

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4546629号
(P4546629)

(45) 発行日 平成22年9月15日(2010.9.15)

(24) 登録日 平成22年7月9日(2010.7.9)

(51) Int. Cl.	F I		
G06F 12/00 (2006.01)	G06F 12/00	533J	
G06F 3/06 (2006.01)	G06F 12/00	518A	
G06F 12/08 (2006.01)	G06F 12/00	531D	
G06F 12/16 (2006.01)	G06F 3/06	304B	
	G06F 3/06	540	
請求項の数 8 (全 21 頁) 最終頁に続く			

(21) 出願番号	特願2000-285237 (P2000-285237)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成12年9月14日(2000.9.14)	(74) 代理人	100093861 弁理士 大賀 真司
(65) 公開番号	特開2002-49517 (P2002-49517A)	(72) 発明者	占部 喜一郎 神奈川県小田原市国府津2880番地 株式会社日立製作所 ストレージシステム事業部内
(43) 公開日	平成14年2月15日(2002.2.15)	(72) 発明者	裏谷 郁夫 神奈川県小田原市国府津2880番地 株式会社日立製作所 ストレージシステム事業部内
審査請求日	平成19年8月6日(2007.8.6)	審査官	桜井 茂行
(31) 優先権主張番号	特願2000-159547 (P2000-159547)		最終頁に続く
(32) 優先日	平成12年5月25日(2000.5.25)		
(33) 優先権主張国	日本国(JP)		

(54) 【発明の名称】 記憶システム、記憶システムの応答方法及び記録媒体

(57) 【特許請求の範囲】

【請求項1】

計算機と、第二のストレージ装置と、に接続された第一のストレージ装置であって、
プロセッサと、

前記プロセッサから参照される制御プログラムやデータが保存されたメモリと、を有し

前記プロセッサは、

前記計算機から送信された複数の書き込み対象のデータに、前記計算機からの書き込み
要求が発行された順に識別子を付与し、

前記書き込み対象のデータを、前記書き込み要求とは非同期に、前記第二のストレージ装
置に送信し、

前記計算機によって実行されるアプリケーションプログラムのトランザクションごと
に対するコミットに基づいて、前記計算機が発行した、前記第二のストレージ装置が前記書
込み対象のデータを格納したかを確認するデータ格納確認コマンドを前記計算機から受信
して、前記第二のストレージ装置のボリュームに書き込まれるデータの識別子を前記
第二のストレージ装置に問い合わせ、

前記データ格納確認コマンドを受信するまでに、前記計算機から書き込まれたデータに
付加した第1の識別子と、前記第二のストレージ装置から得た第2の識別子とを比較し、
この比較結果に基づいて、前記データ格納確認コマンドに対する完了応答を、前記計算機
へ送信する、ストレージ装置。

【請求項 2】

前記識別子は、複数の前記書込み対象のデータの順序性を示すシーケンス番号である、請求項 1 記載のストレージ装置。

【請求項 3】

前記第 2 の識別子は、前記第二のストレージ装置への前記問い合わせの直前に、前記第二のストレージ装置に格納されたデータに付与された識別子であり、

前記第 1 の識別子は、前記データ格納確認コマンドを受信する直前に、前記計算機から書き込まれた最新データに対する識別子である、請求項 2 記載のストレージ装置。

【請求項 4】

前記アプリケーションプログラムが前記トランザクションのコミットポイントでコミットコマンドを送信した直後に、前記計算機から、前記データ格納確認コマンドを受信する、請求項 1 記載のストレージ装置。

10

【請求項 5】

前記計算機は前記アプリケーションプログラムを複数備え、
前記複数のアプリケーションプログラムはグループ分けされており、
前記識別子は、前記複数のアプリケーションプログラムのグループを識別するための情報を含む、請求項 2 記載のストレージ装置。

【請求項 6】

計算機に接続する第一のストレージ装置と、第二のストレージ装置と、の間のデータコピー制御方法であって、

20

前記第一のストレージ装置は、
前記計算機から送信された複数の書き込み対象のデータに、前記計算機からの書き込み要求が発行された順に識別子を付与し、

前記書き込み対象のデータを、前記書き込み要求とは非同期に、前記第二のストレージ装置に送信し、

前記計算機によって実行されるアプリケーションプログラムのトランザクションごとに対するコミットに基づいて、前記計算機が発行した、前記第二のストレージ装置が前記書込み対象のデータを格納したかを確認するデータ格納確認コマンドを前記計算機から受信して、前記第二のストレージ装置のボリュームに書き込まれるデータの識別子を前記第二のストレージ装置に問い合わせ、

30

前記データ格納確認コマンドを受信するまでに、前記計算機から書き込まれたデータに付加した第 1 の識別子と、前記第二のストレージ装置から得た第 2 の識別子とを比較し、この比較結果に基づいて、前記データ格納確認コマンドに対する完了応答を、前記計算機へ送信する、リモートコピー制御方法。

【請求項 7】

計算機に接続する第一のストレージ装置に、第二のストレージ装置へのデータコピーを実行させるためのプログラムが記録されたコンピュータ読取可能な媒体であって、

前記プログラムは、
前記計算機から送信された複数の書き込み対象のデータに、前記計算機からの書き込み要求が発行された順に識別子を付与するステップと、

40

前記書き込み対象のデータを、前記書き込み要求とは非同期に、前記第二のストレージ装置に送信するステップと、

前記計算機によって実行されるアプリケーションプログラムのトランザクションごとに対するコミットに基づいて、前記計算機が発行した、前記第二のストレージ装置が前記書込み対象のデータを格納したかを確認するデータ格納確認コマンドを前記計算機から受信して、前記第二のストレージ装置のボリュームに書き込まれるデータの識別子を前記第二のストレージ装置に問い合わせるステップと、

前記データ格納確認コマンドを受信するまでに、前記計算機から書き込まれたデータに付加した第 1 の識別子と、前記第二のストレージ装置から得た第 2 の識別子とを比較するステップと、

50

この比較結果に基づいて、前記データ格納確認コマンドに対する完了応答を、前記計算機へ送信するステップと、を前記第一のストレージ装置に実行させるものである、媒体。

【請求項 8】

計算機と、

前記計算機に接続する第一のストレージ装置と、

前記第一のストレージ装置に接続する第二のストレージ装置と、

を備える、データコピー制御システムであって、

前記第一のストレージ装置は、

プロセッサと、

前記プロセッサから参照されるプログラムやデータが保存されたメモリと、を有し、

前記プロセッサは、

前記計算機から送信された複数の書き込み対象のデータに、前記計算機からの書き込み要求が発行された順に識別子を付与し、

前記書き込み対象のデータを、前記書き込み要求とは非同期に、前記第二のストレージ装置に送信し、

前記計算機によって実行されるアプリケーションプログラムのトランザクションごとに対するコミットに基づいて、前記計算機が発行した、前記第二のストレージ装置が前記書き込み対象のデータを格納したかを確認するデータ格納確認コマンドを前記計算機から受信して、前記第二のストレージ装置のボリュームに書き込まれるデータの前記識別子を前記第二のストレージ装置に問い合わせ、

前記データ格納確認コマンドを受信するまでに、前記計算機から書き込まれたデータに付加した第 1 の識別子と、前記第二のストレージ装置から得た第 2 の識別子とを比較し、この比較結果に基づいて、前記データ格納確認コマンドに対する完了応答を、前記計算機へ送信する、データコピー制御システム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、ディスク制御装置間の非同期リモートコピーにおけるデータ同期制御に関するものである。

【0002】

【従来の技術】

近年のコンピュータシステムは、銀行及び証券業務の基幹業務を、大型コンピュータによる一括管理から、クライアント・サーバシステムを中心とする分散システムへ移行している。このような分散システム環境では、クライアントからの要求を複数のサーバとディスクアレイ装置を用いてデータを処理する HA (High Availability: 高可用性) クラスタ構成が採られている。このような HA クラスタ構成では、地震などの災害に備えて遠隔地にあるデータセンター間でデータを二重化する方法が採用されてきている。二重化は、通常、2 台のディスクアレイ装置 (記憶システム) を公衆回線や専用回線等を介してリモート接続し、ホストコンピュータ装置からローカルのディスクアレイ装置への書き込みデータを、リモートのディスクアレイ装置へコピーする方法が採られている。

【0003】

ディスクアレイ装置間で二重化する方法は、大別すると同期方式と非同期方式の 2 種類がある。

【0004】

同期方式では、まず、ローカル側のホスト装置からの書き込み要求を、ローカル側のディスクアレイ装置のキャッシュに書込む。続いて、ローカル側のディスクアレイ装置は、キャッシュに書込まれたデータを、リモートのディスクアレイ装置に転送する。更に、ローカル側のディスクアレイ装置は、書き込み要求のデータがリモート側のディスクアレイ装置により受信されたこと示す応答信号を受信すると、この後、ホストに対して、前記書き込み要求に対する応答を返す。つまり、ローカル側のディスクアレイ装置は、リモート側のディ

