

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4452494号
(P4452494)

(45) 発行日 平成22年4月21日 (2010. 4. 21)

(24) 登録日 平成22年2月5日 (2010. 2. 5)

(51) Int. Cl.	F I
G 0 6 F 3/06 (2006.01)	G 0 6 F 3/06 3 0 4 F
G 0 6 F 12/00 (2006.01)	G 0 6 F 3/06 3 0 1 Z
	G 0 6 F 12/00 5 3 1 D

請求項の数 8 外国語出願 (全 26 頁)

(21) 出願番号	特願2003-428627 (P2003-428627)	(73) 特許権者	000005108
(22) 出願日	平成15年12月25日 (2003. 12. 25)		株式会社日立製作所
(65) 公開番号	特開2004-272884 (P2004-272884A)		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9月30日 (2004. 9. 30)	(74) 代理人	100093861
審査請求日	平成18年10月27日 (2006. 10. 27)		弁理士 大賀 真司
(31) 優先権主張番号	10/386277	(72) 発明者	渡辺 直企
(32) 優先日	平成15年3月11日 (2003. 3. 11)		アメリカ合衆国カリフォルニア州サニーベ
(33) 優先権主張国	米国 (US)		イル イーストエルカミノリアル965

審査官 高瀬 勤

最終頁に続く

(54) 【発明の名称】 複数リモートストレージでのリモートコピー停止後のデータ同期化方式

(57) 【特許請求の範囲】

【請求項 1】

第一、第二、第三、第四のストレージボリュームが直列接続され、各ストレージボリュームに保存されたデータを同期化させるストレージシステムにおいて、

前記第二のストレージボリュームと前記第三のストレージボリュームとの間のリモートコピー動作の停止を検出する検出部と、

前記検出部による前記リモートコピー動作の停止が検出されたことを契機に、前記第二のストレージボリュームになされた前記リモートコピー動作の停止までの第一の変更記録を生成するとともに前記第三のストレージボリュームになされた前記リモートコピー動作の停止までの第二の変更記録を生成し、前記第一の変更記録の複写情報を前記第一のストレージボリュームに保存するとともに前記第二の変更記録の複写情報を前記第四のストレージボリュームに保存する制御部と、

前記第一の変更記録の複写情報を用いて、前記第四のストレージボリュームを前記第二のストレージボリュームに同期化させ、又は、前記第二の変更記録の複写情報を用いて、前記第三のストレージボリュームを前記第一のストレージボリュームに同期させるデータ同期化部と、

を備えることを特徴とするストレージシステム。

【請求項 2】

前記第一のストレージボリュームを有するプライマリストレージシステムと、

前記第二、第三及び第四のストレージボリュームを有する少なくとも2台のセカンダリ

10

20

ストレージシステムと、

を備え、

前記プライマリストレージシステムは、

ホストコンピュータ又は前記少なくとも 2 台のセカンダリストレージシステムが発行する I / O 要求に基づいて生成されるコマンドエントリを格納するための第一のキューを有する第一のディスクコントローラを備え、

前記少なくとも 2 台のセカンダリストレージシステムは、それぞれ、

前記ホストコンピュータ又は前記プライマリストレージシステムが発行する I / O 要求に基づいて生成されるコマンドエントリを格納するための第二のキューを有する第二のディスクコントローラを備える

10

ことを特徴とする請求項 1 に記載のストレージシステム。

【請求項 3】

前記第一及び第二のディスクコントローラは、それぞれ、前記第一及び第二のキューとして、ワーキングキュー、中間キュー及びライトヒストリキューを備える

ことを特徴とする請求項 2 に記載のストレージシステム。

【請求項 4】

前記第一及び第二のディスクコントローラは、それぞれ、前記第一及び第二のキューとして、コールバックキュー及びライトヒストリキューを備える

ことを特徴とする請求項 2 に記載のストレージシステム。

【請求項 5】

20

第一、第二、第三、第四のストレージボリュームが直列接続されたストレージシステム内のストレージボリュームに保存されたデータを同期化させるデータ同期化方法において、

前記第二のストレージボリュームと前記第三のストレージボリュームとの間のリモートコピー動作の停止を検出する検出ステップと、

前記リモートコピー動作の停止を検出したことを契機に、前記第二のストレージボリュームにおいて前記第二のストレージボリュームになされた前記リモートコピー動作の停止までの第一の変更記録を生成するとともに前記第三のストレージボリュームにおいて前記第三のストレージボリュームになされた前記リモートコピー動作の停止までの第二の変更記録を生成し、前記第一の変更記録の複写情報を前記第一のストレージボリュームに保存するとともに前記第二の変更記録の複写情報を前記第四のストレージボリュームに保存する制御ステップと、

30

前記第一の変更記録の複写情報を用いて、前記第四のストレージボリュームを前記第二のストレージボリュームに同期化させ、又は、前記第二の変更記録の複写情報を用いて、前記第三のストレージボリュームを前記第一のストレージボリュームに同期させる同期ステップと

を備えることを特徴とするデータ同期化方法。

【請求項 6】

前記第一のストレージボリュームを有するプライマリストレージシステムでは、第一のディスクコントローラが、第一のキューに、ホストコンピュータ又は前記少なくとも 2 台のセカンダリストレージシステムが発行する I / O 要求に基づいて生成されるコマンドエントリを格納させ、

40

前記第二、第三及び第四のストレージボリュームを有する少なくとも 2 台のセカンダリストレージシステムでは、それぞれ、第二のディスクコントローラが、第二のキューに、前記ホストコンピュータ又は前記プライマリストレージシステムが発行する I / O 要求に基づいて生成されるコマンドエントリを格納させる

ことを特徴とする請求項 5 に記載のデータ同期化方法。

【請求項 7】

前記第一及び第二のディスクコントローラは、それぞれ、前記第一及び第二のキューとして、ワーキングキュー、中間キュー及びライトヒストリキューを用いる

50

ことを特徴とする請求項 6 に記載のデータ同期化方法。

【請求項 8】

前記第一及び第二のディスクコントローラは、それぞれ、前記第一及び第二のキューとして、コールバックキュー及びライトヒストリキューを用いる
ことを特徴とする請求項 6 に記載のデータ同期化方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、プライマリ（ローカル）ストレージシステムと、セカンダリ（リモート）ストレージシステムで構成されるデータ処理ストレージシステムに関する。 10

【背景技術】

【0002】

商業、政府、その他の企業体によるデータ処理の使用拡大により、夥しい量のデータが保存されており、この多くのものはそれら企業体の日々のオペレーションに対して極めて重要なものである。

【0003】

例えば、莫大な数の財務上のトランザクションは現在、完全に電子的に処理されている。

【0004】

業務上、例えば、航空会社は将来のチケット予約に関する情報が失われると大混乱になるというリスクを背負っている。信頼できるデータへのニーズの結果として、元データが損傷又は損失してもデータが使用できるように、保持されているデータの一つ以上のコピーを、ローカルサイトから離れた位置に、バックアップすることが普通である。データが重要である程、バックアップ方法は精巧なものでなければならない。 20

【0005】

例えば、重要な又は貴重なデータを守るための一つの方法は、そのデータのコピーをローカルストレージシステムから地理的に離れたサイトにバックアップ保存することである。

【0006】

各リモートストレージシステムは、ローカルストレージシステムで保持されるデータのミラーイメージを維持し、ローカルストレージシステムのローカルデータが更新される度に、その保存されたデータをローカルストレージシステムのローカルデータイメージの“ミラー”変更に更新する。 30

【0007】

ローカルストレージシステムのデータをミラーするためのリモートストレージシステムの一例は、“Method and Apparatus for Mirroring Data in a Remote Storage System”と題される米国特許No. 5,933,653に記載されている。

【0008】

【特許文献 1】米国特許 5 9 3 3 6 5 3 号公報 40

【発明の開示】

【発明が解決しようとする課題】

【0009】

リモートストレージシステムに転送される更新データは、遠隔コピー操作のオーバーヘッドを低減するために、しばしばキューされグループ化されて、インターネットの如き、ネットワーク転送媒体を通して送信される。この結果、リモートサイトにミラーされるデータイメージはローカルサイトのものとは必ずしも同じではない。

【0010】

もし、ローカルデータのミラーの為に、複数のリモートストレージが使用される場合は、少なくとも更新が完了するまでの間は、リモートストレージのデータイメージは互いに 50

異なっている事が有り得る。

【 0 0 1 1 】

これらの異なったデータイメージの存在は、ローカルシステムが損傷した場合には問題となる。ローカルストレージシステムの損傷は、あるリモートストレージシステムに損傷前のローカルストレージシステムに完全ではなくても比較的近いデータイメージを残し、他方、他のシステムでは最終更新操作では決して更新されないより古い“ステール”データイメージが存在する。

【 0 0 1 2 】

かくして、ローカルストレージシステムの損傷時には、リモートストレージシステムのデータをすべて最新データに一致させることを保障する再同期化処理が、システムの再起動以前に必要である。

10

【 0 0 1 3 】

更に解決を要するもう一つの問題は、リモートコピー動作中に“停止 (suspension)”が発生した場合のシステム回復である。予期せぬ事故による途絶、例えば、キャッシュオーバーフロー、コピー中のストレージシステム障害、ネットワーク障害等が発生してリモートコピー動作が停止した場合は、再同期化処理が必要になる。

