

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4378335号
(P4378335)

(45) 発行日 平成21年12月2日(2009.12.2)

(24) 登録日 平成21年9月18日(2009.9.18)

(51) Int. Cl.		F I			
G06F	3/06	(2006.01)	G06F	3/06	301Z
G06F	12/00	(2006.01)	G06F	12/00	518M
G06F	13/10	(2006.01)	G06F	13/10	340A
			G06F	12/00	514A

請求項の数 18 (全 14 頁)

(21) 出願番号	特願2005-262155 (P2005-262155)	(73) 特許権者	390009531
(22) 出願日	平成17年9月9日(2005.9.9)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公開番号	特開2007-72975 (P2007-72975A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公開日	平成19年3月22日(2007.3.22)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
審査請求日	平成20年1月21日(2008.1.21)	(74) 代理人	100086243 弁理士 坂口 博
早期審査対象出願		(74) 代理人	100091568 弁理士 市位 嘉宏
		(74) 代理人	100108501 弁理士 上野 剛史

最終頁に続く

(54) 【発明の名称】 ディスクへのトランザクション・データ書き込みの方式を動的に切り替える装置、切り替える方法、及び切り替えるプログラム

(57) 【特許請求の範囲】

【請求項1】

トランザクション・データをディスクに書き込む装置であって、
前記トランザクション・データを格納するディスク書き込み用キューとメモリ書き込み用キューに関連づける情報を管理するためのキュー管理テーブルを有するメモリであって、前記ディスク書き込み用キューと前記メモリ書き込み用キューがデータの耐久性保証方法を切り替える単位であるリージョンごとに設けられているものと、

先行トランザクションのディスクへの書き込み処理が終了する前に次のトランザクションが発生する頻度として定義されるディスクへの書き込み負荷が所定の閾値を超えたことを条件として前記キュー管理テーブルに前記メモリ書き込み用キューに関連づける情報を登録し、前記負荷が閾値を下回ることを条件として前記キュー管理テーブルから前記メモリ書き込み用キューに関連づける情報を削除する手段と、

前記キュー管理テーブルに前記ディスク書き込み用キューに関連づける情報のみが登録されている場合、前記ディスク書き込み用キューに格納された前記トランザクション・データを受け取り、冗長化メモリを使用せずにディスクに書き込む手段と、

前記キュー管理テーブルに前記メモリ書き込み用キューに関連づける情報が登録されている場合、前記メモリ書き込み用キューに格納された前記トランザクション・データを受け取り冗長化メモリに書き込む手段と、

を有するトランザクション・データをディスクに書き込む装置。

【請求項2】

前記キュー管理テーブルに前記メモリ書き込み用キューに関連づける情報が登録されている場合、前記トランザクション・データを前記冗長化メモリに書き込むために、前記メモリに書き込む手段にスレッドを割り当てる手段をさらに含む請求項 1 に記載のトランザクション・データをディスクに書き込む装置。

【請求項 3】

前記メモリに書き込む手段が、前記メモリ書き込み用キューに格納された前記トランザクション・データをマルチキャストで複数のメモリに書き込むことを特徴とする請求項 1 または 2 に記載のトランザクション・データをディスクに書き込む装置。

【請求項 4】

前記メモリに書き込む手段が、マルチキャストの宛先からの応答により所定の数のメモリへのコピーが完了したことが判明したら、データ書き込みの完了通知を返す手段をさらに含む、請求項 3 に記載のトランザクション・データをディスクに書き込む装置。

10

【請求項 5】

アプリケーションのプロセスの障害が発生した場合にも、前記負荷が閾値を超えたとする、請求項 1 に記載のトランザクション・データをディスクに書き込む装置。

【請求項 6】

データベースの一時的な応答遅延が発生した場合にも、前記負荷が閾値を超えたものとする、請求項 1 に記載のトランザクション・データをディスクに書き込む装置。

【請求項 7】

データベースに障害が発生した場合にも、前記負荷が閾値を超えたものとする、請求項 1 に記載のディスクへの書き込み方式を変更する装置。

20

【請求項 8】

ディスクの障害が発生した場合にも、前記負荷が閾値を超えたものとする、請求項 1 に記載のトランザクション・データをディスクに書き込む装置。

【請求項 9】

前記管理する手段が、さらに、データの耐久性保証方法を切り替える場合、前記ディスク書き込み用キューに排他ロックをかける手段を含む請求項 1 に記載のトランザクション・データをディスクに書き込む装置。

【請求項 10】

ライト・スルー方式とライト・ビハインド方式でトランザクション・データのディスクへの書き込み方式を変更する方法であって、

30

先行トランザクションのディスクへの書き込み処理が終了する前に次のトランザクションが発生する頻度として定義されるディスクへの書き込み負荷が所定の閾値より下から上へ変わったことに応答して、キュー・マネージャがディスク書き込み用キューとメモリ書き込み用キューに関連づける情報をキュー管理テーブルに登録するステップであって、前記ディスク書き込み用キューと前記メモリ書き込み用キューがデータの耐久性保証方法を切り替える単位であるリージョンごとに設けられているものと、

前記負荷が前記閾値より上から下へ変わったことに応答して、前記キュー・マネージャが前記キュー管理テーブルからメモリ書き込み用キューに関連づける情報を削除するステップと、

40

前記キュー管理テーブルに前記ディスク書き込み用キューに関連づける情報のみが登録されている場合、ディスク用データ記録手段が前記ディスク書き込み用キューに格納された前記トランザクション・データを受け取り、冗長化メモリを使用せずにディスクに書き込むステップと、

前記キュー管理テーブルに前記メモリ書き込み用キューに関連づける情報が登録されている場合、メモリ用データ記録手段が前記メモリ書き込み用キューに格納された前記トランザクション・データを受け取り冗長化メモリに書き込むステップと、

