

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4939440号
(P4939440)

(45) 発行日 平成24年5月23日(2012.5.23)

(24) 登録日 平成24年3月2日(2012.3.2)

(51) Int.Cl.

F I

G 0 6 F 12/00 (2006.01)

G 0 6 F 12/00 5 1 8 A

G 0 6 F 12/00 5 3 1 D

G 0 6 F 12/00 5 3 3 J

請求項の数 23 (全 14 頁)

(21) 出願番号 特願2007-556381 (P2007-556381)
 (86) (22) 出願日 平成18年2月17日(2006.2.17)
 (65) 公表番号 特表2008-530716 (P2008-530716A)
 (43) 公表日 平成20年8月7日(2008.8.7)
 (86) 国際出願番号 PCT/US2006/005909
 (87) 国際公開番号 W02006/089263
 (87) 国際公開日 平成18年8月24日(2006.8.24)
 審査請求日 平成21年1月27日(2009.1.27)
 (31) 優先権主張番号 11/061, 152
 (32) 優先日 平成17年2月18日(2005.2.18)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 502303739
 オラクル・インターナショナル・コーポレ
 イション
 アメリカ合衆国、94065 カリフォル
 ニア州、レッドウッド・ショアーズ、オラ
 クル・パークウェイ、500
 (74) 代理人 100064746
 弁理士 深見 久郎
 (74) 代理人 100085132
 弁理士 森田 俊雄
 (74) 代理人 100083703
 弁理士 仲村 義平
 (74) 代理人 100096781
 弁理士 堀井 豊

最終頁に続く

(54) 【発明の名称】 データベースシステムにおいて報告トランザクションを処理する方法およびメカニズム

(57) 【特許請求の範囲】

【請求項 1】

データベースシステムにおいて報告トランザクションを処理するコンピュータによって
 実行される方法であって、プロセッサを用いて以下の処理を実行することを含み、当該
 以下の処理は、

データベースのスナップショットを取得することを備え、前記データベースは、プライ
 マリノードおよび前記データベースの複製されたデータベースを用いないフェイルオーバ
 ノードにリンクされ、前記処理はさらに、

前記プライマリノードで1つ以上の非報告トランザクションを実行することを備え、前
 記1つ以上の非報告トランザクションは、1つ以上の第1のデータベースクエリを、前記
 1つ以上の第1のデータベースクエリの実行において前記データベースの最新の更新を用
 いることなく実行し、前記処理はさらに、

前記プライマリノードで前記1つ以上の非報告トランザクションを実行すると同時に
 、前記フェイルオーバノードで報告トランザクションを実行するために前記スナップショ
 ットを利用することを備え、

前記報告トランザクションは、1つ以上の第2のデータベースクエリを、前記1つ以上
 の第2データベースクエリの実行において前記データベースの最新の更新を用いて実行し

、
 前記プライマリノードのトランザクションは、前記データベースを修正することを許可
 される一方で、前記フェイルオーバノードのトランザクションは、前記データベースを直

10

20

接的に修正することを許可されない、コンピュータによって実行される方法。

【請求項 2】

前記フェイルオーバーノードで 1 つ以上の一時的なテーブルを作成することをさらに備え、前記 1 つ以上の一時的なテーブルは、前記報告トランザクションが前記フェイルオーバーノードで実行されるときに使用される、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 3】

前記 1 つ以上の一時的なテーブルは、前記報告トランザクションにおけるクエリスクリプトを通じて作成される、請求項 2 に記載のコンピュータによって実行される方法。

【請求項 4】

前記 1 つ以上の一時的なテーブルのうち少なくとも 1 つは、前記報告トランザクションにおける 2 つ以上のクエリにアクセス可能である、請求項 2 に記載のコンピュータによって実行される方法。

【請求項 5】

前記データベースにおける 1 つ以上のスキーマを修正することをさらに備え、前記 1 つ以上のスキーマは、前記報告トランザクションが前記フェイルオーバーノードで実行されるときに使用される、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 6】

前記 1 つ以上のスキーマは、前記プライマリノードで実行する前記 1 つ以上の非報告トランザクションにアクセス可能ではない、請求項 5 に記載のコンピュータによって実行される方法。

【請求項 7】

前記 1 つ以上のスキーマのうち少なくとも 1 つは 1 つ以上のテーブルを含む、請求項 5 に記載のコンピュータによって実行される方法。

【請求項 8】

前記プライマリノードで 1 つ以上のユーザ定義プロシージャにアクセスすることをさらに備え、前記 1 つ以上のユーザ定義プロシージャは、前記報告トランザクションが前記フェイルオーバーノードで実行されるときに使用される、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 9】

前記データベースに一時的なスペースを確保することをさらに備え、前記一時的なスペースは、前記報告トランザクションが前記フェイルオーバーノードで実行されるときに使用される、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 10】

前記プライマリノードおよび前記フェイルオーバーノードはクラスタの一部である、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 11】