【 0 0 1 4 】

リモートコピー動作での再同期化の一例は、“Method and Apparatus for Independent Operation of a Remote Data Facility.”と題される米国特許No.6,092,066、に記されている。

20

【 0 0 1 5 】

然しながら、この特許での技術では、ある限定された環境での再同期化処理が可能になるのみである。

【 0 0 1 6 】

ある種のより複雑なシステム停止、例えば、リンク障害、キャッシュオーバーフロー、及び/またはドライブ障害等の2つの障害が組み合わさって発生した場合、システムの再初期化をしないで済む再同期化の方策は用意されていない。

【 0 0 1 7 】

このような状況では、この技術では、少なくとも2つのコピーが存在する構成を保障していない為、再同期化処理のためにはボリューム全体のコピーが通常は必要になる。

30

【課題を解決するための手段】

【 0 0 1 8 】

本発明は、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間でデータのミラー処理を実行する改善されたデータ処理ストレージシステムを提供する。典型的に、各ストレージシステムは、データ記憶の為にボリュームを備え、ボリュームはミラー状態に維持される。

【 0 0 1 9 】

もし、たとえば何れかのストレージボリュームの障害又はネットワーク結合障害により、ボリューム間でのデータ転送が不能になったら、タイムスタンプ付きビットマップがプライマリストレージシステムで生成され、セカンダリストレージシステムの一つに記憶される。これらの記録は、結合リンク確立後にペアを再同期化するために用いられる。

40

【 0 0 2 0 】

もし、ペアの一つのメンバ又は他方で障害が発生した場合には、障害発生時に、記録が障害発生ストレージボリュームとは別のライトオペレーションステータスのストレージボリュームに書き込まれる。この記録は後でストレージボリュームを再同期化するために使用される。

【 0 0 2 1 】

各プライマリ及びセカンダリストレージシステムは追加のボリュームを備え、ミラーリモートペアの動作停止時は、ビットマップがペアの各メンバから追加のストレージボリュームの一つに保存される事が好ましい。

50

【 0 0 2 2 】

既に述べたように、これらのビットマップは、以前のペアのメンバーで保持されていた情報が失われても、新しいペアの再同期化のために使用することが出来る。再同期化処理は、新しいペア間でビットマップを交換して、新しいペアの再同期化に必要なライトオペレーションのセットを決定することにより実行される。

【 0 0 2 3 】

次いで、これらのライトオペレーションが実行され、新しく同期化されたペアが誕生する。

【 0 0 2 4 】

プライマリストレージシステムの第一ストレージボリュームと、第二ストレージボリュームをも備えるセカンダリストレージシステムの第一ストレージボリュームとの間でのデータ転送が途絶した後の、ストレージシステムに保存されたデータの同期化方法は、下記のステップから成り立つ事が好ましい。

【 0 0 2 5 】

第一に、プライマリ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のデータ転送の途絶を検出し、次に、プライマリストレージシステムにプライマリストレージシステムの第一ストレージボリュームに書き込まれたデータの記録を提供し、セカンダリストレージボリュームにセカンダリストレージシステムの第一ストレージボリュームに書き込まれたデータの記録を提供する。

【 0 0 2 6 】

次のステップでは、プライマリストレージシステムの第一ストレージボリュームに書き込まれたデータの記録の少なくとも一部のコピーを第二ストレージボリュームに生成する。

【 0 0 2 7 】

次に、第二ストレージボリュームの少なくとも一部のコピーとセカンダリストレージシステムの第一ストレージボリュームに書き込まれたデータの記録を用いて、セカンダリストレージシステムの第一ストレージボリュームがセカンダリストレージシステムの第二ストレージボリュームに同期化される。

【 0 0 2 8 】

他の実施例として、直列的に結合されたストレージシステムに保存されたデータの同期化方法が提供される。

【 0 0 2 9 】

ストレージシステムでは、第一、第二、第三、及び第四のストレージボリュームが直列的に結合されている。第二ストレージボリュームと第三ストレージボリューム間のデータ転送が途絶された後に、第二ストレージボリュームに第二ストレージボリュームに書き込まれたデータの第一の記録を提供し、第三ストレージボリュームに第三ストレージボリュームに書き込まれたデータの第二の記録を提供することにより、システムが再同期化される。

【 0 0 3 0 】

次に、システムは、第一の記録の少なくとも一部を第一ストレージボリュームにコピーし、第二の記録の少なくとも一部を第四ストレージボリュームにコピーする。コピー終了後のステップでは、コピーされた記録の少なくとも一部を用いて、第二、第三ストレージボリュームの少なくとも一つを、第一、第四ストレージボリュームの少なくとも一つに同期化する。

【 発明の効果 】

【 0 0 3 1 】

何れかのストレージボリュームの障害又はネットワーク結合障害により、ストレージボリューム間でのデータ転送が不能になったら、タイムスタンプ付きビットマップがプライマリストレージシステムで生成され、セカンダリストレージシステムの一つに記憶される。

10

20

30

40

50

【 0 0 3 2 】

ペアの一つのメンバ又は他方で障害が発生した場合には、障害発生時に、記録が障害発生ストレージボリュームとは別のライトオペレーションステータスのストレージボリュームに書き込まれる。この記録は後でストレージボリュームを再同期化するために使用される。

【 発明を実施するための最良の形態 】

【 0 0 3 3 】

図 1 は、プライマリ又はローカルサイト 1 2 と複数のセカンダリ又はリモートサイト 1 4 がデータ通信ネットワーク 1 0 7 で通信的に結合された、データ処理システム 1 0 を構成するストレージシステムを示す。

10

【 0 0 3 4 】

ローカルサイト 1 2 はホストプロセッサ 1 0 1 とローカルストレージシステム 1 0 4 を備えている。

【 0 0 3 5 】

ホストプロセッサ 1 0 1 とストレージシステム 1 0 4 は、ローカルサイト 1 2 で保持されるデータイメージの更新データをリモートサイト 1 4 に通信する為に、リモートサイト 1 4 にデータ通信ネットワーク 1 0 7 を通して結合されている。

【 0 0 3 6 】

かくして、リモートサイト 1 4 はローカルサイト 1 2 のミラーデータイメージを保持する。

20

【 0 0 3 7 】

リモートサイト 1 4 は、各々リモートストレージシステム (セカンダリストレージ 1 , セカンダリストレージ 2) 1 0 5 、 1 0 6 を備えている。各サイト 1 4 は対応するホストプロセッサ (リモートホスト) 1 0 2 、 1 0 3 を含んでも含まなくても良い。

【 0 0 3 8 】

リモートストレージシステム 1 0 5 、 1 0 6 は、ローカルストレージシステム 1 0 4 のミラーデータイメージを保持する為に、ローカルストレージシステム 1 0 4 のストレージボリューム 1 4 2 と極めて類似したストレージボリュームを備えている事が好ましい。

【 0 0 3 9 】

リモートサイト 1 4 は、ローカルストレージシステム 1 0 4 が計画的又は非計画的に停まった場合には、当該データにアクセスされるものの、ローカルサイト 1 2 のデータが破壊されるような災害に見舞われた場合でも、サイトとデータを維持するためには、リモートサイト 1 4 はローカルストレージシステム 1 0 4 から地理的に離れた場所に設置した方がより好ましいと言える。

30

【 0 0 4 0 】

リモートストレージシステム 1 0 5 , 1 0 6 はローカルストレージシステム 1 0 4 と実質的に同じである為に、ローカルストレージシステム 1 0 4 に関する議論はそのまま等しくリモートストレージシステム 1 0 5 , 1 0 6 に当てはまる。

【 0 0 4 1 】

ローカルサイト 1 2 では、ホストプロセッサ 1 0 1 は、ネットワークインターフェース (I / F) 1 1 1 を経由してデータ通信ネットワーク 1 0 7 に結合し、入出力 (I / O) バス 1 0 8 と入出力インターフェース 1 1 0 を経由してローカルストレージシステム 1 0 4 に結合している。

40

【 0 0 4 2 】

ローカルストレージシステム 1 0 4 はディスクコントローラ 1 4 1 を備え、ディスクコントローラ 1 4 1 は、入出力 (I / O) バス 1 0 8 に結合する為の入出力インターフェース 1 3 0 とデータ通信ネットワーク 1 0 7 に結合する為のネットワークインターフェース 1 3 1 とを有する。

【 0 0 4 3 】

ローカルストレージシステム 1 0 4 は更に、ディスクユニット 1 4 0 で構成されるスト

50

レージボリューム 1 4 2 を備え、ディスクコントローラ 1 4 1 が入出力インターフェース 1 3 7 , 1 3 8 を経由してデータバス 1 3 9 によって ストレージボリューム 1 4 2 に結合する。

【 0 0 4 4 】

ディスクコントローラ 1 4 1 自身は、中央処理ユニット (C P U) 1 3 3 を有し、 C P U 1 3 3 は内部バス 1 3 2 経由でメモリ 1 3 4 に結合し、更に本ディスクコントローラ 1 4 1 の各種インターフェース (即ち、入出力インターフェース 1 3 0 , 1 3 8 及びネットワークインターフェース 1 3 1 等) に結合している。

【 0 0 4 5 】

メモリ 1 3 4 は、ホスト プロセッサ 1 0 1 からの I / O 要求に対応して、ストレージボリューム 1 4 2 へのリード、ライトデータを格納する為のキャッシュメモリ 1 3 5 を有する。メモリは更に、コントロール情報 1 3 6 の如きある種のデータ構造や情報を維持する。

