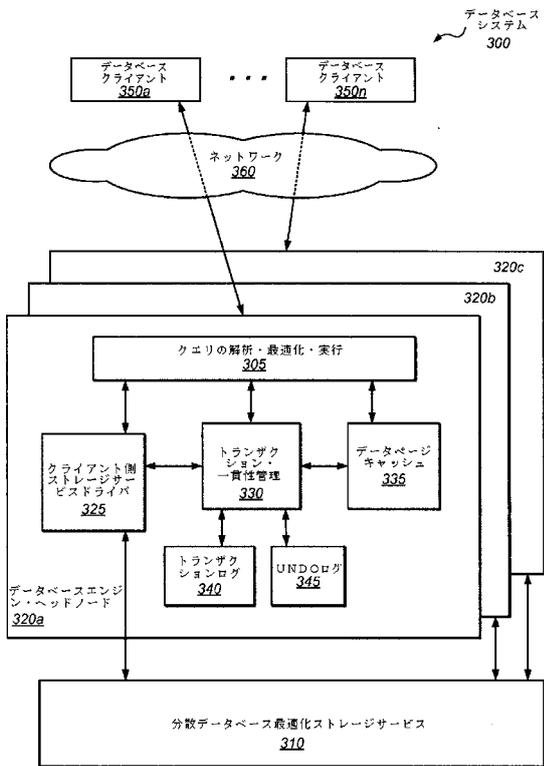
【 図 2 】

図 2



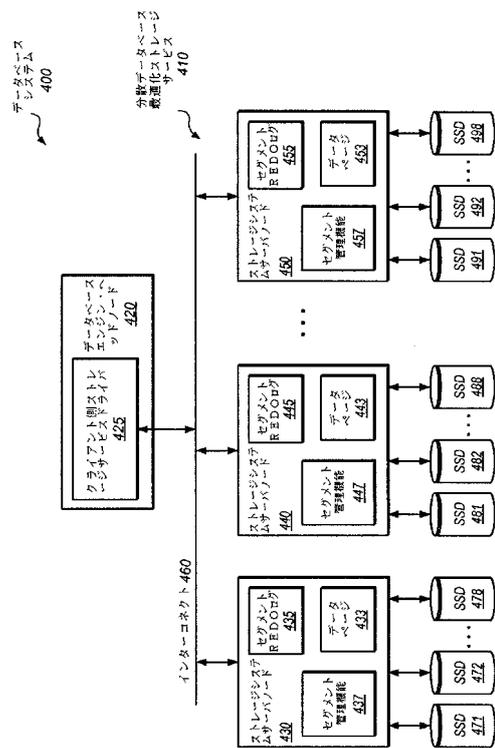
【 図 3 】

図 3



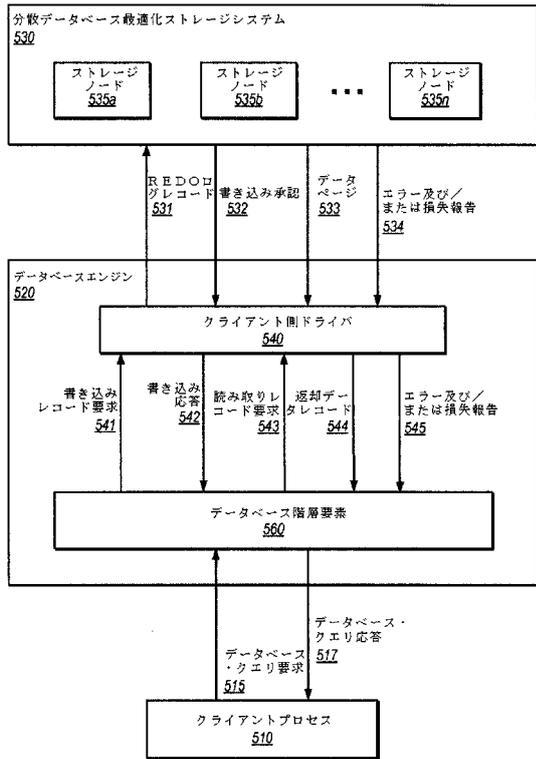
【 図 4 】

図 4



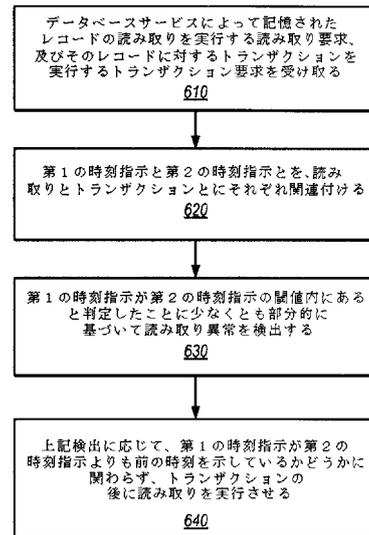
【 図 5 】

図 5



【 図 6 】

図 6



【 図 7 】

図 7 A

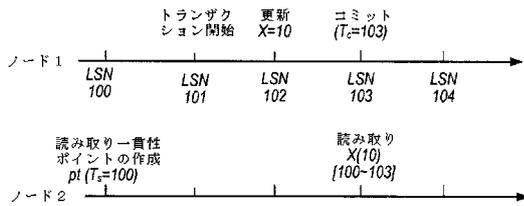


図 7 B

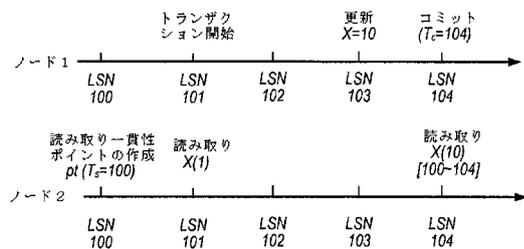
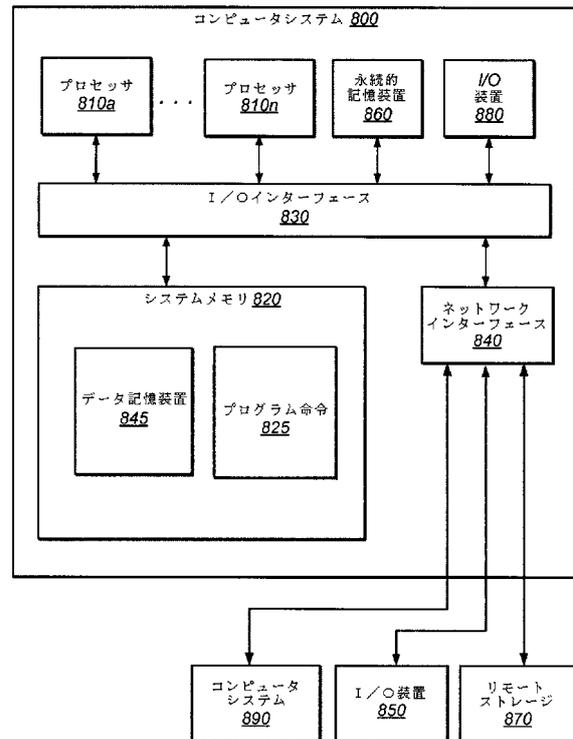


図 7 C



【 図 8 】

図 8



【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US 14/37801

A. CLASSIFICATION OF SUBJECT MATTER IPC(8) - G06F 17/00 (2014.01) CPC - G06F 2201/84 According to International Patent Classification (IPC) or to both national classification and IPC.		
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) CPC: G06F 2201/84; IPC(8): G06F 17/00 (2014.01)		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched CPC: Y10S 707/99954, G06F 17/30306, Y10S 707/99942, G06F 11/1464, G06F 11/1469, G06F 11/1435, G06F 2201/84, G06F 17/30067, G06F 11/1451; IPC(8): G06F 17/00 (2014.01) (keyword limited)		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) PatBase, Google Patents, IEEE; Search Terms: transaction ordering, database, nodes, record read, commit, consistency point, write read request, latch time, accuracy window, potential read anomaly, deadlock detection, latch, unlatch		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 5,452,445 A (Hallmark et al.) 19 September 1995 (19.09.1995), entire document especially abstract, col. 8, ln 24-38, col. 9, ln 15 - col. 10, ln 17, Figure 2, steps 201 - 204, Figure 3, step 301 - 304, Figure 4, steps 401 - 407	1 - 15
A	US 2007/0130238 A1 (Harris et al.) 07 June 2007 (07.06.2007), entire document	1 - 15