前記クラスタは 1 つ以上のさらなるフェイルオーバーノードを含む、請求項 10 に記載のコンピュータによって実行される方法。

【請求項 12】

前記 1 つ以上の非報告トランザクションのうち少なくとも 1 つはリード・ライトトランザクションである、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 13】

前記報告トランザクションおよび前記 1 つ以上の非報告トランザクションはワークロードの一部である、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 14】

前記報告トランザクションはリアルタイムに近い報告を与える、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 15】

前記プライマリノードのみが前記データベースを修正できる、請求項 1 に記載のコンピ

10

20

30

40

50

ユータによって実行される方法。

【請求項 16】

前記スナップショットはユーザコマンドに応答して取得される、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 17】

前記スナップショットはリード・オンリである、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 18】

前記スナップショットは前記プライマリノードによって修正されることができない、請求項 1 に記載のコンピュータによって実行される方法。

10

【請求項 19】

前記スナップショットおよび前記データベースはディスクスペースを共有する、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 20】

前記スナップショットは最新のものである、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 21】

前記スナップショットは、前記フェイルオーバーノードで前記報告トランザクションを実行するために直接に使用される、請求項 1 に記載のコンピュータによって実行される方法。

20

【請求項 22】

コンピュータによって実行されることで、請求項 1 ~ 21 のいずれかに記載の方法を実行する、プログラム。

【請求項 23】

請求項 1 ~ 21 のいずれかに記載の方法を実行する、システム。

【発明の詳細な説明】

【技術分野】

【0001】

背景および概要

この発明はデータベースシステムに関する。より詳細には、この発明は、データベースシステムにおいて報告トランザクション (reporting transaction) を処理する方法およびメカニズムに向けられる。

30

【背景技術】

【0002】

多くのデータベースシステムは、ペーの速い今日の市場において極めて重要である高可用性を保証するためにフェイルオーバークラスタを利用する。フェイルオーバークラスタでは、データベースはプライマリノードおよび少なくとも 1 つのフェイルオーバーノード (スベアノードとしても公知である) にリンクされる。データベースおよびウェブサーバなどのアプリケーションは、誤動作するまでプライマリノードで動作する。誤動作が発生すると、アプリケーションはフェイルオーバーノードで再開される。フェイルオーバーノードおよびプライマリノードが単一のクラスタに属しているので、プライマリノードの故障を検出するために標準的なハートビートメカニズムを使用できる。

40

【発明の開示】

【発明が解決しようとする課題】

【0003】

フェイルオーバークラスタに関する 1 つの問題は、フェイルオーバーノードをプライマリノードと同時に使用できないことである。したがって、プライマリハードウェアが故障したときにのみ使用される追加のハードウェアを購入するコストを正当化することは困難であり得る。ある特定の並列データベースシステムは、2 つ以上のノードがクラスタにおけるデータベースに同時にアクセスできるアクティブ/アクティブクラスタを利用することに

50

よってこの問題を解決する。しかしながら、アクティブ/アクティブクラスタは、クラスタにおけるすべてのノードからの同時の読取および修正が存在する状態でデータベースが確実に一貫性があるようにするために複雑な並行処理制御メカニズムを必要とする。

【 0 0 0 4 】

ユーザが直面する別の問題は、報告トランザクションが他のトランザクションと同時に実行される混合ワークロードを動作させる必要があることである。理想的には、リアルタイムの報告は各報告トランザクションによって与えられる。すなわち、最新の更新からの結果はトランザクションにおけるクエリによって使用される。さらに、ユーザは、非報告トランザクション (non-reporting transaction) と報告トランザクションとの間の (たとえば、CPU またはメモリについての) ハードウェアリソースの競合を回避するために、別個に報告トランザクションを動作させることを好む。

10

【 0 0 0 5 】

アクティブ/アクティブクラスタリングをサポートしないデータベースシステムでは、報告のために複製データベースが作成され、使用され得る。しかしながら、複製データベースがプライマリデータベースの完全なコピーであるので、この解決法は格納コストを2倍にする。さらに、複製データベースはしばしばプライマリデータベースに遅れをとる。なぜなら、プライマリデータベースにおける変更を瞬時に複製することが実現可能でない可能性があるためである。たとえ瞬時の複製が実現可能であったとしても、プライマリデータベースでのすべてのコミットが同期して報告データベースに複製される必要があるだろうという理由で、プライマリデータベースのスループットは大幅に影響を受けるであろう。

20

【 0 0 0 6 】

したがって、フェイルオーバークラスタを利用するデータベースシステムにおいて報告トランザクションを実行することに関するこれらのおよび他の問題に方法およびメカニズムが対処する必要がある。

