

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5286212号  
(P5286212)

(45) 発行日 平成25年9月11日(2013.9.11)

(24) 登録日 平成25年6月7日(2013.6.7)

(51) Int.Cl. F I  
G06F 3/06 (2006.01) G06F 3/06 304F

請求項の数 2 (全 21 頁)

(21) 出願番号	特願2009-223621 (P2009-223621)	(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21年9月29日(2009.9.29)	(74) 代理人	100100310 弁理士 井上 学
(65) 公開番号	特開2011-76130 (P2011-76130A)	(74) 代理人	100098660 弁理士 戸田 裕二
(43) 公開日	平成23年4月14日(2011.4.14)	(72) 発明者	平岩 友理 神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所 システム開発研究 所内
審査請求日	平成24年2月29日(2012.2.29)	(72) 発明者	西 好行 神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所 ソフトウェア事業 部内

最終頁に続く

(54) 【発明の名称】 ストレージクラスタ環境でのリモートコピー制御方法及びシステム

(57) 【特許請求の範囲】

【請求項1】

アプリケーションプログラムと、ストレージクラスタリング制御プログラムと、前記ストレージクラスタリング制御プログラムとは独立して動作するコピー制御プログラムと、  
を実行し、前記アプリケーションプログラムの実行に従って書き込みデータを送信する第一の  
ホスト計算機と、

前記第一のホスト計算機と接続し、第一のボリュームを提供する第一のストレージ装置と、

前記第一のホスト計算機及び前記第一のストレージ装置と接続し、第二のボリュームを提供する第二のストレージ装置と、

前記第一のストレージ装置及び前記第二のストレージ装置と接続し、第三のボリュームを提供する第三のストレージ装置と、

を有する計算機システムであって、

前記第一のストレージ装置は、前記第一のボリュームに前記書き込みデータを格納し、  
前記第一のストレージ装置は、前記第一のボリュームをコピー元とし、前記第二のボリュームをコピー先とする第一のコピーペアについての同期リモートコピー処理によって、  
前記書き込みデータを前記第二のストレージ装置へ送信し、

前記第一のストレージ装置は、前記書き込みデータに基づき、前記書き込みデータが格納された記憶領域のアドレスを含む情報と書き込みデータとの組である第一のジャーナルエントリを作成し、

前記第一のストレージ装置は、前記第一のボリュームをコピー元とし、前記第三のボリュームをコピー先とする第二のコピーペアについての非同期リモートコピー処理によって、前記第一のジャーナルエントリを前記第三のストレージ装置へ送信し、

前記第三のストレージ装置は、前記第一のジャーナルエントリとして送信された前記書き込みデータを前記第三のボリュームに格納し、

前記第二のストレージ装置は、前記第一のストレージ装置より送信された前記書き込みデータを前記第二のボリュームに格納し、前記書き込みデータに基づき第二のジャーナルエントリを作成し、

前記ストレージクラスタリング制御プログラムを実行することで、前記第一のホスト計算機は、前記第一のボリュームへの前記書き込みデータの書き込みが異常終了したことを検知し、前記第一のコピーペアのペア状態を変更することで検知後の前記書き込みデータを前記第二のボリュームに対して書き込み可能とするストレージ装置切り替え指示を前記第二のストレージ装置に送信し、前記ストレージ装置切り替え指示に基づく切り替え処理の完了後に、前記切り替え処理を完了したことを示す情報を出力し、

10

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記切り替え処理を完了したことを示す情報の出力の有無を監視し、前記出力を検知することにより、前記ストレージクラスタリング制御プログラムによって異常終了を検知し前記切り替え処理を完了したことを、前記ストレージ装置切り替え指示の送信より後に検知し、前記第二のボリュームをコピー元とし、前記第三のボリュームをコピー先とする第三のコピーペアについての別な非同期リモートコピー処理によるコピー処理を開始するための差分リシンク指示を前記第二のストレージ装置に送信し、

20

前記第二のストレージ装置は、前記差分リシンク指示を受信すると、前記切り替え処理が行われる前より作成した前記第二のジャーナルエントリを、前記第三のストレージ装置に送信することにより、差分リシンクを実行し、

前記計算機システムは、前記第三のストレージ装置と接続し、前記アプリケーションプログラムを実行することで前記第三のボリュームに格納したデータを読み込む第二のホスト計算機を有し、

前記第一のホスト計算機は、前記第一のコピーペアが前記ストレージクラスタリング制御プログラムの制御対象であることを示すコピー定義情報を有し、

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記第一のコピーペア及び前記第二のコピーペアを対象とするペア状態指示を受付け、前記第一のコピーペアのペア状態又は前記第二のコピーペアのペア状態を表示し、

30

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記コピー定義情報に基づいて、前記第二のコピーペア及び前記第三のコピーペアを対象とする状態変更指示は処理し、前記第一のコピーペアを対象とする状態変更指示の処理は抑止し、

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記差分リシンク指示送信に基く第三のコピーペアを対象とする副ボリューム間差分リシンクが差分コピーで処理できないことを示す情報を前記第二のストレージ装置から取得し、前記ストレージクラスタリング制御プログラムによる前記第一のコピーペアの状態変化を要因とした前記第三のコピーペアの副ボリューム間差分リシンクが失敗したメッセージを表示する

40

ことを特徴とする計算機システム。

#### 【請求項2】

アプリケーションプログラムと、ストレージクラスタリング制御プログラムと、前記ストレージクラスタリング制御プログラムとは独立して動作するコピー制御プログラムと、を実行し、前記アプリケーションプログラムの実行に従って書き込みデータを送信する第一のホスト計算機と、

前記第一のホスト計算機と接続し、第一のボリュームを提供する第一のストレージ装置と、

前記第一のホスト計算機及び前記第一のストレージ装置と接続し、第二のボリュームを

50

提供する第二のストレージ装置と、

前記第一のストレージ装置及び前記第二のストレージ装置と接続し、第三のボリュームを提供する第三のストレージ装置と、

を有する計算機システムであって、

前記第一のストレージ装置は、前記第一のボリュームに前記書き込みデータを格納し、

前記第一のストレージ装置は、前記第一のボリュームをコピー元とし、前記第二のボリュームをコピー先とする第一のコピーペアについての同期リモートコピー処理によって、前記書き込みデータを前記第二のストレージ装置へ送信し、

前記第一のストレージ装置は、前記書き込みデータに基づき、前記書き込みデータが格納された記憶領域のアドレスを含む情報と書き込みデータとの組である第一のジャーナルエントリーを作成し、

前記第一のストレージ装置は、前記第一のボリュームをコピー元とし、前記第三のボリュームをコピー先とする第二のコピーペアについての非同期リモートコピー処理によって、前記第一のジャーナルエントリーを前記第三のストレージ装置へ送信し、

前記第三のストレージ装置は、前記第一のジャーナルエントリーとして送信された前記書き込みデータを前記第三のボリュームに格納し、

前記第二のストレージ装置は、前記第一のストレージ装置より送信された前記書き込みデータを前記第二のボリュームに格納し、前記書き込みデータに基づき第二のジャーナルエントリーを作成し、

前記ストレージクラスタリング制御プログラムを実行することで、前記第一のホスト計算機は、前記第一のボリュームへの前記書き込みデータの書き込みが異常終了したことを検知し、前記第一のコピーペアのペア状態を変更することで検知後の前記書き込みデータを前記第二のボリュームに対して書き込み可能とするストレージ装置切り替え指示を前記第二のストレージ装置に送信し、前記ストレージ装置切り替え指示に基づく切り替え処理の完了後に、前記切り替え処理を完了したことを示す情報を出力し、