を含むディスクへの書き込み方式を変更する方法。

【請求項 11】

前記キュー管理テーブルに前記メモリ書き込み用キューに関連づける情報が登録されて

50

いる場合、前記トランザクション・データを冗長化メモリに書き込むために、書き込みタイミング制御手段が前記メモリに書き込む手段にスレッドを割り当てるステップをさらに含む請求項 10 に記載のディスクへの書き込み方式を変更する方法。

【請求項 12】

前記登録するステップの前に、登録する前記キューに対応するリージョンのディスク書き込み用キューに前記キュー・マネージャが排他ロックをかけるステップを含む請求項 10 に記載のディスクへの書き込み方式を変更する方法。

【請求項 13】

前記削除するステップの前に、削除する前記キューに対応するリージョンのディスク書き込み用キューに前記キュー・マネージャが排他ロックをかけるステップを含む請求項 10 に記載のディスクへの書き込み方式を変更する方法。

10

【請求項 14】

アプリケーションのプロセスのフェイルオーバーが発生した場合にも、前記キュー・マネージャは、前記負荷が閾値を超えたものとみなして処理する、請求項 10 に記載のディスクへの書き込み方式を変更する方法。

【請求項 15】

データベースの一時的な応答遅延が発生した場合にも、前記キュー・マネージャは、前記負荷が閾値を超えたものとみなして処理する、請求項 10 に記載のディスクへの書き込み方式を変更する方法。

【請求項 16】

20

データベースに障害が発生した場合にも、前記キュー・マネージャは、前記負荷が閾値を超えたものとみなして処理する、請求項 10 に記載のディスクへの書き込み方式を変更する方法。

【請求項 17】

ディスクの障害が発生した場合にも、前記キュー・マネージャは、前記負荷が閾値を超えたものとみなして処理する、請求項 10 に記載のディスクへの書き込み方式を変更する方法。

【請求項 18】

請求項 10 から 17 のいずれかに記載の方法を実現するステップを、コンピュータに実行させるためのプログラム。

30

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、オンライン・トランザクション処理（OLTP）等において、動的な負荷の変動、プロセスやハードウェアの障害に応じて、ディスクへのトランザクション・データの書き込みを、ライト・スルー方式とライト・ビハインド方式で動的に切り替えることが可能なトランザクション・データをディスクに書き込む装置、書き込む方法、及び書き込むプログラムに関する。

【背景技術】

【0002】

40

証券取引システムのようなオンライン・トランザクション処理アプリケーションでは、コミット済みのデータが障害によって失われることは許されない。通常、障害に対するデータの耐久性を保証する方法として、コミット時にディスクに直接データを書き込むライト・スルー方式が用いられる。

【0003】

従来は、データのディスク書き込みは一般的にレイテンシが大きくスループットも低いため、それを改善する手法としてライト・ビハインド方式が使われてきた。図 1 は、ライト・ビハインド方式によるデータの書き込み処理を簡単に説明している。ライト・ビハインド方式では、まず、コミット時に異なるノードやプロセス上のメモリ（遠隔メモリ）にデータを一時的にコピーして冗長化しておく（101）。コピー先の遠隔メモリの場所や

50

構成などはどのような障害に耐えたいかによって変わる。コミット完了後、ローカルプロセスのメモリ（ローカルメモリ）上のデータを非同期にディスクに書き込む（102）。もし、ディスクにデータが書き込まれる前に障害が発生しても、遠隔メモリ上からデータを復旧できる（103）。ディスクに書き込みが完了したらローカルメモリおよび遠隔メモリからデータを削除する（104）。メモリはディスクに比べてデータの書き込み速度が高速なため、アプリケーションはレイテンシを小さくしスループットを高められる。

【0004】

しかし、ライト・スルー方式とライト・ビハインド方式を切り替えるというデータの耐久性保証はアプリケーションの起動前に静的に設定するのが一般的であり、実行時の負荷の変化、プロセスやハードウェアの障害に対応して方式を動的に変更することはできなかった。

10

【0005】

また、ライト・スルー方式とライト・ビハインド方式を選択する特許文献1では、マルチプロセッサシステムにおいて、複数のプロセッサが共有する領域へは、CPUのライトアクセスは競合を防止するためライト・スルー方式でアクセスし、1つのCPUに固有のアドレスへのアクセスには、ライト・ビハインド方式でアクセスするのみで、アプリケーションの状況に応じて動的に変更することはできなかった。

【0006】

【特許文献1】特開平9-93265号公報

【発明の開示】

20

【発明が解決しようとする課題】

【0007】

コミット時のデータの保存先を、ディスクではなく遠隔メモリにすることによって、トランザクション・データの耐久性を保証しつつレイテンシを向上できる。しかし、ディスクに比べてメモリは高価であり、結果として容量が小さい。そのため、ライト・ビハインド方式では、ある一定量以上のコミット済みデータが一定時間以内にメモリに蓄積すると、メモリの使用容量が限界に達してしまい、それ以上にコミット済みデータを保存することができなくなる。その結果、メモリが溢れてしまう。一方で、これを防ぐためにトランザクション・データのディスク書き込みを待つと、ディスク書き込み処理がボトルネックとなってしまう。