【課題を解決するための手段】

【 0 0 0 7 】

この発明の実施例は、データベースシステムにおいて報告トランザクションを処理するための改良された方法、システムおよび媒体を提供する。実施例によれば、データベースのスナップショットが取得される。データベースはプライマリノードおよびフェイルオーバーノードにリンクされる。次いで1つ以上の非報告トランザクションがプライマリノードで実行され、プライマリノードで1つ以上の非報告トランザクションを実行するのと同時にフェイルオーバーノードで報告トランザクションを実行するためにスナップショットが利用される。

30

【発明を実施するための最良の形態】

【 0 0 0 8 】

この発明の局面、目的および利点のさらなる詳細について、詳細な説明、図面および特許請求の範囲において以下で説明する。先の一般的な説明および以下の詳細な説明の両方は例示的および説明的なものであり、この発明の範囲に関して限定的であるように意図されるものではない。

40

【 0 0 0 9 】

添付の図面は、この発明をさらに理解できるようにするために含まれ、詳細な説明とともにこの発明の原理を説明するのに役立つ。

【 0 0 1 0 】

詳細な説明

データベースシステムにおける報告トランザクションの処理を開示する。複雑な整合性およびルーティングメカニズムを必要とするアクティブ/アクティブクラスタを利用するか、または追加のハードウェアの購入を必然的に伴い、潜在的にデータが古い別個の複製データベースを有するのではなく、報告トランザクションは、プライマリノードで動作する非報告トランザクションと同時に、データベーススナップショットを使用して、フェイ

50

ルオーバノードで実行される。これは、そうでなければアイドルのままであろうフェイルオーバノードを利用し、最新のスナップショットが使用されるときにリアルタイムに近い報告を与える。

【 0 0 1 1 】

データベースシステムにおいて報告トランザクションを処理する方法を図 1 に示す。102において、データベースのスナップショットが取得される。データベースはプライマリノードおよびフェイルオーバノードにリンクされる。いくつかの実施例では、プライマリノードのみがデータベースを修正することを許可される。クライアント接続部は、すべての報告トランザクションをフェイルオーバノードに向け、すべての他のトランザクションをプライマリノードに向けるよう構成され得るであろう。フェイルオーバノードが、場合によってはデータベースを修正し得るであろうトランザクションを自動的にプライマリノードにルーティングすることも可能であり得る。このルーティングは、セッションがデータベースを修正することになるかどうかを識別するトランザクションにリード・ライト (READ-WRITE) またはリード・オンリ (READ-ONLY) という印をつけることによってなされ得る。

【 0 0 1 2 】

次いで1つ以上の非報告トランザクションがプライマリノードで実行され (1 0 4)、プライマリノードで1つ以上の非報告トランザクションを実行すると同時にフェイルオーバノードで報告トランザクションを実行するためにスナップショットが利用される (1 0 6)。報告トランザクションおよび非報告トランザクションの各々は、1つ以上のクエリを備える。そして、非報告トランザクションはリード・ライトまたはリード・オンリトランザクションであってもよいが、報告トランザクションは通常リード・オンリトランザクションである。

【 0 0 1 3 】

スナップショットは、データベースのある時点のコピーであり、スナップショットが取得された後に修正されるデータベースブロックを除いて、データベースと同一のディスクスペースを共有する。これは、スナップショットが修正されないままであるように、変更されたブロックが新しい場所に書込まれる標準的なコピー・オン・ライトメカニズムによって達成されることができる。スナップショットがリード・オンリであり、プライマリノードによって修正されることができないので、フェイルオーバノードで動作するクエリは、プライマリノードとの調整を必要とすることなく、使用されるスナップショットと一貫性のある結果を返すことになる。そして、スナップショットが一貫性があり、データベース全体のためのものである (すなわち、クエリにおいて参照されるスナップショットおよびテーブルの中の索引がすべて一貫性がある) ので、既存のクエリ実行エンジンは修正される必要がない。さまざまなスナップショット方法論が利用可能であり、ファイル、アプリケーション、システムまたはデータベースレベルで実現され得る。たとえば、ファイルレベルのスナップショットの作成についての説明は、http://www.netapp.com/tech_library/3002.htmlにおいて見ることができる。

【 0 0 1 4 】

スナップショットは、すべての変更されないデータについてデータベースと同一のディスク記憶装置を使用するので、ディスクスペースおよびCPU使用率の両方の点で比較的安価に作成される。したがって、データベースシステムはかなり頻繁に、たとえば10秒ごとにスナップショットを取得するよう構成され得る。しかしながら、ユーザコマンドにตอบสนองして、たとえば報告セッションまたは他のこのようなメトリクスによって所望されるサービスの質に基づいてデータベースシステムがスナップショットを生成することも可能である。最新のスナップショットを使用してフェイルオーバノードで報告トランザクションを実行することにより、リアルタイムに近い報告が与えられることになる。なぜなら、最新の更新は報告トランザクションにおけるクエリによって使用されることになるためである。しかしながら、ユーザは取得された最新のスナップショットよりも古いスナップショットの使用を指定することも許可され得る。