10

20

30

40

50

スクアレイ装置にデータが届けられたことを確認して、ホストへ応答を返す。この応答は、ホスト装置に対しリモート側にデータが渡されたことを保証する。(同期とは、ホストからの書込み要求とリモート側のディスクアレイ装置へのコピーが同期して行われるという意味で用いられる。)この同期方式は、リモート側からの応答信号を待つ遅延が発生するため、ローカルとリモート間でデータ伝送の伝播遅延が少ない比較的近距離(100km以内)に適しているが、例えば公衆回線網等を使用する遠距離転送には適さない。なお、ローカル側およびリモート側のディスクアレイ装置に記録されたデータは、それぞれのドライブ制御回路を介してそれぞれの物理ディスクへ書込まれる。

【0005】

一方、非同期方式は、遠距離転送に向いており、ローカル側のホスト装置からの書込み要求に対するホスト装置への応答(書込み完了)は、ローカル側のキャッシュに前記書込み要求のデータが書かれた時点で、書込み完了をホスト装置へ返される。キャッシュに書かれたデータは、ホスト装置への応答後に、別のタイミングで(この意味で非同期)、リモート側のディスクアレイ装置へコピー(転送)される。この非同期方式では、リモート側のディスクアレイ装置へデータを転送するタイミングとは関係なく、ホスト装置へ上記書込み要求に対する応答を行うので、応答タイミングが同期方式と比べて速くなり、ホストは次の処理に早く移ることができる。

【0006】

なお、リモートコピーに関しては、特表平8-509565号に記載がある。

【0007】

【発明が解決しようとする課題】

上記非同期にリモート側にデータを転送する方法では、ローカル側のディスクアレイ装置が、当該ディスクアレイ装置内にデータを格納した時点で、リモート側にデータがストアされたか否かに関係なく、ホストに対し、書込み要求完了を報告する。このため、ローカル側のホストは、ホスト書込み要求のリモート側への同期完了確認(ホスト書込み要求によるデータがリモート側のディスクアレイ装置へ確実に転送されたかの確認)が困難であった。このホスト書込み要求のリモート側での同期完了確認は、特に、データベースの履歴ログファイル等でのデータベースのトランザクション単位のコミット(データが確実にストレージに記憶されたという保証)が必要とされる。なお、コミットとは、1つのトランザクションに関する複数のデータベースの更新結果を実際の記憶システム上にログファイルと共に書込む一連の処理のことである。

【0008】

また、災害時のデータリカバリの観点から正サイト障害によってディスクアレイ装置内に残存していた未転送のデータは失われ、副サイトに切替り副サイトで運用を開始した時に保証されるデータが分からない問題があった。

【0009】

しかしながら、上記従来技術の非同期転送方式は、非同期転送の特性上から、ホストIOに対する同期確定方法を備えていなかった。即ち、データベース(DB)運用上必要となる、APP(アプリケーションプログラム)からのトランザクションに対するコミットポイントでの書込み要求がリモートサイト(副サイト)に確実に書かれたか否かを確認する方法が提供されていなかった。

【0010】

以下、課題を具体的に説明する。まず、コンピュータが1つの記憶システムと接続される場合を説明し、その後、記憶システムが非同期データ転送(非同期リモートコピー)を行っている場合につき、その課題を明確に説明する。

【0011】

初めに、コンピュータに1つの記憶システムが接続されている場合を説明する。コンピュータのアプリケーションが、書込みコマンド(要求)を実行すると、通常、コミット(commit)コマンド無しの状態では、その書込みコマンドのデータはコンピュータ内のデータバッファ上に書込まれるだけで、データバッファ上のデータと、記憶システム内のデー

10

20

30

40

50

タは一致しない。その後、アプリケーションがコミットコマンドを発行すると、データバッファ上のデータが記憶システムへ実際に書込みコマンドによって書込まれる。その後、記憶システムは、書込みデータをキャッシュメモリに記憶すると（この時点で記憶システム内のデータとコンピュータ内のデータは一致する）、書込みコマンドを発行したコンピュータに対し、書込み要求完了で応答する。コンピュータは、その書込み要求完了を確認すると、前記コミットコマンドに対するリターンをアプリケーションへ発行する。アプリケーションは、このリターンにより、記憶システム内のデータとコンピュータ内のデータが一致していることを知る。

【0012】

次に、非同期リモートコピーを行っている場合について説明する。コンピュータのアプリケーションが、コミットコマンドを発行すると、データバッファ内のデータが書込みコマンドでローカル側の記憶システムのキャッシュへ書込まれる。ローカル側の記憶システムは、その応答として、書込み完了をコンピュータへ返す。コンピュータは、書込み完了を受けると、アプリケーションへコミットに対するリターンを返す。しかしながら、このリターンは、単にローカル側の記憶システムのデータがデータバッファ内のデータと一致したことを示すのみであり、リモート側の記憶システム内のデータとデータバッファ内のデータが一致したことを示すものではない。従って、ローカル側の記憶システムが書込み要求完了を返した後で、かつ、リモート側へのデータコピーが終了する前にローカル側の記憶システム内のデータが消失した場合に、アプリケーションがリモート側のデータを用いて処理を継続しようとする、コミットのリターンを受け、記憶システム上でデータが確定したと確認したにもかかわらず、誤ったデータを用いて処理を継続することになる。即ち、非同期リモートコピー中に、もし、障害等が発生した場合には、従来からのコミット機能ではアプリケーションは満足な結果を得られない場合が発生することになる。

【0013】

このように、従来技術の非同期リモートコピーは、非同期転送の特性上から、ホストIOに対するデータ同期確定方法を備えていない。その結果、従来技術の非同期転送方式では、データベース(DB)運用上必要となる、APP(アプリケーションプログラム)からのトランザクションに対するコミットポイントでの書込み要求がリモートサイト(副サイト)に確実に書かれたか否かを確認できないという課題がある。

【0014】

本発明の目的は、ホストIO(書込み要求)に対するリモートサイトへのデータ同期確定を、ホストアプリケーションのコミット単位または任意な時点での同期確定を可能にすることにある。

【0015】

【課題を解決するための手段】

かかる課題を解決するため本発明においては、ローカル側の上位装置とリモート側の記憶システムとに接続され、前記上位装置から与えられる書込み対象のデータを記憶すると共に、当該書込み対象のデータを前記リモート側の記憶システムに送信する記憶システムであって、前記リモート側の記憶システムが前記書込み対象のデータを受信したか否かを問い合わせる問合せコマンドを前記上位装置から受信する手段と、前記上位装置へ、前記問合せコマンドに対する応答を送信する手段とを備えることを特徴とする。

また本発明においては、データを記憶するディスク装置と、外部上位装置と前記ディスク装置との間で前記データを転送するディスク制御装置とを有する記憶システムであって、前記ディスク制御装置は、前記上位装置から与えられた前記ディスク装置に書き込むべき書込み対象のデータを他の記憶システムへ送信する回路と、前記上位装置から入力されるコマンドを実行する回路及びソフトウェアとを備え、前記コマンドは、前記他の記憶システムが前記書込み対象のデータを受信したか否かに関する情報を問い合わせるものであり、かつ、自記憶システムに対して、前記外部上位装置にその問合せに対する結果を報告することを要求するものであることを特徴とする。

さらに本発明においては、上位装置に対する記憶システムの応答方法であって、前記上

10

20

30

40

50

位装置から送信される書込みコマンドを受信する第1のステップと、前記書込みコマンドに対する応答を前記上位装置に返す第2のステップと、前記書込みコマンド内のデータを他記憶システムへ送信する第3のステップと、前記上位装置から送信される前記他の記憶システムが前記データを受信したか否かを問う問合せコマンドを受信する第4のステップと、前記上位装置へ、前記問合せコマンドに対する応答を送信する第5のステップとを備えることを特徴とする。

さらに本発明においては、ローカル側の上位装置とリモート側の記憶システムとに接続され、前記上位装置から与えられる書込み対象のデータを対応する正ボリュームに格納すると共に、当該書込み対象のデータを副ボリュームが設けられた前記リモート側の記憶システムに送信する記憶システムであって、上位装置との間で情報を送受信するための第1のインタフェース回路と、他の記憶システムとの間で情報を送受信するための第2のインタフェース回路とを備え、前記上位装置から送信される書込み対象のデータを前記第1のインタフェース回路を介して受信し、当該書込み対象のデータを、当該書込みデータに付与した識別情報と共に前記第2のインタフェース回路を介して前記他の記憶システムに送信し、前記上位装置からの前記書込み対象のデータに関する問合せコマンドを前記第1のインタフェース回路を介して受信し、前記問合せコマンドにおいて指定された前記書込み対象データに付与した前記識別情報を前記第1のインタフェース回路を介して前記上位装置に送信すると共に、最後に前記副ボリュームに格納された前記書込み対象のデータに付与された前記識別情報を、前記第2のインタフェース回路を介して前記他の記憶システムに問い合わせ、前記他の記憶システムから送信される、最後に前記副ボリュームに格納された前記書込み対象のデータに付与された前記識別情報を前記第2のインタフェース回路を介して受信し、当該識別情報を、前記第1のインタフェース回路を介して前記上位装置に送信することを特徴とする。