30

【課題を解決するための手段】

【0008】

上記課題を解決するため、本発明においては、ライト・スルー方式とライト・ビハインド方式でトランザクション・データのディスクへの書き込み方式を動的に変更することが可能な、トランザクション・データをディスクに書き込む装置を提供する。該装置は、前記トランザクション・データを格納するディスク書き込み用キューとメモリ書き込み用キューとを管理するためのキュー管理テーブルを有するメモリを備える。また、該装置は、負荷が所定の閾値を超えたことを条件として前記キュー管理テーブルに前記メモリ書き込み用キューを登録し、前記負荷が閾値を下回ることを条件として前記キュー管理テーブルから前記メモリ書き込み用キューを削除する手段を有する。そして、該装置は、前記ディスク書き込み用キューに格納された前記トランザクション・データを受け取りディスクに書き込む手段と、前記メモリ書き込み用キューに格納された前記トランザクション・データを受け取り冗長化メモリに書き込む手段とを有する。該装置により、負荷に応じて、ディスクへのデータ書き込みをライト・スルー方式とライト・ビハインド方式で動的に切り替えることが可能になり、システムのリソースを有効に活用することができる。

40

【発明を実施するための最良の形態】

【0009】

以下、発明の実施の形態を通じて本発明を説明するが、以下の実施形態は特許請求の範囲にかかる発明を限定するものではなく、また実施形態の中で説明されている特徴の組み合わせは、発明の内容を理解しやすくするためのものもあり、その全てが発明の解決手段

50

に必須であるとは限らない。

【0010】

以下の実施の形態では、主に方法またはシステムについて説明するが、当業者であれば明らかとなり、本発明はコンピュータで使用可能なプログラムとしても実施できる。したがって、本発明は、ハードウェアとしての実施形態、ソフトウェアとしての実施形態またはソフトウェアとハードウェアとの組合せの実施形態をとることができる。プログラムは、ハードディスク、CD-ROM、光記憶装置または磁気記憶装置等の任意のコンピュータ可読媒体に記録できる。

【0011】

図2は、負荷の変動、プロセスやハードウェアの障害に応じて、ディスクへのトランザクション・データの書き込みを、ライト・スルー方式とライト・ビハインド方式で動的に切り替えるシステムが作動するハードウェア構成200の概略を示している。中央演算処理装置であるCPU201は各種オペレーティング・システムの制御下で様々なプログラムを実行する。CPU201はバス202を介して、メモリ203、ディスク204、ディスプレイ・アダプタ205、ユーザ・インタフェース206およびネットワーク・インタフェース207と相互接続される。ディスク204には、コンピュータを本発明を実現するためのシステムとして機能させるためのソフトウェア、オペレーティング・システムおよびデータベース・システム等が含まれる。

10

【0012】

ユーザ・インタフェース206を通して、キーボード208およびマウス209に接続され、ディスプレイ・アダプタ205を通してディスプレイ装置208に接続され、ネットワーク・インタフェース207を通してネットワーク211に接続される。ネットワーク211からさらに冗長化用ノード/プロセス212に接続される。冗長化用ノード/プロセスは、遠隔メモリ213を有する。この遠隔メモリ213は、ライト・ビハインド方式でトランザクション・データの耐久性を保証するため、コミット後でディスク204に反映されていないメモリ203のトランザクション・データを冗長化して保持する。このハード構成200は、コンピュータ・システム、バス配置およびネットワーク接続の一実施形態の例にすぎず、本発明の特徴は、さまざまなシステム構成で、同一の構成要素を複数有する形態で、または、ネットワーク上にさらに分散された形態でも実現することができる。

20

30

【0013】

図3は、ディスクへのトランザクション・データの書き込みを、ライト・スルー方式とライト・ビハインド方式で動的に切り替える装置の機能および処理されるデータを、概念的に表したものである。ここで表現される機能ブロックおよびデータの構成は、さらに細かい単位で表現可能であり、また、もっと大きい単位の機能ブロックでも表現可能である。これら機能ブロックの表現単位の違いは、発明を実施する際の本質的問題ではないことは、当業者には明らかである。ディスク301は、ハード・ディスクなどの一般にデータ書き込み速度が遅い装置で、データベース等のデータが最終的に記録される装置である。メモリ302は、ディスク301に比べて、書き込み速度が高速であり、ライト・ビハインド方式での書き込みの場合に、一時的にデータを冗長化するのに使用される。メモリ302は、障害に耐えるために使用されるので、通常は異なるノードやプロセス上に配置される。リージョン303は、データの耐久性保証方法を切り替える単位であり、あるトランザクションでアクセスするデータの集合は、必ず特定のリージョンに含まれる。リージョンの定義は、データ全体を一つのリージョンとすることもできるし、また、証券取引システムなどでは、会社毎の取引データをそれぞれのリージョンとして定義することもできる。ちなみにトランザクションはリージョンの境界を越えてデータにアクセスすることはできない。

40

【0014】

キュー管理テーブル304は、トランザクションのコミット毎のデータが入ったキューを管理するテーブルである。キューはリージョンとデータ記録手段の組み合わせごと

50

意される。図3においては、リージョンaおよびリージョンbと、ディスク用データ記録手段305およびメモリ用データ記録手段306の組み合わせとなるので、キューはディスク用キュー{D、a}と{D、b}、およびメモリ用キュー{M、a}と{M、b}の4つとなる。ディスク用キュー{D、a}と{D、b}は、ディスク301にデータを書き込むためのトランザクション・データを管理し、メモリ用キュー{M、a}と{M、b}は、メモリ302にデータを書き込むためのトランザクション・データを管理する。キュー・マネージャ307が、キュー管理テーブル304に対し、キューの登録、削除および変更等を行う。キュー・マネージャ307の機能は、書き込みタイミング制御手段308などの一部として取り込まれていてもいい。