【 0 0 1 5 】

図 2 は、プライマリノード 2 0 2、フェイルオーバーノード 2 0 4 およびデータベース 2 0 6 を有するクラスタ 2 0 0 を示す。データベース 2 0 6 のスナップショット 2 0 8 が取得されている。複数の非報告トランザクション 2 1 0 a および 2 1 0 b がプライマリノード 2 0 2 で動作している間、スナップショット 2 0 8 はフェイルオーバーノード 2 0 4 で報告トランザクション 2 1 2 を実行するために使用される。いくつかの実施例では、非報告トランザクション 2 1 0 a および 2 1 0 b ならびに報告トランザクション 2 1 2 はワークロードの一部である。

【 0 0 1 6 】

データベースシステムにおいて報告トランザクションを処理するための方法のプロセスフローを図 3 に示す。この実施例によれば、プライマリノードおよびフェイルオーバーノードにリンクされるデータベースのスナップショットが取得される (3 0 2)。3 0 4 において、1 つ以上の非報告トランザクションがプライマリノードで実行される。プライマリノードで 1 つ以上の非報告トランザクションを実行するのと同時にフェイルオーバーノードで報告トランザクションを実行するためにスナップショットが利用される (3 0 6)。次いで、報告トランザクションがフェイルオーバーノードで実行されるときに 1 つ以上の一時的なテーブルが作成され、使用される (3 0 8)。

【 0 0 1 7 】

クラスタ 4 0 0 を図 4 に示す。クラスタ 4 0 0 は、プライマリノード 4 0 2、フェイルオーバーノード 4 0 4 およびデータベース 4 0 6 を含む。この例では、スナップショット 4 0 8 a が取得され、非報告トランザクション 4 1 0 がプライマリノード 4 0 2 で動作している間にフェイルオーバーノード 4 0 4 で報告トランザクション 4 1 2 を実行するために使用される。報告トランザクション 4 1 2 の実行中、一時的な結果を格納するためにトランザクション 4 1 2 におけるクエリスクリプトを通じて一時的なテーブル 4 1 4 a および 4 1 4 b が作成される。これらの一時的なテーブル 4 1 4 a および 4 1 4 b はプライマリノード 4 0 2 に透過的に送られ、プライマリノード 4 0 2 は次いで一時的なテーブル 4 1 4 a および 4 1 4 b のためにデータベース 4 0 6 においてスペースを割当てる。フェイルオーバーノード 4 0 4 において一時的なテーブル 4 1 4 a および 4 1 4 b に後に保存される変更はプライマリノード 4 0 2 に送られる必要はない。

【 0 0 1 8 】

図 4 では、データベース 4 0 6 の新しいスナップショット 4 0 8 b が取得されて、報告トランザクション 4 1 2 における後続のクエリが一時的なテーブル 4 1 4 a および 4 1 4 b にアクセスできるようにする。しかしながら、他の実施例では、作成されるすべての一時的なテーブルよりも少ない一時的なテーブルが後続のクエリによるアクセスのために保有されることになる。したがって、クエリの完了後、フェイルオーバーノードは、テーブルのために割当てられたデータベーススペースを解放するために、一時的なテーブルを削除でき、削除部分をプライマリノードに送ることができる。

【 0 0 1 9 】

一貫性のある結果を保証するために、単一のクエリは通常同一のスナップショットを使用することになる。しかしながら、図 4 の例に見られるように、同一のセッションまたはトランザクション内の後続のクエリは、以前のクエリによって使用されたスナップショットと同一のスナップショットまたはそれよりも最近のスナップショットを使用してもよい。

【 0 0 2 0 】

データベースシステムにおいて報告トランザクションを処理する別の方法を図 5 に示す。5 0 2 において、データベースのスナップショットが取得される。この実施例では、データベースはプライマリノードおよびフェイルオーバーノードにリンクされる。次いで 1 つ以上の非報告トランザクションがプライマリノードで実行され (5 0 4)、プライマリノードで 1 つ以上の非報告トランザクションを実行するのと同時にフェイルオーバーノードで報告トランザクションを実行するためにスナップショットが利用される (5 0 6)。5 0

10

20

30

40

50

8において、報告トランザクションがフェイルオーバーノードで実行されるときにデータベースにおける1つ以上のスキーマが修正され、使用される。1つ以上のスキーマは、プライマリノードで作成されていた可能性があり、フェイルオーバーノードで報告トランザクションが使用するために「印をつけられていた」または「確保されていた」可能性がある。さらに、1つ以上のスキーマへの変更はプライマリノードと調整することなくなされ得る。

