

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第4473612号
(P4473612)

(45) 発行日 平成22年6月2日 (2010.6.2)

(24) 登録日 平成22年3月12日 (2010.3.12)

(51) Int.Cl.

F I

G O 6 F 12/00 (2006.01)

G O 6 F 3/06 (2006.01)

G O 6 F 12/00 5 3 1 D

G O 6 F 12/00 5 3 1 J

G O 6 F 3/06 3 0 4 F

G O 6 F 3/06 3 0 4 P

請求項の数 16 外国語出願 (全 20 頁)

(21) 出願番号	特願2004-81902 (P2004-81902)	(73) 特許権者	000005108
(22) 出願日	平成16年3月22日 (2004.3.22)		株式会社日立製作所
(65) 公開番号	特開2005-18736 (P2005-18736A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成17年1月20日 (2005.1.20)	(74) 代理人	100093861
審査請求日	平成19年2月26日 (2007.2.26)		弁理士 大賀 真司
(31) 優先権主張番号	10/602223	(72) 発明者	山神 憲司
(32) 優先日	平成15年6月23日 (2003.6.23)		アメリカ合衆国カリフォルニア州ロスガト
(33) 優先権主張国	米国 (US)		ス カイルニベル 1 0 8
		審査官	桜井 茂行

最終頁に続く

(54) 【発明の名称】 リモートコピーシステム

(57) 【特許請求の範囲】

【請求項 1】

複数の第1のデータボリュームと、関連する第1のホストから書き込みリクエストを受け取って、前記書き込みリクエストに関連する書き込みデータを前記複数の第1のデータボリュームに記憶するように構成されている第1のストレージコントローラと、を有する第1のストレージシステムと、

第2のデータボリュームと、前記第1のデータボリュームを前記第2のデータボリュームにミラーリングするための前記書き込みデータを含む第1のデータを受信する第2のストレージコントローラと、ジャーナルボリュームと、を有する第2のストレージシステムと、

複数の第3のデータボリュームと、前記複数の第1のデータボリュームを前記複数の第3のデータボリュームに各々ミラーリングするジャーナルを含む第2のデータを受信する第3のストレージコントローラと、を有する第3のストレージシステムと、

を備え、
前記ジャーナルボリュームは、前記複数の第1のデータボリューム及び複数の第3のデータボリュームと対応しており、

前記ジャーナルは、前記書き込みデータと、前記第1のデータボリュームに書き込み順序を提供するシーケンス番号を有し、

前記第1のホストから前記第1のストレージシステムへの書き込みリクエストは、前記第1のデータが前記第2のストレージシステムによって受信された後、完了し、

前記第 1 のホストから前記第 1 のストレージシステムへの書き込みリクエストは、前記第 3 のストレージシステムによって受信された第 2 のデータとは独立して完了し、

前記第 3 のデータボリュームに記憶される書き込みデータは、前記ジャーナルのシーケンス番号によって提供される書き込み順序にしたがって生成され、

前記ジャーナルは、前記ジャーナルボリュームに先入れ先出しで積み重ねられ、

当該ジャーナルは、更新の順序情報及び前記第 1 のデータボリュームの識別情報を含み、当該更新の順序情報及び前記第 1 のデータボリュームの識別情報に基づいて前記第 3 のデータボリュームに非同期で送信されることを特徴とするリモートコピーシステム。

【請求項 2】

前記第 2 のストレージシステムは、前記第 1 のストレージシステムの比較的近くに位置し、前記第 3 のストレージシステムは、前記第 1 のストレージシステムの比較的遠くに位置していることを特徴とする請求項 1 に記載のリモートコピーシステム。

【請求項 3】

前記第 2 のストレージシステムは、前記第 1 のストレージシステムから 100 マイル以内に位置し、前記第 3 のストレージシステムは、前記第 1 のストレージシステムから 100 マイルを超えて位置していることを特徴とする請求項 2 に記載のリモートコピーシステム。

【請求項 4】

前記第 2 のストレージシステムに結合される第 2 のホストをさらに備え、

前記第 2 のストレージシステムは、もし前記第 1 のストレージシステムに障害が発生した場合には、プライマリストレージシステムとして機能するように構成されていることを特徴とする請求項 1 に記載のリモートコピーシステム。

【請求項 5】

前記第 3 のストレージシステムに結合される第 3 のホストをさらに備え、

前記第 3 のストレージシステムは、もし前記第 1 のストレージシステムに障害が発生した場合には、前記第 1 のストレージシステムに代ってプライマリストレージシステムとして機能するように構成されていることを特徴とする請求項 1 に記載のリモートコピーシステム。

【請求項 6】

前記ジャーナルは、前記第 1 のホストからの書き込みリクエストに基づいて前記書き込みデータが前記第 1 のデータボリュームに記憶された時刻に関する情報をさらに含むことを特徴とする請求項 1 に記載のリモートコピーシステム。

【請求項 7】

前記第 1 のストレージシステムに障害が発生した場合、前記ジャーナルは、前記書き込みデータを確保するため、前記書き込みデータが前記第 2 のストレージシステムにコピーされた後、前記第 3 のデータボリュームに受信されることを特徴とする請求項 1 に記載のリモートコピーシステム。

【請求項 8】

前記ジャーナルボリュームが、コントロールデータ領域と、ジャーナルデータ領域と、を含み、当該コントロールデータ領域は、コントロールデータのみを記憶する構成であり、当該ジャーナルデータ領域は、ジャーナルデータのみを記憶する構成であることを特徴とする請求項 1 に記載のリモートコピーシステム。

【請求項 9】

複数の第 1 のデータボリュームと、関連する第 1 のホストから書き込みリクエストを受け取って、前記書き込みリクエストに関連する書き込みデータを前記複数の第 1 のデータボリュームに記憶するように構成されている第 1 のストレージコントローラと、を有する第 1 のストレージシステムと、

前記第 1 のストレージシステムに結合され、且つ、第 2 のデータボリュームと、前記第 1 のデータボリュームに記憶されたデータを前記第 2 のデータボリュームにコピーするための書き込みデータを含む第 1 のデータを同期して受信するよう構成された第 2 のストレ

10

20

30

40

50

ージコントローラと、ジャーナルボリュームと、を有する第2のストレージシステムと、
前記第2のストレージシステムに結合され、且つ、複数の第3のデータボリュームと、
前記ジャーナルボリュームからの第2のデータを非同期で受信するよう構成された第3の
ストレージコントローラと、を有する第3のストレージシステムと、

を備え、

前記ジャーナルボリュームは、前記複数の第1のデータボリューム及び複数の第3のデ
ータボリュームと対応しており、

前記第2のデータは、前記第1のデータボリュームに記憶されたデータを、前記第2の
ストレージシステムを介して、前記第3のデータボリュームにコピーするジャーナルを有
し、

前記ジャーナルは、前記書き込みデータと、前記第1のデータボリュームに書き込み順
序を提供するシーケンス番号を有し、

前記第3のデータボリュームに記憶される書き込みデータは、前記シーケンス番号によ
って提供される書き込み順序にしたがって生成され、

前記ジャーナルは、前記ジャーナルボリュームに先入れ先出しで積み重ねられ、

当該ジャーナルは、更新の順序情報及び前記第1のデータボリュームの識別情報を含み
、当該更新の順序情報及び前記第1のデータボリュームの識別情報に基づいて前記第3の
データボリュームに非同期で送信されることを特徴とするリモートコピーシステム。

【請求項10】

前記第2のストレージシステムは、前記第1のストレージシステムの比較的近くに位置
し、前記第3のストレージシステムは、前記第1のストレージシステムの比較的遠くに位
置していることを特徴とする請求項9に記載のリモートコピーシステム。

【請求項11】

前記第2のストレージシステムは、前記第1のストレージシステムから100マイル以
内に位置し、前記第3のストレージシステムは、前記第1のストレージシステムから10
0マイルを超えて位置していることを特徴とする請求項10に記載のリモートコピーシ
ステム。