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記切り替え処理を完了したことを示す情報の出力の有無を監視し、前記出力を検知することにより、前記ストレージクラスタリング制御プログラムによって異常終了を検知し前記切り替え処理を完了したことを、前記ストレージ装置切り替え指示の送信より後に検知し、前記第二のボリュームをコピー元とし、前記第三のボリュームをコピー先とする第三のコピーペアについての別な非同期リモートコピー処理によるコピー処理を開始するための差分リシンク指示を前記第二のストレージ装置に送信し、

前記第二のストレージ装置は、前記差分リシンク指示を受信すると、前記切り替え処理が行われる前より作成した前記第二のジャーナルエントリーを、前記第三のストレージ装置に送信することにより、差分リシンクを実行し、

前記計算機システムは、前記第三のストレージ装置と接続し、前記アプリケーションプログラムを実行することで前記第三のボリュームに格納したデータを読み込む第二のホスト計算機を有し、

前記第一のホスト計算機は、前記第一のコピーペアが前記ストレージクラスタリング制御プログラムの制御対象であることを示すコピー定義情報を有し、

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記第一のコピーペア及び前記第二のコピーペアを対象とするペア状態指示を受付け、前記第一のコピーペアのペア状態又は前記第二のコピーペアのペア状態を表示し、

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記コピー定義情報に基づいて、前記第二のコピーペア及び前記第三のコピーペアを対象とする状態変更指示は処理し、前記第一のコピーペアを対象とする状態変更指示の処理は抑止し、

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記差分リシンク指示送信に基く第三のコピーペアを対象とする副ボリューム間差分リシンクが差分コピーで処理できないことを示す情報を前記第二のストレージ装置から取得し、前記第二のボリュームからの全コピーによって前記第三のコピーペアのコピーを再開する第二の指

10

20

30

40

50

示を送信する、

ことを特徴とする計算機システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンピュータシステムにおけるデータの保護技術に関する。

【背景技術】

【0002】

社会基盤を担う公共あるいは企業の基盤システムにおいては、当然ながら高い可用性が求められている。そのため、アプリケーションの継続性を実現したストレージクラスタリング技術として特許文献1記載の技術がある。これは通常の業務に使用される第1のストレージ装置と第2のストレージ装置の間のデータの同一性は、同期リモートコピーによって確保される。業務を行うホストにより第1のストレージ装置の障害を検出した際に、ホストアプリケーションの入力/出力が第2のストレージ装置に対して実行できるように、ホストの制御ブロック内の装置アドレス情報を再構成し、ポインタの切り替えを行う。また、同期リモートコピーの制御もあわせて行う。この技術により、ホストアプリケーションを静止させる必要なく、第1のストレージ装置を、第2のストレージ装置に切り替えることができ、業務の可用性を向上させることができる。

10

【0003】

上記に加えて、データストレージ市場では、大量のデータを格納したストレージ装置が災害等で破壊されてもデータを喪失しない、いわゆるディザスタリカバリシステムが要求されている。このような市場の要求に応えるべく、リモートコピー技術を利用してデータをバックアップする計算機システムが提供されている。これは、十分に離れた二つの地点に設置されたストレージ装置に同一のデータを格納するものである。一方のストレージ装置のデータが更新されると、その更新は、リモートコピーによってもう一方のストレージ装置に同期的もしくは非同期的に反映される。このため、二つのストレージ装置のデータの同一性が確保される。

20

【0004】

さらにデータの安全性を高めるため、相互に十分に離れた三つの地点にストレージ装置を設置する計算機システムが特許文献2に開示されている。この計算機システムでは、通常の業務に使用される第1のストレージ装置と、遠隔地の第2のストレージ装置との間のデータの同一性は、同期リモートコピーによって確保される。一方、第1のストレージ装置と、遠隔地の第3のストレージ装置との間のデータの同一性は、非同期リモートコピーによって確保される。

30

【0005】

災害等に起因する障害によって第1のストレージ装置を業務に使用することができなくなった場合、第2のストレージ装置が第1のストレージ装置の業務を引き継ぐ。このとき、第2のストレージ装置も使用することができない場合は、第3のストレージ装置が第1のストレージ装置の業務を引き継ぐ。その結果、深刻な災害が発生した場合でも、データの喪失を防ぐことができる。

40

【0006】

このように三つの地点にストレージ装置を設置する計算機システムでは、通常の運用時には、第2のストレージ装置と第3のストレージ装置との間でデータが複製されない。このため、第2のストレージ装置と第3のストレージ装置の間では、データの同一性が保証されない。したがって、第1のストレージ装置の業務を第2のストレージ装置が引き継いだ後で、さらに第2のストレージ装置に障害が発生した場合、第3のストレージ装置が第2のストレージ装置の業務を引き継ぐことができない。

【0007】

このため、第1のストレージ装置の業務を引き継いだ第2のストレージ装置の運用が開始される前に、第2のストレージ装置と第3のストレージ装置との間のデータの同一性を

50

確保する。第2のストレージ装置の運用が開始された後は、第2のストレージ装置のデータの更新をリモートコピーによって第3のストレージ装置に反映させる。その結果、第2のストレージ装置に障害が発生したときに、第3のストレージ装置が第2のストレージ装置の業務を引き継ぐことができる。

【0008】

この引き継ぎの際に、第2のストレージ装置のデータを全て第3のストレージ装置に複製すれば、これらのストレージ装置のデータの同一性が確保される。しかし、このように全てのデータを複製することは、長い時間を要する。特に、近年の大容量のストレージ装置では、数時間以上を要する場合がある。全てのデータの複製が終了するまで、第2のストレージ装置を業務に使用することができず、長時間のシステム停止によって深刻な経済的損失が発生するおそれがある。この第2のストレージ装置と第3のストレージ装置とのデータの同一性の確保を行う時間を短縮する手段として、特許文献3記載のデータの更新方法がある。この技術は、第2のストレージ装置と第3のストレージ装置の同一性を確保する際に、相互の差異のデータをもう片方に反映することで、データの複製容量を削減し、結果として時間短縮を実現するものである。本明細書では、この特許文献3に示されるようなデータの同一性の確保技術を「副ボリューム間差分リシンク」と呼ぶ。

10

【先行技術文献】

【特許文献】

【0009】

【特許文献1】米国特許第7043665号明細書

20

【特許文献2】米国特許第7167962号明細書

【特許文献3】米国特許第7447855号明細書

【発明の概要】

【発明が解決しようとする課題】

【0010】

特許文献2及び特許文献3の非同期リモートコピーを特許文献1のストレージクラスタリング技術へ適用する場合、管理者が指定した、クラスタを構成する片方のストレージ装置をコピー元として非同期リモートコピーを行うことが考えられ、クラスタ制御とディザスタリカバリ用途の非同期リモートコピーとが連携してより高性能または高信頼なシステムとして動作することができない。

30

【0011】

本発明は、ストレージクラスタリング技術と連携した非同期リモートコピーを提供することを目的とする。

【課題を解決するための手段】

【0012】

本発明では、非同期リモートコピーを制御するホスト計算機上のプログラムは、ストレージクラスタリング環境でのホスト書き込み先ボリュームの切り替えを行うストレージクラスタリング制御プログラムの切り替え指示と非同期に連携して、非同期リモートコピーのペア操作を行う。