【0021】

データベーススキーマはオブジェクトの集まりである。スキーマオブジェクトは、たとえばテーブル、ビュー、シーケンスおよびストアドプロシージャを含むが、それらに限定されない。テーブルは、概してデータベースにおける編成の基本単位であり、それぞれの行および列に格納されたデータを備える。ビューは、1つ以上のテーブルでのデータの特別仕立ての提示である。ビューは、データが基づいているテーブル、すなわちベーステーブルからデータを導き出す。さらには、ベーステーブルがテーブルである可能性もあれば、ベーステーブル自体がビューである可能性もある。ビューの一例は、テーブルからテーブルのデータの列のうち2列をマイナスしたものである。

【0022】

シーケンスは、1つ以上のデータベーステーブルの数値列を識別する固有の数字の連続的なリストである。シーケンスは概して、単一のテーブルまたは複数のテーブルの行について固有の数値を自動的に生成することによってアプリケーションプログラミングを単純化する。シーケンスを使用することによって、2人以上のユーザが概して同時にテーブルにデータを入力できる。ストアドプロシージャは概して、特定のタスクを行なうために実行可能な単位としてともにグループ分けされる1組のコンピュータ命令文である。

【0023】

図6は、プライマリノード602、2つのフェイルオーバーノード604aおよび604b、ならびにデータベース606を有するクラスタ600を示す。データベース606のスナップショット608が取得されている。この実施例では、スナップショット608を介してしかフェイルオーバーノード604aおよび604bに開いていない残りのデータベース606とは異なって、リード・ライトモードでデータベース606内のスキーマ614aおよび614bをフェイルオーバーノード604aおよび604bが利用できる。この状況下で、スキーマ614aおよび614bは、フェイルオーバーノード604aおよび604bで動作する報告トランザクション612aおよび612bによってそれぞれに修正され得る。スキーマ614aおよび614bに含まれるデータがフェイルオーバーノード604a、604bとプライマリノード602との間で共有されないため、プライマリノード602で実行する非報告トランザクション610はデータベース606におけるスキーマ614aおよび614bにアクセスできない。

【0024】

データベースシステムにおいて報告トランザクションを処理するための方法のフローチャートを図7に示す。702において、プライマリノードおよびフェイルオーバーノードにリンクされるデータベースのスナップショットが取得される。704において、1つ以上の非報告トランザクションがプライマリノードで実行される。次いで、プライマリノードで1つ以上の非報告トランザクションを実行すると同時にフェイルオーバーノードで報告トランザクションを実行するためにスナップショットが利用される(706)。

【0025】

この実施例では、報告トランザクションがフェイルオーバーノードで実行されるときに、プライマリノードの1つ以上のユーザ定義プロシージャがアクセスされ、使用される(708)。ユーザ定義プロシージャは、一般に複雑な報告の準備をより容易にするために使用され、通常はプライマリノードで作成され、コンパイルされる。これらのプロシージャには、ちょうど他のデータベースオブジェクトと同様に、フェイルオーバーノードからアクセス可能である。

【0026】

データベースシステム 800 を図 8 に示す。この図はユーザ 802、クライアント 804、プライマリノード 806、フェイルオーバーノード 808 およびデータベース 810 しか示していないが、システム 800 は他のクラスタ、ノード、ユーザ、データベースおよびクライアントを含んでもよい。この例では、ユーザ 802 はクライアント 804 を介してプライマリノード 806 でプロシージャ 818a および 818b を定義している。データベース 810 のスナップショット 812 が取得された後、スナップショット 812 ならびにユーザ定義プロシージャ 818a および 818b を使用して、プライマリノード 806 で非報告トランザクション 814 を動作させるのと同時に報告トランザクション 816 がフェイルオーバーノード 808 で実行される。図 8 に示すように、スナップショット 812 の使用は、ユーザ定義プロシージャ 818a および 818b とは異なって、直接的である。すなわち、スナップショット 812 はプライマリノード 806 を経ることなく使用される。

10

【0027】

データベースシステムにおいて報告トランザクションを処理する別の方法を図 9 に示す。この方法によれば、902 において、データベースのスナップショットが取得される。データベースはプライマリノードおよびセカンダリノードにリンクされる。次いで、904 において、1 つ以上の非報告トランザクションがプライマリノードで実行され、906 において、プライマリノードで 1 つ以上の非報告トランザクションを実行すると同時にフェイルオーバーノードで報告トランザクションを実行するためにスナップショットが利用される。報告トランザクションがフェイルオーバーノードで実行されるときにデータベース

20

【0028】

データベースに一時的なスペースを確保するために、フェイルオーバーノードはメッセージをプライマリノードに送信することができる。なぜなら、確保することは通常、整合性の問題を回避するためにプライマリノードによって行なわれるカタログの変更を必要とするためである。一旦フェイルオーバーノードのためにスクラッチディスクスペースが確保されると、プライマリノードからの介入なしに一時的なスペース自体への書込を行なうことができる。スクラッチスペースによって一時的なファイルを作成することができる。これらの一時的なファイルは時には、メインメモリに適合しない一時的な動作の結果、たとえばもろもろの中間結果、JOIN 法において使用されるハッシュテーブルなどを格納するために必要である。