【請求項12】

前記第2のストレージシステムに結合される第2のホストをさらに備え、

前記第2のストレージシステムは、もし前記第1のストレージシステムに障害が発生し
た場合には、プライマリストレージシステムとして機能するように構成されていることを
特徴とする請求項9に記載のリモートコピーシステム。

【請求項13】

前記第3のストレージシステムに結合される第3のホストをさらに備え、

前記第3のストレージシステムは、もし前記第1のストレージシステムに障害が発生し
た場合には、前記第1のストレージシステムに代ってプライマリストレージシステムとし
て機能するように構成されていることを特徴とする請求項9に記載のリモートコピーシ
ステム。

【請求項14】

前記ジャーナルは、前記第1のホストからの書き込みリクエストに基づいて前記書き込
みデータが前記第1のデータボリュームに記憶された時刻に関する情報をさらに含むこと
を特徴とする請求項9に記載のリモートコピーシステム。

【請求項15】

前記第1のストレージシステムに障害が発生した場合、前記ジャーナルは、前記書き込
みデータを確保するため、前記書き込みデータが前記第2のストレージシステムにコピー
された後、前記第3のデータボリュームに受信されることを特徴とする請求項9に記載の
リモートコピーシステム。

【請求項16】

前記ジャーナルボリュームが、コントロールデータ領域と、ジャーナルデータ領域と、
を含み、当該コントロールデータ領域は、コントロールデータのみを記憶する構成であり

10

20

30

40

50

、当該ジャーナルデータ領域は、ジャーナルデータのみを記憶する構成であることを特徴とする請求項 9 に記載のリモートコピーシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はストレージシステムに係わり、さらに詳細にはリモートコピー機能を行うように構成されたストレージシステムに係わる。

【背景技術】

【0002】

データは全ての計算プロセスの基本となる基礎リソースである。近年のインターネットや e - ビジネスの爆発的な発展に伴い、データ記憶システムの需要は急速に増加している。一般に、ストレージネットワークは二つのアプリケーションや構成を包含するもので、ネットワーク接続ストレージ (NAS) またはストレージ・エリア・ネットワーク (SAN) などである。

【0003】

NAS はストレージサーバとそのクライアントの間でファイルフォーマットのデータを転送するためにイーサネット上で IP を用いる。NAS においては、ディスクアレイやテープアレイなどの総合ストレージシステムは、TCP/IP などのメッセージコミュニケーションプロトコルを用いて、ローカルエリアネットワーク (LAN) を通してメッセージネットワークに直接接続される。クライアントサーバシステムでは、ストレージシステムはサーバとして機能する。

【0004】

一般に、SAN は異種のサーバとストレージリソース間でデータを移動させる専用高性能ネットワークである。従来のメッセージング・ネットワークでは、NAS と異なり、クライアントとサーバ間のトラフィックの輻輳を避けるためにそれぞれに専用のネットワークが用いられる。SAN では、ストレージリソースとプロセッサあるいはサーバ間に直接の接続を設定することが可能である。SAN はサーバの間でシェアすることも、また特定のサーバに専用にすることも可能である。それは一カ所に集中させることも、地理的に離れた場所に拡張することも可能である。SAN インターフェイスは、ファイバーチャネル (FC) やエンタプライズシステムコネクション (ESCON) や小型コンピュータシステムインターフェイス (SCSI) やシリアルストレージアーキテクチャ (SSA) や高性能パラレルインターフェイス (HIPPI) あるいは将来出現するであろう他のプロトコルのような各種のインターフェイスによることが可能である。

【0005】

たとえば、インターネットエンジニアリングタスクフォース (IETF) は、TCP/IP 上でブロックストレージを可能にする新しいプロトコルあるいは標準 iSCSI を開発中であり、一方幾つかの会社は iSCSI を SAN の主たる標準にするためにホストプロセッサから iSCSI - TCP/IP プロトコルスタックを外す努力をしている。用いられているストレージシステムの形式がどうであろうとも、データストレージシステムのユーザは、ストレージユニット (あるいはストレージサブシステム) の障害で貴重なデータが失われることを防ぐためのデータバックアップに強い関心を持っている。

【0006】

したがって、データストレージシステムは一般に、プライマリユニットに障害が発生した時に、エマージェンシーリカバリーのためのデータを記憶するバックアップユニットを含んでいる。しかしながら、障害は、ユニット自身の故障による場合も、たとえばユニットが設置されている場所での地震とか台風による自然災害の発生による場合もある。もしバックアップユニットがプライマリユニットの近くに設置されていると、自然災害が発生した場合に両方とも破壊される可能性がある。したがって、多くのストレージシステムのユーザは、たとえば 100 マイル以上離れた遠い距離を離してプライマリユニットとバック

10

20

30

40

50

クアップユニットを設置する傾向がある。実際に、複数のユーザはプライマリユニットとバックアップユニットを別の大陸に設置することまでやっている。

【発明の開示】

【発明が解決しようとする課題】

【 0 0 0 7 】

現在、ストレージシステムでは、データをバックアップあるいはセカンダリサイトにコピーするために二つの動作モード、すなわち同期モードと非同期モードが用いられている。

【 0 0 0 8 】

同期モードでは、ホストからプライマリストレージシステムへの書き込みリクエストは、書き込みデータがセカンダリストレージシステムにコピーされ、それが確認された後でのみ完了する。したがって、セカンダリストレージシステムから確認を受け取るまで、ホストからの書き込みデータはプライマリシステムのキャッシュに保存されているので、このモードではセカンダリシステムでデータの消失がないことが保証される。

【 0 0 0 9 】

さらに、プライマリストレージシステムにおけるプライマリボリューム（P V O L）とセカンダリストレージシステムにおけるセカンダリボリューム（S V O L）は同一に保持されているので、もしP V O Lに障害が発生しても、速やかにS V O Lを用いてP V O Lを再生することができる。しかしながら、このモードの下では、プライマリストレージシステムとセカンダリストレージシステムを、たとえば100マイル以上離して設置することは出来ない。そのようなことがなければ、ストレージシステムは効率よくホストからの書き込みリクエストを実行できる。

【 0 0 1 0 】

非同期モードでは、ホストからプライマリストレージシステムへの書き込みリクエストは、プライマリシステムのみへ書き込みデータを記憶することで完了する。書き込みデータは、続いてセカンダリストレージシステムにコピーされる。すなわち、プライマリストレージシステムへデータを書き込むことは、セカンダリストレージシステムにデータをコピーすることとは独立した処理である。したがって、プライマリシステムとセカンダリシステムは、たとえば100マイルあるいはそれ以上ずっと離して別々に設置されてもよい。

【 0 0 1 1 】

しかしながら、P V O LとS V O Lは同一に維持されているので、もしプライマリシステムがダウンするとデータが失われる可能性がある。したがって、同期モードと非同期モードの利点を提供する、言い換えればデータ喪失がないことを保証しながらプライマリシステムとセカンダリシステムをずっと離して設置できる、データストレージシステムあるいはリモートコピーシステムを提供することが望まれている。

【課題を解決するための手段】

【 0 0 1 2 】

本発明の実施例は、リモートコピー機能を行うべく構成されたストレージシステムに係わる。一実施例においては、プライマリストレージシステムは、仲介ストレージシステムに書き込みデータをコントロールデータとともに同期して送る。仲介ストレージシステムは書き込みデータとコントロールデータを、たとえばジャーナルボリューム（J N L）などのボリュームに記憶する。仲介ストレージシステムはデータの到着の順序を維持し、コントロールデータの各々に順序情報を割り当てあるいは連結させる。

