

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4289056号  
(P4289056)

(45) 発行日 平成21年7月1日(2009.7.1)

(24) 登録日 平成21年4月10日(2009.4.10)

|                   |              |                  |      |       |      |
|-------------------|--------------|------------------|------|-------|------|
| (51) Int.Cl.      |              | F I              |      |       |      |
| <b>G06F 12/00</b> | <b>12/00</b> | <b>(2006.01)</b> | G06F | 12/00 | 531D |
| <b>G06F 3/06</b>  | <b>3/06</b>  | <b>(2006.01)</b> | G06F | 12/00 | 533J |
|                   |              |                  | G06F | 3/06  | 304F |

請求項の数 9 (全 21 頁)

|              |                               |           |                     |
|--------------|-------------------------------|-----------|---------------------|
| (21) 出願番号    | 特願2003-207782 (P2003-207782)  | (73) 特許権者 | 000005108           |
| (22) 出願日     | 平成15年8月19日 (2003.8.19)        |           | 株式会社日立製作所           |
| (65) 公開番号    | 特開2004-348701 (P2004-348701A) |           | 東京都千代田区丸の内一丁目6番6号   |
| (43) 公開日     | 平成16年12月9日 (2004.12.9)        | (74) 代理人  | 100100310           |
| 審査請求日        | 平成18年3月20日 (2006.3.20)        |           | 弁理士 井上 学            |
| (31) 優先権主張番号 | 特願2003-86920 (P2003-86920)    | (72) 発明者  | 須藤 敦之               |
| (32) 優先日     | 平成15年3月27日 (2003.3.27)        |           | 東京都国分寺市東恋ヶ窪一丁目280番地 |
| (33) 優先権主張国  | 日本国(JP)                       | (72) 発明者  | 馬場 恒彦               |
|              |                               |           | 東京都国分寺市東恋ヶ窪一丁目280番地 |
|              |                               |           | 株式会社日立製作所中央研究所内     |
|              |                               | 審査官       | 田川 泰宏               |

最終頁に続く

(54) 【発明の名称】 計算機システム間のデータ二重化制御方法

(57) 【特許請求の範囲】

【請求項1】

第1の計算機システムと、該第1の計算機システムに接続される第1のストレージシステムとを有する正システムと、第2の計算機システムと、該第2の計算機システムに接続された第2のストレージシステムとを有する副システムとを備え、かつ少なくとも前記第1、第2のストレージシステムの間が相互接続されているシステムにおけるデータ二重化制御方法であって、

前記第1の計算機システムの処理によって前記第1のストレージシステムが保持するデータの更新を再現可能なログを該第1のストレージシステム内の特定記憶装置に登録するステップと、

前記ログの複製のために設定された第2のストレージシステム内の特定記憶装置に前記第1のストレージシステムの特定記憶装置に登録されたログをコピーするステップと、

前記第1計算機システムの処理によって前記第1のストレージシステムが保持するデータを更新するステップと、

前記コピーステップにより前記第2のストレージシステム内の特定記憶装置の保持内容に変更が生じたことを前記第2の計算機システムが検知するステップと、

前記第2の計算機システムが前記第2のストレージシステム内の特定記憶装置の保持内容の変更を読み込むステップと、

前記第2の計算機システムがログを読み込み該ログにしたがって前記第2のストレージに保持する前記データの複製を更新するステップ、

10

20

を有することを特徴とするデータ二重化制御方法。

【請求項 2】

前記第 1 の計算機システムの処理によって前記第 1 のストレージシステムが保持するデータの更新を再現可能なログは、1 つまたは複数のトランザクションからなり該トランザクションの開始と終了を区分する情報を含むことを特徴とする請求項 1 のデータ二重化制御方法。

【請求項 3】

前記第 1 の計算機システムの処理による前記第 1 のストレージシステムが保持するデータの更新を再現可能なログを該第 1 のストレージシステム内の特定記憶装置に登録するステップは、第 1 の計算機システムが指定したログ入出力単位で行うことを特徴とする請求項 1 のデータ二重化制御方法。

10

【請求項 4】

前記第 1 の計算機システムの処理によって前記第 1 のストレージシステムが保持するデータの更新を可能なログを該第 1 ストレージシステム内の特定記憶装置に登録するステップと

前記ログの複製のために設定された第 2 のストレージシステム内の特定記憶装置に前記第 1 のストレージシステムの特定記憶装置のログの変更をコピーするステップとを同期で行うことを特徴とする請求項 1 のデータ二重化方法。

【請求項 5】

請求項 1 記載のデータ二重化制御方法において、前記第 1 の計算機システムが停止したことを検知するステップと、前記第 2 の計算機システムが前記第 1 の計算機システムから業務を引き継ぐステップを更に有することを特徴とするデータ二重化制御方法。

20

【請求項 6】

請求項 1 記載のデータ二重化制御方法において、前記第 1 の計算機システムが停止したことを検知するステップと、前記第 2 の計算機システムが前記第 1 の計算機システムから業務を引き継ぐステップと、前記第 2 の計算機システムが前記第 2 のストレージシステム内の特定記憶装置のログを読み込み該ログにしたがってデータを更新するステップをさらに有することを特徴とするデータ二重化制御方法。

【請求項 7】

データベースを構築する第 1 のデータベースサーバに接続され、前記データベースのデータを格納する記憶領域及び前記データベースの更新履歴を示す更新ログを格納する記憶領域を備える第 1 のストレージシステムと、第 2 のデータベースサーバと、に接続される第二のストレージシステムであって、

30

第 1 の記憶領域と、前記第二のデータベースサーバが用いるデータを格納する第二の記憶領域と、

前記第 1 の記憶領域及び前記第二の記憶領域に対してデータの入出力を制御するストレージ制御部と、を備え、

前記ストレージ制御部は、

前記更新ログを前記第 1 のデータベースサーバから前記第一のストレージシステムを介して、受領する受領部と、

40

前記更新ログを、前記第 1 の記憶領域に格納し、

前記第 1 の記憶領域に格納されたことを前記第二のデータベースサーバに通知し、

前記第二のデータベースサーバから、前記第 1 の記憶領域に格納される更新ログを前記第 2 の記憶領域に反映する反映要求を受け、

前記第 2 の記憶領域に前記更新ログを反映し、

前記反映が完了した場合、完了したことを示す完了報告を前記第 2 のデータベースサーバに通知する、ことを特徴とする第 2 のストレージシステム。

【請求項 8】

データベースを構築する第 1 のデータベースサーバに接続され、前記データベースのデータを格納する記憶領域及び前記データベースの更新履歴を示す更新ログを格納する記憶

50

領域を備える第1のストレージシステムと通信可能な第1の記憶領域及び第2の記憶領域を備える第2のストレージシステムに接続される第2のデータベースサーバであって、前記第1のストレージシステムから送信される前記1の記憶領域に新たな更新ログが格納されたか否かを検出する検出手段と、

前記検出手段により検出した場合は、前記第1の記憶領域から更新ログを読み出す読み出し手段と、

前記読み出された更新ログを前記第2の記憶領域に更新する更新要求を前記第2のストレージシステムに送信する更新要求送信手段と、

前記第2のストレージシステムから前記更新要求に対応する更新処理の完了を示す前記更新完了を受信する手段と、を有することを特徴とする第2のデータベースサーバ。

10

【請求項9】

請求項8記載の第2のデータベースサーバであって、

第一のデータベースサーバの障害を検出した場合、

前記更新ログの前記第1の記憶領域への格納を停止するよう前記第2のストレージシステムに指示する停止指示手段と、

前記第1の記憶領域から更新ログを読み出し、前記第2の記憶領域に未反映の更新ログの有無を確認する確認手段と、

前記未反映の更新ログがあった場合は、前記第2のストレージシステムに前記第2の記憶領域に反映するよう指示する更新指示手段と、

前記更新指示手段に対し、前記第2のストレージシステムから反映完了報告を受信した場合、前記第2のデータベースサーバで、前記第1のデータベースサーバの処理を引き継ぐことを特徴とする第2のデータベースサーバ。

20

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は計算機およびストレージ装置からなる業務システムに関し、特に複数のシステム間でデータを複製するデータ二重化制御方法、およびデータを複製したシステムへの高速な切り替えを実現する方法に関する。