【0015】

ディスク用データ記録手段305は、キュー{D、a}およびキュー{D、b}からトランザクション・データを受け取り、ディスク301に書き込む。メモリ用データ記録手段306は、キュー{M、a}およびキュー{M、b}からトランザクション・データを受け取り、メモリ302に書き込む。キューにはトランザクション・データそのもの、またはトランザクション・データが記録されている場所を示すポインタが入っていて、これらに従って、データ記録手段305および306は、データの書き込み処理を行う。キュー{M、a}およびキュー{M、b}は、通常はシステムのローカルのメモリ上(図2のメモリ203)にあるが、他のシステムに存在することも可能である。書き込みタイミング制御手段308は、ディスク用データ記録手段305およびメモリ用データ記録手段306にスレッドを割り当て、データ記録手段305およびメモリ用データ記録手段306がトランザクション・データを受取って書き込み処理を行う。また、書き込みタイミング制御手段308は、スレッドのプールを管理し、このプールから必要に応じてデータ記録手段305およびメモリ用データ記録手段306にスレッドを割り当て、不要になったスレッドを再びプールに戻す。そして、キューに対しては、それぞれのプロセスがデータの読み取りと書き込みを別々に行うので、データの順番を間違えて処理するのを防止するために、それぞれのプロセスは、キューに対して排他ロックをかける。

【0016】

本発明では、メモリのリソースを最大限に活用するため、ディスクへのトランザクション・データの書き込み処理では、通常の負荷の場合はメモリを使わないようにしておき、低負荷状態から高負荷状態に移行する場合にメモリを効率的に使って一度に大量のデータを処理できるようにする。すなわち、通常の負荷の場合は、メモリへのデータ蓄積を行わないライト・スルー方式でデータ書き込みを処理し、メモリの空き容量を最大に保つことで、高負荷時にライト・ビハインド方式に移行した際のメモリに蓄積できるデータ量を多くする。また、高負荷状態から平均的な負荷に戻る場合は、次に来る高負荷状態に耐えるため、ライト・ビハインド方式からライト・スルー方式に戻してメモリの空き容量を多くし、次のライト・ビハインド方式でメモリに蓄積できるデータ量を多くできるようにする。

【0017】

また、ライト・スルー方式でディスクに障害が発生した場合、アプリケーションを停止させないようにするために、一時的にライト・ビハインド方式に切り替えたい場合がある。逆に、ライト・ビハインド方式で、遠隔メモリに障害が発生し、同じくアプリケーションを停止させないようにするために、ライト・スルー方式に切り替えなければならない場合もある。例えば、ディスク、それを管理するOS(オペレーティング・システム)やデータベースが何らかの原因で一時的に回答が遅くなる応答遅延が発生する場合がある。データベース・プロセスに障害が発生した場合、データベースが一時的に回答できない状況もある。同様に、ライト・ビハインド方式で、遠隔メモリを管理するOSやプロセスの負荷が高くなり応答しなくなる場合がある。これらの障害は短時間で回復することが多いので、ある一定時間以上遠隔メモリやディスクが応答しない場合は、ライト・スルー方式とライト・ビハインド方式を相互に切り替えられると、アプリケーションに対して応答時間を保証できる。

10

20

30

40

50

【 0 0 1 8 】

さらに、ライト・ビハインド方式でアプリケーションが動作中にプロセスに障害が発生し、別プロセスが引き継ぐフェイルオーバーの場合、遠隔メモリとディスク上のデータの差分をディスクに反映する必要がある。しかし、引継ぎ後のプロセスが、差分の反映が完了するのを待ってからでないと新しいトランザクションを実行できないとすると、アプリケーションの可用性が低下してしまう。そこで、アプリケーションの可用性を重視する場合、差分を反映している間はライト・ビハインド方式でディスクに書き込みつつ、同時にアプリケーションの処理を即時に開始できることが望ましい。この場合、差分を反映し終え、さらにフェイルオーバー後のトランザクションによる新たなデータも反映し終えるまでは、ライト・ビハインド方式を使うことになる。さらに、これらの処理が完了した後は、必要に応じてライト・スルー方式に切り替えられることが望ましい。なお、フェイルオーバーも、一時的に応答が遅くなる応答遅延と考えることもできる。

10

【 0 0 1 9 】

さらに、負荷はある特定の種類のデータに集中する場合がある。例えば、証券取引アプリケーションでは、特定の株銘柄の売買が短期間に集中する場合がある。そのため、ライト・スルー方式とライト・ビハインド方式の切り替えは、株取引の銘柄などのデータの種類ごと（例えば、リージョンごと）に行なえるようにしておく必要がある。

【 0 0 2 0 】

図 4 に、図 3 の構成を用いて、ライト・スルー方式を実現した状態の一例を示す。ここでは、リージョン a とリージョン b (4 0 3) のデータは、ディスク用データ記録手段 4 0 5 によって書き込みが行われる。ライト・スルー方式でデータ書き込みをするため、メモリ用データ記録手段 4 0 6 は使用されない。キュー管理テーブル 4 0 4 には、ディスク用データ記録手段 4 0 5 が利用するリージョン a とリージョン b についてのキュー { D、a } およびキュー { D、b } が登録される。ここでは、リージョン a のトランザクションがコミットされる場合を例として説明する。データは関連するキュー { D、a } に格納される。ライト・スルー方式では、コミット完了前にディスクにデータが書き込まれる必要があるため、書き込みタイミング制御手段 4 0 8 は、ディスク用データ記録手段 4 0 5 にスレッドを割り当て、データを渡し、ディスク用データ記録手段 4 0 5 が書き込みを行なう。