【 0 0 1 3 】

つづいて、仲介ストレージシステムは書き込みデータとその対応するコントロールデータを非同期で、あるいはプライマリストレージシステムにおける書き込みデータとコントロールデータの記憶とは関連無しにセカンダリストレージシステムに送る。セカンダリストレージシステムは、コントロールデータとコントロールデータに関連した順序情報に従い、書き込みデータをセカンダリボリューム（S V O L）に記憶する。ここに用いられて

10

20

30

40

50

いる、「ストレージシステム」という用語は、データを記憶するように構成され、たとえばディスクアレイユニットなどの１台ないしはそれ以上のストレージユニットあるいはストレージサブシステムを持つコンピュータシステムを意味する。

【 0 0 1 4 】

したがって、ストレージシステムは、１台ないしはそれ以上のホストと、１ないしはそれ以上のストレージサブシステムを含むコンピュータシステム、あるいは単にストレージサブシステムまたはユニット、あるいはコミュニケーションリンクを介して相互に結合した複数のストレージサブシステムまたはユニットを意味することもある。ここに用いられている、「ストレージサブシステム」という用語は、データを記憶するように構成され、ストレージ領域と１台ないしはそれ以上のホストからのリクエストを処理するストレージコントローラを含むコンピュータシステムを意味する。ストレージサブシステムの例はディスクアレイユニットである。

10

【 0 0 1 5 】

ここに用いられている、「ホスト」という用語は、１台ないしはそれ以上のストレージシステムあるいはストレージサブシステムに結合され、ストレージシステムあるいはストレージサブシステムにリクエストを送るように構成されたコンピュータシステムを意味する。ホストはサーバかあるいはクライアントであるかもしれない。ここに用いられている、「リモートコピーシステム」という用語は、リモートコピー機能を行うように構成されたコンピュータシステムを意味する。リモートコピーシステムは単一のストレージシステム、ストレージサブシステムまたはストレージユニット、あるいはネットワークまたはコミュニケーションリンクで結合された複数のストレージユニット、ストレージシステム、ストレージサブシステムを意味する。

20

【 0 0 1 6 】

したがって、リモートコピーシステムは、プライマリストレージシステム、セカンダリストレージシステム、仲介システムあるいはそれらの組み合わせを意味する。リモートコピーシステムはまた、１台ないしはそれ以上のホストを含むこともある。一実施例においては、プライマリストレージシステム 1 1 0 a のボリュームは、仲介ストレージシステム 1 1 0 c を用いてセカンダリストレージシステム 1 1 0 b にミラーリングされている。仲介ストレージシステム 1 1 0 c は一般にプライマリストレージシステム 1 1 0 a の比較的近くに位置し、一方セカンダリストレージシステム 1 1 0 b は、仲介ストレージシステム 1 1 0 c と / あるいはプライマリストレージシステムから比較的遠方に位置している。

30

【 0 0 1 7 】

プライマリストレージシステムと関連しているホストから書き込みリクエストを受け取ると、書き込みデータはプライマリストレージシステム 1 1 0 a から仲介ストレージシステム 1 1 0 c に同期的にコピーされる。仲介ストレージシステムにコピーされた書き込みデータは、コントロールデータとジャーナルデータを含むジャーナル情報の形式である。書き込みデータはジャーナルデータに対応し、コントロールデータはジャーナルデータについての管理情報を提供する。

【 0 0 1 8 】

次に、ホストの書き込みリクエストに関連して、ジャーナルは仲介ストレージシステム 1 1 0 c からセカンダリストレージシステム 1 1 0 b に非同期的にコピーされる。書き込みデータは、上記のコピー処理の間にプライマリストレージシステムと仲介ストレージシステムが損傷を受けない限り信頼できる。一実施例においては、仲介ストレージシステムは１ないしはそれ以上のジャーナルボリュームを含むがデータボリュームは含まないので、装置コストは低くなる。一実施例では、リモートコピーシステムは第１のストレージコントローラと第１のデータボリュームを含む第１のストレージシステムを含む。第１のストレージコントローラは第１のデータボリュームに対するデータアクセスリクエストを制御するように構成されている。

40

【 0 0 1 9 】

第１のストレージシステムは第１のストレージシステムに関連する第１のホストからの

50

書き込みリクエストを受け取ると第1のデータボリュームに書き込みデータを記憶し、コントロールデータとジャーナルデータを含むジャーナルを生成するように構成されている。第2のストレージシステムはジャーナルボリュームを含み、第1のストレージシステムにより生成されたジャーナルをジャーナルボリュームに受信し記憶するように構成されている。第3のストレージシステムは第2のデータボリュームを含み、第2のストレージシステムからジャーナルを受け取り、コントロールデータにより提供される情報に従ってジャーナルのジャーナルデータを第2のストレージシステムに記憶する。

【0020】

一実施例においては、ストレージシステムは、第1のホストからデータアクセスリクエストを受け取る第1のストレージコントローラと、第1のストレージコントローラと結合し第1のストレージコントローラの制御により情報を記憶するように構成され、プライマリボリュームを含んでいる第1のストレージ領域と、第1のストレージコントローラから少なくとも100マイル離れて設置されている第2のストレージコントローラと、第2のストレージコントローラと結合し第2のストレージコントローラの制御に従い情報を記憶するように構成され、セカンダリボリュームを含んでいる第2のストレージ領域と、を含んでいる。セカンダリボリュームはプライマリボリュームをミラーリングする。

【0021】

第1のストレージコントローラはプライマリボリュームに第1のホストからの書き込みリクエストに関連した書き込みデータを記憶し、書き込みリクエストに応じてコントロールデータとジャーナルデータを含むジャーナルを生成するように構成されている。ジャーナルデータは書き込みデータに対応しており、ジャーナルは第1のストレージ領域の外に設けられたジャーナルボリュームに同期的に転送される。

【0022】

別の実施例においては、リモートコピーシステム動作は、プライマリストレージシステムに関連したプライマリホストから書き込みリクエストを受け取った後にプライマリストレージシステムでコントロールデータとジャーナルデータを生成し、プライマリストレージシステムとセカンダリストレージシステムの間でデータのミラーリングを行うために仲介ストレージシステムにジャーナルを転送することを含んでいる。セカンダリストレージシステムは仲介ストレージシステムから遠く離れて配置されている。

【発明の効果】

【0023】

一実施例においては、仲介ストレージシステムは1ないしはそれ以上のジャーナルボリュームを含むがデータボリュームは含まないので、装置コストは低くなる。

【発明を実施するための最良の形態】

【0024】

図1Aは、本発明の一実施例によるリモートコピーシステム50は複数のストレージシステムを含むことを示している。リモートコピーシステムはプライマリストレージシステム110aとセカンダリストレージシステム110bと仲介ストレージシステム110cを含んでいる。ストレージシステムは相互にコミュニケーションリンク120aと120bを介し結合されている。一実施例においては、以下に説明するようにリンク120aはリンク120bよりかなり短いので、リンク120aはファイバー・チャネルであり、リンク120bは公衆コミュニケーションリンクである。本実施例では、ストレージシステム110aと110bと110cはディスクアレイユニットか、あるいはストレージサブシステムである。

【0025】

図1Bはデータの読み出し/書き込みリクエストを処理するように構成されたストレージコントローラ62と、書き込みリクエストに応じてデータを記憶する記録媒体を含むストレージユニット63を含む典型的なストレージサブシステム60(たとえばプライマリシステム110a)の例を示す。コントローラ62はホストコンピュータ(たとえばホスト100a)と結合したホストチャネルアダプター64と他のサブシステム(たとえばス

10

20

30

40

50

トレージシステム 110c あるいは 110b) と結合したサブシステムチャネルアダプター 66 とストレージサブシステム 60 におけるストレージユニット 63 に結合したディスクアダプター 68 を含んでいる。