【0002】

【従来の技術】

データベースサーバおよびストレージ装置からなる業務システムが複数ある場合のデータ複製方法として、データベースサーバ上で動作するDBMSが実行する方法がある。DBMSがデータ複製する方法については、例えば非特許文献1に記述がある。複数のシステムのデータベースサーバ同士を接続し、一方のシステム上で動作するDBMSの更新情報を別のシステムに転送することでデータ複製する技術である。

30

【0003】

また、同様なシステムのデータ複製方法として、ストレージ装置間のデータコピー機能を使用する方法がある。ストレージ装置間のデータコピー機能については、例えば非特許文献2に記述がある。複数システムのストレージ同士をファイバーチャネルで接続し、一方のストレージ装置のディスクドライブに更新があると、別のストレージ装置のディスクドライブにもデータの更新を反映する技術である。

40

【非特許文献1】

Oracle9i製品カタログ (<http://www.oracle.co.jp/products/catalog/pdf/9iDBr2J07266-01.pdf>)、第6頁。

【非特許文献2】

日立統合ストレージソリューション「Storeplaza」カタログ (<http://www.hitachi.co.jp/Prod/comp/storeplaza/data/stpzclg.pdf>)、第5頁。

【0004】

【発明が解決しようとする課題】

従来のデータ複製方法を実行する場合、通常のデータベース業務を行う以上のコストが必

50

要である。また、複数システム間で同期したデータ複製を行うと業務処理の遅延が発生する。

【 0 0 0 5 】

D B M Sによるデータ複製を行うためには、D B M Sの動作するデータベースサーバが業務処理を行う負荷に加えて、データ複製処理を行う負荷が加わるためにより高性能なデータベースサーバが必要でありコストが増加するという課題がある。また、複製したデータが一致するためには、D B M Sの更新処理を実行するたびにデータベースサーバ間で同期通信を行う必要がある。同期通信中は、D B M Sが次の更新処理を実行できないため、業務が遅延することが課題である。

【 0 0 0 6 】

ストレージ装置でデータ複製を行うためには、D B M Sが扱うデータの更新を全てコピーするため、ストレージ装置間の接続に広帯域の回線を使う必要がある。広帯域の回線を使用することでコストが増大する課題がある。また、複製したデータが一致するためには、ディスクドライブ上のデータが更新されるたびにストレージ装置間で同期通信を行う必要がある。同期通信中はディスクドライブへの次の更新処理が実行できないため、業務が遅延することが課題である。

【 0 0 0 7 】

また、同期通信による遅延を防ぐため、D B M Sやストレージ装置間の通信を非同期で実行する方法が存在するが、障害や災害でデータ複製先のシステムに切り替える場合に、未転送分のデータを複製先で再構築する必要が生じ、システムの切り替えが遅延することが課題である。

【 0 0 0 8 】

【課題を解決するための手段】

サーバが外部から受け付けた要求に応じて業務を実行すると、ストレージ装置に保存されたデータへの更新や追加が必要となる。このストレージ装置のデータ複製を行うために全てのデータを複製するのではなく、複製先としてサーバとストレージ装置を用意し、複製元のサーバで実行された業務を復元可能なログをストレージ装置の特定のディスクドライブに保存し、このディスクドライブが更新されるたびに複製先のストレージ装置にディスクドライブのコピーを行う。複製先のストレージ装置へのディスクドライブのコピーが完了したら、コピーされたストレージ装置からログを保存したディスクドライブが更新されたことを複製先のサーバに通知する。複製先のサーバは、ストレージ装置からログを保存したディスクドライブの変更通知を受信できるようにしておき、通知を受けた後でログをディスクドライブから読み取り、複製元のサーバで行われたのと同じ業務処理を実行する。このログを基にした業務処理の実行後、その結果をストレージ装置に反映することでデータの複製が完了する。

【 0 0 0 9 】

このデータ複製方法を実行しているシステムにおいて、複製元のサーバとストレージ装置が障害や保守操作により停止した場合、複製先のストレージ装置に保存された業務データが最新の状態にあるため、複製元のサーバが受信していた業務を複製先で受信するように変更することで、業務の処理を中止せずにサーバとストレージ装置の切り替えを実行する。

【 0 0 1 0 】

サーバとストレージ装置の切り替え実行後に、複製元と複製先双方のサーバとストレージ装置とがデータ複製のために実行していた処理を交替することで、業務を受信し処理しているシステムが停止した場合、再びサーバとストレージ装置の切り替えを実行する。

【 0 0 1 1 】

【発明の実施の形態】

以下で説明する実施形態では、例として計算機上で動作する業務としてデータベースサーバを取り上げるが、計算機上で実行される業務はデータベースに限定するものではない。計算機上で動作する業務は、複製するデータが、正システムで行った業務によって更新さ

10

20

30

40

50

れるもので、かつ、そのデータの更新を副システムで再現できるログを生成するものであればよい。例えば、ファイルシステムなどでも実施可能である。

〔第1実施形態〕

図1は、本発明が適用されたデータベースサーバとストレージ装置を用いたデータ複製システムの一実施例である。

【0012】

正システムを構成するのは、データベースサーバ2とストレージ装置8である。これらはデータベースサーバ2に内蔵されたストレージ接続装置3とストレージ装置8のディスク制御装置5とがサーバ・ストレージ間接続インタフェース4によって接続される。ストレージ装置8はディスク制御装置5によって読み込み書き込みを行うデータを保存するディスクドライブ6,7を内蔵しており、データベースサーバ2が業務ネットワーク1を通じて業務要求を受け取って処理したデータや、その処理に必要なデータおよびデータベースサーバ2内部で実行された業務データを保持する。

10

【0013】

データベースサーバ2とストレージ装置8とは、サーバ・ストレージ間接続インタフェース4を通してデータの読み込み・書き込みを行うだけでなく、データベースサーバ2が要求したディスクドライブ6,7の変更があった場合、ストレージ装置8からデータベースサーバ2に通知を行う方法を有している。

【0014】

副システムを構成するのは、データベースサーバ12とストレージ装置18である。これらはデータベースサーバ12に内蔵されたストレージ接続装置13とストレージ装置18のディスク制御装置15とがサーバ・ストレージ間接続インタフェース14によって接続される。ストレージ装置18はディスク制御装置15によって読み込み書き込みを行うデータを保存するディスクドライブ16,17を内蔵しており、データベースサーバ12が業務ネットワーク1を通じて業務要求を受け取って処理したデータや、その処理に必要なデータおよびデータベースサーバ12内部で実行された業務データを保持する。

20

18とは、サーバ・ストレージ間接続インタフェース14を通してデータの読み込み・書き込みを行うだけでなく、データベースサーバ12が要求したディスクドライブ16,17の変更があった場合、ストレージ装置18からデータベースサーバ12に通知を行う方法を有している。

30

【0015】

ディスク制御装置5とディスク制御装置15とはストレージ装置間接続インタフェース20により接続される。これにより、正システムのストレージ装置8と副システムのストレージ装置18は互いに接続される。ストレージ装置8とストレージ装置18は、一方のディスクドライブの一つを複製元に、他方のディスクドライブの一つを複製先にあらかじめ設定しておくことで、ストレージ装置間接続インタフェース20を通して内容を複製する方法を有している。

【0016】

以下、本実施形態のデータ複製方法およびシステム切り替え方法の動作を説明する。本実施形態では業務を通常実行している正システムと、正システムが何らかの理由で稼働不可能になった時に業務を引き継ぐ副システムとの間でデータ複製を行うものとする。

40

【0017】

まず、データ複製方法を実現するための初期設定を正システム、副システム双方について行う。

【0018】

正システムの初期設定は、業務システムに応じたデータベースを構築することから始める。ストレージ装置8のディスク制御装置5で、データベースサーバ2が使用可能なディスクドライブ6,7を割り当てる。データベースサーバ2はデータベースのデータを保持するディスクドライブ6とデータベースのログを保持するディスクドライブ7とを設定する。ここで言うログとは、データベースの更新作業を逐一表すもので、ログを再実行するこ

50

とでデータベースの再構築が可能なものである。例えば、データベースが実行したトランザクションログやデータベースサーバが受け取った業務要求全てのSQLコマンドである。トランザクションは、データベースの処理を複数個まとめた処理単位で、その処理が全て成功するか失敗するかのいずれかになる。そのため、多数の業務要求を処理しなければならない業務システムは、トランザクション単位で処理を行うことでデータベースに不整合を発生させないために使用する。