【0016】

さらに本発明においては、コンピュータ上で動作するアプリケーションから複数の書込み要求を受けた際に、前記複数の書込み要求内のデータを他記憶システムに非同期にコピーする記憶システムと接続される前記コンピュータに、インストールされるプログラムが記録された記録媒体であって、前記プログラムは、前記アプリケーションから問合せコマンドを受ける第1ステップと、前記記憶システムから、前記記憶システムに記憶された前記アプリケーションに関する書込み対象のデータを識別する識別情報を得る第2ステップと、前記記憶システムから、前記他の記憶システムに記憶された前記アプリケーションに関する書込み対象のデータを識別する識別情報を得る第3ステップと、前記第2ステップで得た識別情報と前記第3ステップで得た識別情報を基に、前記問合せに回答する第4のステップとを備える処理を前記コンピュータに実行させることを特徴とする。

さらに本発明においては、ローカル側ホストからの書込み対象のデータを、他の記憶システムへリモート非同期データ転送を行う記憶システムであって、前記記憶システムは、前記書込み対象のデータが前記他の記憶システムへ記憶されたか否かを前記ローカル側ホストへ報告することを特徴とする。

さらに本発明においては、他の記憶システムへデータの非同期リモートコピーを行う記憶システムに接続されるコンピュータにインストールされるプログラムを記録する記録媒体であって、前記プログラムは、前記コンピュータ内で動作するアプリケーションからコミットコマンドを受けると前記コンピュータ内のバッファにある前記アプリケーションが発行した書込み対象のデータを実際に前記記憶システムへ書込む書込み要求を発行する第1のステップと、前記記憶システムから前記書込み要求に対する書込み完了を受信する第2のステップと、前記書込み対象のデータが前記他の記憶システムにコピーされたか否かを問い合わせる問合せコマンドを前記記憶システムへ発行する第3のステップと、前記記憶システムから前記問い合わせコマンドに対する応答を受ける第4のステップと、前記書込み完了を受信し、かつ、前記応答を受けた後で、前記コミットコマンドに対する応答を返す第5のステップとを備える処理を前記コンピュータに実行させることを特徴とする。

【0017】

10

20

30

40

50

【発明の実施の形態】

(実施例)

以下、本発明の一実施例を図1～図3により詳細に説明する。

【0018】

最初に、図3を用いて本発明を適用するシステム構成例を示す。図3は、一般的なH A クラスタの全体構成を示したシステムブロック図である。以下、ディスクアレイ装置13aをローカル(又は正サイト)、ディスクアレイ装置13bをリモート(又は副サイト)とし説明する。また、参照番号に「a」を添付した方をローカル側、参照番号に「b」を添付した方をリモート側として説明する。

【0019】

H A クラスタを構成するホストコンピュータ装置1a、1bはデータベース等のA P P (アプリケーション・プログラム)2a、2b、ミドルソフト3a、3b、A P P やミドルソフトを制御するO S (オペレーティング・システム)4a、4b、及び、ディスクアレイ装置との間でコマンド等の入出力を制御する“ I O I / F ” (入出力インタフェース)23a、23bで構成される。

【0020】

ミドルソフト3a、3bは、リモート側(副サイト)では、正サイト障害での副サイト運用時に、副サイト側のホスト装置に対し、副サイトのディスクアレイ装置への書込み禁止状態を解除し、一方、正サイト側では、正サイトと副サイト間でデータを一致させる初期ペア状態の生成、ペア状態のサスペンド(論理的なペアを切断した状態)等のペア制御指令を大型ディスクアレイ装置13a、13bに送るコマンド制御ソフトである。本発明の一実施例では、このミドルソフト3aにより、本願特有のS y n c コマンドを提供する(詳細は後述する)。

【0021】

ホストコンピュータ装置1a、1bからのI O 要求は“ I O I / F ” (23a、23b)を介して大型ディスクアレイ装置13a、13bに送られる。

【0022】

大型ディスクアレイ装置13a、13bは、H O S T I / F 制御回路17a、17bでホストからのI O 要求を受け、キャッシュ21a、21bに、書込み要求のデータ(W R I T E データ)を書込む。このキャッシュ21a、21bに書かれたデータはドライブ制御22a、22bによって実際の物理ディスク9a、9bに書込まれる。ここで、キャッシュに書込まれたデータに関しては、大型ディスクアレイ装置は、そのデータの記憶を保証するものとする。また、大型ディスクアレイ装置は、複数の物理ディスク9をR A I D (R e d u n d a n c y A r r a y o f I n e x p e n s i v e D r i v e s) 1、R A I D 5 等で保護し物理ディスクの障害に備える。なお、R A I D に関しては、“A C a s e f o r R e d u n d a n t A r r a y s o f I n e x p e n s i v e D i s k s (R A I D) ” b y D a v i d A . P a t t e r s o n , e t a l . , D e c 1987, U. C. B e r k l e y に記載がある。

【0023】

大型ディスクアレイ装置間でのデータの二重化では、まず、ローカル大型ディスクアレイ装置13aのキャッシュ21aに書かれたデータが、リモートI / F 制御回路18a、18bおよび公衆回線や専用線等のデータ転送路20a、20bを介してリモート大型ディスクアレイ装置13bに送られる。リモート側の大型ディスクアレイ装置13bは、キャッシュ21bに、この送られたデータを書き込む。その後、リモート側の大型ディスクアレイ装置13bは、ローカル側での物理ディスク装置9aへの書込みと同様に、物理ディスク9bへ受信したデータを書込む。

【0024】

この二重化された状態では、リモート側のH O S T I / F 制御17bはリモートホスト1bからの書込みを禁止する。

【0025】

ここで、リモート側への未転送データは、B I T M A P メモリ11a、11bによって差

10

20

30

40

50

分管理される。このメモリ内の差分情報によって管理される差分データ（未転送データ）が、ホスト装置 1 a からの I O 要求に同期しないで非同期にリモート側に転送される。この差分データを管理する差分 B I T M A P メモリについては後述する。

【 0 0 2 6 】

なお、図 3 において、本願発明で定義する S y n c コマンド等のディスクアレイ装置内での処理は、H O S T / I F 制御回路やリモート I / F 制御回路等に、マイクロプロセッサおよびマイクロプロセッサから参照されるプログラムやデータが保持される一般的なメモリ（D R A M 等）を設け、マイクロプロセッサがプログラムを実行することによって実現される。また、制御回路内に専用回路を設けることで実現しても良い。一方、当該装置を制御するマイクロプロセッサ、プログラムを記憶するメモリ等を、H O S T / I F 制御回路やリモート I / F 制御回路等以外の場所に設け、これにより、H O S T / I F 制御回路やリモート I / F 制御回路を制御しながら実行しても良い。

10

【 0 0 2 7 】

図 1 は、図 3 のシステム構成において本願発明を実現する論理ブロックを示す。

【 0 0 2 8 】

ローカル側のホスト装置 1 a の論理ブロックは、アプリケーションソフトである A P P （アプリケーションプログラム）2 a、大型ディスクアレイ装置 1 3 a、1 3 b 内のボリュームの初期ペア状態の生成、ペア状態のサスペンド等のペア制御を実行するミドルソフト 3 a、及び O S 4 a で構成される。

【 0 0 2 9 】

ローカルサイトである大型ディスクアレイ装置 1 3 a と、リモートサイトである大型ディスクアレイ装置 1 3 b は、図示しない光ファイバ、広域回線網等のデータ転送路（図 3 の 2 0 a、2 0 b）で接続される。大型ディスクアレイ装置 1 3 a、1 3 b 内のボリュームは、P V O L（正ボリューム）1 0 a と S V O L（副ボリューム）1 0 b で構成され、ホストからのデータは P V O L 1 0 a から S V O L 1 0 b にコピーされ二重化される。なお、これらボリュームは物理ディスク（図 3 の 9 a、9 b）の中に保持される。

20

【 0 0 3 0 】

B I T M A P 1 1 0（図 3 の B I T M A P メモリ 1 3 a に保持される）、B I T M A P 1 2 0（図 3 の B I T M A P メモリ 1 3 b に保持される）は、それぞれ P V O L（1 0 a）と S V O L（1 0 b）間のデータ差分管理テーブル（ボリューム単位に管理可能）であり、P V O L と S V O L の全データブロックを数 1 0 K B（キロバイト）単位でビットマップ化したものである。B I T M A P 1 1 0 は、P V O L（正ボリューム）に格納されているがリモート側へ未転送であるデータを表し、B I T M A P 1 2 0 はリモート側へ転送されたが S V O L（副ボリューム）へまだ格納されていないデータを表す。通常、ペア（P A I R）状態（二重化状態）がサスペンド（P S U S）した状態となると、ホストからの新たなデータは P V O L 1 0 a にのみ書込まれるので、このデータ分が P V O L（1 0 a）と S V O L（1 0 b）の不一致として、この B I T M A P 1 1 0、1 2 0 によって差分管理されることになる。