【0026】

本発明の実施例では、これらのアダプターの各々はデータを送受信するポート（図示せず）とポートを経由したデータ転送を制御するマイクロプロセッサ（図示せず）を含んでいる。コントローラ 62 はストレージユニット 63 から読み出したあるいは書き込むデータを一時的に記憶するために用いるキャッシュメモリ 70 も含んでいる。一実施例では、ストレージユニットは複数の磁気ディスクドライブ（図示せず）である。サブシステムはホストコンピュータに複数の論理ボリュームをストレージ領域として提供する。ホストコンピュータはこれらの論理ボリュームの識別子をストレージサブシステムからデータを読み出すあるいはそこにデータを書き込むために用いる。論理ボリュームの識別子は論理ユニット番号（LUN）と呼ばれる。論理ボリュームは単一の物理ストレージデバイスあるいは複数のストレージデバイスに含むことができる。同様に、複数の論理ボリュームを単一の物理ストレージデバイスに関連させることが出来る。

【0027】

より詳細なストレージサブシステムの説明が、参考として援用されている、現在の特許権者が権利を有する、2002年6月5日に出願された日本国特許出願2002-163705により優先権が主張されており、2003年5月21日に出願された「データストレージサブシステム」という名称の米国特許出願_____により提供されている。

【0028】

図1Aに戻ると、本発明の実施例においては、書き込みデータはプライマリシステム110aと仲介システム110cとの間では同期して、仲介システム110cとセカンダリシステム110bとの間では非同期で送られる。したがって、仲介システム110cは一般にセカンダリシステム110bからの距離に比較してプライマリシステム110aの比較的近くに位置している。たとえば、仲介システムはプライマリシステムから約5マイル以下、あるいは10マイル以下もしくは20マイル以下に位置している。一実施例においては、仲介システムはプライマリシステムから100マイルを超えて離れることは無い。比較してみると、本発明の一実施例においては、仲介システムはセカンダリシステムから50マイル以上あるいは100マイル以上あるいは200マイル以上あるいは異なった大陸に位置している。

【0029】

リモートコピーシステムはプライマリストレージシステム110aにコミュニケーションリンク130aを経由して結合しているプライマリホスト100aとコミュニケーションリンク130bを経由してセカンダリストレージシステム110bに結合しているセカンダリホスト100bを含んでいる。プライマリホストはプライマリストレージシステム110aのストレージ領域あるいはボリュームにアクセス（読み出し及び書き込み）するためのアプリケーションプログラム（APP）102を含んでいる。一実施例においては、APP102はもしプライマリホスト100aあるいは/およびプライマリストレージシステム110aが使えなくなった時には（言い換えれば障害が発生した時には）、特に企業などのビジネスユーザに支障なくデータにアクセス出来るようにセカンダリホスト100bにフェイルオーバーを動作させることが出来る。

【0030】

リモートコピーマネージメントソフトウェア（RCM）101は、リモートコピーシステムを制御するためのインターフェイスをユーザおよび/あるいはアプリケーションに提供するためにホスト100aと100bの双方のうえで走る。システム50は三つの形式のボリュームを含んでいる。プライマリストレージシステムは1ないしはそれ以上のプライマリボリューム（PVOL）111aを含んでいる。PVOL111aはAPP101が読み出しあるいは書き込むプロダクションデータを含んでいる。セカンダリストレージシステムは1ないしはそれ以上のセカンダリボリューム（SVOL）111bを含んでい

10

20

30

40

50

る。S V O L 1 1 1 b は P V O L 1 1 1 a のコピーデータを含んでいる。仲介ストレージシステムは 1 ないしはそれ以上のジャーナルボリューム (J V O L あるいは J N L) 1 1 2 を含んでいる。J V O L 1 1 2 は P V O L 1 1 1 a に書き込まれた書き込みデータとそれに対応するコントロールデータを記憶する。

【 0 0 3 1 】

図 2 は本発明の一実施例による J N L ボリュームまたは J V O L 1 1 2 の例を示す。J V O L はリモートコピーを実行するに際し用いられるジャーナルあるいはジャーナル情報を記憶するように構成されている。ジャーナルはジャーナルデータとそのコントロールデータを対として持っている。ジャーナルデータは P V O L 1 1 1 a に記憶されたデータあるいはホスト 1 0 0 から P V O L に書き込まれたデータに対応している。コントロールデータは対応するジャーナルデータに関係した管理情報を含む。

10

【 0 0 3 2 】

本発明の実施例においては、コントロールデータとジャーナルデータは両方ともシーケンシャルに、言い換えれば受け取ったジャーナルは F I F O メモリーに順番に積み重ねられるように、同じジャーナルボリュームに記憶される。仲介ストレージシステムは複数のこのような F I F O メモリーを含むことがある。一実施例では、最初のコントロールデータはボリューム 1 1 2 において最初に記憶され、それに対応する最初のジャーナルデータは同じボリュームに記憶される。

【 0 0 3 3 】

次に第 2 のコントロールデータは最初のジャーナルデータの隣に記憶され、第 2 のコントロールデータに対応する第 2 のジャーナルデータはその隣に記憶され、以下同様になる。コントロールデータは、ジャーナルデータが得られる P V O L 1 1 1 a に対する識別子であるインデックス (I D X) 2 1 1、たとえばプライマリストレージシステム 1 1 0 a あるいはジャーナルグループ 2 0 0 (図 4) における P V O L に割り当てられた固有の番号を、含んでいる。アドレス 2 1 2 は書き込みデータがそこから書かれる P V O L のオフセットアドレス、たとえば書き込みデータのスターティング論理ブロックアドレス (L B A)、を提供する。長さ 2 1 3 は書き込みデータの長さ、たとえば論理ブロックの数あるいは書き込みデータの総バイト数、を提供する。

20

【 0 0 3 4 】

時刻 2 1 4 はホストが P V O L 1 1 1 a にデータを書き込んだ時を示す。シーケンス番号 (S E Q #) 2 1 5 は書き込みのシーケンス情報を提供する。すなわち、シーケンス番号はプライマリストレージシステム 1 1 0 a 内の書き込み順序を提供する。J V O L アイデンティフィケーション (J V O L _ I D) 2 1 6 は対応するジャーナルデータ、たとえばプライマリストレージシステムあるいはジャーナルグループ 2 0 0 におけるジャーナルボリュームに割り当てられた固有の番号、を含むジャーナルボリュームを識別する。ジャーナルオフセット (J F S) 2 1 7 は、ジャーナルデータの記憶を開始する、あるいはジャーナルデータのスターティングアドレスである、ジャーナルボリュームにおけるオフセットアドレスを提供する。あるいは、コントロールデータが対応するジャーナルデータに近接して記憶されているので、コントロールデータは J V O L _ I D 2 1 6 と J O F S 2 1 7 を含まないかもしれない。

30

40

【 0 0 3 5 】

仲介ストレージシステム 1 1 0 c は、第 1 のポインター (J O P T R) 2 1 8 と第 2 のポインター (J I P T R) 2 1 9 を保持している。J O P T R 2 1 8 はリモートストレージシステム 1 1 0 b に送られるジャーナルを指す。J I P T R 2 1 9 はプライマリシステムから受け取った次のジャーナルを記憶するアドレスを指す。したがって、J I P T R 2 1 9 は、リモートストレージシステム 1 1 0 b に未だ送られていないジャーナルを新しいジャーナルが重ね書きしないように、J O P T R 2 1 8 より前に進んではならない。

【 0 0 3 6 】

図 8 に関連して説明されるが、仲介ストレージシステム 1 1 0 c はジャーナルボリューム上で空間を割り付ける。この空間の割付は J I P T R 8 2 0 にコントロールデータとジ

50

ジャーナルデータの長さを加える、言い換えれば次の $JIPTR = JIPTR + (\text{コントロールデータとジャーナルデータの長さ})$ にすることで実行される。前述のように、仲介ストレージシステム 110c がプライマリストレージシステム 110a から JNLWR コマンドを受け取ると、仲介ストレージシステム 110c はコントロールデータとジャーナルデータを記憶するために JNL ボリューム 112 上に、キャッシュメモリと併行してスペースを割り付ける。キャッシュメモリにジャーナルが記憶されると、仲介ストレージシステム 110c は JNLWR コマンドの完了を送る。続いてジャーナルは JNL ボリューム 112' に記憶される。