【0019】

副システムにも正システムと同様のデータベースを構築する。ストレージ装置18において、ストレージ装置8でデータベースサーバ2が使用可能としたディスクドライブ6と同様なディスクドライブ16と、ディスクドライブ7と同様なディスクドライブ17をディスク制御装置15でデータベースサーバ12が使用可能となるように割り当てる。データベースサーバ12は、データベースサーバ2同様に、データベースのデータを保持するディスクドライブ16とデータベースのログを保持するディスクドライブ17とを設定する。

10

【0020】

次に、正システムのストレージ装置8と副システムのストレージ装置18との間で、ストレージ装置間接続インタフェース20を通じてデータベースのログを保持するディスクドライブ7をディスクドライブ17にコピーするように設定する。このディスクドライブコピーは、同期コピー、非同期コピーいずれとも可能である。例えば、同期コピーとは、データベースサーバ2からのディスクドライブ7への情報書き込み要求に対して、ディスクドライブ7への情報書き込みとディスクドライブ17への情報書き込みとの両方が終了してから、データベースサーバ2に書き込み完了報告を行うものを指す。一方、非同期コピーは同期コピーと異なり、ディスクドライブ7への情報書き込みが完了した時点でデータベースサーバ2に書き込み完了報告を行う。非同期コピーを行う場合、正システムと副システムのログディスクが常に一致するとは限らず、システム切り替え時にデータが欠損することがある。

20

【0021】

そして、副システムのデータベースサーバ12からストレージ装置18のディスク制御装置15に対して、ディスクドライブ17の更新が行われたらデータベースサーバ12に通知を行うように設定する。

30

【0022】

正システムに障害・災害などが発生した場合に、副システムに切り替えるため正システムの停止を迅速に検知する必要がある。そのため、正システムのデータベースサーバ2と副システムのデータベースサーバ12の間で正システムが稼働していることを通知するための通信設定を行う。例えば、正システムのデータベースサーバ2から副システムのデータベースサーバ12に業務ネットワーク1を経由して一定時間間隔で通知を行う方法がある。また、正システムの稼働状態を監視する外部のサーバから副システムへの切り替えを指示する方法や、副システムから一定時間間隔で正システムに稼働状態を問い合わせる方法もある。

40

【0023】

以上のような設定がデータベースサーバ2, 12とストレージ装置8, 18で完了した後、正システムのデータベースサーバ2で業務処理を開始する。以下では、データ複製の手順について説明する。

【0024】

データ複製第1ステップ101：業務処理要求は、業務ネットワーク1を通じてデータベースサーバ2に到着する。業務処理要求は業務ネットワーク上のプロトコルに応じて送付され、データベースサーバ2の管理するデータ内容を参照するものや更新するものからなる。例えば、TCP/IPプロトコルによって送付される、SQLコマンドの組み合わせからなる。

【0025】

50

業務処理要求を受信したデータベースサーバ2は、ネットワークプロトコル層の解析を行い、データベースへの業務処理内容を取り出し、業務処理内容の解析を行った後、業務処理を実行する。例えば、TCP/IPプロトコルの解析を行い、SQLコマンドを取り出し、その処理をデータベースで実行する処理がある。

【0026】

データ複製第2ステップ102：業務処理の内容が、データベースの更新処理を伴う場合にはストレージ装置内に保持しているデータを更新する必要がある。その場合には、ストレージ装置接続装置3からサーバ・ストレージ接続インタフェース4を通じて、ディスク制御装置5に対してディスクドライブ7への更新ログの書き込みをストレージ装置8に指示する。例えば、データベースサーバ2にホストバスアダプタを装着し、ファイバーチャネルケーブルを通じてSCSIコマンドをディスクコントローラに送信することに当たる。また、この更新ログの書き込みはデータ更新の度に常に実行するだけでなく、トランザクション単位で実行が完了したときに行う方法、データベースサーバ2で用意したログ用バッファの空きが無くなった場合や一定時間が経過した後などのようにデータベースサーバ2が指定する入出力の単位で条件に応じて実行する方法なども用いることができる。また、本実施形態では簡単のため1回の書き込み要求のように図示したが、通常は別のディスクドライブへの書き込み要求は複数の要求に分けて送信される。

10

【0027】

データ複製第3ステップ103：更新ログの書き込み要求を受けたディスク制御装置5は、ディスクドライブ7へと情報の書き込みを行う。

20

【0028】

データ複製第4ステップ104：ディスクドライブ7はその情報の書き込みが終了すると副システムのストレージ装置18内のディスクドライブ17にコピーするように設定されているため、ディスク制御装置5はストレージ装置間接続インタフェース20を通じて副システムのストレージ装置18にあるディスク制御装置15にディスクドライブ7の更新内容を送信し、ディスクドライブ17に書き込むよう指示し、ディスク制御装置15はディスクドライブ17へと書き込みを行う。この書き込みが完了したら、ディスク制御装置5からデータベースサーバ2に更新ログ書き込み完了報告を行う。例えば、ストレージ装置間接続インタフェース20としてファイバーチャネルケーブルを用い、ストレージ装置の管理ソフトウェアでディスクドライブ7、17のコピーを設定することで実現できる。また、本実施形態では、ディスクドライブ7の更新直後にディスクドライブ17へのコピーを行う同期コピー方法としているが、一定時間間隔でコピーを実行する非同期コピー方法を用いることも可能である。ただし、非同期コピー方法を用いた場合、ディスクドライブ17への書き込みが完了しなくてもストレージ装置8内のディスクドライブ7への書き込みが完了したらデータベースサーバ2に完了報告を行うため、システム切り替え時にディスクドライブのデータがコピーされていない事態も発生しうる。

30

【0029】

データ複製第5ステップ105：更新ログの書き込み完了報告を受けたデータベースサーバ2は、ディスク制御装置5にディスクドライブ6へ更新データの書き込みを指示し、ディスク制御装置5はディスクドライブ6への情報の書き込み処理を行う。

40

【0030】

データ複製第6ステップ106：ディスクドライブ17への更新を実行後、ディスク制御装置15はあらかじめデータベースサーバ12から更新を通知するように指定されているため、更新が発生したことをデータベースサーバ12に通知する。この更新通知要求と更新通知のインタフェースは、例えば、データベースサーバ12からストレージ装置18内の特殊なディスクドライブへの読み込み要求の応答としてディスク制御装置15が通知する方法や、データベースサーバ12から更新通知を要求するディスクドライブ17への専用コマンドに対する応答としてディスク制御装置15が通知する方法、また、ディスク制御装置15からディスクの更新を通知する専用の割り込みインタフェースをデータベースサーバ12内のストレージ接続装置13に設ける方法などがある。また、ディスク制御装

50

置 1 5 からデータベースサーバ 1 2 への通知は、更新が発生する度に実行する方法に限定するわけではなく、一定時間間隔ごとに通知する方法や、データベースサーバ 2 から指示をストレージ装置 8 に発行したものをストレージ装置 1 8 に伝えてデータベースサーバ 1 2 への通知を実行させる方法などがある。

【 0 0 3 1 】

データ複製第 7 ステップ 1 0 7 : ディスクドライブ 1 7 の更新通知を受けたデータベースサーバ 1 2 は、ディスクドライブ 1 7 の更新分の情報を読み込み、その更新ログにしたがって、ディスクドライブ 1 6 上のデータを更新するようにストレージ接続装置 1 3 からサーバ・ストレージ間接続インタフェース 1 4 を通じてディスク制御装置 1 5 に情報書き込み要求を通知する。例えば、データベースサーバ 1 2 にホストバスアダプタを装着し、ファイバーチャネルケーブルを通じて S C S I コマンドをディスクコントローラに送信する方法がある。

10

【 0 0 3 2 】

正システムのデータベースサーバ 2 が業務処理要求を実行するたびに、このように正システムのデータを副システムに複製することで、正システムのデータベースサーバ 1 2 にデータ複製のための負荷をかけることなく、また業務ネットワーク 1 上にデータ複製のためのデータ送信を行うことなく、ストレージ装置 8 , 1 8 間でのデータ転送量を小さくして、データ複製のコストを抑え、業務の遅延を小さくすることができる。

【 0 0 3 3 】

正システムが災害や機器の障害などにより停止した場合、副システムに業務処理を切り替える。正システムの保守作業を行うため必要がある場合でも、正システムを停止させ、副システムに業務処理を切り替えることがある。図 2 は正システム停止後に副システムが業務処理を引き継ぐ処理を行う手順を示した。装置構成は図 1 と同様であるため詳細は省略する。以下では、業務処理引き継ぎの手順を説明する。