20

【 0 0 2 1 】

ディスク用データ記録手段 4 0 5 によるデータの書き込みが完了するとコミットが完了する。リージョンの異なるトランザクションは互いに関連がないため、書き込みタイミング制御手段 4 0 8 は、データ記録手段に対してキューごとにスレッドを割り当てる。なお、図 4 ではライト・スルー方式でデータ書き込みを行うので、キュー管理テーブル 4 0 4 を経由しないでディスク用データ記録手段 4 0 5 がデータを受取ることも可能である。しかし、ライト・スルー方式でのデータ書き込みの際にも、データがキュー管理テーブル 4 0 4 を経由するようにすれば、書き込みタイミング制御手段 4 0 8 や、ディスク用データ記録手段 4 0 5 の機能は、ライト・ビハインド方式の際の機能と共通化できるため、機能構成を単純なものとするのが可能となる。また、必ずキューを経由してデータ書き込み処理をすれば、ライト・ビハインド方式からライト・スルー方式に切り替えた場合にも、データ順にキューの残存データの書き込み処理をしやすくなる。

30

40

【 0 0 2 2 】

図 5 に、図 3 の構成を用いて、ライト・ビハインド方式を実現した状態の一例を示す。ここでは、リージョン a とリージョン b (5 0 3) のデータは、それぞれディスク用データ記録手段 5 0 5 によりディスク 5 0 1 に書き込まれ、メモリ用データ記録手段 5 0 6 によりメモリ 5 0 2 に書き込まれる。従って、キュー管理テーブルには、ディスク用データ記録手段 5 0 5 に渡すデータを管理するためのキュー { D、a } およびキュー { D、b } と、メモリ用データ記録手段 5 0 6 に渡すデータを管理するためのキュー { M、a } およびキュー { M、b } の 4 つが登録される。ライト・ビハインド方式でのデータ書き込みでは、コミット処理の前に冗長化用の遠隔メモリ (メモリ 5 0 2) に書き込む必要がある。

50

書き込みタイミング制御手段508は、メモリ用データ記録手段506にスレッドを割り当てる。

【0023】

アプリケーションからのトランザクション・データはキュー{M、a}および{M、b}に格納される。メモリ用データ記録手段506により、メモリ502へのデータの書き込みがなされると、アプリケーションのコミット処理が完了する。その後、書き込みタイミング制御手段508は、ディスク用データ記録手段505にスレッドを割り当て、ディスク用データ記録手段505はデータを受取って書き込み処理を行う。なお、ディスク用データ記録手段505が、ディスク501にデータの書き込みを完了すると、冗長化したデータは不要となるので、書き込みタイミング制御手段508は、メモリ502のデータの削除を指示する。ライト・ビハインド方式の場合、アプリケーション側のコミット処理は、メモリ502への書き込み処理が完了していればいいので、ディスク用データ記録手段505に対するスレッドの割り当てはコミット完了前であってもよい。また、メモリからのデータの削除処理は、個々のトランザクション毎に行なってもいいし、性能向上のために複数トランザクション分の削除をまとめてバッチ処理してもよい。

10

【0024】

図6に、リージョンaについて、ディスクへのデータ書き込み方式を、ライト・スルー方式からライト・ビハインド方式へ変更する処理のフロー600を例示する。ステップ601で処理が開始される。ステップ602で負荷の検出を行う。ステップ603で負荷が所定の閾値を超えたか否か判断する。例えば、所定の閾値は、先行トランザクションの処理のディスク書き込みが終了する前に次のトランザクションが発生する程度のトランザクション量および頻度を設定する。これは、アプリケーションのスループットやレイテンシに対する要求を満たす範囲内で負荷に対する閾値を設定することになる。ステップ603で、負荷が閾値を超えていないと判断される場合(No)は、ステップ606に進んで、ディスクへのデータ書き込み方式を変更することなく処理を終了する。

20

【0025】

一方、ステップ603で、負荷が閾値を超えていると判断された場合(Yes)は、ステップ604に進む。ステップ604では、キュー{D、a}について、リージョンaの識別子を鍵とした排他ロックを取得する。排他ロックの鍵はリージョンの識別子以外のもので可能である。この排他ロックの取得は、データ書き込みをライト・スルー方式からライト・ビハインド方式に変更している間に、アプリケーションやデータ書き込み手段がキューにアクセスしてエラーとなることを防ぐためである。ステップ605では、キュー管理テーブルにキュー{M、a}を登録する。ステップ606で排他ロックを解除する。ステップ607で、キュー{M、a}について、メモリ用データ記録手段スレッドを割り当てる。但し、ステップ607では、当該スレッドのインスタンスが残っている場合は、スレッドの割り当ては行う必要はない。この処理は不要となる。ステップ608で変更処理を終了する。ステップ606が完了した後、メモリ用データ記録手段にスレッドが割り当てられ、ライト・ビハインド方式によるデータの書き込み処理が始まる。

30

【0026】

ちなみに、ステップ602では、アプリケーションの動作中にプロセスに障害が発生し、別のプロセスに引き継ぐフェイルオーバー、データベースの一時的な応答遅延や、データベースのフェイルオーバー、また、ディスクのフェイルオーバーを検出した場合にも、負荷が閾値を超えたとみなして、ステップ603からステップ604に進んで、ディスクへのデータ書き込みをライト・スルー方式からライト・ビハインド方式に変更するようにしてもよい。

40

【0027】

ちなみに、ステップ602で負荷の検出できる例として、アプリケーションのフェイルオーバー時の差分反映が挙げられる。ライト・ビハインド方式でアプリケーションが動作中にプロセスに障害が発生し、別プロセスが引き継ぐフェイルオーバーの場合、遠隔メモリとディスク上のデータの差分をディスクに反映する必要がある。しかし、これが完了す

50

るのを待っていたのでは、アプリケーションの可用性が低下してしまう。そこで、本発明を用いて、差分を反映している間はライト・ビハインド方式でディスクに書き込みつつ、同時にアプリケーションの処理を即時に開始することができる。その後、差分を反映し終え、さらにフェイルオーバー後のトランザクションによる新たなデータも反映し終えた時点で、必要ならライト・スルー方式に切り替えてもよい。