30

【 0 0 3 1 】

ローカル側の F I F O 7、リモート側の F I F O 8 は、それぞれ、ローカル側の大型ディスクアレイ装置 1 3 a、リモート側の大型ディスクアレイ装置 1 3 b 間の非同期転送用のバッファキューであり、ペア（P A I R）状態の時に使用される。

40

【 0 0 3 2 】

ローカル側のホスト 1 a からの I O（入出力）の書込みデータは、大型ディスクアレイ装置 1 3 a の P V O L（1 0 a）対応のキャッシュに置かれる。その後、データは、物理ディスク 9 a に書込まれると同時に、その I O 単位に、その I O を識別するための識別子であるシーケンス番号が付加され、即ち、ホスト I O データは順序化され、ホストコンピュータからの書込み要求の発行順に、一旦 F I F O 7 にキューされる。なお、識別子は番号の他、タイムスタンプ等、データをホストコンピュータからの書込み要求の発行順に一意に識別できる識別子ならその種類は問わない。また、F I F O バッファ 7 にキャッシュ上

50

のデータをキューイングする際は、このF I F Oバッファに、キャッシュ上のデータの複製を作り実際にキューイングする方法でも良いが、データ量の増加を防ぐため、キャッシュ上のデータを管理するポインタ・アドレス等を管理するキューであることが望ましい。この場合、実際のデータをリモート側へ転送する場合は、キャッシュから実際のデータが転送されることになる。

【 0 0 3 3 】

このシーケンス番号を付加したデータは、ホストへのI O完了応答報告と非同期に、ローカル側の大型ディスクアレイ装置 1 3 a からリモート側の大型ディスクアレイ装置 1 3 b へ転送される。リモート側の大型ディスクアレイ装置 1 3 b は、このデータを受信すると、このデータをシーケンス番号順にF I F O 8 にキューイングする。このF I F O 8 にキューされたデータは、シーケンス番号順にS V O L (副ボリューム) 1 0 b 対応のキャッシュに置かれ、その後、物理ディスク 9 b に書込まれる。なお、受信データを一旦F I F O 8 にキューイングした後にシーケンス番号順にキャッシュに置く理由は、ローカル側の大型ディスクアレイ装置 1 3 a の制御の都合や転送路の通信状態の関係で、リモート側の大型ディスクアレイ装置 1 3 b は、書込み要求のデータを必ずしも、ホストコンピュータからの書込み要求の発行順序で、受信しないからである。本実施例では、前述のように、F I F O 8 に受信データをキューイングし、その後、受信データに割り付けられているシーケンス番号が番号順に揃った段階で、こんどは、F I F O からデータをシーケンス番号順に読み出してS V O L 対応のキャッシュへ正式に書込むことで、その順序性を保証している。当然のことながら、リモート側の大型ディスクアレイ装置 1 3 b がデータを受信した順序がそのまま、ホストコンピュータから発行される書込み要求の発行順序と一致することが保証されるなら、大型ディスクアレイ装置 1 3 b は、受信データをキャッシュへS V O L データとして直接書込める。

【 0 0 3 4 】

大型ディスクアレイ装置 1 3 a 、 1 3 b 間の伝送路等の転送障害等によって非同期転送が出来ない場合、大型ディスクアレイ装置 1 3 a 、 1 3 b は、F I F O 7 、 8 にキューされた未転送データをB I T M A P 1 1 0 、 1 2 0 に差分データとして、ボリューム単位に、管理し、二重化を障害サスペンド(P S U E)状態とする。

【 0 0 3 5 】

ローカル側からリモート側へのデータの転送が進むと、F I F O 7 内にはデータが無くなる。一方、F I F O 8 には、データが受信され、例えば、F I F O 8 内にシーケンス番号「 1 」～「 5 」を持つデータが記憶される。その後、F I F O 8 からS V O L 対応のキャッシュへ受信データが格納されると、F I F O 8 内のデータは減る。この際、S V O L へ新たに書込まれたデータに割り付けられていたシーケンス番号が新たなS V O L 書込みシーケンス番号となる。従って、F I F O 8 から実際にS V O L 対応のキャッシュへ書込まれたデータに割り当てられているシーケンス番号をリモート側の大型記憶システム内に設けられるメモリにS V O L 書込みシーケンス番号として記憶しておき、ローカル側から問合せがあった場合に、このメモリ内のシーケンス番号をS V O L 書込みシーケンス番号として報告すれば良い。そして、リモート側がS V O L 書込みシーケンス番号として「 5 」を報告すると、ミドルウェア 3 a は、(P V O L シーケンス番号「 5 」) < = (S V O L 書込みシーケンス番号「 5 」) が成立したことから、データ同期が確立したと判断する。

【 0 0 3 6 】

次に、本願発明である、ホスト 1 a のA P P 2 a とミドルソフト 3 a と大型ディスクアレイ装置 1 3 a の連携によってA P P 2 a のC o m m i t 時のリモートサイトへの同期確認がどのように実行されるか説明する。

【 0 0 3 7 】

これは、ホスト 1 a のミドルソフト 3 a が、A P P 2 a から同期確認の要求を受け、大型ディスクアレイ装置 1 3 a 、 1 3 b の状態とF I F O 7 、 8 キュー等から得られる、P V O L とS V O L に格納されるデータに関するシーケンス番号を比較し、コミットコマンド

発行直前の“WRITE I/O”(書込み要求)のデータがリモートサイトに同期したかを確認し、APP2aへリターンを返すことで、可能になる。

【0038】

以下、ホスト1aのAPP2aとミドルソフト3aと大型ディスクアレイ装置13aの連携によってAPP2aのCommitのリモートサイトへの同期要求がどのように行われるか具体的に説明する。

【0039】

ホスト1aのAPP2aからのCommitが発行されAPP2aのCommitをリモートサイトへ同期要求する場合、まず、APP2aはトランザクションとしてデータベースへWRITEし、最後にCommitを発行する。通常はこのCommitで完了である。

10

【0040】

本願発明では、大型ディスクアレイ装置間でデータが二重化構成の場合、更に、Commitの直後に、本発明の特徴である“Syncコマンド”を発行する。Syncコマンドは、ホスト装置上で動作するライブラリ及びホストコマンドあり、ソフトウェア製品としてCDROM、フロッピーディスク等のソフトウェア記憶媒体を介して提供される。このSyncコマンドはミドルソフト3によってアプリケーションに提供される。また、パラメータとしてgroup(後で説明する)と最大の同期完了待ち時間を指定するtimeoutによって定義される。

【0041】

20

なお、このSyncコマンドを発行するタイミングは、Commit直後に限定されず、アプリケーションソフトが、リモート側とローカル側でデータの同期が必要と判断した時点で任意に発行可能である。言い換えれば、同期を取る/取らないの選択が、Syncコマンドの発行の有/無で任意にできるので、Syncコマンドを使用するアプリケーションソフトに対する自由度がある。但し、非同期リモートコピーにおいて、トランザクション単位でコミットを行う場合、リモート側の大型ディスクアレイ装置内に書込みデータが確実に記憶されているか否かを確認するという意味では、アプリケーションにおける、コミットコマンドの発行、引き続いての本発明のSyncコマンドの発行という2つのコマンドの流れは必然である。

【0042】

30

このミドルソフト3aは、Syncコマンドを受けると、大型ディスクアレイ装置のペア状態をチェックし、PAIR状態であればFIFO7の最新のPVOLシーケンス番号(最も新しくホストから受領したデータであって、ローカル側のキャッシュ21aに書かれており、リモートへ未コピーのデータに付与されているシーケンス番号)をローカル側のディスクアレイ装置13aから取得し、このPVOLの最新シーケンス番号をCommit直後のシーケンス番号として保持する。次に、ミドルソフト3aは、ローカル側の大型ディスクアレイ装置13aを介して、リモート側の大型ディスクアレイ装置13b内のSVOL書込みシーケンス番号(最も新しくSVOL(リモート側のキャッシュ21b)へ書き込まれたデータに付与されているシーケンス番号)を取得する。更に、PVOLシーケンス番号とSVOL書込みシーケンス番号を比較し、

40

(PVOLシーケンス番号)

$< =$ (SVOL書込みシーケンス番号)

が成立するまでSVOL書込みシーケンス番号を繰返し取得し、テストを繰り返す。条件が成立すると、ミドルソフト3aは、同期完了の応答としてこのSyncコマンドの呼び出しもとであるAPP2aにリターンを返す。APP2aは、このリターンを受けると、つまり、このSyncコマンドが完了したことで同期完了とみなす。

【0043】

ここで、このPVOLシーケンス番号とSVOL書込みシーケンス番号の取得およびデータ同期について、以下、詳細に説明する。PVOLシーケンス番号は、最も新しくローカル側のキャッシュに書込まれたデータに付与されたシーケンス番号であるから、例えば、