20

【 0 0 3 4 】

正システムが停止すると、システムを切り替えて副システムのデータベースサーバ 1 2 とストレージ装置 1 8 で業務処理を引き継ぐ。正システムの停止は、例えばデータベースサーバ 2 とデータベースサーバ 1 2 との間で一定時間間隔で通信を行うハートビート通信や、データベースサーバ 2 , 1 2 以外の監視サーバを業務ネットワーク 1 に接続してハートビート通信を行う方法で検出可能である。業務処理の引き継ぎは、データベースサーバ 2 で受け取っていた業務処理要求をデータベースサーバ 1 2 が受け取れるように設定を変更することで可能である。例えば、データベースサーバ 2 が業務要求受信に用いていたネットワークアドレスを引き継ぐ方法がある。

30

【 0 0 3 5 】

システム切り替え第 1 ステップ 2 0 1 : まず、データベースサーバ 1 2 からディスクドライブ 7 からディスクドライブ 1 7 への情報の書き込みを停止するようにストレージ装置 1 8 に指示を出してから、データベースサーバ 1 2 がログの書き込まれたディスクドライブ 1 7 を参照し、未実行の業務がディスクドライブ 1 7 に存在するか確認する。ディスクドライブ 7 からディスクドライブ 1 7 への書き込みを停止する理由は、データベースサーバ 2 が正常に稼働しているか不明でありディスクドライブ 7 への書き込みが正常に行われるかも明らかでないため、業務を引き継ぐ副システムにそのような書き込みを反映させないためである。

40

【 0 0 3 6 】

システム切り替え第 2 ステップ 2 0 2 : 未実行の業務処理があればその業務を実行してディスクドライブ 1 6 のデータ更新を行うようストレージ接続装置 1 3 からサーバ・ストレージ間接続インタフェース 1 4 を通じてディスク制御装置 1 5 に情報の書き込み要求を通知する。ただし、正システムのデータベースサーバ 2 からディスクドライブ 7 へのログ書き込みが、トランザクション単位で行われなかった場合、ディスクドライブ 1 7 上のログがトランザクションの途中で途切れていることも起こりうる。そのような場合は、データベースサーバ 1 2 が途切れたトランザクションによるデータ更新を取り消すためのデータ

50



更新が必要になる。

【0037】

システム切り替え第3ステップ203：ディスク制御装置15は、要求を受けた情報の書き込みをディスクドライブ16に行う。

【0038】

システム切り替え第4ステップ204：ディスクドライブ16の情報書き込みが完了したら、データベースサーバ12で業務要求を受け付けて業務処理を開始する。

【0039】

さらに、正システムが障害・災害から回復し再び動作するようになった場合や保守作業完了で正システムが動作可能になった場合、本実施形態で説明してきたデータ複製方法を、副システムから正システムに複製する方向に適用することで、正システムが停止中に副システムで実行した業務処理によるデータやログの更新を正システムに反映させることができる。

10

【0040】

例えば、ストレージ装置8, 18の間でディスクドライブ17の更新部分をディスクドライブ7にコピーする設定を行い、データベースサーバ2でディスクドライブ7の更新ログにしたがってディスクドライブ6のデータを更新することで、副システムのデータ複製を実行可能である。このように、正システムと副システムが同時に停止することがなければ、交互に本発明のデータ複製方法を適用することで業務停止時間を小さくできる。

【0041】

また、本実施形態では正システムと副システムが一對一の形態を説明したが、正システムから複数の副システムへのデータ複製を行う方法や、正システムから副システムへ複製したデータをさらに別の副システムへデータ複製を行なう方法も容易に構築可能である。

20

【0042】

図3から図6に、本実施形態の主な構成要素であるデータベースサーバおよびストレージ装置の処理手順をフローチャートで示した。以下で各図のフローチャートについて説明する。

【0043】

図3に示した正システムデータベースサーバ処理手順のフローチャートについて説明する。

30

【0044】

まず、データベースサーバの初期設定を行なう(301)。例えば、初期設定には、データベースの構築やディスクドライブの割り当てなどがある。

【0045】

次に、データ複製システムを構築するまで、ストレージ装置と副システムの初期設定完了を待つ(302)。例えば、ストレージ装置間のディスクコピーの設定や副システムのデータベース構築の完了を待つことになる。

【0046】

データ複製システムの初期設定が一通り完了すると、業務処理要求受付を開始する(303)。例えば、インターネット経由で行う商取引の商品管理などが業務処理にあたる。

40

【0047】

業務処理要求の受付開始後が、業務処理要求が到着したか(304)判定する。業務処理要求が到着していれば、業務処理を実行する(305)。業務処理要求が到着していなければ業務処理要求を待ち、業務処理要求の到着を判定を続ける。業務処理要求が到着し、業務処理がデータ更新を伴う場合、業務処理のログをストレージ装置に書き込む(306)。そして、ログの書き込み完了報告を受信する(307)。その後でデータの更新要求をストレージ装置に送る(308)。

【0048】

ストレージ装置に対して行なったデータの更新書き込み要求について、ストレージ装置からの書き込み完了報告を受信する(310)ことで業務処理要求が完了する。

50

## 【 0 0 4 9 】

正システムは一度稼働すると、ここで説明したように、業務処理要求の実行とそれに伴うログの更新およびデータの更新を繰り返し行なう。

## 【 0 0 5 0 】

図 4 に示した副システムデータベースサーバ処理手順のフローチャートについて説明する。

## 【 0 0 5 1 】

正システム停止後に交替して業務処理を実行するのがデータ複製の目的であるため、副システムのデータベースサーバには正システムの設定にあわせた初期設定を行う(401)。例えば、コピーするログディスクやデータディスクの用意などがデータ複製のためには必要となる。

10

## 【 0 0 5 2 】

次に、正システムからコピーされたログを参照するため、ストレージ装置のログディスク更新を検出する設定をする(402)。例えば、副システムのデータベースサーバから副システムのストレージ装置にログディスクの更新を通知するように指定する方法や、副システムのデータベースサーバから定期的にストレージ装置内のログディスクを読み込んで更新が行われたかを判定する方法などが考えられる。これらにより、ログの更新を検出できる。

## 【 0 0 5 3 】

そして、ログが更新されるとその内容を副システムで実行してデータベースのデータも更新する、データ複製処理を開始する(403)。

20

## 【 0 0 5 4 】

システム切り替えが必要かを判定するため、正システムが正常に稼働しているか(404)判定する。例えば、正システムから副システムに対して10秒間稼働状態の通知がない場合にシステムを切り替えるとか、正システムと副システムのデータベースサーバ以外の稼働状態監視サーバを業務ネットワークに接続して状態監視をさせることによって正システムが稼働していないと判断した場合は業務処理を副システムが引き継ぐというように方針決めておき、判定を実行することになる。なお、このように正システムの稼働状態を判定する方法は複数考えられ、ここに示した方法に限定されるものではない。

## 【 0 0 5 5 】

もし正システムが正常に稼働していないと判定した場合には、システム切り替えの処理を行なう。まず、正システムの業務引継処理を実行する(410)ことで、業務ネットワークと接続可能としてから、ログディスクの更新分で未実行の業務処理を実行する(411)。そして実行の結果、データの更新をストレージ装置に送る(412)。ストレージ装置のデータ更新完了報告を受信する(413)ことで、正システムのデータ複製が完了したとみなす。そして、ログディスクの未実行業務処理のデータがすべてストレージ装置に反映されたら、業務処理要求受付を開始する(414)。

30

## 【 0 0 5 6 】

もし稼働状態が正常であると判定した場合は、ストレージ装置からログディスクの更新を検出したか(405)判定する。ログディスクの更新を検出する方法は、ストレージ装置からデータベースサーバに割り込みを発生させる方法や、データベースサーバからストレージ装置へ発行する特殊なI/Oコマンドへの応答を返す方法、あるいはデータベースサーバが一定時間間隔でログディスクの内容を読み込み、その情報を解析する方法などがある。ログディスクの更新を検出したら、ログディスクの更新分を読み込む(406)。そして、そのログ更新分を適用してデータ更新を実行する(407)。さらに、データの更新実行によって発生するデータの更新をストレージ装置に送る(408)ことで、データが正システムの最新のものと一致するようにする。ストレージ装置のデータ更新完了報告を受信する(409)と再び状態通知の受信やストレージ装置の更新通知待ちの処理を繰り返してデータ複製をしながら、システム切り替えの準備をする。