30

【0029】

図 10 は、プライマリノード 1002 ならびに 3 つのフェイルオーバーノード 1004a、1004b および 1004c を有するクラスタ 1000 を示し、それらはすべてデータベース 1006 にリンクされる。この図では、ユーザ定義プロシージャ 1012 は、リード・ライトトランザクション 1010a およびリード・オンリトランザクション 1010b とともに、プライマリノード 1002 で見られることができる。報告トランザクション 1014a および 1014b はフェイルオーバーノード 1004a で動作している。さらに、報告トランザクション 1014d、1014e および 1014f がフェイルオーバーノード 1004c で動作している間、報告トランザクション 1014c はフェイルオーバーノード 1004b で動作している。データベース 1006 の 3 つのスナップショット 1008a、1008b および 1008c は異なるときに取得された。報告トランザクションの各々はスナップショットのうちの 1 つを使用して実行され得る。しかしながら、同一のフェイルオーバーノードでの報告トランザクションは同一のスナップショットを利用する必要はない。たとえば、フェイルオーバーノード 1004c での報告トランザクション 1014d、1014e および 1014f は各々が異なるスナップショット 1008 を使用できる。

40

【0030】

図 10 に示すように、3 つの一時的なスペース 1016a、1016b および 1016c は、フェイルオーバーノード 1004a、1004b および 1004c のためにそれぞれにデータベース 1006 に確保されている。フェイルオーバーノード 1004a、1004

50

bおよび1004cの各々は要求をプライマリノード1002に送信して、それぞれのクラッチスペースを確保する。他の実施例では、フェイルオーバーノード1004a、1004bおよび1004cは1つ以上の一時的なスペースを共有してもよい。

【0031】

システムアーキテクチャの概要

図11は、この発明の実施例を実現するのに好適なコンピュータシステム1100のブロック図である。コンピュータシステム1100は、プロセッサ1104、システムメモリ1106（たとえばRAM）、静的記憶装置1108（たとえばROM）、ディスクドライブ1110（たとえば磁気もしくは光学）、通信インターフェイス1112（たとえばモデムもしくはイーサネット（登録商標）カード）、ディスプレイ1114（たとえばCRTもしくはLCD）、入力装置1116（たとえばキーボード）およびカーソル制御装置1118（たとえばマウスもしくはトラックボール）などのサブシステムおよび装置を相互接続する、情報を通信するためのバス1102または他の通信メカニズムを含む。

【0032】

この発明の一実施例によれば、コンピュータシステム1100は、システムメモリ1106に含まれる1つ以上の命令の1つ以上のシーケンスを実行するプロセッサ1104によって特定の動作を行なう。このような命令は、静的記憶装置1108またはディスクドライブ1110などの別のコンピュータ可読媒体からシステムメモリ1106に読取られることができる。代替的な実施例では、この発明を実現するためにソフトウェア命令の代わりにまたはソフトウェア命令と組合せられてハードワイヤード回路が使用されてもよい。

【0033】

本明細書において使用される「コンピュータ可読媒体」という用語は、実行のためにプロセッサ1104に命令を与えることに関与する任意の媒体を指す。このような媒体は、不揮発性媒体、揮発性媒体および伝送媒体を含むがそれらに限定されない多くの形態を取り得る。不揮発性媒体はたとえばディスクドライブ1110などの光学または磁気ディスクを含む。揮発性媒体はシステムメモリ1106などのダイナミックメモリを含む。伝送媒体はバス1102を備えるワイヤを含む同軸ケーブル、銅線および光ファイバを含む。伝送媒体は、電波および赤外線データ通信中に発生するものなどの音波または光波の形態も取り得る。

【0034】

コンピュータ可読媒体の一般的な形態は、たとえばフロッピー（登録商標）ディスク、フレキシブルディスク、ハードディスク、磁気テープ、他の磁気媒体、CD-ROM、他の光学媒体、パンチカード、紙テープ、穴のパターンを有する他の物理的な媒体、RAM、PROM、EPROM、FLASH-EPROM、他のメモリチップもしくはカートリッジ、搬送波、またはコンピュータが読取ることができる他の媒体を含む。

【0035】

この発明の実施例では、この発明を実施するための命令のシーケンスの実行は単一のコンピュータシステム1100によって行なわれる。この発明の他の実施例によれば、通信リンク1120（たとえばLAN、PTSNまたはワイヤレスネットワーク）によって結合される2つ以上のコンピュータシステム1100が互いに連携してこの発明を実施するのに必要な命令のシーケンスを実行してもよい。

【0036】

コンピュータシステム1100は、通信リンク1120および通信インターフェイス1112を介して、プログラムすなわちアプリケーションコードを含むメッセージ、データおよび命令を送信および受信できる。受信されたプログラムコードは、受信したままでプロセッサ1104によって実行されてもよく、および/または後に実行するためにディスクドライブ1110もしくは他の不揮発性記憶装置に格納されてもよい。