【発明の効果】

40

【0013】

本発明によれば、クラスタ制御と連携した高性能または高信頼なディザスタリカバリ用途の非同期リモートコピーを実現することができる。

【図面の簡単な説明】

【0014】

【図1】図1は、本実施形態の計算機システムの構成を表している。

【図2】図2は、ストレージ装置の構成の詳細を表している。

【図3】図3は、コピー制御情報の構成の詳細を表している。

【図4】図4は、コピーグループ情報の構成の詳細を表している。

【図5】図5は、ペア情報の構成の詳細を表している。

50

【図6】図6は、ストレージクラスタリング及び副ボリューム間差分リシンクを適用したコピーグループの関係を示した模式図である。

【図7】図7は、ストレージクラスタリング制御プログラムの動作を示したフローである。

【図8】図8は、コピー制御プログラムのストレージクラスタリング制御に対応するための動作を示したフローである。

【図9】図9は、コピー制御プログラムの実際にコマンド発行処理を行う部分の動作を示したフローである。

【図10】図10は、コマンドラインインタフェースの説明図である。

【発明を実施するための形態】

【0015】

本発明の実施例を、図1乃至図10を用いて説明する。

【0016】

計算機からの入出力制御や、ストレージ装置内でのボリュームに基づく記憶領域使用方法は、従来技術に示したものと同様である。

【実施例1】

【0017】

図1は、本実施形態の計算機システムの構成を表している。本システムは、ホスト計算機1000と、ホスト計算機1000からアクセスするストレージ装置2000、ホスト計算機1000とストレージ装置2000をつなぐデータネットワーク1009、管理計算機1100、ホスト計算機1000とストレージ装置2000と管理計算機をつなぐ管理ネットワーク1010で構成される。ホスト計算機1000には、ホスト計算機1000の設定など、運用を管理するためのホスト計算機入出力装置1003が付属している。管理計算機には、管理計算機の設定など、運用を管理するための管理計算機入出力装置が付属している。

【0018】

データネットワーク1009は、データ通信のネットワークであって、本実施形態では、SAN (Storage Area Network) である。なお、データネットワーク1009は、データ通信のネットワークであればSAN以外のネットワークでもよく、例えばIPネットワークでもよい。

【0019】

また、管理ネットワーク1010は、データ通信のネットワークであって、本実施形態では、IPネットワークである。なお、管理ネットワーク1010は、データ通信のネットワークであればIPネットワーク以外のネットワークでもよく、例えばSAN (Storage Area Network) でもよい。

【0020】

また、データネットワーク1009及び管理ネットワーク1010は、同一のネットワークであってもよい。また、管理計算機とホスト計算機1000とは、一つの計算機によって実現されてもよい。

【0021】

なお、説明の都合上、本実施形態では、ストレージ装置2000を3台、ホスト計算機1000を1台としたが、これらは複数あってもよい。

【0022】

ストレージ装置2000は、データを格納する領域であるボリューム2301を有する。また、ボリューム2301も、データをホスト計算機1000からの書き込みにより格納するデータボリューム(「A0001」「A0002」と図示)と、リモートコピーを行う際にコピー中のデータを格納するジャーナルボリューム(「JNL-A」と図示)という用途がある。本実施形態のリモートコピーは、「A0001」から「C0001」がペア構成でデータのコピーが行われる場合、「A0001」の書き込みのデータが「JNL-A」にも格納され、「JNL-C」に転送され、「JNL-C」をもとに「C000

10

20

30

40

50

1」に反映されるものとする。ここではストレージ装置をボリュームについてのみ記載している。ストレージ装置の詳細は後述する。

【0023】

なお、以後の説明ではクラスタリング対象の二つのストレージ装置を指す場合はストレージ装置Aやストレージ装置Bと呼び、遠隔地のストレージ装置を指す場合はストレージ装置Cと呼び、特定のストレージ装置を指さない場合はストレージ装置またはストレージ装置2000と呼ぶ。

【0024】

ボリューム2301は、コピーペアを構成できる。なお、コピーペアは、複製元のボリューム2301と、当該複製元のボリューム2301に記憶されたデータの複製を記憶する複製先のボリューム2301とから構成される。

10

【0025】

ホスト計算機1000は、CPU1001、メモリ1002、キーボードやマウスやディスプレイ等の入出力装置1003、ストレージI/F1004、管理I/F1005を備え、各々が接続されている。

【0026】

ストレージI/F1004は、ホスト計算機1000をデータネットワーク1009に接続するネットワークインタフェースである。ストレージI/F1004は、データネットワーク1009を介してストレージ装置2000とデータ及び制御指示を送受信する。

20

【0027】

管理I/F1005は、ホスト計算機1000を管理ネットワーク1010に接続するネットワークインタフェースである。管理I/F1005は、管理ネットワーク1010を介してストレージ装置2000及び管理計算機とデータ及び制御指示を送受信する。

【0028】

メモリ1002には、アプリケーション1006、ストレージクラスタリング制御プログラム3000、コピー制御プログラム4000、OS1007、入出力構成制御情報1008、及び、コピー制御情報5000が格納される。

【0029】

アプリケーション1006は、OS1007に要求を出すことにより、ストレージ装置2000に備わるボリューム2301にデータを読み書きする。例えば、アプリケーション1006は、DBMS(Data Base Management System)又はファイルシステム等である。

30

【0030】

OS1007は、このシステムを管理するオペレーティングシステムである。アプリケーション1006などからの要求に従い、入出力の対象となる機器を入出力構成制御情報1008から特定し、ストレージ装置とのデータの入力/出力を行う機能をその一部として有する。

【0031】

ストレージクラスタリング制御プログラム3000は、従来技術記載のストレージクラスタリングを制御する。すなわち、ストレージクラスタリングの対象となるボリューム間の同期コピーの制御、及び、入出力対象となっているボリュームの障害が発生した場合に、入出力を一旦停止させ、同期コピーのコピー方向を逆転させ、入出力構成制御情報内のボリューム情報を書き換えることで入出力が元の同期コピー先に切り替わるようにした上で入出力を再開させるという処理を行う。

40

【0032】

コピー制御プログラム4000は、コピー運用の手順に基づきコピー制御を行う。また、コピー制御のため、ストレージマイクロプログラム2206にコピーペアを制御させる要求及びコピーペアの状態を取得させる要求を送信する。

【0033】

コピー制御情報5000は、コピー制御プログラム4000がコピー制御を行うための

50

情報を格納した領域である。詳細は後述する。

アプリケーション 1006、ストレージクラスタリング制御プログラム 3000、コピー制御プログラム 4000、OS 1007はCPU 1001によって実行される。

【0034】

なお、説明の都合上、図1では、アプリケーション 1006を一つとしたが、複数あってもよい。

【0035】

なお、ストレージ装置Cは非同期リモートコピーによって複製されたデータを用い、ホスト計算機 1000が実行していた業務処理を実行するホスト計算機 1000Bが接続される。

10

【0036】

なお、典型的にはホスト計算機 1000とストレージ装置Aとストレージ装置Bとは一つのサイト(ローカルサイト)に存在し、ホスト計算機 1000Bとストレージ装置Cとは別なサイト(リモートサイト)に属することで、ローカルサイトでの災害を想定したディザスタリカバリ環境を構築する。しかし、サイトと各装置・計算機の関係はこれ以外の場合もありうる。