10

【 0 0 4 6 】

C P U 1 3 3 は、ストレージボリューム 1 4 2 に対する (例えばホストプロセッサ 1 0 1 により齎される) あらゆるデータイメージの変更をリモートストレージシステム 1 0 5 , 1 0 6 に送信する、リモートコピー処理を通常通りに実行する事が好ましい。

【 0 0 4 7 】

かくして、リモートストレージシステム 1 0 5 , 1 0 6 は、ローカルストレージシステム 1 0 4 で維持されるデータをミラーする。

20

【 0 0 4 8 】

要約すると、リモートコピー処理は以下の如く動作する。

【 0 0 4 9 】

ディスクコントローラ 1 4 1 は、ストレージボリューム 1 4 2 のデータ (データイメージ) に対する I / O ライト要求、即ち、追加、変更、削除、あるいはその他の更新を受け取ると、そのデータをストレージボリューム 1 4 2 に書き込む。ストレージボリューム 1 4 2 へのこの様な各書き込み又は少なくともデータイメージの一部分はリモートサイト 1 4 にミラーされる為に、データメッセージが生成され、リモートコピーキュー (ここでは示されていない) に格納される。

【 0 0 5 0 】

30

リモートコピーキューは、C P U 1 3 3 上で稼動しているリモートコピープロセスにより周期的に監視されている。一つ以上のデータメッセージを有するキューが検出されると、データメッセージは取り出されて、リモートストレージシステム 1 0 5 , 1 0 6 の各々に送信され、データメッセージ中のデータが書き込まれ、結果として、ローカルサイト 1 2 のデータに対するデータイメージが更新されることになる。

【 0 0 5 1 】

データメッセージに関する情報は、ストレージシステム 1 0 4 、1 0 5 、1 0 6 でヒストリカル情報として保持され、ローカルストレージシステム 1 0 4 から送信済み、送信途上、又はリモートストレージシステム 1 0 5 , 1 0 6 での受信完了を示している。各ストレージシステムは、ヒストリカル情報をキュー構造で管理する事が好ましい。

40

【 0 0 5 2 】

図 2 はキュー構造を示している。図 2 で示す如く、ディスクコントローラ 1 4 1 は、ワーキングキュー 1 1 0、ロールバックキュー 1 1 1、中間キュー 1 1 2、およびライトヒストリキュー 1 1 3 を有する。

【 0 0 5 3 】

ワーキングキュー、ロールバックキュー、中間キュー、及びライトヒストリキュー 1 1 0 ~ 1 1 3 は、リモートストレージシステム 1 0 5 (各々キュー 1 2 0、1 2 1、1 2 2 及び 1 2 3)、1 0 6 (各々キュー 1 3 0、1 3 1、1 3 2 及び 1 3 3) でミラーされている。これらのキューは、先入れ先出し (F I F O) 構成で実装される。

【 0 0 5 4 】

50

I/Oリードライト要求は典型的に、コマンドエン트리とこれに伴う又は継続するデータ（ライト要求の場合）により構成される。

【0055】

コマンドエントリは、データが書き込まれるストレージボリューム142上（ライト要求の場合）、又はデータが読み出される（リード要求の場合）ボリューム上のロケーションを特定し、その他実装によって必要とされる情報を特定する。

【0056】

ホストプロセッサ101より、リモートストレージシステム105, 106でミラーされているデータイメージの変更を伴うI/O要求を受けとった場合には、コマンドエントリには、シーケンス番号が割り当てられる。コマンドエントリはコマンドエントリと割り当てられたシーケンス番号で形成される事になる。

10

【0057】

コマンドエントリは次いで、ワーキングキュー110に格納される。キューは、リモートストレージシステム105, 106にデータを送信するに先立つ、データとデータメッセージに対するヒストリカル情報を構成する。

【0058】

コマンドエントリがワーキングキュー110に滞在している間に、対応するライト要求が処理される。処理には、ライト要求に対応するデータ（要求と共に受信するか、要求に継続して受信するかは、ホストプロセッサ101とローカルストレージシステム104に用いられる通信プロトコルによる）を受信し、データの為にキャッシュメモリ135のエリアを割り当て、当該エリアに受信データを書き込む事が含まれる。

20

【0059】

キャッシュメモリ135内のデータへのポインターは、対応するコマンドエントリに関連付けられる。I/O要求データが正常に受信されたか、エラーが検出されたかのステータスを示すステータスメッセージが、I/O要求の発行元に返される。

【0060】

図2は、ディスクコントローラ141がシーケンス番号“15”を割り当てられたライト要求を受領中の状態を示す。コマンドエントリは割り当てられたシーケンス番号とライト要求のコマンドエントリで形成される。

【0061】

30

コマンドエントリは次いで、上述されたように更なる処理の為にワーキングキュー110に格納される。コマンドエントリは、I/O要求が正常に処理された場合には、ワーキングキュー110からロールバックキュー111に移される。

【0062】

ロールバックキュー111は、後に十分に説明されるように、ローカルストレージシステム104とリモートストレージシステム105, 106でのロールバック同期化の為に用いられる一時的な保存エリアである。

【0063】

リモートストレージシステム105, 106も基本的には同じ目的の為に、ロールバックキュー111を含む同一のキュー構造を有する。

40

【0064】

ロールバックキュー111に保持されるコマンドエントリに対応するデータは、一つのストレージシステムで障害が検出された場合に廃棄処理をする場合、又は生存ストレージシステムの間での巡回処理をする場合を含めて、本発明での同期化処理に於いて使用される。

【0065】

コマンドエントリは、ロールバックキュー111から関連データがストレージボリューム142への書き込み待ちか書き込み中であるエントリを保存する中間キュー112に移動される。

【0066】

50

ローカルストレージボリューム 1 4 2への書き込みが完了すると、コマンドエントリは、プライマリストレージシステム 1 0 4をミラーするリモートストレージシステムの保存するデータイメージの更新のためにリモートストレージシステム 1 0 5 , 1 0 6に送信される、データメッセージ形成のためのリモートコピー要求を形成するために使用される。

【 0 0 6 7 】

次いで、コマンドエントリからポインターが抽出され、エントリはライトヒストリキュー 1 1 3に格納される。

【 0 0 6 8 】

上述したように、図 2 は、ホストプロセッサ 1 0 1より受信しシーケンス番号 1 5 が割り当てられ、従ってワーキングキュー 1 1 0に格納される、I/O 要求からのコマンドエントリを示す。

【 0 0 6 9 】

シーケンス番号 1 3 と 1 4 のコマンドエントリはこの時点でロールバックキュー 1 1 1に保持され、関連データのストレージボリューム 1 4 2への書き込み待ちの状態になっている。

【 0 0 7 0 】

中間キュー 1 1 2 は、シーケンス番号 1 0 , 1 1、及び 1 2を持つコマンドエントリを保持している。シーケンス番号 1 0 のコマンドエントリは書き込み中か、次の書き込み待ちの何れかである。シーケンス番号 7 , 8 及び 9 のコマンドエントリの関連データはストレージボリューム 1 4 2に書き込み済みで、従ってライトヒストリキュー 1 1 3に格納されている。

【 0 0 7 1 】

リモートストレージシステム 1 0 5 , 1 0 6 は実質的に同様なキュー構造を保持し、同様の方法で動作する。かくして、例えば、リモートストレージ 1 0 5 はシーケンス番号 1 0 を割り当てられたコマンドエントリを持つデータメッセージを受信中であり、全データが受信中の間は、ワーキングキューに格納されている。

【 0 0 7 2 】

シーケンス番号 1 0 はローカルストレージシステム 1 0 4によって割り当てられたものである。ひとたび、データメッセージが受信されると、コマンドエントリは、ワーキングキューから、シーケンス番号 6 ~ 9 のコマンドエントリを現在保持していることを図 2 が示しているロールバックキュー 1 2 1に移動される。

【 0 0 7 3 】

中間キュー 1 2 2 はシーケンス番号 5 が割り当てられたデータメッセージのコマンドエントリを有し、エントリはリモートストレージシステム 1 0 5のストレージボリューム 1 4 0に書き込み中である。

【 0 0 7 4 】

書き込みが完了すると、エントリはシーケンス番号 1 ~ 4 のデータメッセージのコマンドエントリと共に、ライトヒストリキュー 1 2 3に移動される。ライトヒストリキュー 1 2 3の深さによっては、最も早いエントリ、例えばシーケンス番号 1のエントリは、シーケンス番号 5 のコマンドエントリを格納させるために、排出されても良い。

【 0 0 7 5 】

他のリモートストレージシステム 1 0 6も同様なデータキュー (1 3 0、1 3 1、1 3 2 及び 1 3 3) を有する。図 2 では、リモートストレージ 1 0 6は現在シーケンス番号 1 2に関するデータメッセージを受信中であり、コマンドエントリはワーキングキュー 1 3 0に格納されている。

【 0 0 7 6 】

ロールバックキュー 1 3 1は現在シーケンス番号 8 ~ 1 1の制御情報を保有している。ストレージシステムがヒストリカル情報を追跡するためのキューは、メモリ 1 3 4のみならず、ストレージボリューム 1 4 2にも保持される事が好ましい。