40

## 【 0 0 5 7 】

50

図5に示した正システムストレージ装置処理手順のフローチャートについて説明する。

【0058】

まず、ストレージ装置内のディスクドライブをデータベースサーバに割り当てるなどの、初期設定を行う(501)。

【0059】

そして、本実施形態のデータ複製方法を行なうため、正システムのログディスクを副システムのログディスクに対応付けしコピーの設定をする(502)。この設定を行なう前に、副システムのデータベースサーバとストレージ装置の初期設定を完了しておく必要がある。

【0060】

設定が完了したら、読み込み・書き込み処理を開始(503)し、データベースサーバからのデータ更新要求などを受け付ける状態になる。

【0061】

処理要求受信(504)を待つ状態から要求を受信すると、まず書き込み要求か(505)判定する。書き込み要求でなければ、読み込み要求された情報をデータベースサーバに転送(511)し、データベースサーバに情報読み込み完了報告を送信する(512)。実際はディスクドライブのコントロールなどの要求も受信するが、ここでは読み込み要求と同じものとみなしている。書き込み要求を受信した場合には、要求された情報をディスクに書き込む(506)処理を行い、そのディスクがコピーを設定したディスクか(507)判定する。コピー設定されていなければ、データベースサーバにデータ書き込み完了報告を送信する(510)。コピーを設定したディスクであれば、副システムのストレージ装置に書き込み要求と情報を転送(508)し、副システムのストレージ装置から書き込み完了報告を待つ(509)。副システムから完了報告を受け取るとデータベースサーバに情報書き込み完了報告を送信する(510)。ここでは、正システムと副システムのストレージ装置間で同期コピーを行なう方法としている。

【0062】

このように、正システムのストレージ装置はデータベースサーバからの処理要求を待ち、ディスクの情報読み込み・書き込み処理と副システムへのディスクコピー処理を繰り返す。

【0063】

図6に示した副システムストレージ装置処理手順のフローチャートについて説明する。

【0064】

まず、ストレージ装置内のディスクドライブをデータベースサーバに割り当てや外部ストレージ装置からのディスクコピー設定などの、初期設定を行う(601)。

【0065】

そして、読み込み・書き込み処理開始(602)後、データベースサーバからの要求を受信可能な状態となる。さらに、データベースサーバから更新通知するディスクを指定される(603)ことでデータ複製の準備が整う。

【0066】

処理要求受信(604)を開始し、処理要求を受けるとその要求が書き込み要求か(605)判定する。書き込み要求でなければ、読み込み要求された情報を要求元に転送(610)し、要求元に情報転送完了報告を送信する(611)ことで処理要求の実行が完了する。一方、書き込み要求であった場合、まず、要求された情報をディスクに書き込む(606)。そして、要求元に情報書き込み完了報告を送信する(607)。通常の処理要求はこれで処理が完了するが、データ複製方法を実施するためにディスクへの書き込みが発生した場合は、そのディスクがデータベースサーバに更新通知を指定されたディスクか(608)判定し、指定されたディスクであれば、データベースサーバに更新を通知する(609)。指定されていなければ、通知せずに処理を完了する。

【0067】

このように、副システムのストレージ装置は正システムからのログディスク書き込みと副

10

20

30

40

50

システムのデータベースサーバからの読み込み・書き込み要求を処理し、更新を通知するディスクであればその通知を送信する処理を繰り返して、データ複製処理を実現する。

【0068】

本実施形態のように正システムから副システムに切り替えるが、上記のようなデータ複製方法を用い、特にストレージ装置8, 18間のディスクコピーが同期コピー方法で行うとデータ複製に必要なログを漏れなくコピーできるため、システム切り替えに伴う遅延を小さくすることができる。

[第2実施形態]

第1実施形態では、ディスクドライブの更新をデータベースサーバに通知するために通常のディスクドライブ読み込み書き込み以外のインタフェースを必要としたが、以下で説明する第2実施形態においてはディスクドライブの読み込み書き込みインタフェースのみでデータ複製を実現する。データ複製システムの構成は、図1に示す第1実施形態と同様に構成される。

10

【0069】

第1実施形態では、副システムにおいてディスクドライブ17の更新をデータベースサーバ12に通知するための設定をした。これに対し、本実施形態ではデータベースサーバ12からディスクドライブ17をポーリングで監視し、更新を検知する。

【0070】

更新の検知は以下のような手順で可能である。データベースサーバ12で、ディスクドライブ17にログが書き込まれる位置を保持し、その位置のデータを定期的に読み込んで更新されたかを判定する。更新されていれば、ログに従って業務処理をデータベースサーバ12で行う。処理が完了したら、ログが書き込まれる位置の更新を行い、再び定期的に読み込んで更新されたかの判定処理を繰り返す。

20

【0071】

データベースの更新ログは、通常一定の領域に順次上書きされないようにディスクに書き込む。そして、領域の終端まで書き込むと再び領域の先頭から順次書き込む。そのため、更新ログの書き込みが上書きされる前にデータベースサーバで内容を読み込むことができることと、更新ログを1つずつ区別することができるようになっていたことが保証される場合、ポーリングで監視してデータのディスクドライブを更新することでデータの複製が可能である。

30

【0072】

このように、データベースサーバからディスクドライブの更新をポーリングによって監視する方法でデータ複製する場合、ポーリング間隔を十分小さくすることによって、システム切り替えによる遅延を小さくすることができる。また、第1実施形態と同様のシステム構成であり、データ複製にかかるコストを小さくできる。さらに、ストレージ装置8、18間のディスクコピーを同期コピー方法で行うとデータ複製に必要なログを漏れなくコピーできるため、システム切替に伴う遅延を小さくすることができる。

[第3実施形態]

第1実施形態、第2実施形態では、データベースサーバ2, 12とストレージ装置8, 18が直接接続された場合やストレージエリアネットワークで接続されたことを前提としていたが、本実施形態では、ストレージ装置としてネットワークアタッチトストレージ(NAS)装置を使用して実現する。本実施形態のデータ複製システムの構成は、図1に示す第1実施形態と同様の構成である。

40

【0073】

本実施形態では、ストレージ装置間のディスクコピー方法が第1実施形態、第2実施形態と異なる。NAS装置はファイルシステムでのアクセス要求を受信する。そのため、ディスク制御装置5, 15間のストレージ間接続インタフェースもファイル単位でのアクセスを実行する。そのため、ディスク制御装置5内でディスクドライブ7の変更を検知するのではなく、データベースサーバ2が操作するログファイルの更新を検知する必要がある。更新の検知には、ディスク制御装置5でログファイルの更新を定期的に監視するデーモン

50

を実行しておき、更新が起こったらストレージ間接続インタフェースを通じてファイルのコピーをディスクドライブ17に書き込む。また、副システムでログの更新を通知するインタフェースも、前述の第1実施例、第2実施例のものとは異なる。ディスク制御装置15にはログファイル更新を検知するデーモンを備える。データベースサーバ12にはこのデーモンと通信を行うプロセスを生成しておき、つまり、ログファイルの更新があったら通知される機構を構築する。あるいは、第2実施形態で示したデータベースサーバ12からポーリングで監視する方法をログファイルに適用してもよい。

【0074】

また、ログデータのコピー方法としてディスク制御装置5と別のディスクドライブ単位で変更を検知可能なディスク制御装置をストレージ装置8, 18に設け、その間でディスクの更新を行う方法も可能である。この場合、別途設けたディスク制御装置により、第1実施形態のようにディスクドライブ7からディスクドライブ17へのデータコピーを実行する。ログデータの更新をデータベースサーバ12に通知する方法は、上記のようなディスク制御装置15のデーモンとデータベースサーバ12のプロセスで通信を行う方法やデータベースサーバ12からポーリングで監視する方法によって可能である。