【0037】

図2は、図1におけるストレージ装置の構成の詳細を表している。

【0038】

ストレージ装置 2000は、ディスク装置 2100及びディスクコントローラ 2200を備える。ディスク装置 2100は、ホスト計算機 1000に書き込み要求されたデータを格納し、ホスト計算機 1000から読み込み要求された格納データをホスト計算機に送信する。ディスクコントローラ 2200は、ストレージ装置 2000の処理を制御する。

20

【0039】

ディスク装置 2100は、複数のボリューム 2301を備える。ボリューム 2301は、物理的な記憶領域であるハードディスクドライブ(HDD)、フラッシュドライブ(SSD)、又は、論理的な記憶領域である論理デバイス(Logical Device)のいずれであってもよく、本発明では、ボリュームの種類を問わない。なお、説明の都合上、図2では、3つのボリューム 2301を示したが、ボリューム 2301の数はいくつであってもよい。

30

【0040】

ディスクコントローラ 2200は、ホストI/F 2201、管理I/F 2202、ディスクI/F 2203、メモリ 2204、及び、CPU 2205を備える。

【0041】

メモリ 2204には、ストレージマイクロプログラム 2206及びコピーペア情報 2207が格納される。

【0042】

ストレージマイクロプログラム 2206は、CPU 2205によって実行される。ストレージマイクロプログラム 2206は、ホスト計算機 1000からの要求に応じて、コピーペアを制御し、コピーペアの状態を取得し報告する。

40

【0043】

ストレージマイクロプログラム 2206によるコピーペアの制御については後ほど説明する。

【0044】

コピーペア情報 2207は、ストレージ装置 2000に備わるボリューム 2301のうち、コピーペアを構成するボリューム 2301の情報が格納される。

【0045】

なお、本実施形態では、ストレージマイクロプログラム 2206及びコピーペア情報 2207は、ディスクコントローラ 2200のメモリ 2204に格納されるとしたが、本発明はこれに限定されない。例えば、ストレージマイクロプログラム 2206及びコピーペ

50

ア情報 2207 は、ディスクコントローラ 2200 に接続されるフラッシュメモリに格納されてもよいし、ディスク装置 2100 に備わるボリューム 2301 に格納されてもよい。

【0046】

ホスト I/F 2201 は、ストレージ装置 2000 をデータネットワーク 1009 に接続するネットワークインタフェースである。ホスト I/F 2201 は、データネットワーク 1009 を介してホスト計算機 1000 とデータ及び制御指示を送受信する。

【0047】

管理 I/F 2202 は、ストレージ装置 2000 を管理ネットワーク 1010 に接続するネットワークインタフェースである。管理 I/F 2202 は、管理ネットワーク 1010 を介してホスト計算機 1000 及び管理計算機 1200 とデータ及び制御指示を送受信する。

10

【0048】

ディスク I/F 2203 は、ディスクコントローラ 2200 をディスク装置 2100 に接続するインタフェースである。

【0049】

次に、ストレージ装置によるリモートコピーについて説明する。

【0050】

同期リモートコピーは、主ボリュームに対してホストから書き込み要求を受け付けた場合、当該書き込みデータを副ボリュームにデータコピーした後に、ホストに対して書き込み完了を返すコピー方式である。

20

【0051】

同期リモートコピーが実行される際、正と副のボリュームのデータおよびコピーの状況を示したり操作するために、ディスクコントローラはペア状態 (Simplex、Initial-Copying、Duplex、Suspend 及び Duplex-Pending) と呼ばれる情報を管理する。以下、各ペア状態毎に同期リモートコピーの処理内容について説明する。

【0052】

< Simplex 状態 >

Simplex 状態は、正と副のボリューム間でコピーが開始されていない状態である。

30

【0053】

< Duplex 状態 >

Duplex 状態は、同期リモートコピーが開始され、後述する初期化コピーも完了して正副のボリュームのデータ内容が同一となった状態である。同期リモートコピーの場合、主ボリュームに対して行われた書き込みの内容が副ボリュームに対してコピーされた後に、書き込みを行ったホストに対して正常完了のメッセージが返される。従って、書き込み途中の領域を除けば、主ボリュームのデータと副ボリュームのデータの内容は同じとなる。

【0054】

40

< Initial-Copying 状態 >

Initial-Copying 状態は、Simplex 状態から Duplex 状態へ遷移するまでの中間状態であり、この期間中に、必要ならば主ボリュームから副ボリュームへの初期化コピー (主ボリュームに既に格納されていたデータのコピー) が行われる。初期化コピーが完了し、Duplex 状態へ遷移するために必要な処理が終わったら、ペア状態は Duplex となる。なお、当該状態には Simplex 状態にてホスト計算機から「作成」指示を受信した契機で遷移する。

【0055】

< Suspend 状態 >

Suspend 状態は、主ボリュームに対する書き込みの内容を副ボリュームに反映さ

50

せない状態である。この状態では、正副のボリュームのデータは同じでない。オペレータやホストからの「停止」指示を契機に、ペア状態は他の状態から S u s p e n d 状態へ遷移する。

【 0 0 5 6 】

それ以外に、同期リモートコピーを行うことが出来なくなった場合に自動的にペア状態が S u s p e n d 状態に遷移することが考えられる。

【 0 0 5 7 】

以下の説明では、後者の場合を障害 S u s p e n d 状態と呼ぶことにする。障害 S u s p e n d 状態となる代表的な原因としては、正副のボリュームの障害、正副のディスクコントローラの障害、正副間の通信障害が考えられる。S u s p e n d 状態となった正副のディスクコントローラは、当該状態となった以降の正と副のボリュームに対する書き込み位置を記憶する。

10

【 0 0 5 8 】

< Duplex - Pending 状態 >

Duplex - Pending 状態は、Suspend 状態から Duplex 状態に遷移するまでの中間状態である。この状態では、主ボリュームと副ボリュームのデータの内容を一致させるために、主ボリュームから副ボリュームへのデータのコピーが実行される。正と副のボリュームのデータが同一になった後、ペア状態は Duplex となる。なお、Duplex - Pending 状態におけるデータのコピーは、Suspend 状態で正と副のディスクコントローラが記録した書き込み位置を利用して更新が必要な部分だけをコピーする差分コピーが用いられる。なお、当該状態には Suspend 状態（障害 S u s p e n d 状態を含む）にてホスト計算機から「リシンク」指示を受信した契機で遷移する。

20

【 0 0 5 9 】

なお、以上の説明では Initial - Copying 状態と Duplex - Pending 状態は別々な状態としたが、これらをまとめて一つの状態として管理装置の画面に表示したり、状態を遷移させても良い。

【 0 0 6 0 】

次に非同期リモートコピーについて説明する。

【 0 0 6 1 】

非同期リモートコピーの場合、Duplex 状態の副記憶領域に対する書き込みデータの反映は、記憶装置 1 2 0 のホスト 1 1 0 に対する書き込みの正常完了メッセージの送信とは無関係（非同期）に行われる。

30

【 0 0 6 2 】

非同期リモートコピーの場合、正記憶領域から副記憶領域へのデータコピーの方法として、以下の方法がある。

【 0 0 6 3 】

例えば、正記憶装置が、データが書き込まれた記憶領域のアドレスを含んだ制御情報及び書き込まれたデータの組（以下「ジャーナルエントリ」と呼ぶ）を、データの書き込みの度に作成し、これを副記憶装置へ転送し、副記憶領域へ反映させる方法がある（記憶先はキャッシュメモリ又は / 及びジャーナルボリュームである）。さらにこの発展形として、ジャーナルエントリの制御情報に書き込みの時間順序を示す情報を含め、副記憶領域へジャーナルエントリを反映させる際にはこの時間順序を示す情報を利用して時間順序どおりに反映する方法がある。