【 0 0 7 7 】

10

20

30

40

50

ローカルストレージシステム 104 は更に、メモリ 134 中にリモートコピーステータステーブル 114 を保持し、送信済みのシーケンス番号、リモートストレージシステム 105、106 にて受信され応答が得られたシーケンス番号を特定できるようになっている。

【0078】

例えば、リモートストレージシステム 105 (テーブル 114 では“S1”で表示)で受信済みの最新のデータメッセージはシーケンス番号9で、リモートストレージシステム 106 (“S2”)で受信済みの最新のデータメッセージはシーケンス番号11である事が判る。リモートコピーステータステーブル 114 は更に、リモートストレージシステムのロールバックキュー 121、131 とライトヒストリキュー 123、133 に関する情報も保持している。

10

【0079】

かくして、リモートコピーステータステーブル 114 では、リモートストレージシステム 105、106 のロールバックキュー 121、131 は各々4データメッセージエントリ分の“長さ”を有し、10MBのデータを保持できる事がわかる。

【0080】

リモートストレージシステム 105、106 各々のライトヒストリキュー 123、133 は5個のデータメッセージに5個のエントリを有する。ライトヒストリキュー 123、133 のサイズはバイト単位の容量又は他の形でリモートコピーステータステーブル 114 に表示されても良い。

20

【0081】

図3及び図4は、本発明が解決しようとしている問題点を説明するダイアグラムである。

【0082】

図3はプライマリストレージボリューム20に対応して、少なくとも2つのセカンダリストレージボリューム21、22を備える、多重リモートコピー動作の為の構成を説明する。

【0083】

23、24の矢印は、プライマリストレージボリューム20とセカンダリストレージボリューム21、22との間のリモートコピーリンクを示す。

30

【0084】

矢印25が示す図の次の部分では、リモートコピーリンクの一つ23が停止(suspend)している。リンクはネットワーク障害又はストレージボリューム障害で停止され得る。ここで用いられる“停止(suspend)”の用語は、原因の如何によらず、二つのユニットの間でのデータ転送が途絶していることを示す。

【0085】

停止が発生すると、停止中のペアの各ストレージボリューム(プライマリストレージボリューム20とセカンダリストレージボリューム21)は自らのビットマップ情報を生成する。

【0086】

40

これらのビットマップ情報は、図3ではBp、Bsで表示されている。このビットマップ情報は、停止時のストレージボリュームの状態を反映している。ビットマップ情報生成のための方法については、後に記すが、関連米国特許No. 10/042,376に記されている。

【0087】

図の次の部分は、停止に加えて発生する三つの可能なシナリオを図で表示したものである。

【0088】

第一のシナリオは矢印27により図の左側に示され、第二のシナリオは矢印28で示され、第三のシナリオは矢印29で示されている。矢印27で示され、A-3のラベルで示される第一のシナリオでは、セカンダリストレージボリューム21が障害になった場合で

50

ある。この場合は、矢印 30 で示されるように、プライマリストレージボリューム 20 とセカンダリストレージボリューム 22 のペアが残っている。

【0089】

図 3 の矢印 28 で示される第二のシナリオでは、セカンダリストレージ 22 が障害になった場合である。この状況に対する対応は、矢印 31 と A - 6 のラベルで示される。この状況では、プライマリストレージ 20 と稼動可能なセカンダリストレージ 21 間のリンクが回復されるのを待つ。

【0090】

リンクが回復すると、プライマリストレージボリューム 20 とセカンダリストレージボリューム 21 は 2 つのビットマップ B p と B s の情報を用いて再同期化を行う事が出来る。

10

【0091】

図 3 の右側部分に、矢印 29 で第三の状態が示される。この状態では、プライマリストレージボリューム 20 が障害になり、ビットマップ B p はもはや使用不能である。この状態では、図 3 の右側下の A - 8 でラベルされる部分で示されるように、データの全コピー動作がセカンダリストレージボリューム (S 1) 20 とセカンダリストレージボリューム (S 2) 21 を再同期化するための唯一の手段である。

【0092】

図 4 は、直列接続多重コピーを実現するプライマリとセカンダリストレージの異なった構成を示す図である。

20

【0093】

図の最上段に示すとおり、プライマリストレージボリューム (P) 40 は、セカンダリストレージボリューム (S 1) 42, (S 2) 44, (S 3) 46 との間でリモートコピー機能を実現する為に、直列的に構成されている。

【0094】

矢印 41 は、セカンダリストレージボリューム 42 とセカンダリストレージボリューム 44 と間のリモートコピーリンクが停止している状況を示している。この状態では、更に図示されているように、ビットマップ B p と B s がセカンダリストレージボリューム 42 と 44 によって生成されている。この状況において、追加して発生する 4 種類の障害が存在する。これらは、図の下側の部分に矢印 1, 2, 3 及び 4 で示されている。

30

【0095】

矢印 1 では、プライマリストレージボリューム 40 が障害になっている場合を示している。この状況では、ケース B - 4 に示すように、構成は変更され、リモートコピーペア S 1 と S 2 及び S 2 と S 3 は同じままである。

【0096】

ケース 2 では、停止ペアのメンバではないセカンダリストレージボリュームが障害になっている。この状況では、ケース B - 6 に示されるように、リモートコピーペアは不変のままで、構成のみが変更されている。

【0097】

図示されているように、プライマリストレージボリューム 40 はセカンダリストレージボリューム 42 及び 44 との関連して動作を継続する。

40

【0098】

矢印 3 は停止ペアのメンバであるセカンダリストレージボリュームが障害になった場合を反映している。このケースでは、ビットマップ B p からの変更情報は失われる。

【0099】

この結果、B - 8 に示されているように、プライマリストレージボリューム 40 とセカンダリストレージボリューム 44 は再同期化不能である。プライマリストレージボリューム 40 とセカンダリストレージボリュームのペアは、新しいデータのコピーを全て転送することによってのみ再同期が可能になる。

【0100】

50

矢印 4 は最後のケースで、セカンダリストレージボリュームが障害になっている。前述した状況では、セカンダリストレージボリュームの障害でビットマップ B s が損失して、変更情報が失われる結果となる。この状況では、ケース B - 10 に示されるように、例えばセカンダリストレージボリューム 46 をターゲットに全データの再コピーを実施することが唯一の解である。

【0101】

図 3、図 4 に示したビットマップ喪失の状況は、本発明の技術により克服する事が出来る。このことについては次で説明する。

【0102】

図 3、図 4 の幾つかの状況で説明したように、図 3 の A - 3 と A - 7 及び図 4 の B - 7 と B - 9 の場合には、ミラーペアのビットマップの一つが失われている事に注意が必要である。

【0103】

本発明の実施例によれば、ペアの停止を検出した場合、停止されたりリモートコピーペアのメンバ外のストレージボリューム上でビットマップのコピーを生成する事により、この問題は克服される。離れたロケーションにビットマップをコピーすることにより、ビットマップの一つ (B p 又は B s) が失われた場合でも、容易に再同期を実施することが出来る。図 5 に本発明の好ましい実施例を説明する。

【0104】

図 5 の上部で示す通り、プライマリストレージボリューム 50 はセカンダリストレージボリューム 52、54 とペア状態になっている。矢印 56 で示す図の次の部分では、プライマリストレージボリューム 50 とセカンダリストレージボリューム 52 の間での停止の発生を示している。

【0105】

ケース A 2' で示す通り、この停止の発生を契機に、プライマリストレージボリューム 50 は異なるセカンダリストレージボリューム 54 上にビットマップのコピーを生成する。セカンダリストレージボリューム 54 は、プライマリストレージボリューム 50 をミラーするボリュームとは異なるボリュームである。ビットマップ B p' は、ビットマップ B p と必ずしも同一でなくても良い。セカンダリストレージボリューム 54 でのビットマップコピーに対する最小限の要求は、その開始時間がオリジナルコピー B p の開始時間に比べて同じかより古い事である。

【0106】

ビットマップの開始時刻がより古い場合は、ビットマップコピーが示す全ての変更は、単に同じデータを同じアドレスに上書きする事になるだけである。本発明の親出願では、ストレージが停止発生時と同時刻からのビットマップコピーを生成する技術を開示しているが、本発明では、ストレージが同時刻のビットマップコピーを生成するのが不可能な場合でも、代わりにより古いビットマップコピーを用いる事により十分に満足できる解を得られることを述べている。

【0107】

図 5 のケース A - 7 で示す通り、プライマリビットマップ B p が損失しても、セカンダリストレージボリューム 52 と 54 はプライマリビットマップ B p のコピー B p' を用いて再同期することが可能である。

【0108】

本再同期化については、図 5 の下側部分 A - 8 に示されている。

【0109】

図 6 は本発明のもう一つの好ましい実施例を示す。図 6 は図 4 と同様に、図の上部での障害と同じ状態を示す。

【0110】

B - 2 でリモートコピー動作が停止する。本状態の検出を契機に、ストレージシステムはビットマップ B p のコピーを B - 2' に示されるように他のボリュームに作る。

【 0 1 1 1 】

図示されているケースでは、プライマリストレージボリューム 6 0 はプライマリビットマップ B p のコピー B p ' を受け取り、セカンダリストレージボリューム 6 6 はビットマップ B s のコピー B s ' を受け取る。