【0037】

図3は本発明の他の実施例における JNL ボリュームあるいは JVOL 112' を示す。上記にて説明したように、JVOL はリモートコピーを行うために用いられるジャーナルあるいはジャーナル情報を記憶するように構成されている。ジャーナルはジャーナルデータとそのコントロールデータの一对を含んでいる。コントロールデータは JVOL 112' のコントロールデータ領域 230 に記憶され、ジャーナルデータは JVOL 112' のジャーナルデータ領域 240 に記憶される。一実施例においては、JVOL は先入れ先出し (FIFO) ストレージであって、コントロールデータとジャーナルデータは受け取った順序で読み出される。一実施例では、コントロールデータは、ジャーナルデータが引き出される PVOL 111a の識別子であるインデックス (IDX) 211' を含んでおり、たとえば固有の数値がプライマリストレージシステム 110a あるいはジャーナルグループ 200 (図4) の PVOL に割り当てられる。

【0038】

アドレス 212' は書き込みデータがそこから書かれる PVOL のオフセットアドレス、たとえば書き込みデータのスターティング論理ブロックアドレス (LBA)、を提供する。長さ 213' は書き込みデータの長さ、たとえば論理ブロックの数あるいは書き込みデータの総バイト数、を提供する。時刻 214' はホストが PVOL 111a にデータを書き込む時を示す。シーケンス番号 (SEQ#) 215' は書き込みのシーケンス情報を提供する。すなわち、シーケンス番号は、プライマリストレージシステム 110a 内における書き込み順序を与える。JVOL アイデンティフィケーション (JVOL_ID) 216' は対応するジャーナルデータを含むジャーナルボリュームを識別し、たとえばプライマリストレージシステムあるいはジャーナルグループ 200 において固有の数値がジャーナルボリュームに割り当てられる。ジャーナルオフセット (JOS) 217' はジャーナルデータの記憶が始まるジャーナルボリュームのオフセットアドレスあるいはジャーナルデータのスターティングアドレスを提供する。コントロールは、本発明の実施例ではコントロールデータとジャーナルデータは異なった領域に記憶されているので、さらに JOS 217' と JVOL_ID 216' を含んでいる。

【0039】

ジャーナルには二つの形式があり、更新ジャーナルとベースジャーナルである。更新ジャーナルはホストから書き込まれたデータに対するジャーナルである。ジャーナルはホストがデータを PVOL 111a に書き込む際に取りられる。ベースジャーナルは、組み合わせを行う前に PVOL 111a に存在している、先に存在しているデータに対するジャーナルである。ベースジャーナルは PVOL の新しいコピーが作成された時、あるいは再同期化が必要な際に取りられる。

【0040】

図3は本発明の一実施例によるジャーナルグループを示す。ジャーナルグループはジャーナルが生成されるべきボリューム (1 ないしはそれ以上のボリューム) のセットである。ボリュームは只一つのジャーナルグループ 300 の一部であることも可能である。ジャーナルグループ 300 は 1 ないしはそれ以上のデータボリューム 111 あるいは 1 ないしはそれ以上のジャーナルボリュームを含んでいる。ジャーナルグループ 300 はマスターあるいはリストアの何れかの属性を持っている。マスタージャーナルグループ 300a はジャーナルを生成し、一方リストアジャーナルグループ 300b はジャーナルを S VOL

10

20

30

40

50

1 1 1 bに更新する。

【0041】

マスタージャーナルグループはプライマリストレージシステム110aに関連し、1ないしはそれ以上のPVOL111aを含み、随意に1ないしはそれ以上のジャーナルボリュームを含む。リストアジャーナルグループはセカンダリストレージシステム110bに関連し、1ないしはそれ以上のPVOL111bを含み、随意に1ないしはそれ以上のジャーナルボリュームを含む。ジャーナルグループはたとえばJVOL112のような仲介的な属性(図示せず)を持つこともある。本発明の実施例では、仲介的なジャーナルグループはマスタージャーナルグループやリストアグループとグループ化している。このような仲介的なジャーナルグループ(図示せず)は仲介ストレージシステム110cと連携し、1ないしはそれ以上JVOLを含み、随意に1ないしはそれ以上のSVOL111bを含む。

10

【0042】

図5は本発明の一実施例によるジャーナルグループ(JNLG)テーブル400を示している。ストレージシステムは、対応するジャーナルグループが生成されると、JNLGテーブル400を生成し保持する。

【0043】

図4はプライマリストレージシステム110aとセカンダリストレージシステム110bと仲介ストレージシステム110cにより保持されているJNLGテーブル400の内容を示す。ジャーナルグループ番号(GRNUM)410はストレージシステム110内でジャーナルグループに個別に割り当てられた番号を示す。ジャーナルグループ名(GRNAME)420は一般にユーザが割り当てるものであるが、ジャーナルグループに与えられた名称を示す。もし、2ないしはそれ以上のグループが同じ名称GRNAME420を持つ場合は、それらはリモートミラーリング関係にある。ジャーナルグループの属性(GRATTR)430はジャーナルグループに割り当てられた属性、たとえばマスター(MASTER)や仲介(INTERMEDIARY)やリストア(RESTORE)などを示す。

20

【0044】

上記で説明したように、マスタージャーナルグループはジャーナルグループにおけるデータボリューム(PVOL)からジャーナルを生成する。仲介ジャーナルグループはマスタージャーナルグループとリストアジャーナルグループ間に於いて過渡的なジャーナルグループである。リストアジャーナルグループはジャーナルボリュームからデータボリューム(SVOL)に対してジャーナルを再生する。グループステータス(GRSTS)440はジャーナルグループのステータスを示す。ジャーナルグループは以下のCOPYやPAIRやSUSPやSMP Lなどのステータスをもつこともある。

30

【0045】

COPYステータスはベースジャーナルを取り出す先であるジャーナルグループにデータボリュームがあることを示す。すなわち、組み合わせ(言い換えればベースジャーナルに記憶されている事前に存在するデータ)前にPVOLに記憶されているデータがSVOLにコピーするために検索されることである。PAIRステータスは全ての予め存在したデータはSVOLにコピーされ、SVOLへコピーするためにジャーナルグループは更新ジャーナルから更新されたデータを検索されあるいは検索してきた。SUSPステータスすなわちサスペンドステータスは、ジャーナルグループが、更新ジャーナルから更新したデータを取ることをあるいは検索することを中断していることを示す。SMP Lステータスは、ジャーナルグループのどのボリュームもベースジャーナルから予め存在するデータを取得することを開始していない、言い換えればリモートコピーを開始しようとすることを示している。

40

【0046】

アレイ(DVOL__INFO)450はジャーナルグループの全てのデータボリュームについての情報を記憶する。アレイの各エントリーは次の情報を含んでいる。ストレージ

50

システム 110 においてボリュームに個別に指定されるボリューム識別 (VOL ID) 451 (VOL ID 451 は一般に整数の値を持つ) と、ユーザによりボリュームに指定されるボリューム名称 (VOL NAME) 452 と、たとえば COPY や PAIR や SUSP や SMPLE やその他のボリュームのステータスを示すボリュームステータス (VOL STS) 453 を含んでおり、ポインター (BJPtr) 454 はベースジャーナルを取得する経過を保持する。ジャーナルグループにおけるデータボリュームにはグループにおいて DVOL__INFO 450 の索引のために個別のインデックスが指定される。

【0047】

アレイ (JVOL__INFO) 460 はジャーナルグループにおける全てのジャーナルボリュームに関連する情報を記憶する。JVOL__INFO 460 はジャーナルボリューム 112 の識別子である VOL ID 461 を含んでいる。ジャーナルグループにおけるジャーナルボリュームにはグループにおいて JVOL__INFO 450 アレイの索引のために個別のインデックスが指定される。