40

【 0 0 6 4 】

また、本方法の効率的な方法として、正記憶領域の同一領域に対する書き込みが連続して発生した場合には、正記憶装置は、途中の書き込みに対するジャーナルエントリは副記憶装置へ転送せず、最後の書き込みに対するジャーナルエントリのみを転送する方法がある。

【 0 0 6 5 】

50

さらに、もう一つの例としては、ある一定時間に主ボリュームに対して書き込まれたデータを差分データとして保持し、副ボリュームへコピーする方法がある。この方法の場合、差分データを副ディスクコントローラへ全て転送してからボリュームへコピーする。

【0066】

非同期リモートコピーも、ペア状態 (Simplex、Initial-Copying、Duplex、Suspend、Duplex-Pending及びSuspending) を用いて操作と管理を行う。Simplex、Initial-Copying、Suspend及びDuplex-Pending状態については同期リモートコピーと同様である。

【0067】

< Duplex状態 >

Duplex状態も基本的には同期リモートコピーの場合と同じであるが、書き込みデータの副ボリュームへのコピーが非同期に行われるため、副ボリュームは主ボリュームより少し遅れながら追従する。

【0068】

< Suspending状態 >

Suspending状態とは、Duplex状態からSuspend状態へ遷移するまでの中間状態である。非同期リモートコピーの場合は、Suspending状態を経由してSuspend状態へ遷移する。同期リモートコピーのSuspend状態で述べた正副のボリュームに対する書き込み位置の記録には、コピーできなかったジャーナルエントリを書き込み位置の記録に加える。

【0069】

図3は、図1におけるコピー制御情報の構成の詳細を表している。

【0070】

コピー制御情報5000は、コピーグループ情報5100、ペア情報5200、ペア状態情報5300、コピー環境情報5400から構成される。

【0071】

コピーグループ情報5100は、複数のコピーペアをグループ化したコピーグループに関する情報を示す。なお、コピーグループ情報5100については、図4で詳細を説明する。

【0072】

ペア情報5200は、コピーグループ毎に、グループとしてまとめて操作するコピーペアの情報を有する。なお、ペア情報5200については、図5で詳細を説明する。

【0073】

ペア状態情報5300は、ペア情報5200に記載されたペア各々のペア状態の情報を保持する領域であり、コピー制御プログラム4000がストレージ装置2000に発行したペアの状態取得指示の結果得られた状態が格納される。

【0074】

コピー環境情報5400は、コピー制御を行うパラメタを格納する領域である。本実施の形態では、オプション5401を有し、ここには副ボリューム間差分リシンクを実行する場合に、差分リシンクができなかった場合はリシンク処理を中断する「中断」と、差分リシンクができなかった場合は差分リシンクのかわりにボリューム全体のコピーを行ってリシンク処理を行う「続行」の、2つのいずれかの値が設定される。

ペア状態情報5300に格納される指示の結果得られた情報以外のこれらの値は、入出力装置1003もしくは管理計算機を介してユーザが設定する。

【0075】

図4は、図3におけるコピーグループ情報5100の構成の詳細を表している。

【0076】

コピーグループ情報5100は、グループ識別子5101、第一属性5102、第二属性5103のセットからなる。グループ識別子5101は、コピーグループの識別子を格

10

20

30

40

50

納する領域である。第一属性は、同期コピーや非同期コピーなど、コピー種別を格納する領域である。第二属性5103は、ストレージクラスタリング機能の適用対象であるか否かという、第一属性5102とは別の機能での属性を格納する領域である。

【0077】

図5は、図3におけるペア情報5200の構成の詳細を表している。

【0078】

ペア情報5200は、グループ識別子5201、プライマリデバイス識別子5202、セカンダリデバイス識別子5203を含む。図3はコピーグループがG1からG3まで3つある場合の例である。例えばG1は、デバイス識別子が「A0001」を主ボリューム、「B0001」を副ボリュームとするペアと、「A0002」を主ボリューム、「B0002」を副ボリュームとするペアの2つのペアから構成されている。このデバイス識別子により、対象となる装置、及び、その装置内のいずれのボリュームかが特定できる。

10

【0079】

図6は、図4、及び、図5における情報の例示を元にした、本実施形態におけるストレージクラスタリング及び副ボリューム間差分リシンクを適用したコピーグループの関係を示した模式図である。

【0080】

G1は、ストレージクラスタリングが適用中であるコピーグループで、ストレージ装置上は同期コピーが設定され、「A0001」ボリュームへの書き込みは、同期的に「B0001」ボリュームに反映される。同様に「A0002」ボリュームへの書き込みは、同期的に「B0002」ボリュームに反映される。すなわちG1は、「A0001」と「B0001」、及び、「A0002」と「B0002」の2つのコピーペアから構成されている。

20

【0081】

更に、「A0001」ボリュームと、「A0002」ボリュームは、G2として非同期コピーが設定されている。すなわちG2は、「A0001」と「C0001」、及び、「A0002」と「C0002」の2つのコピーペアから構成され、書き込みが非同期的にコピーされている。

【0082】

このG1とG2で示された構成は、マルチターゲット構成と呼ばれる。

30

【0083】

また、G1とG2のコピーグループを対象に、副ボリューム間差分リシンクを適用するためのG3が設定されている。G1の副ボリュームとG2の副ボリュームを対象に、「B0001」と「C0001」、及び、「B0002」と「C0002」の2つのコピーペアから構成される。通常G3ではボリュームへの書き込みの反映はおこなわれない。これを「差分リシンク待機」と呼び、図6ではこれを破線で示している。G3は、副ボリューム間差分リシンク指示がされた時に、差分のデータの調整を行い、G2のかわりにペアとして非同期コピーを開始し、ペア状態になる。その際、G2が「差分リシンク待機」に変わる。指示の仕方の詳細は後述する。

【0084】

なお、差分リシンク待機状態のペアの正ボリュームを提供するストレージ装置は、クラスタを構成するもう一方のストレージ装置から同期リモートコピーを介して受信した書き込みデータをジャーナルエントリ化して、ジャーナルボリュームに格納している。

40

【0085】

図7は、ストレージクラスタリング制御プログラム3000の動作を示したフローである。

【0086】

ストレージクラスタリング制御プログラムは、クラスタリング適用中のボリュームへの読み込み又は/及び書き込みの成否を監視する(ステップ3001)。

【0087】

50

次に読み込み又は / 及び書き込みが失敗してスワップが必要になったか、もしくは、ユーザからの指示があってスワップが必要になったか、いずれかの状態を検知する（ステップ3002）。すなわち非計画切り替えと計画切り替えのいずれかによる処理の開始を判断する。なお、「スワップ」とは読み込み又は書き込み要求を出すストレージ装置を切り替えることを指す。

【0088】

上記ステップ3002での判定が偽であった場合、ステップ3001に戻り読み込み又は / 及び書き込み要求の監視を継続する。

【0089】

上記ステップ3002での判定が真であった場合、以下の切り替え処理を行う。

10

【0090】

まず、対象となるボリュームへの読み込み又は / 及び書き込み要求の処理を一時的に停止する（ステップ3003）。

【0091】

ストレージクラスタリング対象のボリュームを有するペア状態の変更を行う（ステップ3004）。すなわち、図6の例で述べると、「A0001」のデータが「B0001」に向けて同期的にコピーされていたものを、「停止」指示をストレージ装置へ送り、「B0001」への書き込みができる状態に変更する。「A0002」と「B0002」のペアも同様に処理する。