【 0 1 1 2 】

かくして、B - 7 で示されるセカンダリストレージボリューム 6 2 の障害では、ビットマップ B p が失われても、ビットマップを用いて、再同期化が不能になることを避ける事が出来る。B - 8 に示す通り、セカンダリストレージボリューム 6 4 はプライマリストレージボリューム 6 0 に再同期化する事が可能である。

【 0 1 1 3 】

もし停止がセカンダリストレージボリューム 6 2 の代わりにセカンダリストレージボリューム 6 4 の障害で発生しても、セカンダリストレージボリューム 6 6 はセカンダリストレージボリューム 6 2 に再同期化可能である。

【 0 1 1 4 】

図 7 は本発明の好ましい一つの実現方法を説明するためのフローチャートである。同期型リモートコピーでは、各ホスト I / O 動作がリモート I / O 動作を引き起こす。

【 0 1 1 5 】

かくして、図 7 に記す処理は、ホスト I / O 毎に実行される。ローカルストレージは、図 7 の処理が完了するまでは、ホストに対するステータス報告を行わない。この準備動作を図 7 a に示す。

【 0 1 1 6 】

図 7 に、ローカルストレージシステム 1 0 4 (図 1) がホストプロセッサ 1 0 1 (図 1) から受け取った更新データをリモートストレージシステム 1 0 5 , 1 0 6 にコピーする為のリモートコピー動作の基本的ステップを説明する。

【 0 1 1 7 】

ホストプロセッサ 1 0 1 からのローカルストレージシステム 1 0 4 で保存されているデータイメージを変更する I / O ライト要求は、リモートストレージシステム 1 0 5 , 1 0 6 で保存されているミラーデータイメージに対しても同様の変更を必要とする。

【 0 1 1 8 】

I / O ライト要求は、割り当てられたシーケンス番号と要求に対するデータのポインターを備えたコマンドエントリを生成する。コマンドエントリは、全データが受信され、受領応答がホストプロセッサ 1 0 1 に返送されるまでの間は、ワーキングキュー 1 1 0 に格納されている。

【 0 1 1 9 】

次いで、コマンドエントリはロールバックキュー 1 1 1 に移動される。ロールバックキュー 1 1 1 が満杯になるか、フラッシュコマンドが受信されると、コマンドエントリは中間キュー 1 1 2 に移動される。中間キューに滞在する間に、要求された対応データはストレージボリューム 1 4 2 に書き込まれる。

【 0 1 2 0 】

C P U 1 3 3 上で実行しているリモートコピー処理は周期的に中間キューの内容を調査して、各リモートストレージシステム 1 0 5 , 1 0 6 にコピーすべき更新データコマンドのエントリが存在しないかを判定する。

【 0 1 2 1 】

次いで図 7 を参照すると、ローカルストレージシステム 1 0 4 は、各リモートストレージ 1 0 5 , 1 0 6 が次のデータメッセージを受信できるかをチェックする。チェックは R C ステータステーブル 1 1 4 を参照して行う。R C ステータステーブル 1 1 4 により、各リモートストレージシステム 1 0 5 , 1 0 6 がどのデータメッセージを受信して応答済みか、又リモートストレージシステム 1 0 5 , 1 0 6 が保持する各種キューのサイズ情報等がローカルストレージシステム 1 0 4 に伝わる。

【 0 1 2 2 】

10

20

30

40

50

これにより、ローカルストレージシステム 104 は、特定のリモートストレージシステムが別のデータメッセージや関連データを受信できる余裕があるか否かを判定する。もしなければ、ステップ 501 から抜け、プロセス 6001 にて停止が発生しているかの判定が開始される。(あとで議論される。)

更に、リモートストレージシステム自身が、例えば SCSI の “ ビジー (0 x 0 8) ” 又は “ キュー満杯 (0 x 2 8) ” の何れかを用いて、更なるデータメッセージ受信の不可を返答する事も可能である。同期型リモートコピー動作の場合は、リモート (セカンダリ) ストレージシステムのキュー状態のチェックは不要である。

【 0 1 2 3 】

もし対象のリモートストレージシステムにデータメッセージを受信するのに十分な余裕がある場合には、ローカルストレージシステム 104 はリモートコピー (RC) コマンドを、ステップ 503 にてデータが後続するデータメッセージの形で、リモートストレージシステム (たとえばリモートストレージシステム 105) に発行する。

【 0 1 2 4 】

次いで、ローカルストレージシステム 104 は、データメッセージが受信されたか否かを示すステータス報告の受信を待つ (ステップ 504)。

【 0 1 2 5 】

受信後には、ローカルストレージシステム 104 は、ステップ 505 にて、全てのリモートストレージシステムがデータメッセージでの更新を完了したかをチェックする。もし完了していなければ、処理はステップ 507 に移動し、RC ステータステーブル 114 を更新して、対象リモートストレージシステムはデータメッセージを受信した事を反映させて、次のリモートストレージシステムがデータメッセージを受信するためにステップ 501 に戻る。

【 0 1 2 6 】

しかし、もしステップ 505 で全てのリモートストレージシステムがデータメッセージを受信した事が判明したら、データメッセージ (コマンドエントリ) は、ステップ 506 にてライトヒストリキューに移動され、RC ステータステーブル 114 はステップ 507 で更新され、このデータメッセージに対する処理は 508 で完了する。

【 0 1 2 7 】

ステップ 501 の判定にて、リモートストレージがデータメッセージを受信できない場合には、処理フローはステップ 6001 に移動し、停止状態が発生しているかを判定する。この動作での最初の判定は、リモートコピーペアのステータスについてである。リモートコピーペアの継続を妨げる何らかの異常事態が存在すれば、リモートコピーペアは停止しなければならないと判定される。

【 0 1 2 8 】

本判定の結果として、停止が発生し、主要な処理はステップ 6002 に示すビットマップのコピーを生成する事である。ビットマップを生成する方法は、以下に議論される。

【 0 1 2 9 】

しかしながら、以下でも議論されるように、ひとたびビットマップが生成されると、制御フローはステップ 508 に移動し、多重リモートコピー処理を完了し、再同期処理へ移行する。本発明の親出願の技術を用いることにより、プライマリストレージは停止ポイントを検出し、本ポイントをシーケンス番号で示し、ビットマップのコピーを有する他のストレージに提供することができる。

【 0 1 3 0 】

図 8 は、リモートストレージシステムがリモートコピー処理のデータメッセージを受信するためのステップを広範囲に説明するフローダイアグラムである。最初のステップ 7001 は停止が発生しているかを判定する。この判定は、図 7 に関連して説明したものと同一技術を用いて実行される。

【 0 1 3 1 】

もし停止が発生していたら、図 9 に関連して以下に説明する技術を用いて、ビットマッ

10

20

30

40

50

プのコピーがステップ 7 0 0 2 で生成される。

【 0 1 3 2 】

もし停止が発生していなければ、本処理フローはステップ 6 0 1 に移動する。ステップ 6 0 1 にて、リモートストレージシステムはデータメッセージを受信して、ステップ 6 0 2 にて、このデータメッセージのためのキュー資源の使用状態をチェックする。

【 0 1 3 3 】

言い換えると、利用できる余裕があるかを判定する。ステップ 6 0 2 での判定は、他のリモートストレージシステムのキュー内容に比較して、各キュー（ロールバック、中間、及びライトヒストリキュー）の整合状態も考慮して判定される。

【 0 1 3 4 】

もしステップ 6 0 2 での判定の結果、リモートストレージシステムがデータをこの時点では受領できない場合は、ステップ 6 0 2 を抜け、処理はステップ 6 0 6 に移動して、リモートストレージシステムは“ ビジー ”ステータスメッセージをローカルストレージシステムに返送し、受信処理より抜ける。

【 0 1 3 5 】

次いで、ローカルストレージシステム 1 0 4 は、後で当該トランザクションを再度試行しなければならないと判定する。反対に、ステップ 6 0 2 にて、データが受信可能と判定された場合には、データはステップ 6 0 3 にて受信される。

【 0 1 3 6 】

ステップ 6 0 4 にてリモートストレージシステムはデータ転送完了のステータスを返送して、ステップ 6 0 5 にて、データメッセージはメッセージ受信のために用いられたワーキングキューよりロールバックキューに移動される。ステップ 6 0 7 にてこの処理は完了する。

【 0 1 3 7 】

図 9 は、例えば図 7 のステップ 6 0 0 2 又は図 8 のステップ 7 0 0 2 で実行される、ビットマップのコピーを生成する手順の好ましい実施例を説明するフローチャートである。

【 0 1 3 8 】

図 1 0 は、図 9 のフローチャートの説明で用いられる構成例を説明するダイアグラムである。

【 0 1 3 9 】

図 1 0 のストレージ 7 2 は、ボリューム 7 2 と 7 4 で構成される停止中のリモートコピーペアのプライマリストレージである。

【 0 1 4 0 】

停止ポイントからのビットマップを生成するのに用いることができる 3 つの異なる技術が存在する。停止ポイントからのビットマップを生成する一つの技術は、本発明の親出願で述べられている技術を用いることである。この技術は図 9 に要約されている。

【 0 1 4 1 】

ストレージは、変更情報のすべてを反映する、ロールバックキューとライトヒストリキューを持っているため、多重リモートコピー環境での各種ストレージボリュームは同一時刻又は同一シーケンス番号からのビットマップを生成する事が出来る。