【 0 0 2 8 】

また、ステップ 6 0 2 で負荷の検出できる例として、データベースの一時的な応答遅延への対応が挙げられる。ライト・スルー方式でデータベースに書き込み中、データベースが何らかの原因で一時的に応答が遅くなる場合がある。これに備えて、ある一定時間以上データベースが応答しない場合は、ライト・ビハインド方式に切り替えることによって、アプリケーションに対して応答時間を保証できる。その後データベースが応答したら、ライト・スルー方式に戻してもよい。動的な切り替えにより、耐えられる応答遅延時間を長くできる。

10

【 0 0 2 9 】

さらに、ステップ 6 0 2 で負荷の検出できる例として、データベースのフェイルオーバー時間の隠蔽する場合が挙げられる。ライト・スルー方式でデータベースに書き込み中、データベースに障害が発生しデータベース・プロセスがフェイルオーバーする場合がある。通常、データベースのフェイルオーバーには時間がかかり、その間アプリケーションへの応答を返せなくなるが、これを応答遅延の1つととらえ、同様にライト・ビハインド方式に切り替えれば、データベースのフェイルオーバー時間をアプリケーションに対して隠蔽することができる。これによりアプリケーションの可用性を上げられる。動的な切り替えにより、耐えられるデータベースのフェイルオーバー時間を長くできる。

20

【 0 0 3 0 】

そして、ステップ 6 0 2 で負荷の検出できる例として、ディスクのフェイルオーバーも挙げられる。ライト・スルー方式またはライト・ビハインド方式で使用中のファイルやデータベースまたはディスクそのものに障害が発生した場合、書き込み先を他のファイルやデータベースに切り替える必要がある場合もある。これは単にディスク用データ記録部の書き込み先を切り替えるだけで実現できる。

【 0 0 3 1 】

また、ステップ 6 0 2 の負荷の検出は、リージョン毎に行うことも可能である。例えば、リージョン b について負荷が閾値を超えたと判断した場合は、ステップ 6 0 4 で、リージョン b の識別子を鍵とした排他ロックを取得し、ステップ 6 0 5 でキュー管理テーブルにキュー { M、b } を登録する。ステップ

30

【 0 0 3 2 】

図 7 に、リージョン a について、ディスクへのデータ書き込み方式を、ライト・ビハインド方式からライト・スルー方式へ変更する処理のフロー 7 0 0 を例示する。ステップ 7 0 1 で処理が開始される。ステップ 7 0 2 で負荷の検出を行う。ステップ 7 0 3 で負荷が所定の閾値を下回るか否か判断する。ステップ 7 0 3 で、負荷が所定の閾値を下回ると判断される場合 (N o) は、ステップ 7 0 7 に進んで、ディスクへのデータ書き込み方式を変更することなく処理を終了する。一方、ステップ 7 0 3 で、負荷が閾値を下回ると判断された場合 (Y e s) は、ステップ 7 0 4 に進む。ステップ 7 0 4 では、キュー { D、a } およびキュー { M、a } につき、リージョン a の識別子を鍵とした排他ロックを取得する。ステップ 7 0 5 では、キュー管理表からキュー { M、a } を削除する。ステップ 7 0 6 で排他ロックを解除する。ステップ 7 0 7 で変更処理を終了する。

40

【 0 0 3 3 】

また、ステップ 7 0 3 では、アプリケーションの動作中のフェイルオーバーの完了、データベースの一時的な応答遅延の回復や、データベースのフェイルオーバーの完了、またはディスクのフェイルオーバーの完了を検出した場合にも、ステップ 7 0 4 に進んで、ディスクへのデータ書き込みをライト・ビハインド方式からライト・スルー方式に変更するようにしてもよい。なお、ステップ 7 0 5 では、キュー { M、a } のインスタンスそのも

50

のを削除してしまうのではなく、図8のように単にキュー管理テーブルから登録を抹消しておく。キュー{M、a}のインスタンスが残っていれば、キュー管理テーブルから抹消された後も、メモリ用データ記録手段に割り当てられたスレッドは、キュー{M、a}から残存するライト・ビハインド方式のデータを取得し書き込み処理を継続できる。但し、キュー{M、a}のデータ書き込みが終われば、リソースに応じて、インスタンスを削除することが望ましい。

【0034】

図9は、図6と図7のそれぞれの処理を、同一フローで行う場合の一例を示したフロー900である。ステップ902の負荷の検出はステップ602およびステップ702と共通化できるので、二つの処理フローは図9のように融合することができる。ステップ903とステップ913は負荷が閾値よりも上か下かを判断する処理である。ステップ903では、負荷が閾値より下から上に変化があったか否かを判断する。ステップ913では、負荷が閾値より上から下に変化があったか否かを判断する。ステップ904からステップ906は、図6のステップ604からステップ607と同じである。ステップ914からステップ916は、図7のステップ704からステップ706と同じである。ステップ908では処理は終了となるが、負荷の変化に対応するため、定期的に処理は再開され、ステップ902に戻って再び負荷の検出をし、ライト・スルー方式とライト・ビハインド方式の切り替え処理を始める。

【0035】

本発明は、様々な形態で実施可能である。例えば図10に示すように、データベースとマルチキャスト通信を用いたライト・ビハインド方式に適用できる。ディスクをデータベース1001として、ディスク用データ記録手段1005がデータを書き込み、また、マルチキャスト通信1010を用いて、ディスク用データ記録手段1005が複数の遠隔メモリ1002にコピーすることで冗長度を上げられる。メモリ用データ記録手段1006がデータ書き込み完了を通知するタイミングを変えることで遠隔メモリの障害に対する耐久度を変えることもできる。例えば、コピー先の遠隔メモリが3つある場合、そのうち所定の数(例えば、最初の2つ)の遠隔メモリからのコピー完了応答が届いた時点で書き込み完了を通知すれば、最大1つの遠隔メモリの障害に耐えることができる。