【0092】

20

次に、入出力構成制御情報1008の交換書き換えを行う（ステップ3005）。通常、アプリケーションは、読み込み又は / 及び書き込み要求の発行先として指定する情報とともに読み込み又は / 及び書き込み要求がOS1007に渡る。OS1007は渡された読み込み又は / 及び書き込み要求を処理する時に、入出力構成制御情報1008に含まれる読み込み又は / 及び書き込み要求発行先の情報に対応するデバイス識別子を求め、そのデバイス識別子を対象にした読み込み又は / 及び書き込み要求処理を行う。本ステップでは、クラスタリング対象のボリューム情報と、そのペアとされているボリュームの情報を交換する。「A0001」と「B0001」の例で述べると、交換前に「A0001」を指し示していた情報は、交換後に「B0001」を指し示すようにし、逆に交換前に「B0001」を指し示していた情報は、交換後に「A0001」を指し示すようになる。

30

【0093】

それから、上記ステップ3003で停止させた読み込み又は / 及び書き込み要求処理を再開する（ステップ3006）。

【0094】

最後に、交換処理を完了したことを示すメッセージを出力し（ステップ3007）、ステップ3001の読み込み又は / 及び書き込み要求監視に戻る。

【0095】

なお、ステップ3005ではコピー方向の反転を伴った「リシンク」指示をストレージ装置に送信してもよい。仮に当該処理を行う要因が切り替え元ストレージ装置2000とホスト計算機1000との間のネットワーク障害の場合は、逆方向の同期リモートコピーが開始され、データの冗長性が確保される。

40

【0096】

図8は、コピー制御プログラム4000の、ストレージクラスタリング制御に対応するための動作を示したフローである。

【0097】

コピー制御プログラム4000は、ストレージクラスタリング制御プログラム3000が行う交換の完了メッセージの出力の有無を監視する（ステップ4101）。

【0098】

メッセージが検知されたかを判定する（ステップ4102）。

【0099】

50

上記ステップ4102での判定が偽であった場合、ステップ4101に戻りメッセージの監視を継続する。

【0100】

上記ステップ4102での判定が真であった場合、以下の処理を行う。

【0101】

まず、G1の状態を確認する指示を発行し、ストレージ装置2000からG1の状態情報を取得し、ストレージクラスタリング制御プログラム3000による交換処理後の状態であるか否かを確認する(ステップ4103)。

【0102】

次に、G2及びG3の状態を確認する指示を発行し、ストレージ装置2000からG2及びG3の状態情報を取得し、障害でペア状態が停止しているなど差分リシンクができない状態あるか否かを確認する(ステップ4104)。

【0103】

上記ステップ4103、及び、ステップ4104の結果から、G1が交換処理後の状態であり、かつ、G2及びG3の状態に異常がない、副ボリューム間差分リシンクをすべき状態か否かを判定する(ステップ4105)。

【0104】

上記ステップ4105での判定が真であった場合、以下の処理を行う。

【0105】

差分リシンク指示を発行する(ステップ4106)。ここでG2がペア状態で、G3が差分リシンク待機状態であるのを、G3がペア状態で、G2が差分リシンク待機状態に切り替える指示である。この切り替えを、ボリュームの情報を全てコピーするのではなく、G2の副ボリューム(すなわちG3の副ボリューム)と、G3の主ボリュームの差分だけをコピーするように設定しなおすことで、全てコピーするよりも短期間にコピーできる副ボリューム間差分リシンクを行う。

【0106】

次に、ストレージ装置2000から情報を取得し、ステップでの副ボリューム間差分リシンクが実行できたか否かを判定する(ステップ4107)。ストレージクラスタリング制御プログラムへの作りこみを抑止するため、本フローでの差分リシンクがストレージクラスタリングのボリューム交換と非同期的に実行される。そのため、ボリューム交換実行と差分リシンク実行の間に時間差が生じることがある。通常その差は非同期コピーのジャーナルボリュームの容量を大きくとることで解決されるが、ジャーナルボリュームの容量を少なく運用していた場合など、差分リシンクが失敗に終わることがある。この事象を本ステップで検知する。

【0107】

上記ステップ4107での判定が真であった場合、差分リシンクが成功したため本フローの処理を終了する。

【0108】

上記ステップ4107での判定が偽であった場合、差分リシンクできなかったため、以下の処理を行う。

【0109】

コピー環境情報5400に設定されたオプション情報が、「続行」であるか「中断」であるかを判定する(ステップ4108)。

【0110】

上記ステップ4108での判定が「続行」であった場合、G3の主ボリュームのデータをG3の副ボリュームに全てコピーすることでペア状態を確立させる全コピーをすることで、G2を待機状態に、G3をペア状態に切り替える指示を行い、処理を終了する(ステップ4109)。

【0111】

上記ステップ4108での判定が「中断」であった場合、又は、上記ステップ4105

10

20

30

40

50

での判定が偽であった場合、切り替えができなかったことを示すメッセージを出力し、処理を終了する（ステップ4110）。

【0112】

図9は、コピー制御プログラム4000の実際にコマンド発行処理を行う部分の動作を示したフローである。

【0113】

コピー制御プログラム4000は、「作成」、「停止」、「リシンク」指示やペア状態取得指示など、コピー制御に関する指示を受け取る（ステップ4201）。

【0114】

指示で指定されているコピーグループに対応するコピーグループ情報5100の第二属性5103を参照し、ストレージクラスタリングの対象となっているか否かを判定する（ステップ4202）。

10

【0115】

上記ステップ4202での判定が真であった場合、以下の処理を行う。

【0116】

指示で指定された処理内容が、停止やリシンクなどのペアの状態を変える制御系の指示であるか否かを判定する（ステップ4203）。

【0117】

上記ステップ4203での判定が真であった場合、その指示は許可されていないものとしてエラーとして処理し、終了する（ステップ4204）。

20

【0118】

上記ステップ4203での判定が偽であり指示で指定された処理内容がペアの状態情報取得などの参照系の指示であった場合、又は、上記ステップ4202での判定が偽であった場合、要求通りコマンド発行を行い、処理を終了する（ステップ4205）。

【0119】

図10は、コピー制御プログラム4000がユーザもしくはスクリプトプログラム等に提供するコマンドラインインタフェースの説明図である。

【0120】

コマンドラインインタフェースは、コマンド4301、対象4302、パラメタ4303からなる。

30

【0121】

コマンド4301は、ペア作成や状態取得などの指示を記述する。

【0122】

対象4302は、対象となるグループやペアを記述する。

【0123】

パラメタ4303は、指示に必要なパラメタがあればそれを記述する。

【0124】

コピーグループG1に対して状態を取得する場合は「状態取得，G1」と指示する。コピーグループG3を差分リシンク待機状態に設定する場合は、「ペア作成，G3，差分リシンク待機状態」と指示する。

40

【0125】

これらの指示やパラメタなどは、図10の記載以外にもコピーの機能に従い様々な形態がありうる。

【0126】

以上、本実施形態での図8を用いたコピー制御プログラム4000の、ストレージクラスタリング制御に対応するための動作において、副ボリューム間差分リシンク指示の前に、差分リシンクが可能な状態かを予め確認した上で差分リシンクを実施する形態を説明したが、これをストレージクラスタリングの動作を検知した後に、差分リシンクを行い、その結果成功しなかった場合に、改めて関連するコピーグループの状態を確認し、その結果ペアの状態異常などストレージによる差分リシンクが可能なエラーだった場合に、ペア状