【 0 1 4 2 】

図 9 のステップ 8 0 0 1 で示すとおり、ストレージは自分自身が停止中のリモートコピーペアのプライマリストレージボリュームかセカンダリストレージボリュームかをチェックする。

【 0 1 4 3 】

図 1 0 で示されるケースでは、ストレージボリューム 7 2 が本ペアのプライマリストレージボリュームで、ストレージボリューム 7 4 がセカンダリストレージボリュームである。言い換えれば、“ プライマリ ” 及び “ セカンダリ ” の用語は、相対的な用語である。プライマリストレージボリューム 7 2 が実行する処理は、ステップ 8 0 0 2 にて示す。

【 0 1 4 4 】

10

20

30

40

50

ステップ 8 0 0 3 にて、ストレージボリューム 7 2 は、時刻又はリモートコピーコマンドからのシーケンス番号で表示される停止ポイントを検出する。これを契機に、ストレージボリューム 7 2 は、“ビットマップコピー生成コマンド”をプライマリストレージボリューム 7 0 に対して発行する。

【 0 1 4 5 】

次のステップ 8 0 0 4 で示す通り、停止中のリモートコピーペアのプライマリストレージボリューム 7 2 が例えば図 1 0 のストレージボリューム 7 8 のような、もう一つのセカンダリストレージを持っている場合も、同様の生成ステップが実行される。

【 0 1 4 6 】

かくして、ステップ 8 0 0 5 で示す通り、ビットマップはセカンダリストレージボリューム 7 8 にも生成される。

【 0 1 4 7 】

ストレージボリューム 7 4 で実行される処理は、既に述べられたものと同様である。停止中のリモートコピーペアのセカンダリストレージが、ストレージボリューム 7 6、8 0 のように、一つ以上のセカンダリストレージ自身を有する場合には、ストレージ 7 4 上のセカンドビットマップのコピーは、ストレージボリューム 7 6、8 0 の何れかに作成される。この処理はステップ 8 0 0 7 に示される。

【 0 1 4 8 】

図 1 1 は本発明の他の実施例を示す。図 1 1 に於いて、行は各種ストレージを示し、図の番号のついたブロックは各種データ転送動作を示す。

【 0 1 4 9 】

例えば、図 1 1 の最上段は 1 4 個のデータがホストシステムからリモートシステムに転送された事を示す。このデータが受け取られると、セカンダリストレージボリューム S 1、S 2、及び S 3 に転送される。

【 0 1 5 0 】

データブロック 1 の開始時刻の相違は、データを一つのロケーションから他のロケーションに転送するのに必然的に発生する処理及び転送遅延を反映している。

【 0 1 5 1 】

図 1 1 では、動作 1 0 0 0 で、データブロック 4 をセカンダリストレージボリューム S 1 から S 2 に送信する間に、停止が発生した事を想定している。これを時刻 t 0 とする。実際、第四番目のリモートコピー要求は失敗している。

【 0 1 5 2 】

結果として、少なくともビットマップ B p と B p ' は変更情報を含まなければならない。

【 0 1 5 3 】

時刻 t 1 でセカンダリストレージボリューム S 1 は動作 1 0 0 1 に示すように、ビットマップのコピーを持っているプライマリストレージボリューム P に“停止検出”メッセージを発行する。

【 0 1 5 4 】

時刻 t 2 にて、プライマリストレージボリューム P は“停止検出”メッセージを受信し、リモートコピー動作をストップする。

【 0 1 5 5 】

次いで、プライマリストレージボリューム P は、応答 1 0 0 2 を要求の送信者に送信する。プライマリストレージボリューム P が停止発生以後の変更を記録できれば、リモートコピー動作の凍結中にも、プライマリストレージボリューム P はホスト結合を維持する。

【 0 1 5 6 】

もしプライマリストレージボリューム P がこの情報の維持が不可能なら、ビットマップのコピーが生成される迄、プライマリストレージボリューム P はホストからディスクコネクトしなければならない。説明の実施例では、この変更情報はキャッシュメモリ又はディス

10

20

30

40

50

クに格納される。

【 0 1 5 7 】

時刻 t_3 で、セカンダリストレージボリューム S 1 はプライマリストレージボリューム P からの応答を受信する。セカンダリストレージボリューム S 1 は変更情報と共に “ビットマップ生成” メッセージを送信する（動作 1 0 0 3 で示す）。変更情報は、停止ポイント t_0 から “応答” 受信ポイント（時刻 t_3 ）までカバーする。

【 0 1 5 8 】

時刻 t_4 で、“ビットマップ生成” メッセージが変更情報と共に受信される。これを契機に、プライマリストレージボリューム P は、ビットマップ B p に対応するビットマップのコピー B p' の生成を開始する。

10

【 0 1 5 9 】

この処理が終了次第、プライマリストレージボリューム P と セカンダリストレージボリューム S 1 との間のリモートコピー動作を再開する。これは動作 1 0 0 4 として示される。

【 0 1 6 0 】

時刻 t_5 にて、セカンダリストレージボリューム S 1 は同様にリモートコピー動作を開始する。

【 0 1 6 1 】

ビットマップ B s のコピー生成については、動作 1 0 0 5 に関連して、後に示す。停止が 1 0 0 5 で（又はより早く）発生したとすると、セカンダリストレージボリューム S 1 と S 2 のみが、停止を認識している。リモートコピー要求 No. 4 は失敗し、セカンダリストレージボリューム S 2 には届かない。

20

【 0 1 6 2 】

かくして、少なくともビットマップ B s と B s' はリモートコピー要求 No. 3 の処理後の変更情報を含まなければならない。従って、セカンダリストレージボリューム S 2 は停止中のペアのビットマップ B s の生成を開始する。

【 0 1 6 3 】

時刻 t_2 にて、セカンダリストレージボリューム S 3 は“停止検出”メッセージを受信して、これを契機に、ビットマップのコピー B s' の生成を行う。セカンダリストレージボリューム S 3 は、その後、要求の送信者（セカンダリストレージボリューム S 2）に対して、“応答”を送信する。これは、図 1 1 の動作 1 0 0 6 に示されている。

30

【 0 1 6 4 】

セカンダリストレージボリューム S 2 が“応答”を受信すると、セカンダリストレージボリューム S 2 と S 3 との間のリモートコピーリンクを再開する。

【 0 1 6 5 】

この結果、セカンダリストレージボリューム S 2 は、例えば図 1 1 に示すような要求 2 0 ~ 2 5 の実行等の他の用途に用いられる事が可能となる。

【 0 1 6 6 】

図 1 1 に関連して述べられた動作を要約すると、その動作の共通の手順は；

（１）第一のストレージボリュームで停止を検出して、ビットマップを生成する。

40

【 0 1 6 7 】

（２）本ビットマップのコピーを有するべき第二のストレージボリュームを選択する。

【 0 1 6 8 】

（３）第一と第二のストレージボリュームとの間のリモートコピー動作を凍結する。

【 0 1 6 9 】

もしこれら 2 つのストレージボリュームの間のリモートコピー動作が第二のストレージボリュームから第一の方向であれば、第二のストレージボリュームが、第一と第二の間のリモートコピー動作を凍結するために、第一ストレージボリュームに“凍結”メッセージを送信しなければならない。

【 0 1 7 0 】

50

(4) 第一のストレージボリュームでの変更情報(停止ポイントからの変更を含む)を第二のストレージボリュームに送信する。

【0171】

(5) 第二のストレージボリュームでビットマップのコピーを生成する。

【0172】

(6) 第一と第二のストレージボリューム間のリモートコピー動作を再開する。

【0173】

図12では、ローカルストレージシステム(又はデージーチェーン構成での中間的サイトになっている場合のリモートストレージシステム)の障害が検出された場合の同期化処理で行われるステップを示す。ステップ701で障害が検出される。障害検出は多数の方法でなされる。

10

【0174】

例えば、ストレージシステム同士はハートビートメッセージを交換しあって、もし未検出なら障害を示す。又は、ストレージシステムは自ら障害を検し、しかもその障害をシステムの他の要素に対して報告出来ることもある。又は他の全ての慣用的障害検出技術が適用されても良い。

【0175】

障害が検出されると、残りのストレージシステム、即ち図2の構成に示すストレージシステム105, 106などは、互いに交信しあって、新しいリモートコピーマネージャを互選する。

20

【0176】

選出されたストレージシステムは、残り全てのストレージシステムのデータイメージが同一状態に同期化されるように、リモートコピー動作を制御する。あるいは、リモートコピーマネージャは、システム管理者により事前に指定されても良い。

【0177】

しかしながら、リモートコピーマネージャが決まると、選択されたマネージャはローカルストレージシステムとなり、ステップ703にて、残りのストレージシステムに維持される各種キューの内容と構造に関する情報収集を行う。この情報には、ロールバックキューとライトヒストリキューのエントリ範囲(エントリ数)即ち、データメッセージを含むエントリ数等、が含まれる。

30

【0178】

例えば、図2において、ローカルストレージシステム104が障害になり、ステップ702にて、リモートストレージシステム105が選択又は他の判定によりリモートコピーマネージャになると、リモートストレージシステム106は、ロールバックキュー131内にデータメッセージ8、9、10、及び11がありライトヒストリキュー133内にデータメッセージ3、4、5、及び6があることを報告する。