50

図1においては、FIFO7から最新のPVOLシーケンス番号「5」が獲得され、これがミドルソフト3aへ報告される。一方、リモート側のSVOL書込みシーケンス番号に関しては、FIFO8にはデータ「1」「2」がキューイングされた状態であることから、SVOLには「0」のシーケンス番号を持つデータまでが書込まれていることが分かる。つまり、FIFO8を用いて、番号「0」が、SVOL書込みシーケンス番号であることが分かり、このSVOL書込みシーケンス番号「0」がミドルソフト3aへ報告される。この状態では、(PVOLシーケンス番号「5」) > (SVOL書込みシーケンス番号「0」)であるから、ミドルソフト3aは、同期がとれていないと判断する。

【0044】

ローカル側からリモート側へのデータの転送が進むと、FIFO7内にはデータが無くなる。一方、FIFO8には、データが受信され、例えば、FIFO8内にシーケンス番号「1」～「5」を持つデータが記憶される。その後、FIFO8からSVOL対応のキャッシュへ受信データが格納されると、FIFO8内のデータは減る。この際、SVOLへ新たに書込まれたデータに割り付けられていたシーケンス番号が新たなSVOL書込みシーケンス番号となる。従って、FIFO8から実際にSVOL対応のキャッシュへ書込まれたデータに割り当てられているシーケンス番号をローカル側の大型記憶システム内に設けられるメモリにSVOL書込みシーケンス番号として記憶しておき、ローカル側から問合せがあった場合に、このメモリ内のシーケンス番号をSVOL書込みシーケンス番号として報告すれば良い。そして、リモート側がSVOL書込みシーケンス番号として「5」を報告すると、ミドルウェア3aは、(PVOLシーケンス番号「5」) <= (SVOL書込みシーケンス番号「5」)が成立したことから、データ同期が確立したと判断する。

【0045】

なお、上記において、PVOLシーケンス番号、SVOL書込みシーケンス番号をFIFOから求めるのではなく、ローカル側およびリモート側のRAM等の記憶手段を別途設け、2つの番号を随時これらに記憶するようにしておき、必要な時は、これを読むようにしても良い。

【0046】

また、上記説明では、ホスト1a内のAPP2aはOS4aからは1つのプロセスであるが、複数のプロセスとしてAPP2aは存在することもできる。この場合、この複数のプロセスからそれぞれSyncコマンドが発行されるが、それぞれのSyncコマンド受領時点でのPVOL最新シーケンス番号をそれぞれのSyncコマンド対応に取得し、これをSVOL最新シーケンス番号と比較することで、複数のプロセス単位でも同期確認することが可能である。

【0047】

また、ホスト1a内のアプリケーションソフト単位にグループ化して、このグループ単位に同期確認をとることも可能である。この場合、このグループ単位に非同期転送のシーケンス番号を維持する。大型ディスクアレイ装置内に複数のグループを定義しておき、アプリケーションソフトが同期確認を行う場合、そのアプリケーションソフトが同期完了要求にグループ名を指定することでグループ単位に独立して同期確認が可能になる。

【0048】

図2は、図1を用いて説明した事項の全体制御フローを示す図である。以下制御フローに基づき詳細に説明する。ここで、図2の説明で用いるSVOLシーケンス番号は、図1の説明のSVOL書込みシーケンス番号と同意味で使用する。

【0049】

まず、制御フローは、ホスト1aのアプリケーションソフトであるAPP2aと、ペア制御及びSyncコマンドを実行するミドルソフト3aと、大型ディスクアレイ装置13aホスト1aと物理ディスク9a間でデータの転送制御を行うDKC130a(図3参照)の制御を示す。

【0050】

APP2aは、コミットポイントになると、コミットコマンドを発行し、ホストコンピュ

10

20

30

40

50

ータ1 aのデータバッファ内に記憶していたデータを大型ディスクアレイ装置1 3 aへ書込む(図2には示していない)。

【0051】

次に、APP 2 aは、データベースへのコミット完了直後にSyncコマンドを発行する。ここで、Syncコマンドは、Syncコマンドであることを示すコマンドコード(図示せず)の他、2つの引数を有する。第1引数であるgroupは、前述したグループ名を指定する。第2引数であるtimeoutは、最大の同期完了待ち時間を指定する。APP 2 aは、Syncコマンドを発行した後、同期完了の報告をミドルソフト3 aから受けるか、若しくは同期失敗の報告を受けるまで待つ。

【0052】

ミドルソフト3 aは、Sync(group、timeout)コマンド3 1を実行する。Sync(group、timeout)コマンド3 1は、このコマンド内でまずローカル側の大型ディスクアレイ装置1 3 a内のPVOL(正ボリューム)1 0 aのペア状態を調べるために、PVOL状態取得コマンド2 0 1を大型ディスクアレイ装置1 3 aに発行する(PVOL状態取得ステップ3 2)。大型ディスクアレイ装置1 3 aの制御部分であるディスク制御部(DKC)1 3 0 aは、このコマンド応答としてPVOLペア状態2 0 2を返す(PVOL状態応答ステップ3 9)。なお、groupが指定されていた場合は、このgroup単位にPVOLのシーケンス番号を管理することになる。以下、groupがある特定グループに限定されていると考え、groupの違いを省いて説明する。

【0053】

ミドルソフト3 aは、PVOLペア状態2 0 2を大型ディスクアレイ装置1 3 aから受けると、PVOL状態のチェックを行い、状態がPAIR以外(P S U S , P S U E)であれば二重化はサスペンドであるとして同期失敗をAPP 2 aへ返す(PVOL状態のチェックステップ3 3)。なお、大型ディスクアレイ装置はこれらの状態を管理する機能を有するものとする。

【0054】

一方、PVOL状態のチェックステップ3 3は、状態がPAIRであれば二重化状態と判断し、FIFO7にキューされている書き込みデータの最新のPVOLシーケンス番号を調べるため、PVOLシーケンス番号取得コマンド2 0 3を、大型ディスクアレイ装置1 3 aに発行する(PVOLシーケンス番号を取得するステップ3 4)。

【0055】

大型ディスクアレイ装置1 3 aは、このコマンド応答として、FIFO7上にキューされている最新のデータに付加されている最新のPVOLシーケンス番号をPVOL最新シーケンス番号2 0 4として、ミドルソフト3 aへ返す(PVOL最新シーケンス応答ステップ4 0)。

【0056】

本発明でのミドルソフト3 aは、この取得したPVOL(正ボリューム)シーケンス番号を同期確認の間保持し、後で取得するSVOL(副ボリューム)シーケンス番号との比較に使用する。

【0057】

次に、ミドルソフト3 aは、リモートサイトのSVOLシーケンス番号を入手するために、SVOLシーケンス番号取得コマンド2 0 5を大型ディスクアレイ装置1 3 aに発行する(SVOLシーケンス番号を取得するステップ3 5)。大型ディスクアレイ装置1 3 aは、リモートサイトの大型ディスクアレイ装置1 3 bからSVOL1 0 bに書き込まれたデータに対応する最新のシーケンス番号を取得し、このSVOLシーケンス番号取得コマンドの応答として、最新のSVOL書き込みシーケンス番号2 0 6をミドルソフト3 aへ返す(SVOL書き込みシーケンス応答ステップ4 1)。

【0058】

なお、リモートサイトの大型ディスクアレイ装置1 3 bからSVOL1 0 bに書き込まれたデータに対応する最新のシーケンス番号を取得するには、まず、ローカル側の大型ディ

10

20

30

40

50

スクアレイ装置 13 a が、リモート側の大型ディスク装置 13 b に対して、S V O L へ書込まれた最新のデータに関する最新の S V O L 書込みシーケンス番号を問い合わせるコマンドを発行する。本実施例では、このコマンドは、通常のコピーデータを転送する通信路を用いて、リモート I / F 制御部 18 a、18 b を介して問い合わせられるものとする。これを受信したリモートの大型ディスクアレイ装置は、この問い合わせコマンドを解析し、S V O L へ書込まれた最新のデータに関する最新の S V O L 書込みシーケンス番号を取得する。

【 0 0 5 9 】

ここで、最新の S V O L 書込みシーケンス番号は、大型ディスクアレイ装置 13 b が有する特定のメモリに記録され、リモート側がコピーデータを確実に保証できる形で保持した段階で、更新されているものとする。本実施例では、先に説明したように、リモート側のデータ受信順序がホストコンピュータが発行したデータ順序と必ずしも一致していない場合を考慮し、F I F O 8 を使用する場合を示している。そのため、リモート側がコピーデータを確実に保証できる形とは、本実施例では、受信データを、シーケンス番号に抜けが無い状態で、F I F O 8 から S V O L に対応するキャッシュメモリに管理を移した状態である（なお、本システムは、キャッシュに記憶されたデータは物理ディスクに書き込むことを保証するものとする）。そして、この状態のデータに付与されたシーケンス番号が S V O L 書込みシーケンス番号となる。例えば、F I F O 8 内にシーケンス番号「2」、「3」、「5」を有する受信データがあるとすると、シーケンス番号「2」「3」を有する受信データは S V O L に対応するキャッシュへ書込まれるがシーケンス番号「5」を有する受信データは S V O L 対応のキャッシュへは書込まれない。この段階での S V O L 書込みシーケンス番号は「3」である。その後、シーケンス番号「4」を有するデータを受信すると、シーケンス番号「4」および「5」の受信データが順次 S V O L 対応のキャッシュに格納される。ここで初めて、リモート側がコピーデータを確実に保証できる形で保持した段階の最新の S V O L 書込みシーケンス番号は、「5」ということになる。一方、ローカル側の記憶システムや、ローカル側とリモート側間の通信路が、リモート側での受信データの順序性を保証できるものであれば、リモート側の記憶システムは受信データを受信した時点で、例えば、リモート I / F 制御回路がデータを受信した時点で、その受信したデータに付与されているシーケンス番号を、本実施例でいう最新の S V O L 書込みシーケンス番号として用いることができる。