50

態の確認などの回復処理を行い、再度、副ボリューム間差分リシンクを指示する形態としてもよい。この場合、G 1などのコピーグループの状態確認を完了する時間を待たずに差分リシンクを行うことで、ディザスタリカバリ可能な状態により早く戻ることができる。

【 0 1 2 7 】

また、本実施形態での図 8 を用いたコピー制御プログラム 4 0 0 0 の、ストレージクラスタリング制御に対応するための動作において、ステップ 4 1 0 4 において、関連するコピーグループ G 2 及び G 3 の状態の取得を行った。これを、状態取得の前に、副ボリューム間差分リシンクを実行できる状態に変更するための G 2 に対するペアの停止指示を実施する形態としてもよい。ストレージクラスタリングによるボリュームの切り替えは、ユーザによるメンテナンスの準備として計画的に実施、もしくは、障害から非計画的に実施される。この場合、ストレージ障害の場合は、G 2 の状態が停止状態になるが、データネットワーク 1 0 0 9 は正常に動作しているがホスト計算機 1 0 0 0 とストレージ装置 2 0 0 0 の間のなんらかの回線障害の場合は、G 2 の状態はペア状態のままである。すなわちストレージクラスタリングの動作に伴って差分リシンクを実施すべき場合でも G 2 の状態が複数とりうる可能性がある。差分リシンクができるようにペア状態を確定させるステップを行うことで、計画的もしくは非計画的なボリュームの切り替えに対しても、コピー制御プログラム 4 0 0 0 は一つの手順で実施できる。それによりエラー時のユーザ操作の必要性を削減することができる。

【 0 1 2 8 】

以上の説明により、

アプリケーションプログラムと、ストレージクラスタリング制御プログラムと、コピー制御プログラムと、を実行する第一のホスト計算機と、

前記第一のホスト計算機と接続し、第一のボリュームを提供する第一のストレージ装置と、

前記第一のホスト計算機及び前記第一のストレージ装置と接続し、第二のボリュームを提供する第二のストレージ装置と、

前記第一のストレージ装置及び前記第二のストレージ装置と接続し、第三のボリュームを提供する第三のストレージ装置と、

を有する計算機システムについて説明した。

【 0 1 2 9 】

そして、

前記第一のストレージ装置は、前記第一のボリュームをコピー元とし、前記第二のボリュームをコピー先とする第一のコピーペアについての同期リモートコピー処理によって、前記アプリケーションプログラムの実行に従ってホスト計算機が送信した書き込みデータを前記第二のストレージ装置へ送信し、

前記第一のストレージ装置は、前記第一のボリュームをコピー元とし、前記第三のボリュームをコピー先とする第二のコピーペアについての非同期リモートコピー処理によって、前記アプリケーションプログラムの実行に従ってホスト計算機が送信した書き込みデータを前記第三のストレージ装置へ送信し、

前記ストレージクラスタリング制御プログラムを実行することで、前記第一のホスト計算機は、前記書き込みデータの書き込みが異常終了したことを検知し、前記第一のコピーペアのペア状態を変更することで検知後の前記書き込みデータを前記第二のボリュームに対して書き込み可能とする第一の指示を前記第二のストレージ装置に送信し、

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記ストレージクラスタリング制御プログラムによって異常終了を検知したことを、前記第一の指示の送信より後に検知し、前記第二のボリュームをコピー元とし、前記第三のボリュームをコピー先とする第三のコピーペアについての別な非同期リモートコピー処理によるコピー処理を開始するための差分リシンク指示を前記第二のストレージ装置に送信することを説明した。

【 0 1 3 0 】

また、前記計算機システムは、前記第三のストレージ装置と接続し、前記アプリケーションプログラムを実行することで前記第三のボリュームに格納したデータを読み込む第二のホスト計算機を有してもよいことを説明した。

【0131】

また、

前記第一のホスト計算機は、前記第一のコピーペアが前記ストレージクラスタリング制御プログラムの制御対象であることを示すコピー定義情報を有し、

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記第一のコピーペア及び前記第二のコピーペアを対象とするペア状態指示を受付け、前記第一のコピーペアのペア状態又は前記第二のコピーペアのペア状態を表示し、

10

前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記コピー定義情報に基づいて、前記第二のコピーペア及び前記第三のコピーペアを対象とする状態変更指示は処理し、前記第一のコピーペアを対象とする状態変更指示の処理は抑止してもよいことも説明した。

【0132】

また、前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記差分リシンク指示送信に基づく第三のコピーペアを対象とする副ボリューム間差分リシンクが差分コピーで処理できないことを示す情報を前記第二のストレージ装置から取得し、前記ストレージクラスタリング制御プログラムによる前記第一のコピーペアの状態変化を要因とした前記第三のコピーペアの副ボリューム間差分リシンクが失敗したメッセージを表示してもよいことを説明した。

20

【0133】

また、前記コピー制御プログラムを実行することで、前記第一のホスト計算機は、前記差分リシンク指示送信に基づく第三のコピーペアを対象とする副ボリューム間差分リシンクが差分コピーで処理できないことを示す情報を前記第二のストレージ装置から取得し、前記第二のボリュームからの全コピーによって前記第三のコピーペアのコピーを再開する第二の指示を送信してもよいことを説明した。

【0134】

このようなクラスタリング制御と非同期リモートコピー制御との連携を取ることで、仮に第一のストレージ装置が機能停止した場合でも非同期RCを継続できることから計算機システムの信頼性が向上する。また、コピー制御プログラムは、障害検知のための高度な処理を含むストレージクラスタリング制御プログラムとは別なプログラムであるため、処理のモジュール性が向上し、安定した動作が可能となる他、計算機システムが非同期リモートコピーを実施しない状況から実施する状況にシステム構成を変更する場合（または逆の場合）にも柔軟な設定変更が可能となる。

30

【0135】

また、非同期リモートコピーの場合、対象となるストレージ装置の通信距離が遠く、そして一時的な通信不良が発生する可能性があり、また同期リモートコピーと比較して複雑な処理を必要とするため、ペア操作指示に時間を必要とすることがある。そのため、ストレージクラスタ制御プログラムからコピー制御プログラムへのストレージ装置切り替えの通知を、切り替え指示よりも後に行う点は、アプリケーションプログラムからの読み込み又は/及び書き込み要求の停止時間短縮、又はストレージクラスタリング制御プログラムの処理をプログラム設計者が想定した以上に停止させないことによる安定動作につながる。

40

【0136】

しかし、本実施例はこれら以外の内容についても開示している。

【符号の説明】

【0137】

1000...ホスト計算機 1100...管理計算機 2000...ストレージ装置 3000...ストレージクラスタリング制御プログラム 4000...コピー制御プログラム