【0179】

障害検出時には、リモートストレージシステム105は、好都合にも、中間キューの内容を空にしてストレージボリュームに対するデータ書き込みを維持できるため、リモートストレージシステム106の中間キューの内容は、短時間の間にライトヒストリキューに追加される。

40

【0180】

ステップ1101にて、次にペアは停止中か否か判定される。もし、停止中でなければ、以下に示す通り、処理はステップ704に移る。

【0181】

これに対して、停止中の場合は、停止中のペアは再同期化されなければならない。これはステップ1102に示され、処理は図13で示され、そこで説明される。

【0182】

残りのストレージシステムについてのリモートコピー環境に関する情報を受領し、その後一つに障害が発生した場合、選択されたリモートコピーマネージャは、ロールバック処

50

理を行うかロールフォワード処理を行うかの決定を行う。

【0183】

典型的には、この決定は、リモートコピー環境の管理者がユーザによって、障害時での同期化処理を最善に行う方法として、フラグをセットしておくことにより、事前に決定されている。

【0184】

ステップ704にて、ロールバック処理が用いられることを決定した場合は、ステップ705にて、リモートコピーマネージャは、全てのストレージシステムが共有している番号中で、最大のシーケンス番号を持つデータメッセージを決定する。

【0185】

例えば、図2の例では、シーケンス番号9を持つデータメッセージが条件に該当する。それ故に、ステップ706にて、選択されたリモートコピーマネージャは、ロールバックコマンドを他のすべてのストレージシステムに対して発行し、受領時にシーケンス番号9を超えるデータメッセージを破棄させる。

【0186】

かくして、再び図2に戻って、リモートストレージシステム106は、ロールバックコマンドを受領すると、シーケンス番号10と11のデータメッセージを破棄して、本処理はステップ715で完結する。

【0187】

もしこれに対して、ステップ704にて、ロールフォワードが決定されたなら、ステップ707にて、ロールフォワードシーケンス番号が決定される。

【0188】

本処理は、リモートコピー処理が各種のロールバックキューとライトヒストリキューの内容を比較して、最新のデータメッセージ番号、もしある場合、を有するストレージシステムを決定する事よりなされる。

【0189】

かくして、図2にて、選出されたリモートコピーマネージャがストレージシステム105の場合は、リモートストレージシステム106がシーケンス番号10と11を有するデータメッセージを持っており、自分自身は持っていない事を悟る。したがって、データメッセージ10と11をリモートストレージシステム106からコピーさせることによって、各システムで維持されていたデータイメージが同期化される。

【0190】

かくして、ステップ708にて、選出されたリモートコピーマネージャが最新のデータストレージシステムを決定すれば、ステップ708より709に移り、選択されたRCマネージャは、最新データメッセージを有するストレージシステムから、更新データを受け取る。これは選択されたRCマネージャ自身であることも、又他のストレージシステムのひとつである場合もある。

【0191】

いずれにしても、選択されたRCマネージャは更新データを受け取り、ステップ710にて、更新を必要としている他のストレージシステムに更新データを選択的に又は部分的に転送して、全てのリモートストレージシステムのデータイメージの同期化処理を行う。処理はステップ715で完結する。

【0192】

反対に、ステップ708にて、同期化処理は最新の更新データメッセージを持つストレージシステムが実施すべきと判定され、そのストレージシステムが選択されたリモートコピーマネージャではないなら（又は他のストレージシステムが同期化システムと前もって決定されたら）、RCマネージャは、ストレージシステムが必要としている更新情報を比較して、最新の更新データを持つストレージシステムに対して、更新情報を送付する。

【0193】

次いで、ステップ712にて、データメッセージ形式の更新データが、データイメージ

10

20

30

40

50

の同期化を必要としている全てのストレージシステムに送られ、本処理はステップ 7 1 5 で完結する。

【 0 1 9 4 】

もしロールバック処理が選択される場合は、本発明の親出願で議論されたように、例えば、リモートストレージシステム 1 0 5 , 1 0 6 (図 2) のロールバックキュー 1 2 1 , 1 3 1 は整合されている事が好ましい。

【 0 1 9 5 】

整合処理は、ローカルストレージシステム 1 0 4 にて、CPU 1 3 3 上で走行するリモートコピー処理とリモートコピーステータステーブル 1 1 4 (リモートストレージシステム 1 0 5 , 1 0 6 が維持している各種キューの内容の情報が提供されている) を用いて実行される。

10

【 0 1 9 6 】

例えば、あるリモートストレージシステムはローカルストレージシステムからのリモートコピーデータメッセージを受信不能で、他のリモートストレージシステムは可能と言う事態が発生する。

【 0 1 9 7 】

このような場合は、同期処理が必要なときには確実に実行されるように、他のリモートストレージシステムと比較時に、ひとつのリモートストレージシステムのキューの中のデータメッセージ間に少なくとも一個の共通データメッセージが存在することが必要であることを、同期化の目的のためには、考慮しなければならない。

20

【 0 1 9 8 】

図 1 3 は再同期化処理の詳細を説明するフローダイアグラムである。最初に、ビットマップのコピーを有する各ストレージボリュームは、リモートコピーペアのステータスをチェックし、ステータスに関する情報を交換する。これはステップ 1 2 0 1 に示す。ビットマップのコピーを用いて再同期化を必要とする環境では、再同期化処理は停止状態にあるペア間で必要であるのみならず、ビットマップコピー即ち、 B_p , B_p' , B_s , B_s' を有するストレージボリュームのペアの組合せでも必要である。

【 0 1 9 9 】

ステップ 1 2 0 2 にて、ビットマップコピーを用いた再同期化が不必要と判定されたなら、処理は終了する。一方、再同期化処理が必要なら、次いでビットマップが交換され (1 2 0 3) 、マージされ (1 2 0 4) 、ステップ 1 2 0 5 で示されるように、マージしたビットマップを用いた再同期化処理が実行される。

30

【 0 2 0 0 】

以上で説明した方法と装置により、多重リモートコピー環境でのコピーを用いた再同期化が可能になる。例えば、リモートコピーの停止処理を採用したシステムのケースでは、特に、再同期化が必要になり次第、システム管理者の介入なしに、自動的に実行される。

【 0 2 0 1 】

本発明の好ましい実施例について説明してきたが、本発明の範囲を逸脱する事なく、多数の改変が可能である事に注意願いたい。

【 図面の簡単な説明 】

40

【 0 2 0 2 】

【 図 1 】 は、本発明の一実施例を実装するローカル及び複数のリモートストレージシステムを含むデータ処理システムを示すブロックダイアグラムである。

【 図 2 】 は、ローカル及びリモートストレージシステムの各々で実装され、ローカルストレージシステムによってリモートストレージシステムへ転送されるデータ更新の履歴情報を保持する為のキューの構造を示すダイアグラムである。