【0075】

また、NAS装置ではログファイルをコピーする際に、ファイルを全てストレージ装置8, 18間で転送する必要がある。データ転送量を削減するため、ログをおくためのディレクトリを作成し、更新ログを1個ずつファイルとしてそのディレクトリに置いていくことで更新ログの転送量を削減可能である。データベースは更新ログを作成した日時をファイル名に使用して書き込む。これによりログの一意の識別が可能になる。また、一定以上の時間が経過したログを削除することで、ログの複製を保証することができ、ディスクドライブを使い尽くすことがなくなる。ログデータの更新のデータベースサーバ12への通知は、上記同様であるが、ディスク制御装置15で実行するログ更新を監視するデーモンは、ログディレクトリの下にあるファイルの監視を行い、新たなログファイルが作成された場合にデータベースサーバ12のプロセスに通知を行う。データベースサーバ12でファイル更新をポーリングで監視するデーモンを実行する場合も、同様にログディレクトリの下に新たなファイルが作成されたかを監視する。

【0076】

このようにシステムを構築することで、データ複製が実現される。システム構成は第1実施形態と同様であり、コストを削減することが可能である。また、システム切り替え時の遅延も小さくすることが可能である。

[第4実施形態]

図7は本発明が適用されたデータベースサーバとストレージ装置およびそれらの上で動作するソフトウェアを用いたデータ複製システムの実施例である。本実施例では、サーバ上で動作するアプリケーションプログラムをDBMSとしているが、データの更新をログとして出力するプログラムであれば本発明は適用可能である。例えばトランザクションモニタであってもよい。

【0077】

データ複製システムの構成は、図1に示す第1実施形態と同様に構成されるが、図7においてはデータベースサーバ2、12上のストレージ接続装置3、13は省略した。データベースサーバ2、12上では、それぞれOS72、75とDBMS71、74とが動作する。OS72、75はデータベース2、12のハードウェア制御や他のアプリケーションプログラムの動作環境として実行される。また、DBMS71、74はシステムの業務を実行するアプリケーションプログラムである。さらに、ストレージ装置8、18上のディスク制御装置5、15上ではデータベースサーバ2、12からの要求を受付けてディスクドライブ6、7、16、17の更新を行う制御ソフトウェア73、76が動作する。第1実施形態でストレージ接続装置3、13を通して行っていた処理は、DBMS71、74からはOS72、75の機能によって内部で処理されているものとする。

【0078】

本実施形態では第2実施形態と同様のシステム構成で、各装置とそれら装置の上で動作するソフトウェアによってデータ複製を実現する。つまり、ストレージ装置18からデータベースサーバ12へと、ディスクドライブ16、17の情報が変更されたことを通知するインタフェースは持たない。

#### 【0079】

正システムを構成するデータベースサーバ2上で動作するOS72はストレージ装置8内のディスクドライブ6、7に情報を書き込んだり、あるいはその情報を読み込んだりする操作を実行できる。制御ソフトウェア73はサーバ・ストレージ間接続インタフェース4を通じて送られるOS72からの要求に応じて情報のディスクドライブ6、7への反映や情報のOS72への転送などの処理をする。そして、DBMS71は、OS72の上で動作し、システムの業務を行う。副システムはデータベースサーバ12とストレージ装置18、およびこれらの上で動作するOS75、DBMS74そして制御ソフトウェア76で正システムと同様の処理を行うことが可能な構成とする。

10

#### 【0080】

制御ソフトウェア73、76の間では、ストレージ間接続インタフェース20を通じて互いの保持するディスクドライブの内容を転送し、指定されたディスクドライブに情報を反映することが出来る。例えば、ストレージ装置間のリモートコピーや同期リモートコピーを用いることで実現できる。

#### 【0081】

まず、正システムと副システムの初期設定を行う。正システムでOS72が認識したストレージ装置8内のディスクドライブ6、7をDBMS71のデータディスクとログディスクに割り当てる。副システムでは、正システムのディスク構成にあわせて、OS75が認識したストレージ装置18内のディスクドライブ16をデータディスクに、ディスクドライブ17をログディスクにDBMS74が割り当てる。また、DBMS71とDBMS74との間で互いの稼動状態をチェックする。正システムが停止したときに、迅速に副システム業務を引き継ぐためである。

20

#### 【0082】

ストレージ装置8、18の間の初期設定は、ディスクドライブ7の情報の変更を、ディスクドライブ17にコピーするようにする。これによって、DBMS71で実行した処理のログを副システムにコピーし、DBMS74から参照できるようにする。このコピーの方式は同期、非同期いずれとも可能であるが、非同期の場合、正システムが停止した時点までの全てのログを転送できる保証はない。一方、同期とした場合、正システムが停止するまでのログが完全に副システムで参照可能となり、高速で正確なデータ複製が可能である。

30

#### 【0083】

ここまでの設定が終わったら、以下で説明する手順を実行することで正システムから副システムへのデータ複製が可能である。そこで、正システムで業務を開始する。DBMS71は業務ネットワーク1からの業務処理要求を受付けて処理する。このときの処理の内容を逐一再生可能なものがログにあたり、DBMS71からOS72の機能を經由して制御ソフトウェア73にディスクドライブ7に書き込みように要求する(702)。ただし、業務処理要求は業務ネットワークから受け取るだけでなく、データベースの構築や保守のためにデータベースサーバ2上で実行した処理についても業務処理要求に含まれる。そして、完全なデータ複製を行うためには、これら全ての業務処理要求についてログを作成しディスクドライブ7に書き込みを行う必要がある。

40

#### 【0084】

DBMS71がディスクドライブ7にログの書き込みを要求する単位には、例えばトランザクションをコミットした単位とする方法がある。あるいは、DBMS71が持つログ用のバッファが一杯になった場合や一定時間経過した場合にDBMS71が持つログ用のバッファの内容をディスクドライブ7に書き込み要求する方法などがある。いずれの場合にも、DBMS71が実行した業務のログが副システムのDBMS74に欠損なく、またD

50

B M S 7 1、7 4それぞれが管理するデータの間で不整合が生じることなく複製するためには、D B M S 7 1がディスクドライブ7に要求した単位でディスクドライブ7とディスクドライブ17間の情報が同期していることが必要である。以下では、ディスクドライブ7の情報とディスクドライブ17の情報とが同期しているものとする。しかし、必ずしも同期している必要は無く、非同期でディスクドライブ7からディスクドライブ17へのコピーを行ってもよい。但し、非同期コピーを行った場合、正システムと副システムのログが一致することは保証されず、正システムから副システムへ切り替える時にD B M S 7 1とD B M S 7 4が扱うデータに不整合があり得るため、正システムと副システム間でデータ内容を確認したり修正したりする作業が必要になることがある。

【0085】

要求を受けた制御ソフトウェア73は、要求をログ用のディスクドライブ7に書き込む(703)。このディスクドライブ7は初期設定でストレージ装置18のディスクドライブ17にコピーすると指定しているため、制御ソフトウェア73はストレージ装置間接続インタフェース20を通して制御ソフトウェア76に変更分の書き込み要求を転送し、これを受けた制御ソフトウェア76がディスクドライブ17に要求された情報を書き込む(704)。ここで、ディスクドライブ7とディスクドライブ17の間で同期コピーを行うためには、ディスクドライブ17への情報書き込みが完了してからデータベースサーバ2に完了報告を行い、その報告を受けた後でD B M S 7 1が業務処理の続きを行うようにすることが必要である。こうすると、D B M S 7 1の実行した業務のログと副システムに複製しているログとが一致する。一方、ディスクドライブ7への書き込みが完了した時点でデータベースサーバ2に書き込み完了報告を行い、ディスクドライブ17へのコピーをD B M S 7 1からの要求と無関係の任意の時点で行うと、非同期コピーとなり副システムで参照するログが正システムのものとは一致しないことが起こりうる。

【0086】

さらに、D B M S 7 1は業務処理によって生じるデータの変更をO S 7 2の機能を通じて制御ソフトウェア73に送る(704)。そして制御ソフトウェア73は要求された情報書き込み処理や読み込み処理をディスクドライブ6に対して行う(705)。ただし、D B M S 7 1からのディスクドライブ7へのログ書き込み要求とディスクドライブ6へのデータ書き込み要求は順序が決められているわけではない。データ書き込みはO S 7 2上のデータバッファに蓄積して定期的に要求する方法や一定時間以上業務処理で使用しなかったら書き込む方法などがある。正システムの実行した業務を追うためには、実行した処理のログは副システムに必ず転送しなければならないが、データについては正システムが業務に使用する範囲で整合性が保たれていればよい。