【 0 0 6 0 】

リモート側のディスクアレイ装置 13 b は、取得した最新の S V O L 書込みシーケンス番号を、前記問い合わせコマンドの応答として、ローカル側の大型ディスクアレイ装置 13 a へ転送する。

【 0 0 6 1 】

また、リモート側の最新のシーケンス番号を取得する別の方法として、リモート側の大型ディスクアレイ自体が、定期的に、又は、S V O L へのデータ書込み更新があった際に、最新の S V O L 書込みシーケンス番号をローカル側へ通知し、これをローカル側で記録するようにし、S y n c コマンドが発生した時点で、ローカル側に記録したこの最新の S V O L シーケンス番号を読み出して用いるようにしても良い。

【 0 0 6 2 】

ミドルソフト 3 a は、次のステップ 36 で、保持していた P V O L (最新) シーケンス番号と S V O L (書込み) シーケンス番号を比較し、
P V O L シーケンス番号 < = S V O L 書込みシーケンス番号
であれば当該 P V O L シーケンス番号に対応するデータは S V O L 側に書き込み済みであるとして、A P P 2 a に同期完了を返す (P V O L と S V O L シーケンス比較ステップ 36)。

【 0 0 6 3 】

一方、

P V O L シーケンス番号 > S V O L シーケンス番号

10

20

30

40

50

であれば同期が完了していない（リモート側へデータが書き込まれていない）として、次のステップに進みタイムアウトのチェックを行い指定 `time out` 値を超えていれば同期完了タイムアウトとして同期失敗を `APP 2 a` へ返す（タイムアウトのチェックステップ 37）。ここで、指定 `time out` 値を超えていなければ一定時間待ち（`WAIT` ステップ 38）、その後、ステップ 35 からステップ 38 を同期が完了するまで繰り返す。繰り返しの中で、ステップ 36 で “ `PVOL` シーケンス番号 \leq `SVOL` シーケンス番号 ” の条件が成立した時点で、同期完了し、`Sync` コマンドを発行した呼出し元である `APP 2` に制御が戻る。`APP 2 a` は、`Sync` コマンドの制御が戻った時点で、戻り値をチェックし同期完了を確認する。

【0064】

以上説明したように、ホスト 1 a の `APP 2 a` は、ミドルソフト 3 a と連携して、データベース等のトランザクションのコミットを行った直後に本願で示した `Sync` コマンドを実行するので、本 `Sync` コマンド発行前にディスクアレイ装置へ書込まれたデータが、リモート側のディスクアレイ装置にデータが確実に格納されているか否かを `APP 2` は、知ることが出来る。見方を変えれば、本実施例を使用するホストコンピュータ内のアプリケーションからは、非同期コピーが、`Sync` コマンドを発行することによってアプリケーション毎に同期処理（データがコンピュータ内とリモートの記憶装置内で一致する処理）ができるリモートコピーとして見えることになる。

【0065】

（他の実施例）

図 4 および図 5 は、ホスト装置 1 a のミドルソフト制御の負荷を軽減する目的で、ミドルソフト 3 a の制御をローカル側の大型ディスクアレイ装置 13 a 内で実施した本発明の他の実施例を示したものである。

【0066】

図 4 は、図 1、図 2 のところで説明したミドルソフト制御を、大型ディスクアレイ装置内で実施した場合の、本発明の論理的なブロック図を示す。図 4 において、図 1 のミドルソフト 3 a で実現していた機能を図 4 の `Sync` 制御 300（後述）で置換えた以外は、論理的構成はほぼ同じである。`Sync` 制御 300 の機能は、“ `HOST I/F` 制御回路（17 a）” または “ `リモート I/F` 制御回路（18 a）” 上で本機能を実現するプログラムを実行することで実現される。

【0067】

ホスト 1 a 内の `APP 2 a` は、前記実施例と同様にトランザクションとしてデータベースへ `WRITE` し、最後に `Commit` を発行し、その後、本発明の特徴コマンドである `Sync` コマンドを発行する。この `Sync` コマンドは、OS を介してディスクアレイ装置 13 a に入力される。なお、ここでは、`Commit` 動作を確認するという意味で `Sync` コマンドを `Commit` の発行の後に入れているが、`Sync` コマンドを発行するタイミングは、この時点に限定されるものでなく、アプリケーションプログラムが必要と認識した時点で発行できる。

また、以降説明するディスクアレイ装置側での `Sync` コマンドの処理も、当該 `Sync` コマンドの発行時点に左右されるものではない。

【0068】

この `Sync` コマンドは、大型ディスクアレイ装置によって、これに接続されるホストに対して提供され、大型ディスクアレイ装置に内蔵されたマイクロプログラム等の機能（後述の `Sync` 制御 300 が相当する）である。例えば、図 1 の `HOST I/F` 制御回路内のマイクロプロセッサ等が、`リモート I/F` 制御回路等と連絡しながら実行する。

【0069】

ここで、`Sync` (`appID`、`group`、`timeout`) コマンド 51 は、パラメータとして、ローカル側の大型ディスクアレイ装置 13 a の `Sync` 制御 300 が `Sync` 要求プロセスを識別するための `appID` と、前記の `group`（前述したグループ名）と、および、最大の同期完了待ち時間を指定する `Timeout` を有する。

10

20

30

40

50

【 0 0 7 0 】

この Sync コマンド 5 1 は、OS 4 a を介し大型ディスクアレイ装置 1 3 a の Sync 制御 3 0 0 に渡される。この Sync 制御 3 0 0 は Sync コマンドを受けると、前記実施例と同様にボリュームのペア状態をチェックし、PAIR 状態であれば、FIFO 7 の “最新の PVOL シーケンス番号” と Sync コマンドパラメータとして渡される “app ID” をディスクアレイ装置内のメモリに登録し、この PVOL 最新シーケンス番号を Commit の直後のシーケンス番号として保持する。

【 0 0 7 1 】

次に Sync 制御 3 0 0 は、リモート側の SVOL 書込みシーケンス番号を取得し、“PVOL シーケンス番号”

10

< = “SVOL 書込みシーケンス番号”

が成立するまで SVOL 書込みシーケンス番号を繰返しテストし、この Sync 制御 3 0 0 内で待つ。条件が成立するとこの Sync 制御 3 0 0 は、同期完了の応答として、OS 4 a を介し “app ID” を、呼び出しもとである APP 2 a に戻す。APP 2 a は、この “app ID” を受けると、Sync コマンドが完了したことで同期完了とみなす。なお、PVOL シーケンス番号および SVOL 書込みシーケンス番号の具体的な取得方法は、先の実施例と同様である。

【 0 0 7 2 】

上記では、ホスト 1 a 内の APP 2 a は、OS 4 a からは 1 つのプロセスとして説明したが、複数のプロセスとして APP 2 a は存在することもできる。この場合、複数のプロセスから Sync コマンドが発行されると、Sync 制御 3 0 0 は、app ID で識別し、プロセス対応のそれぞれの Sync コマンド受領時点での app ID と PVOL 最新シーケンス番号を登録して、後で得る SVOL シーケンス番号と比較するので複数のプロセス単位で同期確認することが可能である。結果は、登録した app ID を有するアプリケーションへ報告される。

20

【 0 0 7 3 】

また、ホスト 1 a 内のアプリケーションソフトを複数合わせてグループ化して、このグループ単位に同期確認をとることも可能である。この場合、このグループ単位に非同期転送のシーケンス番号を維持する。大型ディスクアレイ装置内に複数のグループを定義しおき、アプリケーションソフトが同期確認を行う場合、そのアプリケーションソフトが同期完了要求にグループ名を指定することでグループ単位に独立して同期確認が可能になる。

30

【 0 0 7 4 】

図 5 は、図 4 で説明した事項の全体制御フローを示す図である。以下制御フローに基づき詳細に説明する。ここで、図 5 の説明で用いる SVOL シーケンス番号は、図 4 の説明の SVOL 書込みシーケンス番号と同意味で使用する。制御フローは、ホスト 1 a のアプリケーションソフトである APP 2 a と、ペア制御及び Sync コマンドを実行するローカル側の大型ディスクアレイ装置 1 3 a 内の Sync 制御 3 0 0 と、及び、リモート側大型ディスクアレイ装置 1 3 b のディスク制御 DKC 1 3 0 b の制御を示す。