【0037】

先の明細書では、具体的な実施例を参照してこの発明について説明してきた。しかしな

10

20

30

40

50

がら、この発明のより広い精神および範囲から逸脱することなくさまざまな修正および変更がなされ得ることは明白であろう。たとえば、プロセスアクションの特定の順序付けを参照して上述のプロセスフローを説明する。しかしながら、説明するプロセスアクションの多くの順序付けはこの発明の範囲または動作に影響を及ぼすことなく変更され得る。したがって、明細書および図面は限定的な意味ではなく例示的な意味で考えられるべきである。

【図面の簡単な説明】

【 0 0 3 8 】

【図 1】この発明の実施例に従ってデータベースシステムにおいて報告トランザクションを処理する方法のフローチャートである。

10

【図 2】この発明の一実施例に従うフェイルオーバークラスタにおける報告トランザクションの実行を示す。

【図 3】この発明の別の実施例に従ってデータベースシステムにおいて報告トランザクションを処理するための方法のプロセスフローを示す。

【図 4】この発明の別の実施例に従ってクラスタにおいて報告トランザクションがいかに処理されるかの一例である。

【図 5】データベースシステムにおいて報告トランザクションを処理する方法の一実施例を示す。

【図 6】複数のフェイルオーバーノードを有するクラスタを示す。

【図 7】データベースシステムにおいて報告トランザクションを処理するための方法の別の実施例を示す。

20

【図 8】サンプルのデータベースシステムを示す。

【図 9】この発明のさらなる実施例に従ってデータベースシステムにおいて報告トランザクションを処理するための方法のプロセスフローである。

【図 10】この発明のさらなる実施例に従うフェイルオーバークラスタにおける複数の報告および非報告トランザクションの実行を示す。

【図 11】この発明の実施例が実現され得るシステムアーキテクチャの図である。

【図 1】

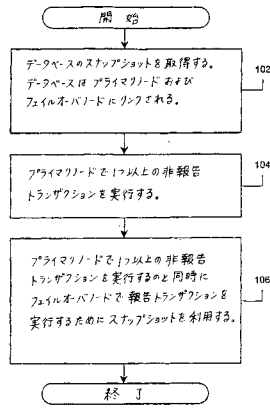


FIG. 1

【図 2】

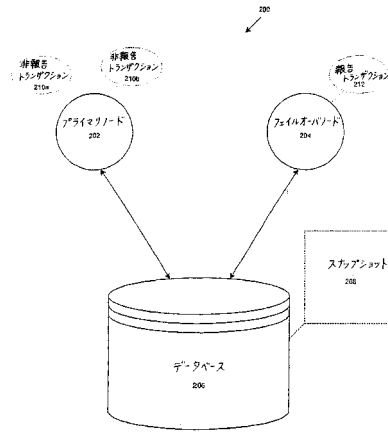


FIG. 2

【図 3】