【0048】

図 6 は本発明の一実施例によるベースジャーナル生成のためのプロセス 500 を示す。プロセス 500 はまた初期コピープロセスとも呼ばれる。ベースジャーナルはユーザが PAIR__CREATE または PAIR__SYNC コマンドを発行すると取得される。PAIR__CREATE コマンドによりプライマリストレージシステムの第 1 のボリュームがセカンダリストレージシステムの第 2 のボリュームと組み合わせられる。本発明の実施例では、もし以下の条件、すなわち、(1) プライマリストレージシステム 110a とセカンダリストレージシステム 110c に保持されている二つの JNL グループ 300 が同じ GRNAME 420 を持っている、(2) プライマリストレージシステム 110a の二つの JNL グループ 300 の一つが GRATR 430 において MASTER 属性を持ち、セカンダリストレージシステム 110c における他方が RESTORE 属性を持っている、(3) 二つの JNL グループ 300 からの二つのデータボリュームが同じインデックス番号を持っている、が満たされる場合には、二つのデータボリュームはペアの関係にある。

【0049】

PAIR__SYNC コマンドは、ペアの関係にあるボリュームに対して、双方が同じデータを保有するように同期をさせるかあるいはミラーリングをさせる。プライマリストレージシステムはこれらの二つのコマンドの一つを受け取ると、プロセス 500 を実行する。ステップ 510 においては、ベースジャーナルポインター (BJPtr) 454 は初期化されて、データボリュームにおける第 1 のデータ (たとえばデータボリューム上の第 1 のブロック、トラック、ブロックの一塊やアドレス指定可能な任意のデータなど) からベースジャーナルの取得を開始する。

【0050】

次のターゲットが検査される (ステップ 515)。次のターゲットは BJPtr 454 の値により得られる。たとえばもし BJPtr が j を持っている、次のターゲットは j 番目のブロックである。一実施例においては、より効率の良い処理のために、ジャーナルについて一度に数ブロックのデータが取得される。したがって、本例では j 番目のブロックの後、次のターゲットは j 番目のブロックから n ブロックである。プロセスは追加のターゲットが有るか否かを定める (ステップ 520)。プロセス 500 はもうそれ以上のターゲットが無い、言い換えれば全てのベースジャーナルが取得された場合には終了する。しかしながら、もしさらにターゲットが存在する場合には、ターゲットブロックに関するコントロールデータが生成される (ステップ 525)。生成されたコントロールデータはプライマリストレージシステム 110a のキャッシュメモリに記憶される。ターゲットデータは PVOL からキャッシュメモリに読み込まれる (ステップ 530)。ターゲットデータを読んだ後で、コントロールデータは確認される。

【0051】

コントロールデータは、以下の情報、すなわち、IDX 211 やアドレス 212 や長さ 213 を含んでいる。他の情報も同様に含まれている可能性がある。ジャーナルデータと

コントロールデータは仲介ストレージシステム 110c に送られる (ステップ 535)。一般に、ステップ 530 で確認されたジャーナル (コントロールデータとジャーナルデータの一对) のみが仲介システムに転送される。一実施例では、複数の確認されたジャーナルが単一のコマンドで共に送られる。ジャーナルを無事に転送し終わると B J P t r 4 5 4 は次のターゲットへとインクリメントされる。すなわち、B J P t r 4 5 4 は $j + n$ にインクリメントされる。ステップ 515 と 540 はターゲットが無くなるまで繰り返される。

【0052】

図 7 は本発明の一実施例による更新ジャーナルを生成するためのプロセス 600 を示す。プロセス 600 はまた更新コピープロセスとも呼ばれる。ユーザにより P A I R _ C R E A T E コマンドあるいは P A I R _ S Y N C コマンドが発行された後でプロセス 600 は開始される (ステップ 602)。すなわち、プライマリストレージシステムが更新ジャーナルの取得を開始する。ステップ 602 はもし P V O L が予め存在するデータを持っている場合は、プロセス 500 の後で実行される。プライマリストレージシステム 110a は書き込みコマンドを受け取ったか否かを判定する (ステップ 610)。この時には全ての読み取りコマンドが無視される。

【0053】

ストレージはまた P V O L のペアステータスが C O P Y であるか P A I R であるかを判定する。もしこれらの条件が満足されると、プライマリストレージシステムはペアステータスが C O P Y であるかを判定する (ステップ 615)。もしそうであるならば、ベースジャーナルが書き込みターゲットアドレスに対して既に取得されているかを判定するためにステータスがチェックされる (ステップ 620)。これはポインター B J P t r 4 5 4 を判定することで行われる。すなわち、もし (書き込みターゲットアドレス) \leq B J P t r 4 5 4 であるならば、プロセス 600 はステップ 625 に進む。もしステップ 620 が真であるかステップ 615 が偽である、言い換えれば、ベースジャーナルを取得するプロセスが完了している場合は、更新ジャーナルは書き込みのために取得される。この目的のために、コントロールデータは最初に生成される。

【0054】

コントロールデータは I D X 2 1 1 とアドレス 2 1 2 と長さ 2 1 3 を含んでいる。書き込みコマンドはアドレス 2 1 2 と長さ 2 1 3 を含んでいる。コントロールデータには他の情報が含まれている可能性がある。書き込みデータはホストから受け取り、キャッシュメモリに記憶される (ステップ 630)。書き込みデータはステップ 625 で生成されたコントロールデータに係わるジャーナルデータに対応する。コントロールデータとジャーナルデータは仲介ストレージシステム 110c に転送される (ステップ 635)。プロセス 600 は仲介ストレージシステム 110c からの受け取り通知を待つ (ステップ 640)。受け取り通知を受け取ると、書き込み完了がホストに送られる (ステップ 645)。仲介ストレージシステムからの受け取り通知が受けとられるまでは書き込み完了がホストに通知されないで、書き込みデータのプライマリシステムと仲介システムへの記憶は保証される。

【0055】

図 8 は本発明の一実施例による、コントロールデータとジャーナルデータを含むジャーナル情報を仲介ストレージシステム 110c に転送するプロセス 700 を示している。プライマリストレージシステム 110a は、仲介ストレージシステム 110c にジャーナルデータを送るための J N L 書き込みコマンド (J N L W R コマンド) を発行する (ステップ 702) 一実施例においては、コマンドは 1 ないしはそれ以上のパラメータ、たとえばジャーナルデータの長さなどを含んでいる。コントロールデータの長さは、固定データ長、たとえば 64 バイト、が本発明の実施例のコントロールデータに用いられているのでコマンドパラメータには含まれていない。

【0056】

もう一つの方法としては、可変長のコントロールデータが用いられることもあり、その

10

20

30

40

50

場合はその長さについての情報がパラメータに含まれる必要がある。仲介ストレージシステム 110c は、コマンドパラメータにおいて提供される情報に従い JNL ボリューム上に記憶空間を割り付ける（ステップ 710）。以下でより詳細に説明するように、空間の割り付けは、書き込み性能を向上させるために、キャッシュメモリーバッファ上でも実行される。割り付けられたバッファはボリュームにおいて割り付けられたストレージ空間に関係している。

【0057】

一旦ストレージ割り付けが実行されると、転送の準備が出来たパケットあるいはメッセージがプライマリストレージシステム 110a へ送られる。プライマリシステムは、転送の準備が出来たパケットを受け取ると、ジャーナル情報を仲介ストレージシステムに転送する（ステップ 720）。一実施例においては、コントロールデータが最初に送られ、次いでジャーナルデータが送られる。仲介ストレージシステムはキャッシュメモリー上のバッファにコントロールデータとジャーナルデータを記憶させる（ステップ 730）。ジャーナルは最終的には、ステップ 720 におけるバッファストレージ結合によりプライマリストレージシステム 110b がアイドル状態である時に、仲介システムにおける割り付けられた JNL ボリュームに記憶される。