【 0 0 7 5 】

APP 2 a は、データベースへのコミット完了直後に Sync (app ID、group、timeout) コマンド 5 1 を発行する。Sync コマンド 5 1 の第 1 引数である app ID は、大型ディスクアレイ装置 1 3 a の Sync 制御 3 0 0 が、ホスト 1 a 内の Sync 要求プロセス (アプリケーション) を識別するための引数である。これにより、Sync 要求を発行するプロセス (アプリケーション) 毎に同期確定を行うことが可能となる。第 2 引数である group は、前述したグループ名を指定する場合に用いる。第 3 引数である timeout は、最大の同期完了待ち時間を指定する。

40

【 0 0 7 6 】

大型ディスクアレイ装置 1 3 a の Sync 制御 3 0 0 は、Sync (app ID、group、timeout) コマンド 5 1 を実行する。まず Sync 制御 3 0 0 は、大型ディスクアレイ装置 1 3 a 内の PVOL 9 のペア状態を調べる (PVOL 状態のチェックステ

50

ップ533)。

【0077】

S y n c制御300は、P V O L状態のチェック33によって状態がP A I R以外(P S U S、P S U E)であれば二重化はサスペンドであるとして同期失敗をA P P 2 aに返す。状態がP A I Rであれば二重化状態であり、書き込みデータはF I F O 7にキューされているので、最新のP V O Lシーケンス番号とa p p I Dをディスクアレイ装置内のメモリに登録する(a p p I DとP V O Lシーケンス番号を登録するステップ534)。この登録したP V O Lシーケンス番号とa p p I Dは、同期確認の間保持され、後で、S V O Lシーケンス番号との比較に使用される。

【0078】

次に、S y n c制御300は、P V O Lシーケンス番号とリモートサイトのS V O Lシーケンス番号の比較のために、S V O Lシーケンス番号取得コマンド501をリモート側の大型ディスクアレイ装置13bのディスク制御部130b(図3のリモート側のディスクアレイ装置参照)に発行する(S V O Lシーケンス番号を取得するステップ535)。

【0079】

リモート側の大型ディスクアレイ装置13bは、このコマンド501の応答としてS V O L書き込みシーケンス応答として最新のS V O L書き込みシーケンス番号502を、ローカル側のディスクアレイ装置13aへ返す(S V O L書き込みシーケンス応答ステップ541)。なお、S V O L書き込みシーケンス番号の具体的な取得方法は先の実施例と同様である。

【0080】

ローカル側のディスクアレイ装置13aのS y n c制御300は、次のステップ36で、保持していたa p p I Dに対応するP V O Lシーケンス番号とS V O Lシーケンス番号を比較し、

“ P V O Lシーケンス番号 ” < = “ S V O Lシーケンス番号 ”

であれば当該P V O Lシーケンス番号はS V O L側に書き込み済みであるとして、該当a p p I Dと合せて同期完了を、アプリケーションソフトA P P 2 aを返す。一方、

“ P V O Lシーケンス番号 ” > “ S V O Lシーケンス番号 ”

であれば同期が完了していないので次の待ちステップ537に進む(a p p I DのP V O LとS V O Lシーケンスを比較するステップ536)。

【0081】

次のタイムアウトをチェックするステップ537では、タイムアウトのチェック537を行い、指定t i m e o u t値を超えていれば同期完了タイムアウトとして同期失敗をA P P 2 aへ返す。一方、指定t i m e o u t値を超えていなければ一定時間W A I T 538する(W A I Tステップ538)。

【0082】

その後、ステップ535から同期が完了するまでステップ536を繰り返し、

“ P V O Lシーケンス番号 ” < = “ S V O Lシーケンス番号 ”

の条件が成立した時点で、同期が完了したとして、S y n cコマンド51の呼び出し元であるA P P 2 aに、S y n c制御から、制御を返す。A P P 2 aは、S y n cコマンドの制御が戻った時点で、a p p I Dをチェックすると共に戻り値(同期が完了したか否かを示す値)をチェックし同期完了を確認する。

【0083】

以上説明したようにディスクアレイ装置が、S y n cコマンドを受信すると、ホストからの書き込み要求のデータが確実にリモート側のディスクアレイ装置にコピーされたか否かを確認した上で、ホストに対し報告することが可能となる。従って、ホストはディスクアレイ装置が非同期リモートコピーを行っている場合でも、S y n cコマンドを発行することによって、S y n cコマンドの発行前にディスクアレイ装置へ書込まれているデータがリモート側へ転送されているかを正確に知ることができ、データベース等で必要となるC o m m i t制御を確実なものとするができる。見方を変えれば、本実施例を使用するホストコンピュータ内のアプリケーションからは、非同期コピーが、S y n cコマンドを発

10

20

30

40

50

行することによってアプリケーション毎に同期処理（データがコンピュータ内とリモートの記憶装置内で一致する処理）ができるリモートコピーとして見えることになる。

【0084】

上記実施例では、コミットコマンドの発行に続けて Sync コマンドの発行という形で発明を説明したが、次のような変形も可能である。

【0085】

一つ目は、最初の実施例においてミドルソフト 3 a を機能拡張し、これにより、非同期データ転送を行っている場合は、アプリケーションから受けた従来のコミットコマンドを「従来のコミットコマンド + Sync コマンド」と認識して、処理する方法である。具体的には、機能拡張されたミドルソフト 3 a は、コミットコマンドを受けると、まず、コンピュータ内のデータバッファ内にあるデータを記憶システムへ書込む。それと合わせて、コミットの対象となるデータが記憶されるボリュームが非同期コピー中か否かを判断する。対象ボリュームが非同期コピー中の場合は、ミドルソフト 3 a は本発明の最初の実施例に対応した問合せコマンドをローカル側の記憶システムへ発行する。その後、ミドルソフト 3 a は、ローカル側の記憶システムから書込み完了を受け、この時点ではアプリケーションに対してコミットコマンドに対するリターンを返さず、先の問合せコマンドに対する応答が有るのを待つ。そして、ミドルソフト 3 a は、記憶システムから問合せコマンドに対する応答を受け、リモート側の記憶システムで書込みデータの確定が取れた時点で初めて、前記アプリケーションへ、コミットコマンドに対するリターンを返すようにする。この方法では、アプリケーションは Sync コマンドを発行する必要が無いので、非同期コピー中であることを意識しなくてすむというメリットがある。

【0086】

2 つ目は、ローカル側の記憶システムが非同期リモートコピーを行っている際、書込みコマンドを発行するホストコンピュータに対して、リモート側の記憶システムが書込みコマンドのデータをどの時点まで記憶したかを、定期的に、報告する方法である。例えば、ローカル側の記憶システムは、定期的に、ホストからのある書込み要求が発行された時刻を記憶しておき、この書込み要求以前に発行された全ての書込み要求に対応するデータがリモート側の記憶システムへ確実に記憶された段階で、先に記録した時刻をホストコンピュータへ返すようにする。

このようにすれば、ホストコンピュータは定期的にどの時刻までの書込み要求が確実にリモート側の記憶システムへ記憶されたかが分かるようになる。一方、ホスト側のアプリケーションは、従来と同様にコミットコマンドを発行する。このコミットコマンドを受けた時点でミドルソフトは、ローカル側の記憶システムへデータバッファ内のデータを書込むため、書込み要求を発行するが、この際、書込み要求を発行した最終時刻を記憶するようにする。そして、ミドルソフトは、この書込み要求に対する全ての書込み完了を受け、かつ、当該最終時刻に等しいか又は遅い、先に述べたローカル側の記憶システムが定期的に報告する時刻を受けた場合に、アプリケーションに対し、先のコミットコマンドに対するリターンを返すようにする。この 2 つ目の方法では、ホスト側のアプリケーションは Sync コマンドを発行する必要が無いので、非同期コピー中であることを意識しなくてすむというメリットがある。また、記憶システム側は自発的にリモート側の状況をホストコンピュータに通知するだけなので、ホストコンピュータと記憶システム間のインタフェースが簡単なものとなる。

【0087】

なお、この 2 つ目の方法において、ローカル側の記憶システムが定期的にどの時刻までの書込み要求をリモート側の記憶システムに確実に書込んだかをホストコンピュータに報告する例をあげている。このような機能を有する記憶システムは、本実施例で説明するようなアプリケーションのためのデータ同期確認にももちろん利用可能であるが、これ以外にも、ホストコンピュータ側から記憶システム状態を知る上で重要な情報を提供しているわけであり、本記憶システムはホスト側から管理しやすい、又は、制御しやすい、又は、使いやすい記憶システムということができる。

10

20

30

40

50

【 0 0 8 8 】

【 発明の効果 】

本発明では、非同期リモートコピーを行っている際、ホストに対し、書込み要求のデータがリモート側に記憶されたか否かを判断するのに必要なリモート側の記憶システムに関する情報をローカル側の記憶システム経由でホストに提供するので、アプリケーションは、ホスト内データとリモート側記憶システム内のデータとの一致を確認できる。

【 0 0 8 9 】

また、本発明では、非同期リモートコピーを行う記憶システムが、ホストに対し、書込み要求のデータがリモート側に記憶されたか否かを通知するので、ホスト内のアプリケーションは、ホスト内データとリモート側記憶システム内のデータとの一致を確認できる。