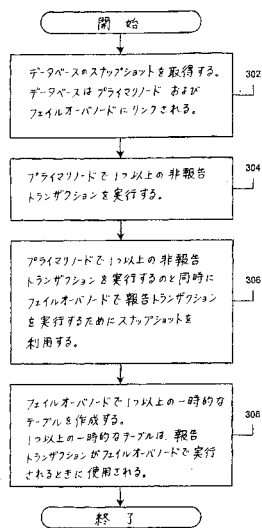


FIG. 3

【図 4】

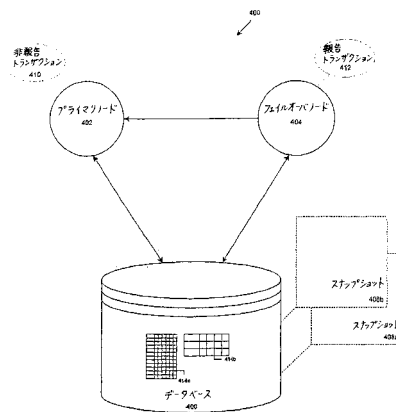


FIG. 4

【図 5】

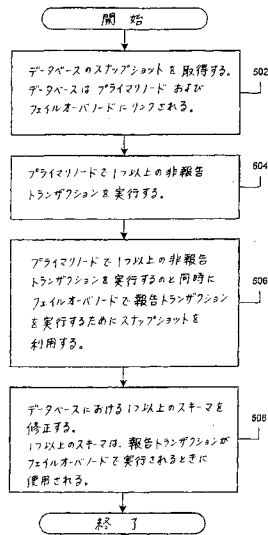


FIG. 5

【図 6】

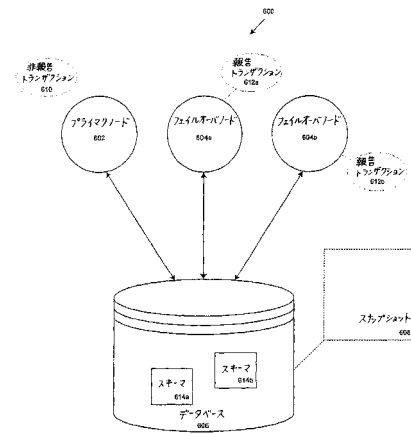


FIG. 6

【図 7】

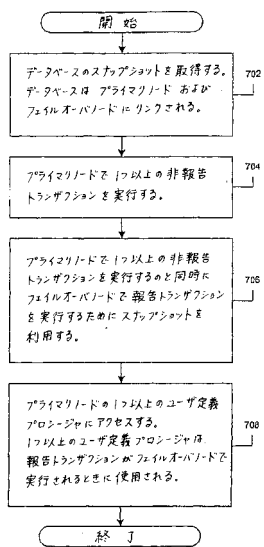


FIG. 7

【図 8】

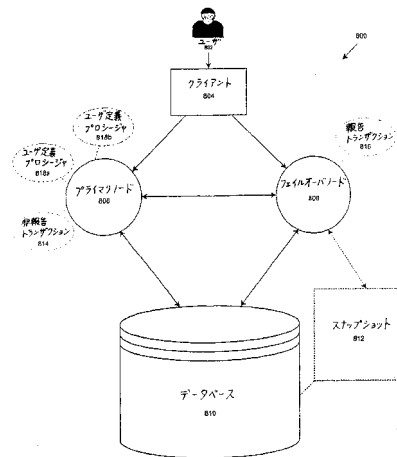


FIG. 8

【図 9】

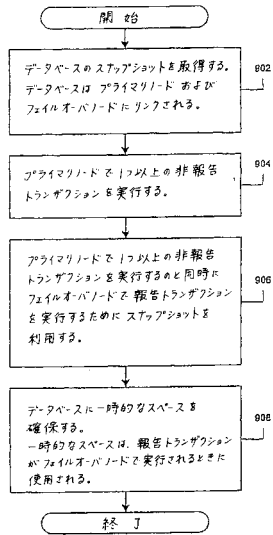


FIG. 9

【図 10】

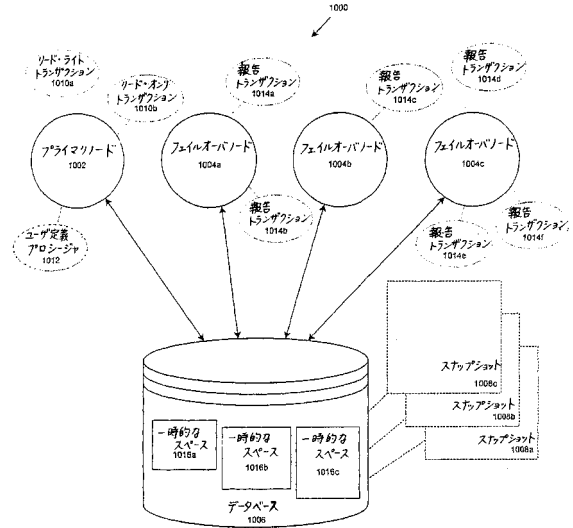


FIG. 10

【図 11】

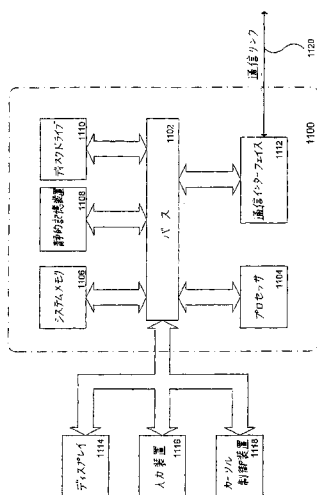


FIG. 11

フロントページの続き

(74)代理人 100098316

弁理士 野田 久登

(74)代理人 100109162

弁理士 酒井 將行

(74)代理人 100111246

弁理士 荒川 伸夫

(72)発明者 チャンドラセカラン, サシカンス

アメリカ合衆国、 9 5 1 3 4 カリフォルニア州、サン・ノゼ、ルネサンス・ドライブ、 4 3 2 5
、ナンバー・ 2 1 3

(72)発明者 ブルシーノ, アンジェロ

アメリカ合衆国、 9 4 0 2 2 カリフォルニア州、ロス・アルトス、ディステル・ドライブ、 4 3
6

審査官 桜井 茂行

(56)参考文献 特開 2 0 0 5 - 2 0 2 9 1 5 (J P , A)

特開 2 0 0 1 - 1 5 9 9 8 5 (J P , A)

米国特許出願公開第 2 0 0 5 / 0 1 3 8 3 1 2 (U S , A 1)

米国特許第 6 5 2 9 9 1 7 (U S , B 1)

(58)調査した分野(Int.Cl. , D B 名)

G06F 12/00

G06F 11/14

G06F 17/30