【0058】

さらに、シーケンス番号および/あるいはその時のタイムスタンプがジャーナルに割り付けられ、言い換えればコントロールデータに付加される。シーケンス番号は順番に（プライマリシステムから）受け取られたジャーナルに割り当てられ、JNL ボリュームに記憶される。シーケンス番号によりプライマリシステムから受け取ったジャーナルの順番が分かるので、データリカバリープロセスに役立つ。本発明の実施例では、仲介システム 110c はジャーナルにシーケンス番号を付加し、さもなければ、シーケンス番号を管理する。仲介システムにおいてはカウンター 152 は、プライマリシステム 110a（図 1 を参照）により転送されたジャーナルにシーケンス番号を付加するために用いられる。タイムスタンプも、仲介ストレージシステム 110c がジャーナルを受け取った時刻を示すために、ジャーナルに付加される。

【0059】

他の実施例においては、シーケンス情報は、ジャーナルが仲介システムに転送される前に、プライマリシステムにおいてジャーナルに付加される。同様に、タイムスタンプは、それが仲介システムへ転送された時刻を示すために、プライマリシステムによってジャーナルに付加することもある。一旦ジャーナルが正常に受信され記憶されると、仲介システムはジャーナルを問題なく受信したことの認証をプライマリシステムに送る。その後は、プライマリストレージシステム 110a は書き込みジャーナルコマンドの完了を発行する。

【0060】

図 9 は本発明の一実施例による、仲介ストレージシステム 110c からセカンダリストレージシステム 110b へジャーナルを送るプロセス 900 を示す。本発明の実施例においては、仲介ストレージシステム 110c の JNL ボリューム 112 に記憶されたジャーナルは、プライマリストレージシステム 110a の書き込みコマンドとは非同期で、言い換えればプライマリシステムからの書き込みコマンドの合間に、セカンダリストレージシステム 110b へ送られる。プロセス 900 はプライマリシステムから仲介システムへジャーナルを転送することに関してはプロセス 700 と同様である。パラメータが付いた JNL WR コマンドが仲介システムからセカンダリストレージシステムへ発行される（ステップ 902）。セカンダリシステムはパラメータにより規定されたデータ長に従いキャッシュメモリー上にバッファ空間を割り付け、転送準備完の通知を仲介システムに送り返す（ステップ 910）。仲介システムはコントロールデータとそれに対応するジャーナルデータを含んでいるジャーナルを送る（ステップ 920）。一実施例においては、コントロールデータは最初に転送され、続いてジャーナルデータが転送される。セカンダリシステムは割り当てられたバッファ空間にジャーナルを記憶し、仲介システムにジャーナ

10

20

30

40

50

ル受信の受け取り通知を送る（ステップ 930）。

【0061】

ジャーナルデータはプロセス 700 のステップ 730 において割り当てられたシーケンス番号とタイムスタンプに基づき S V O L に記憶される。たとえば、下位のシーケンス番号を持つジャーナルは上位のシーケンス番号を持つジャーナルの前に再記憶される。受け取り通知を受信すると、仲介システムはデータ書き込みの完了を示す W R J N L コマンドを発行する（ステップ 940）。仲介システムのジャーナルボリュームに関連したポインター、たとえば J O P T R 810 は、次のバッチのデータまで前進しセカンダリシステムにコピーされる。

【0062】

図 10 は、本発明の一実施例にしたがって、セカンダリストレージシステム 110b の S V O L へ、その対応するコントロールデータを用いてジャーナルデータを記憶するプロセス 1000 を示している。セカンダリストレージシステム 110b は、R E S T O R E の属性を持つ J N L グループ 300 上で定期的にプロセス 1000（ステップ 1002）を実行する。一実施例では、プロセス 1000 は 10 秒ごとに実行される。S V O L に記憶する予定の、コントロールデータとジャーナルデータを含む次のジャーナルは、シーケンス番号を用いて選択される（ステップ 1005）。この目的のために、セカンダリストレージシステム 110b は再記憶した、言い換えればそのジャーナルを S V O L に記憶した、ジャーナルのシーケンス番号の追跡をしている。セカンダリストレージシステムは、最も最近に再記憶されたジャーナルのシーケンス番号とキャッシュメモリーに一時的に記憶されたジャーナルに係わるシーケンス番号とを比較し、再記憶する次のジャーナルを決定する。

【0063】

ステップ 1005 で選択されたジャーナルのコントロールデータは、ジャーナルデータの記憶領域、言い換えれば特定の S V O L とその中の場所、を決めるために用いられる（ステップ 1010）。たとえば、コントロールデータの以下の場所、すなわち I D X 211、アドレス 212、長さ 213、が調べられる。I D X 211 は M A S T E R J N L グループの P V O L 111a、言い換えればプライマリシステムにおけるプライマリストレージボリュームのインデックスを示す。ジャーナルデータは同じインデックスを持つ S V O L に記憶される（ステップ 1015）。すなわち、ジャーナルデータは I D X 211 で特定され、アドレス 212 で示されるアドレスで、長さ 213 に対応する長さにより S V O L に記憶される。本発明の実施例では、コントロールデータは、P V O L と S V O L におけるジャーナルデータの記憶場所はミラーリングされているので、S V O L には記憶されない。

【0064】

図 11 は本発明の一実施例の場合に於いて、プライマリシステム 110a に障害が発生した場合に、リモートコピーシステム 50 がフェイルオーバーを実行するところを示している。フェイルオーバーはプライマリストレージシステム 110a またはホスト 100a またはその両者が停止し、セカンダリホスト 100b が適当なアプリケーションを走らせる状態あるいはプロセスを言い、ここでセカンダリストレージシステムは新しい「プライマリ」ストレージシステムとして機能する。もしプライマリストレージシステム 110a が未だ動作可能であるか、あるいは障害が発生した後で再起動した場合は、セカンダリストレージシステム 110b のデータボリューム 111b を P V O L と設定し、二つのサイトの間でミラーリングを継続することが求められる。

【0065】

仲介ストレージシステム 110c は従来と同様に仲介ストレージとして用いられる。しかしながら、仲介ストレージシステム 110c はセカンダリストレージシステム 110b から遠く離れている可能性があるので、新しいプライマリシステム（言い換えればセカンダリシステム 110b）はホスト 100b の書き込みリクエストに対して非同期でジャーナルを送信する。この目的のために、プロセス 600 はステップ 635 と 640 無しに実

10

20

30

40

50

行される。すなわち、ジャーナルは同期して生成されるが、ホストの書き込みリクエストに対しては非同期で送られる。コントロールデータには、仲介システム 110c へ送信する前に、新しいプライマリシステム 110b でタイムスタンプとシーケンス番号が付加される。したがって、フェイルオーバーの際には、仲介システムでは、このようなステップを実行する必要はない。

【0066】

図12は本発明の他の実施例における、リモートコピーシステム50'を示す。システム50'はプライマリストレージシステム110a'とセカンダリストレージシステム110b'と仲介ストレージシステム110c'を含む。プライマリストレージシステムは複数のボリューム111a'を含み、コミュニケーションリンク130a'を経由してプライマリホスト100a'と結合されている。プライマリホスト100a'はアプリケーション102'とRCM101a'を含む。セカンダリシステムは複数のボリューム111b'を含み、コミュニケーションリンク130b'を経由してセカンダリホスト100b'と結合されている。セカンダリシステムはRCM101b'を含んでいる。仲介システムはジャーナルボリューム112'と複数のデータボリューム111c'を含んでいる。一実施例では、データボリューム111c'は、たとえば他のストレージシステムにあるジャーナルボリューム112'から遠く離れた場所に位置しているかもしれない。