【0036】

また、本発明はリージョンが1つしかない場合にも適用可能である。リージョンがシステムに1つしかない場合も許される。この場合はシステム全体でライト・スルー方式またはライト・ビハインド方式の耐久性保証方法を統一的に切り替えて利用することになる。

【0037】

本手法により、データの耐久性保証方法であるライト・スルー方式とライト・ビハインド方式を動的な負荷の変動に応じて切り替えることが可能になる。これにより、ディスク・ボトルネックになるような高負荷時に閾値の範囲でメモリ使用量を節約し、閾値を超えるような高負荷に備えることが可能になる。また、ライト・スルー方式とライト・ビハインド方式の両方における対故障性を向上でき、フェイルオーバーにかかる時間を小さくできて、データの耐久性保証方法の切り替えの範囲を柔軟に定義できる。

【0038】

以上、本発明を実施の形態を用いて説明したが、本発明の技術的範囲は上記実施の形態に記載の範囲には限定されない。上記実施の形態に、多様な変更または改良を加えることが可能であることが当業者に明らかである。その様な変更または改良を加えた形態も本発明の技術的範囲に含まれ得ることが、特許請求の範囲の記載から明らかである。

【図面の簡単な説明】

【0039】

【図1】ライト・ビハインド方式によるデータの書き込み処理を簡単に示す。

【図2】ディスクへのトランザクション・データの書き込みを、ライト・スルー方式とライト・ビハインド方式で動的に切り替えるシステムが作動するハードウェア構成図の一例を示す。

10

20

30

40

50

【図3】ディスクへのトランザクション・データの書き込みを、ライト・スルー方式とライト・ビハインド方式で動的に切り替える装置の機能および処理されるデータの一例を示す。

【図4】ライト・スルー方式を実現した状態の一例を示す。

【図5】ライト・ビハインド方式を実現した状態の一例を示す。

【図6】ディスクへのデータ書き込み方式を、ライト・スルー方式からライト・ビハインド方式へ変更する処理のフローを例示する。

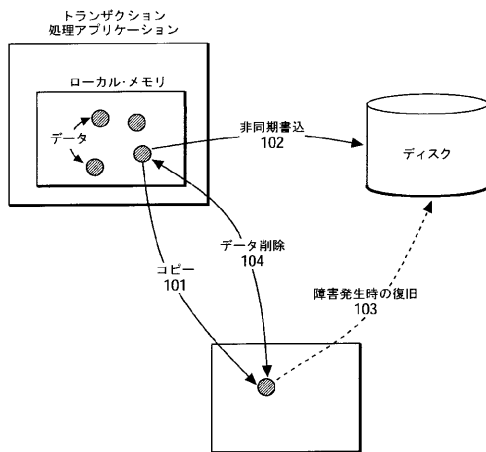
【図7】ディスクへのデータ書き込み方式を、ライト・ビハインド方式からライト・スルー方式へ変更する処理のフロー700を例示する。

【図8】単にキュー管理テーブルから登録を抹消しておく状態を示す。

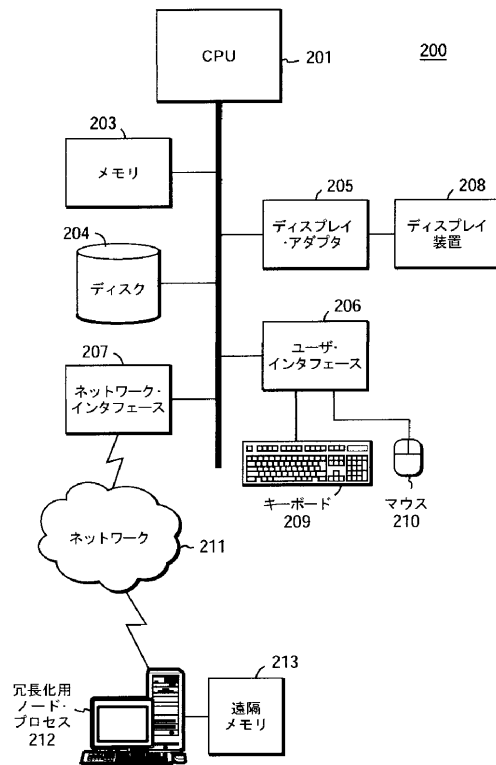
【図9】図6と図7のそれぞれの処理を、同一フローで行う場合のフローの一例を示す。