10

【 図面の簡単な説明 】

【 図 1 】 本発明の論理的なブロック図を示す。

【 図 2 】 本発明の全体制御フローを示す。

【 図 3 】 本発明を適用する一般的な H A クラスタの全体構成図である。

【 図 4 】 本発明の第 2 の実施例の論理的なブロック図を示す。

【 図 5 】 本発明の第 2 の実施例に相当する全体フローを示す。

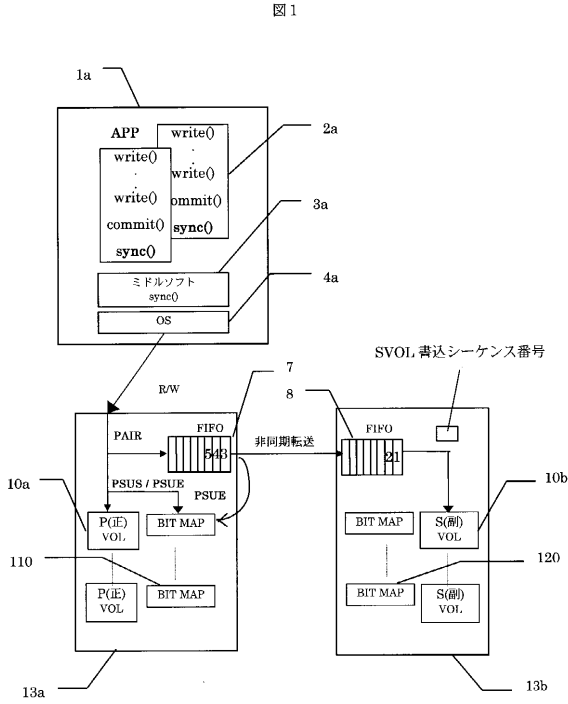
【 符号の説明 】

- 1 a、1 b ホスト装置
- 2 a、2 b アプリケーションソフト (A P P)
- 3 a、3 b ミドルソフト
- 3 0 0 S y n c 制御
- 4 a、4 b オペレーティングシステム (O S)
- 5 P V O L 最新シーケンス番号
- 6 S V O L 書込みシーケンス番号
- 7 F I F O (P V O L 対応)
- 8 F I F O (S V O L 対応)
- 1 0 a P (正) V O L
- 1 0 b S (副) V O L
- 1 1 0 B I T M A P (P V O L 対応)
- 1 2 0 B I T M A P (S V O L 対応)
- 1 3 a ローカル側大型ディスクアレイ装置
- 1 3 b リモート側大型ディスクアレイ装置
- 3 1 実施例の S y n c コマンド
- 5 1 他の実施例の S y n c コマンド

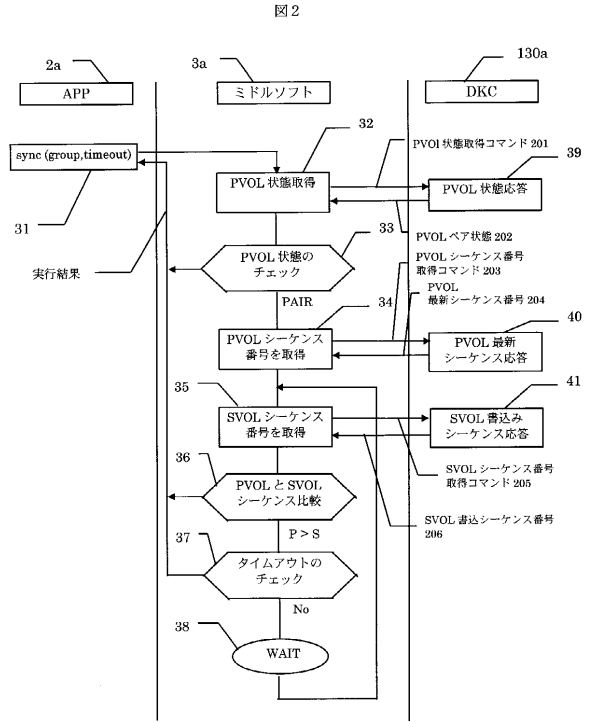
20

30

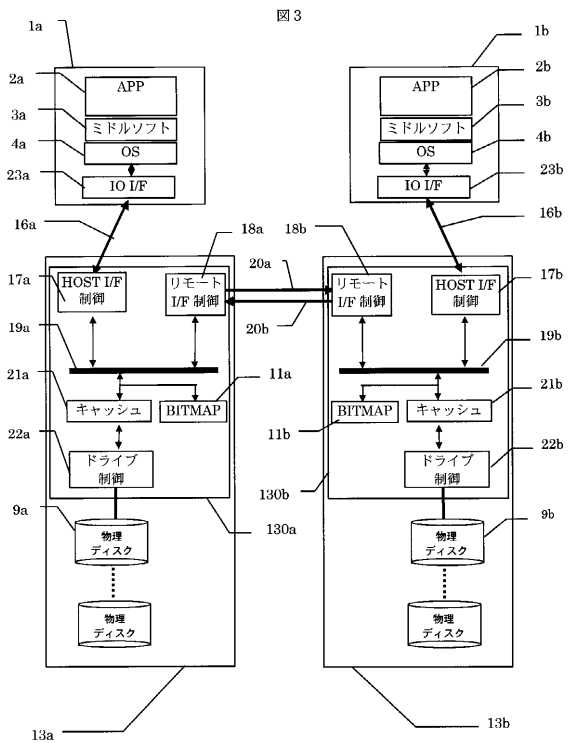
【図1】



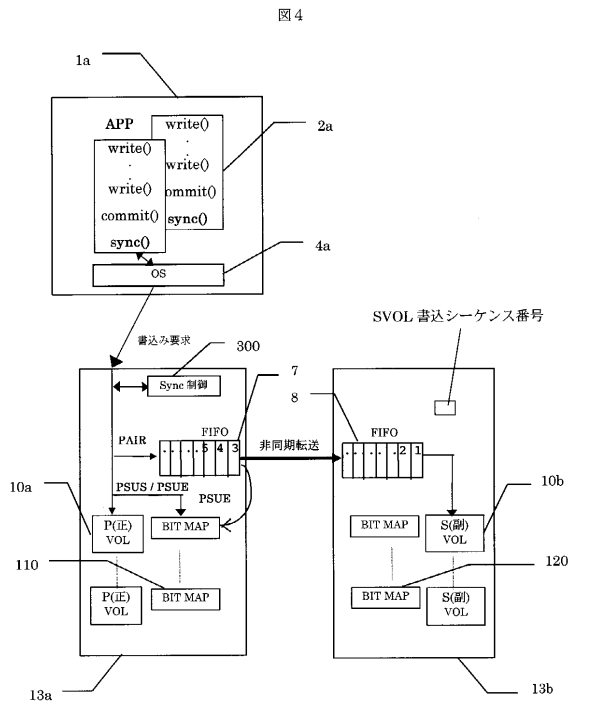
【図2】



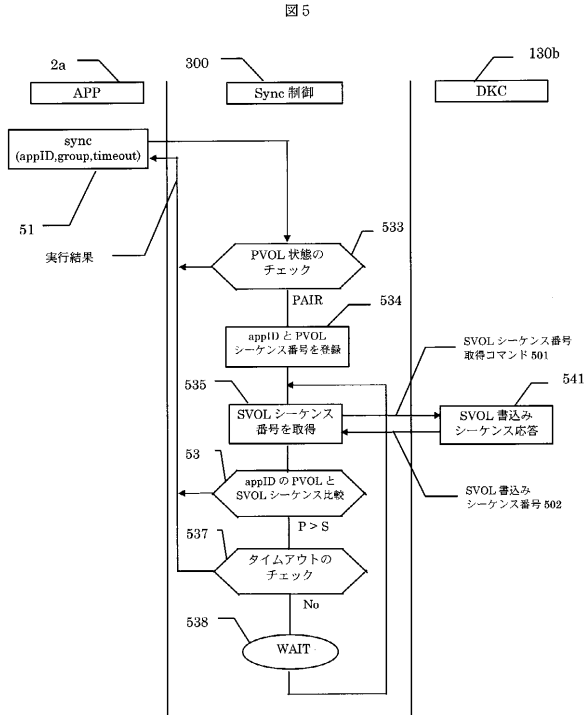
【図3】



【図4】



【 図 5 】



フロントページの続き

(51)Int.Cl.

F I

G 0 6 F 12/08 5 4 1

G 0 6 F 12/08 5 5 7

G 0 6 F 12/16 3 1 0 J

(56)参考文献 特開平 0 7 - 2 3 4 8 1 1 (J P , A)
特開平 1 0 - 0 7 4 1 5 7 (J P , A)
特開平 0 8 - 2 6 3 3 5 9 (J P , A)
特表平 0 8 - 5 0 9 5 6 5 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)

G06F 12/00

G06F 13/00

G06F 3/06

G06F 12/08

G06F 12/16