A	US 2006/0112222 A1 (Barrall) 25 May 2006 (25.05.2006), entire document	1 - 15
A	US 7,707,219 B1 (Bruso et al.) 27 April 2010 (27.04.2010), entire document	1 - 15
A	US 2005/0066095 A1 (Mullick et al.) 24 March 2005 (24.03.2005), entire document	1 - 15
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/>		
* Special categories of cited documents:		
"A" document defining the general state of the art which is not considered to be of particular relevance		"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
"E" earlier application or patent but published on or after the international filing date		"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)		"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
"O" document referring to an oral disclosure, use, exhibition or other means		"&" document member of the same patent family
"P" document published prior to the international filing date but later than the priority date claimed		
Date of the actual completion of the international search 28 August 2014 (28.08.2014)		Date of mailing of the international search report 03 OCT 2014
Name and mailing address of the ISA/US Mail Stop PCT, Attn: ISA/US, Commissioner for Patents P.O. Box 1450, Alexandria, Virginia 22313-1450 Facsimile No. 571-273-3201		Authorized officer: Lee W. Young PCT Helpdesk: 571-272-4300 PCT OSP: 571-272-7774

フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US

(72)発明者 マダヴァラプ , 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 ニューコム , クリストファー・リチャード
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 グプタ , アヌラグ・ウィンドラス
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