【図10】データベースとマルチキャスト通信をライト・ビハインド方式に適用した一例を示す。

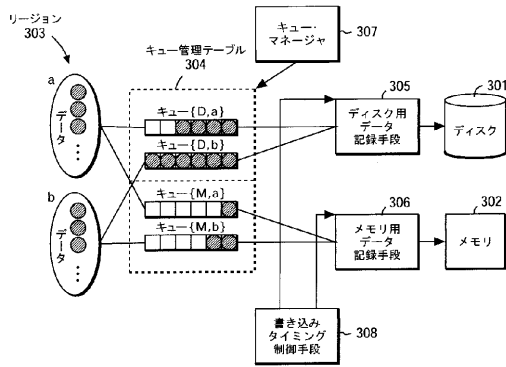
【図1】



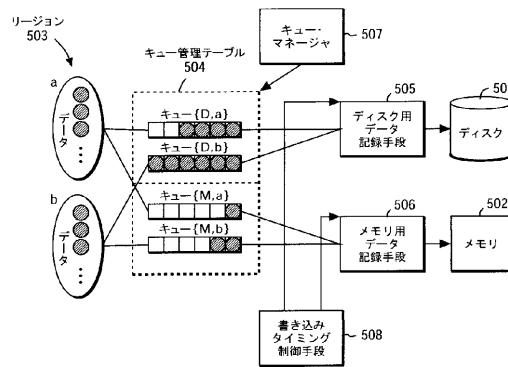
【図2】



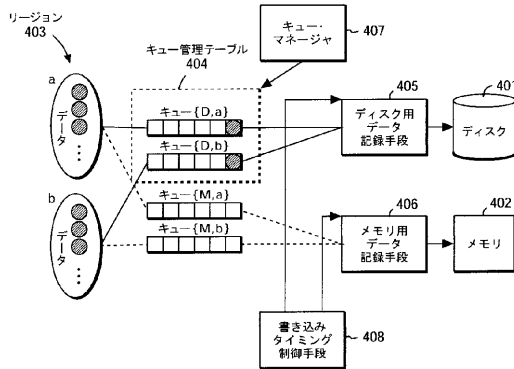
【図3】



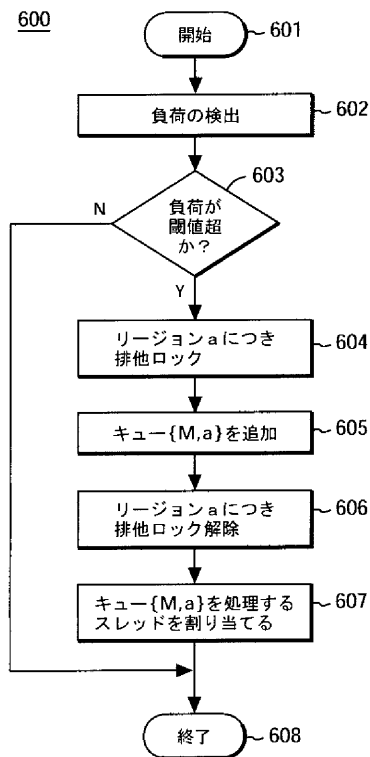
【図5】



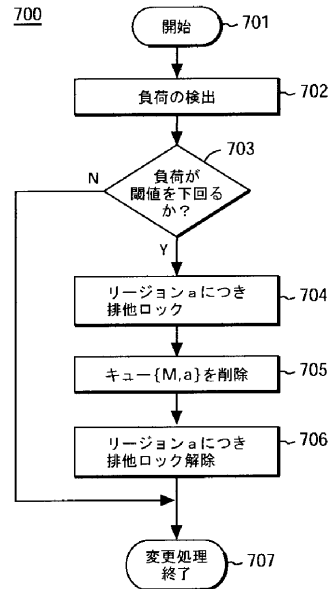
【図4】



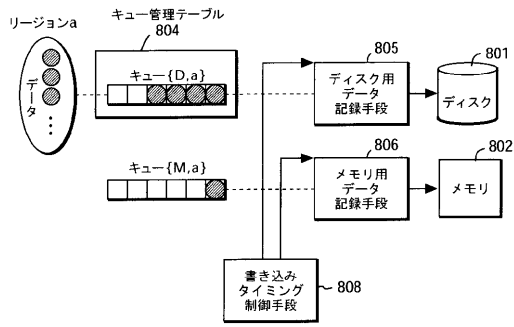
【図6】



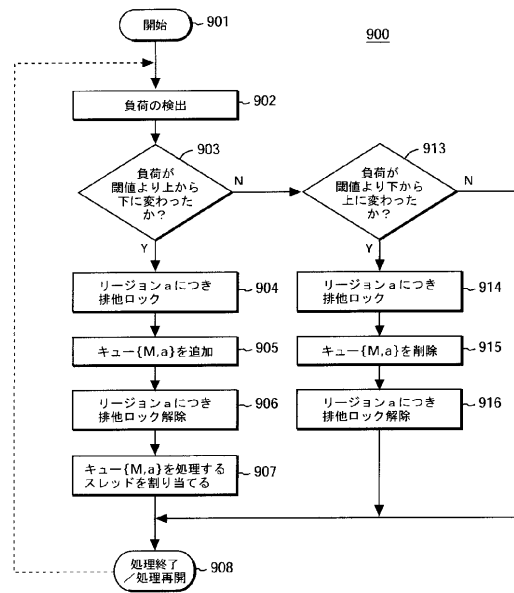
【図7】



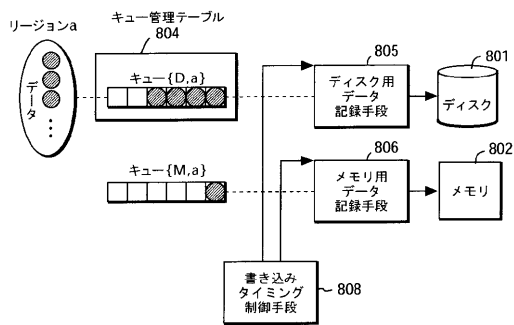
【 図 8 】



【 図 9 】



【 図 10 】



フロントページの続き

(74)代理人 100112690

弁理士 太佐 種一

(72)発明者 根山 亮

神奈川県大和市下鶴間1623番地14 日本アイ・ビー・エム株式会社 東京基礎研究所内

(72)発明者 山本 学

神奈川県大和市下鶴間1623番地14 日本アイ・ビー・エム株式会社 東京基礎研究所内

(72)発明者 小澤 陽介

神奈川県大和市下鶴間1623番地14 日本アイ・ビー・エム株式会社 東京基礎研究所内

審査官 木村 雅也

(56)参考文献 特開平07-072980(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06

G06F 12/00

G06F 13/10