【0087】

一方、副システムのD B M S 7 4はO S 7 5の機能を使用して制御ソフトウェア76に要求を出し(706)、ディスクドライブ17の情報が変更されていないか情報を読み込む(707)。この結果、ディスクドライブ17上のログ内容に変更があったら、そのログの業務を実行し、副システムのディスクドライブ16のデータを更新するため、制御ソフトウェア76に情報更新要求を発行する(708)。この要求を受けた制御ソフトウェア76はディスクドライブ16に情報変更の内容を反映する(709)。これらの処理を一定時間ごとに行い、ログの処理を実行してデータを更新して正システムのデータを副システムに複製する。あるいは、実施例1のように、サーバ12からストレージ制御装置15あるいは制御ソフトウェア76に指示しておき、ディスクドライブ17の情報が更新されたらそれをサーバ12に通知するインタフェースによりログの変更を認識する方法でも可能である。

【0088】

D B M S 7 4がディスクドライブ17からログの変更を読み込んでディスクドライブ16にあるデータに適用し、データを複製する方法は複数ある。例えば、D B M S 7 1からディスクドライブ7への書き込みをトランザクションのコミットした単位で行っていれば、D B M S 7 4はディスクドライブ17で新たに書き込まれたログをそのまま実行すればよ

10

20

30

40

50

い。また、DBMS 7.1からディスクドライブ7への書き込みがトランザクションのコミットと無関係に実行されている場合は、書き込まれたログをそのまま実行する方法と、DBMS 7.4でトランザクションの管理テーブル上に確保しておき、そこに読み込んだログを保持しながらトランザクションのコミットを受け取ってから管理テーブルにあるログをデータに反映する方法などがある。

#### 【0089】

もしも、DBMS 7.1とDBMS 7.4との間で行っている稼働状態監視で、DBMS 7.4がDBMS 7.1の停止を検出した時には、DBMS 7.1が行っていた業務を引き継ぐように設定を変更する(710)。ただし、DBMS 7.4が業務処理を引き継ぐ前に、それまでにDBMS 7.1で実行された正システムの業務データと副システムのデータとの間で整合性を保証するための処理、つまりリカバリを行う必要がある。ディスクドライブ7とディスクドライブ17の間のログ同期がトランザクションのコミット単位で行われていれば、ディスクドライブ17に副システムで未実行のログがあるかを確認し、あればそのログの処理を実行して制御ソフトウェア76に情報更新要求を発行する。(708)そして、制御ソフトウェア76は要求にディスクドライブ16上の情報を更新する(709)。ディスクドライブ17の未実行のログ全てについてディスクドライブ16上のデータ更新を実行したらリカバリが完了し、業務ネットワーク1からの業務処理要求をサーバ12で受けてDBMS 7.4が業務処理を開始する。

10

#### 【0090】

上記の、システム間で業務処理を引き継ぐ前に業務データの整合性を保証するために実行する処理、つまり、リカバリはディスクドライブ17へのログ書き込みがコミット単位であることが保証される場合と、保証されない場合とで異なる。ここで、コミット単位での書き込みが保証されない場合とは、例えば、チェックポイント処理によりコミット前の情報を書き込む場合や、DBMS 7.1の持つバッファがあふれてコミット前の情報が書き出される場合などがある。以下、コミット単位での書き込みが保証される場合とされない場合、それぞれについてリカバリの方法を説明する。

20

#### 【0091】

まず、コミット単位でのログ書き込みが保証される場合のリカバリ処理について説明する。コミット単位での書き込みが保証される場合は、ディスクドライブ17に書き込まれたログを全てデータに反映すればよい。したがって、ディスクドライブ17で未実行であるログを認識し、該ログをロールフォワードすることによりリカバリが実現する。同期コピーを用いた場合は、正システムでコミットした情報は副システムに全てコピーされたことが保証され、かつ、該リカバリ処理によりコピーされたログを漏れなく適用できるため、本システムにおいては、トランザクション欠損が無いことが保証される。

30

#### 【0092】

コミット単位での書き込みが保証されない場合は、ディスクドライブ7にはコミット済みのトランザクションログとコミット未済みのトランザクションログとが含まれる。ここで、整合性を保証すべきなのは、コミット済みトランザクションだけである。なぜなら、コミット未済みの処理による変更は、正システムのデータに反映されていないからである。そのため、リカバリはコミット済みのトランザクションについてのみ実行し、未済みのトランザクションは実行しない。このようリカバリ方法として、以下の3つの方法が考えられる。

40

#### 【0093】

例えば、ディスクドライブ7に書かれたログを全てロールフォワードした後、未コミット分を検出し、検出した未済みのトランザクションのみをロールバックする方法が考えられる。例えば、ログを端末から先頭に向かってスキャンし未コミットのトランザクションを検出する方法がある。あるいは、あらかじめトランザクションの状態を管理するテーブル(以降、トランザクション管理テーブルとする)をDBMS 7.4が用意し、正システムのDBMS 7.1の停止を検出する前からディスクドライブ17へのログ書き込みが行われてロールフォワードをするたびにその状態を変更しておく方法がある。後者の場合、リカバ

50



りするときにトランザクション管理テーブルに登録されている未コミットのトランザクションをロールバックすればよい。

【0094】

また、正システムのDBMS71が停止したことを検出した時点以降にログの変更を検出した場合については、ロールフォワードを行う前に、トランザクションの状態を調査してからロールフォワードを行う方法も考えられる。コミット済みのトランザクションのロールフォワードが終了したらDBMS71が停止したことを検出した時点より以前について、コミット未済みのトランザクションを調査して、コミット未済みのトランザクションが存在したらロールバックを行う。コミット未済みのトランザクションの調査はDBMS71が停止したことを検出した時点からログの終端までスキャンしたり、管理テーブルの情報を用いるなどして行えばよい。

10

【0095】

また、コミットが確定したトランザクションのみロールフォワードする方法も考えられる。例えば、上記管理テーブルにおいてコミット済み状態になった時点で、該トランザクションのロールフォワードをする方法が考えられる。

【0096】

以上、コミット単位でのログ書き込みが保証されない場合の3つのリカバリ方法を同期コピーとともに用いると、正システムでディスクドライブに書き込まれた情報は副システムにコピーされたことが保証され、かつ、該リカバリ方法により、コミット済みのトランザクションはもれなくロールフォワードされるため、本システムにおいてはトランザクション欠損がないことが保証される。

20

【0097】

このように、正システムで実行した業務のログを副システムにコピーして、そのログから業務データを復元する処理を行うことで、正システムのデータベースサーバの処理負荷を上げることなく、また正システムと副システムの間で情報を転送するための通信帯域を抑制して低コストのデータ複製が可能となる。また、正システムから副システムへログを同期コピーで転送すると、業務処理の欠損を起こすことなくデータ複製を完了し、正システムの障害発生時や保守により正システムを停止しなければならない時、迅速に副システムに業務の実行を切り替えることが可能である。

30

【0098】

【発明の効果】

本発明によれば、データベースサーバとストレージ装置からなる複数のシステム間で、低コストかつ通常業務の遅延が小さいデータ複製を実現する。また、システム切り替え時の遅延を小さくすることができる。

【図面の簡単な説明】

【図1】正システムと副システムとからなるデータ複製システムおよび正・副システム間のデータ複製方法の概念図である。

【図2】正システムと副システムとからなるデータ複製システムおよび正・副システム間のシステム切り替え方法の概念図である。

【図3】正システムのデータベースサーバが行う処理手順のフローチャートである。

40

【図4】副システムのデータベースサーバが行う処理手順のフローチャートである。

【図5】正システムのストレージ装置が行う処理手順のフローチャートである。

【図6】副システムのストレージ装置が行う処理手順のフローチャートである。

【図7】正システムと副システムとからなるデータ複製システムおよび正・副システム間のシステム切り替え方法の概念図である。

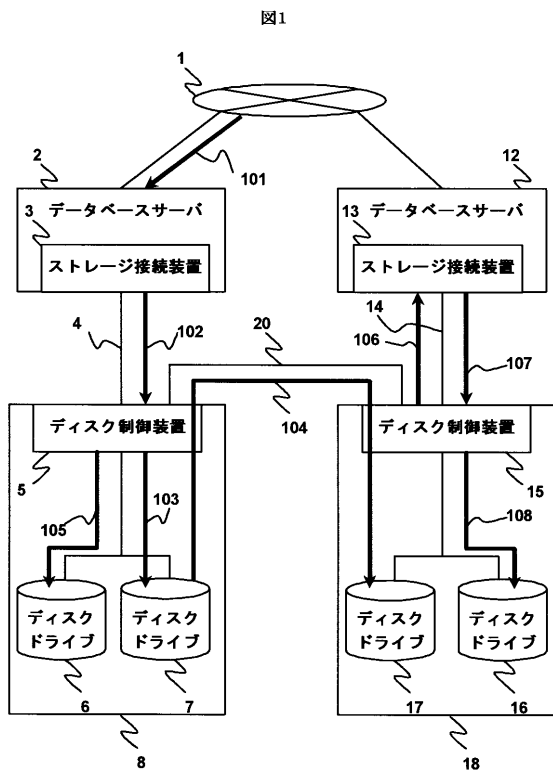
【符号の説明】

1：業務ネットワーク、2：正システムのデータベースサーバ、3：ストレージ接続装置、4：サーバ・ストレージ間接続インタフェース、5：ディスク制御装置、6：データディスクドライブ、7：ログディスクドライブ、8：正システムのストレージ装置、12：副システムのデータベースサーバ、13：ストレージ接続装置、14：サーバ・ストレ