【0067】

仲介ホスト100c'はコミュニケーションリンク130c'を経由して仲介システム110c'に結合している。仲介ホストはRCM101c'を含んでいる。システム50'では、仲介システムとセカンダリシステムは共にプライマリシステムとミラーリングを行っている。仲介システムにおけるデータミラーリングは、上述のプロセス700、900および1000を用いて行われる。システム50'は、プライマリストレージシステム110a'が停止した場合は、仲介ストレージシステム110c'あるいはセカンダリストレージシステム110b'にフェイルオーバーすることができる。このような構成に於いては、仲介ストレージシステムは、セカンダリストレージシステムよりもプライマリシステムやユーザに近い可能性があるので、セカンダリシステムよりも効率の良いストレージセンターの役を果たすものである。フェイルオーバーの際には、仲介ホスト100c'は単独であるいはホスト100a'と協調してプライマリホストとして機能する。

【0068】

上記の詳細説明は本発明の特定の実施例を説明するために提供されたものであって、制約を目的としたものではない。多くの改造や変更が本発明の範囲内に於いて可能である。したがって、本発明は特許請求の範囲により定義されるものである。

【図面の簡単な説明】

【0069】

【図1A】は、本発明の一実施例による3箇所のデータセンタを有するリモートコピーシステムを示す。

【図1B】は、本発明の一実施例による典型的なストレージシステムの例を示す。

【図2】は、本発明の一実施例による仲介ストレージシステムにおいて提供されるジャーナルボリュームを示す。

【図3】は、本発明の他の実施例による仲介ストレージシステムにおいて提供されるジャーナルボリュームを示す。

【図4】は、本発明の一実施例によるマスターとリストアの属性を持つジャーナルグループを示す。

【図5】は、図1のリモートコピーシステムにおいてストレージシステムにより維持されているジャーナルグループテーブルを示す。

【図6】は、本発明の一実施例によるベースジャーナルを生成するプロセスを示す。

【図7】は、本発明の一実施例による更新ジャーナルを生成するプロセスを示す。

【図8】は、発明の一実施例によるジャーナルをプライマリストレージシステムから仲介ストレージシステムに転送するプロセスを示す。

【図 9】は、本発明の一実施例によるジャーナルを仲介ストレージシステムからセカンダリストレージシステムへ送るプロセスを示す。

【図 10】は、本発明の一実施例によるセカンダリストレージシステムにおいてジャーナルを復元するプロセスを示す。

【図 11】は、本発明の一実施例によるリモートコピーシステムにおいて実行されるフェイルオーバーあるいはフォールバックプロセスを示す。

【図 12】は、本発明の他の実施例によるリモートコピーシステムを示す。

【符号の説明】

【 0 0 7 0 】

1 0 0 a . . . ホスト

1 0 1 a . . . R C M、ユーザコマンド

1 0 2 . . . A P P

1 1 1 a . . . P V O L

1 1 1 b . . . S V O L

1 1 0 a . . . ストレージシステム

1 2 0 a . . . J N L を取得 (同期)

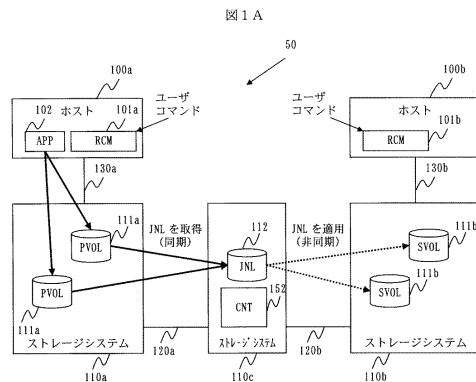
1 2 0 b . . . J N L を適用 (非同期)

2 3 0 . . . コントロールデータ領域、コントロールデータ

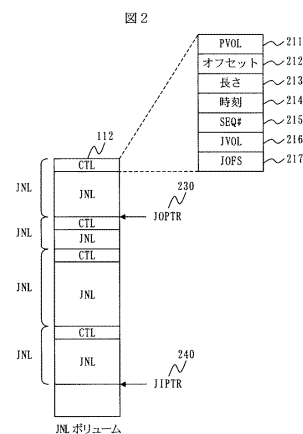
2 4 0 . . . ジャーナルデータ領域、ジャーナルデータ

10

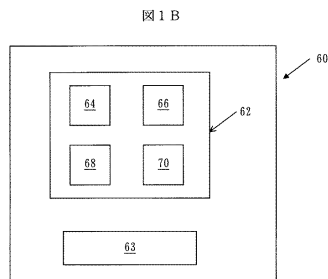
【図 1 A】



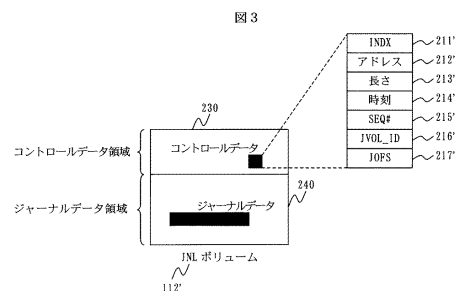
【図 2】



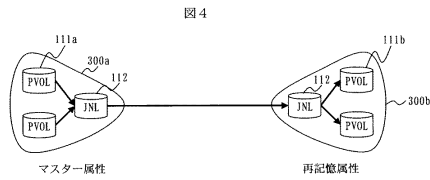
【図 1 B】



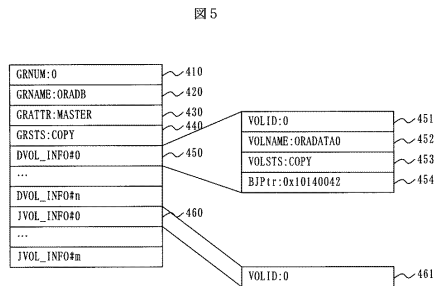
【図 3】



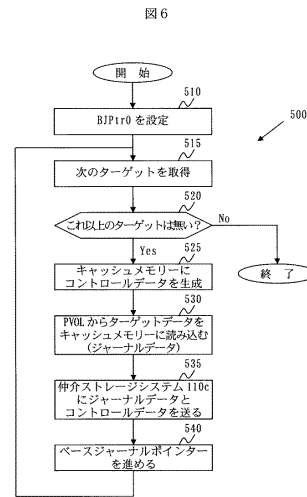
【図 4】



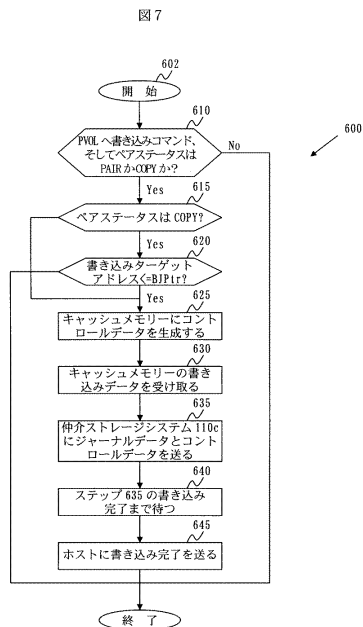
【図 5】



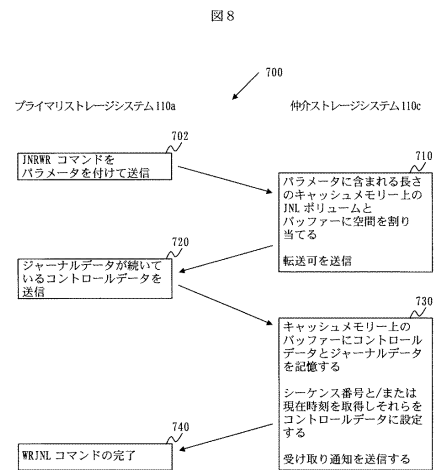
【図 6】



【図 7】



【図 8】



フロントページの続き

(56)参考文献 特開2003-122509(JP,A)
特開平10-049418(JP,A)
特開2000-181634(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F 12/00
G06F 3/06