【 図 3 】 は、多重リモートコピー環境で停止を引き起こす各種の障害形態を示すダイアグラムである。

【 図 4 】 は、異なった多重リモートコピーアーキテクチャにおいて停止を引き起こす各種の障害形態を示すダイアグラムである。

50

【図 5】は、本発明の方法の一実施例を示すダイアグラムである。

【図 6】は、本発明の他の実施例を示すダイアグラムである。

【図 7】は、停止発生時の多重リモートコピープロセスを示すフローチャートである。

【図 7 a】は、図 7 に関連する準備用ステップを示すフローチャートである。

【図 8】は、停止発生時の他の多重リモートコピープロセスを示すフローチャートである。

。

【図 9】は、停止発生時のビットマップのコピーを生成する為の好ましい方法を示すフローチャートである。

【図 10】は、図 9 の説明に用いられるダイアグラムである。

【図 11】は、ビットマップのコピーを生成するための動作の詳細シーケンスを示す。

10

【図 12】は、再同期化の一方法を示すフローチャートである。

【図 13】は、再同期化の他の方法を示すフローチャートである。

【符号の説明】

【0203】

101・・・プライマリ ホスト（ホストプロセッサ）

102・・・リモートホスト、

103・・・リモートホスト

104・・・ローカルストレージシステム

105・・・リモートストレージシステム

106・・・リモートストレージシステム

20

107・・・データ通信ネットワーク

110, 130, 137, 138・・・I/O I/F（インターフェース）

111, 131・・・ネットワーク I/F（インターフェース）

132・・・内部バス

133・・・CPU

134・・・メモリ

135・・・キャッシュメモリ

136・・・コントロール情報

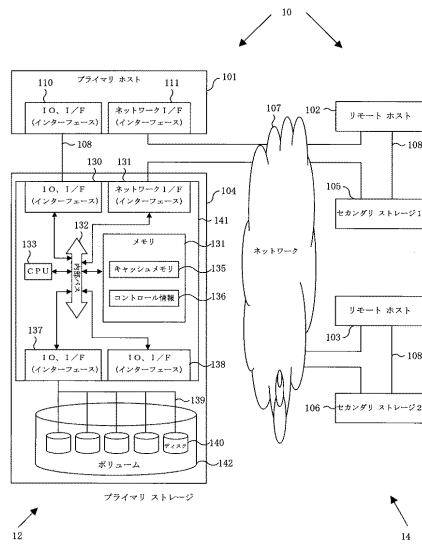
140・・・ディスクコントローラ

142・・・ストレージボリューム

30

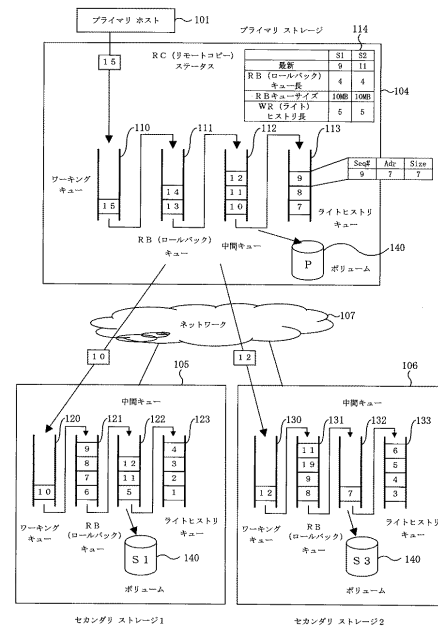
【図 1】

図 1



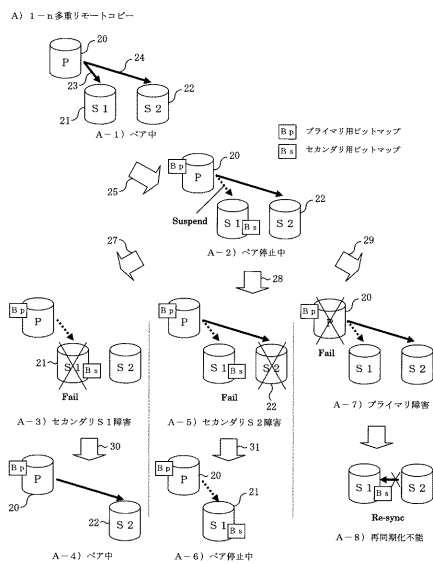
【図 2】

図 2



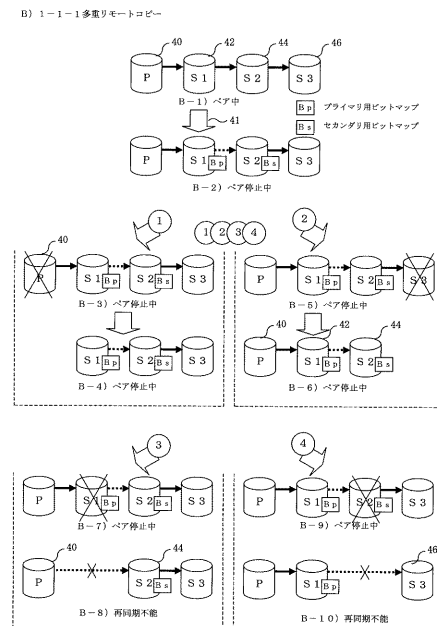
【図 3】

図 3

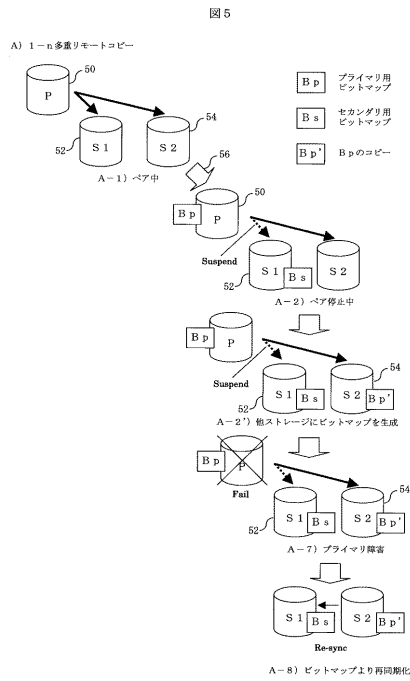


【図 4】

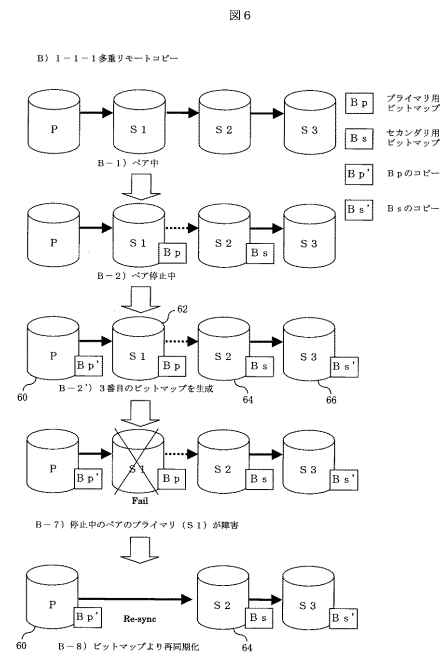
図 4



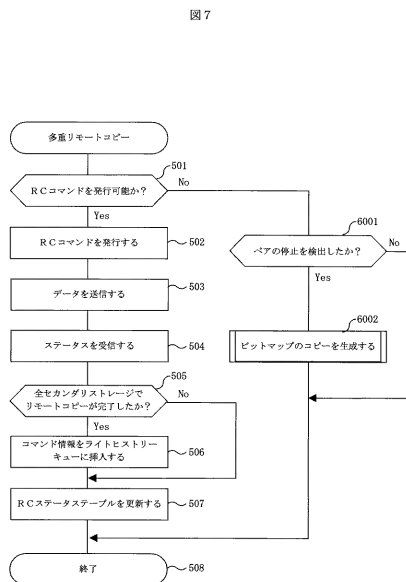
【図 5】



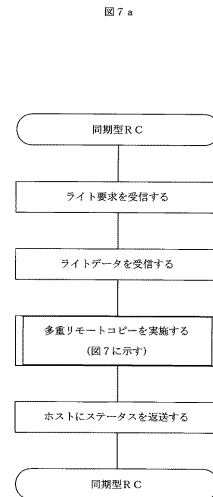
【図 6】



【図 7】

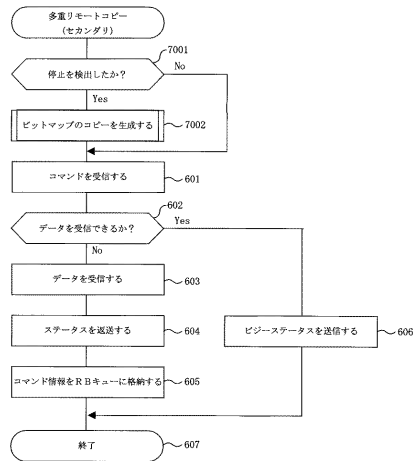


【図 7 a】



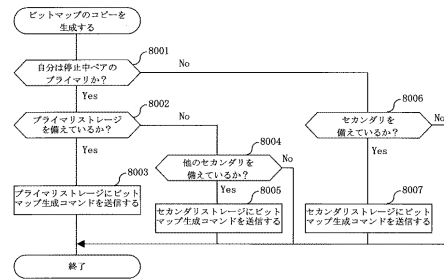
【図 8】

図 8



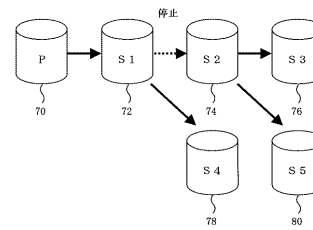
【図 9】

図 9



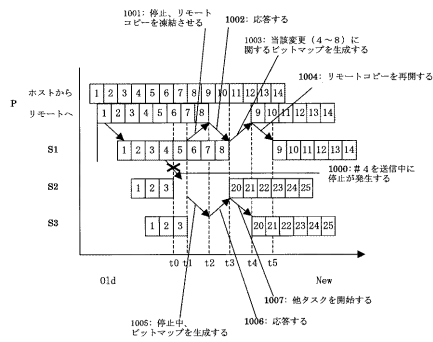
【図 10】

図 10



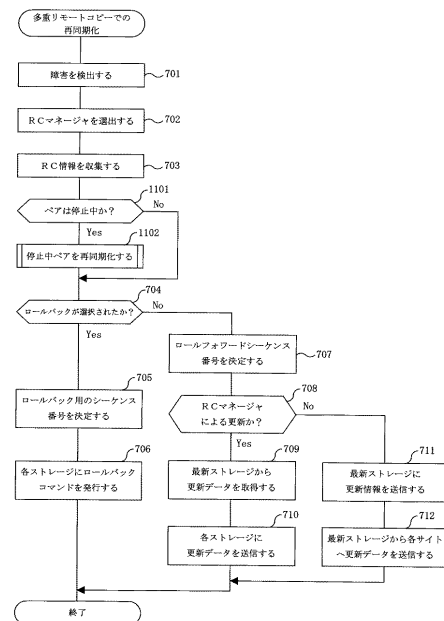
【図 11】

図 11



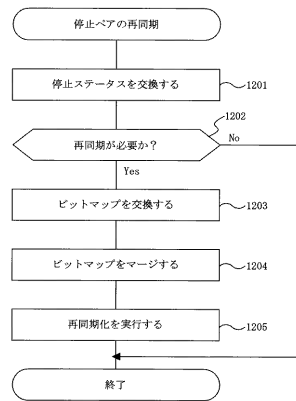
【図 12】

図 12



【図 13】

図 13



フロントページの続き

(56)参考文献 特開平10-133926(JP,A)
特開2001-290687(JP,A)
特開2000-305856(JP,A)
特開平11-338647(JP,A)
特開2003-263280(JP,A)
特開2003-122509(JP,A)
特開2003-131917(JP,A)

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