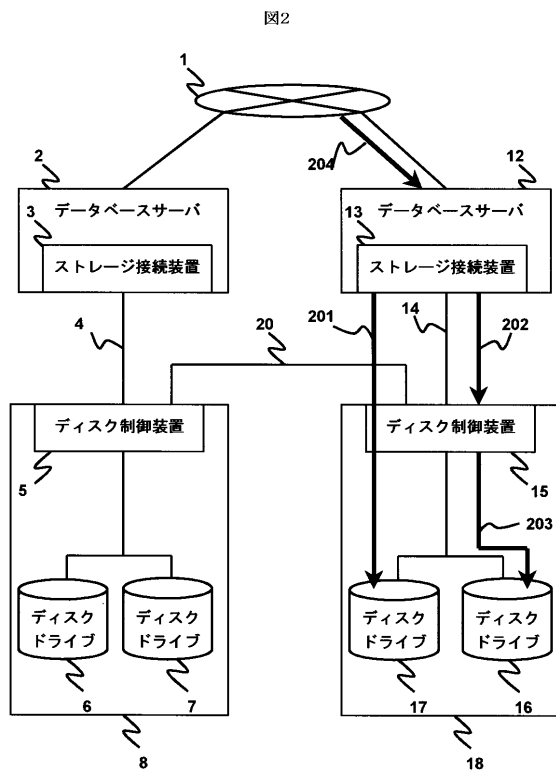
50

ジ間接続インタフェース、15：ディスク制御装置、16：データディスクドライブ、17：ログディスクドライブ、18：副システムのストレージ装置、20：ストレージ装置間接続インタフェース。

【図1】



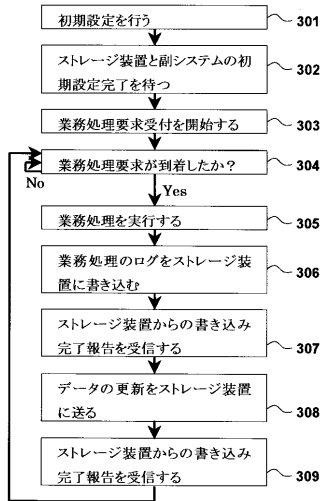
【図2】



【 図 3 】

図3

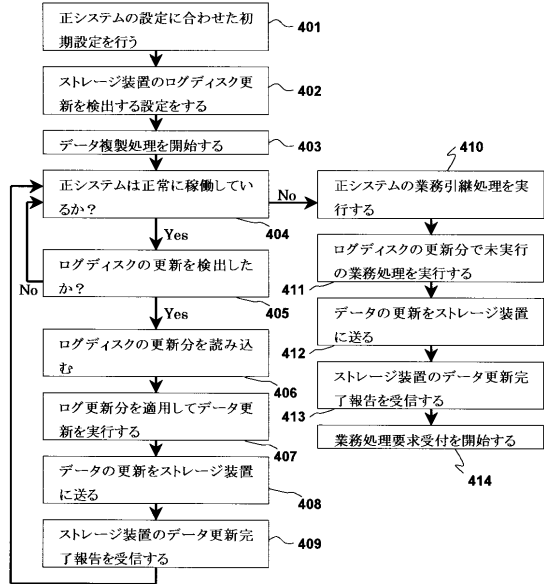
正システムデータベースサーバ処理手順



【 図 4 】

図4

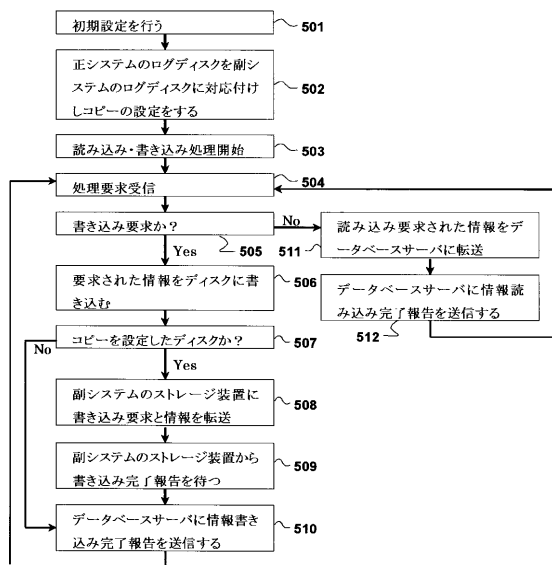
副システムデータベースサーバ処理手順



【 図 5 】

図5

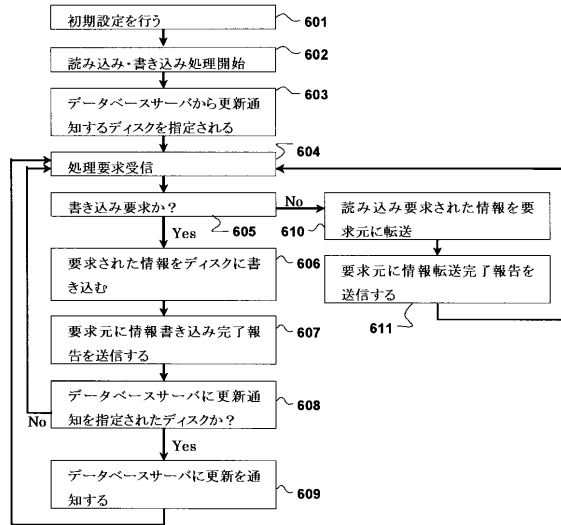
正システムストレージ装置処理手順



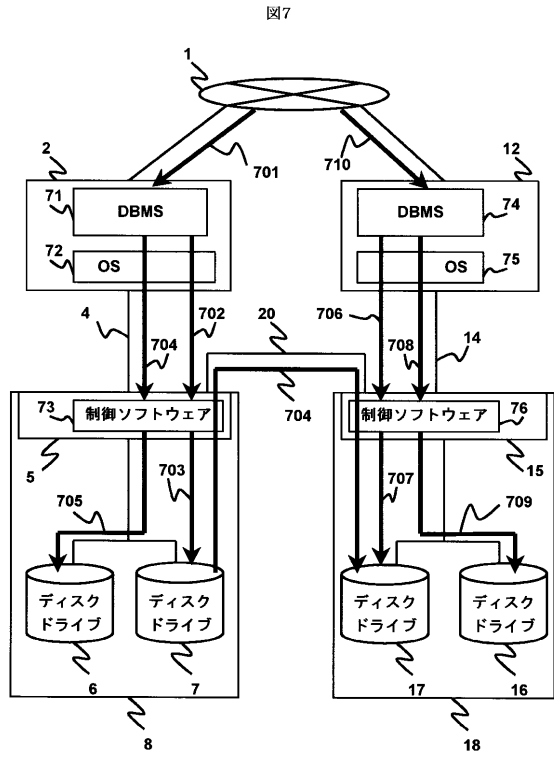
【 図 6 】

図6

副システムストレージ装置処理手順



【図7】



---

フロントページの続き

(56)参考文献 森側 真一, [動向]多様化する災害対策, 日経オープンシステム 第105号, 日本, 日経B P  
社, 2001年12月15日, p.104-107

(58)調査した分野(Int.Cl., D B名)

G06F 12/00

G06F 3/06