50

【図1】

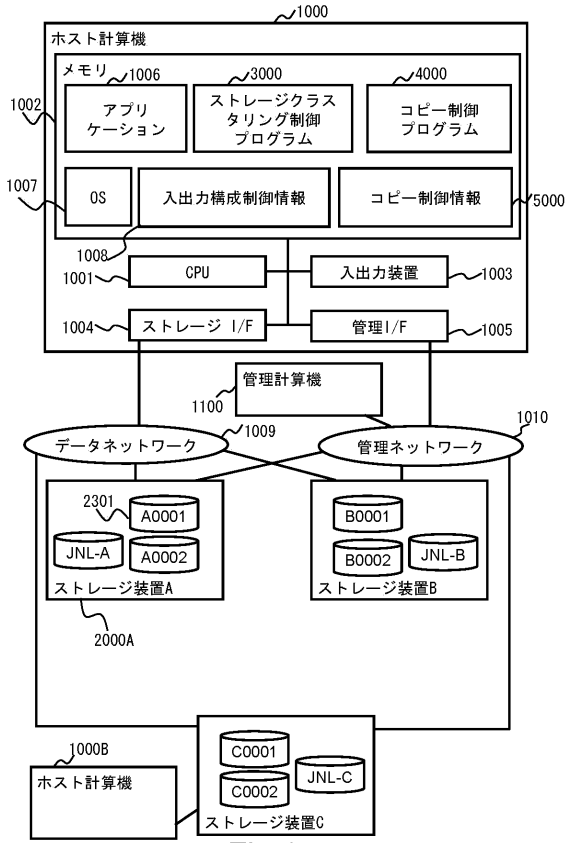


Fig.1

【図2】

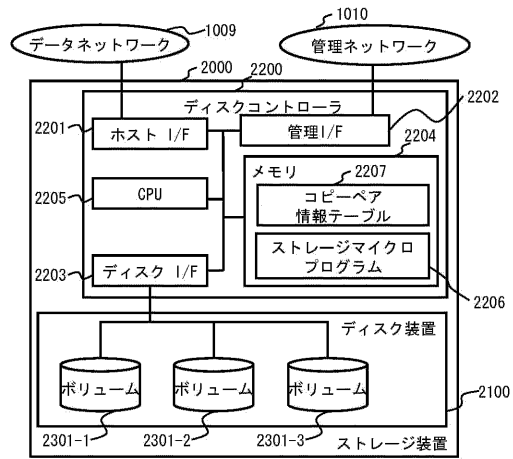


Fig.2

【図3】

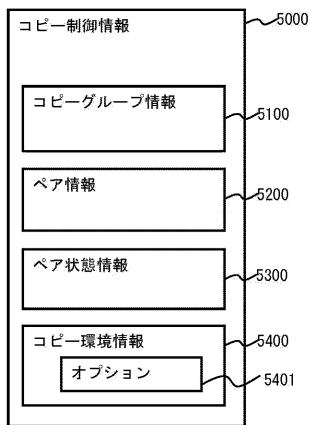


Fig.3

【図4】

5100		
5101	5102	5103
グループ識別子	第一属性	第二属性
G1	同期コピー	クラスタリング機能適用中
G2	非同期コピー	—
G3	非同期コピー	—

Fig.4

【図5】

グループ識別子		G1	5201
5202	プライマリデバイス識別子	セカンダリデバイス識別子	5203
	A0001	B0001	
	A0002	B0002	
グループ識別子		G2	
	プライマリデバイス識別子	セカンダリデバイス識別子	
	A0001	C0001	
	A0002	C0002	
グループ識別子		G3	
	プライマリデバイス識別子	セカンダリデバイス識別子	
	B0001	C0001	
	B0002	C0002	

Fig.5

【図6】

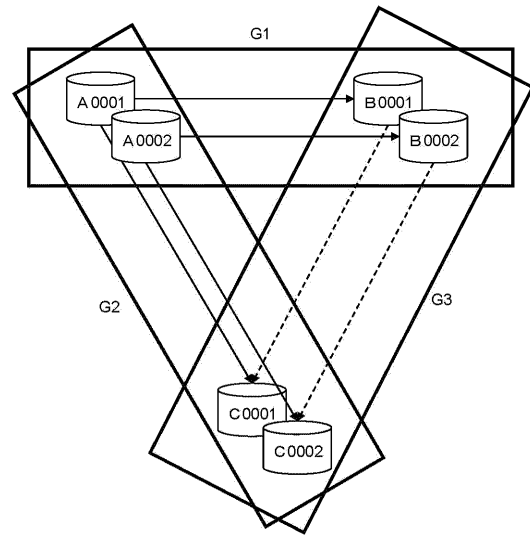


Fig.6

【図7】

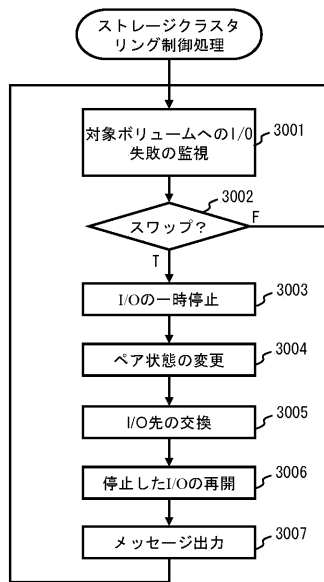


Fig.7

【図8】

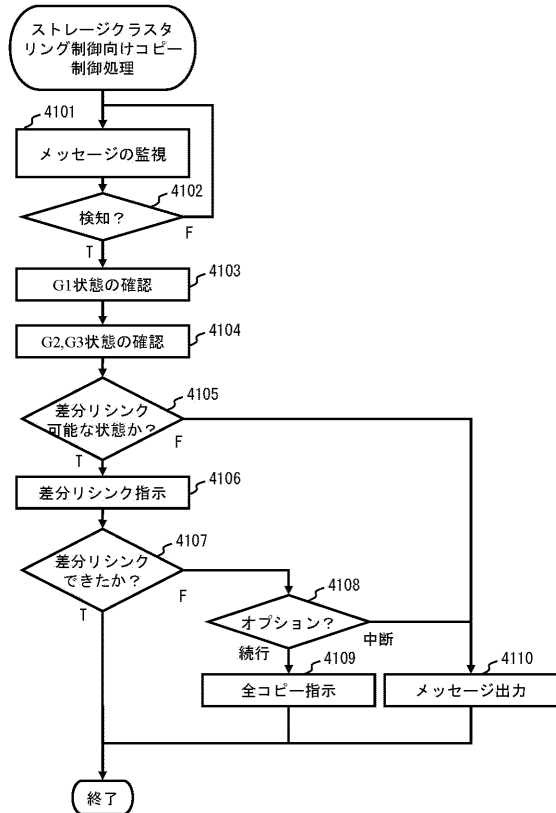


Fig.8

【図9】

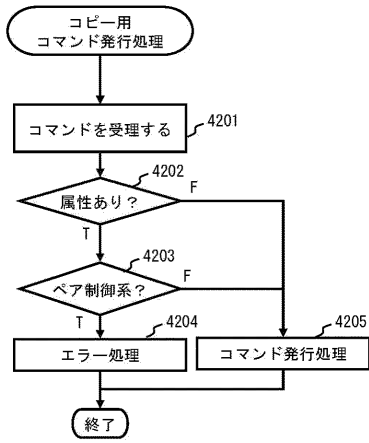


Fig.9

【図10】

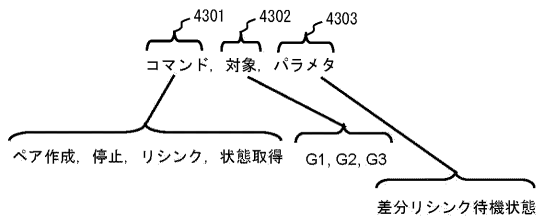


Fig.10

---

フロントページの続き

(72)発明者 小山田 健一

神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所 ソフトウェア事業部内

(72)発明者 松井 義典

神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所 ソフトウェア事業部内

審査官 坂東 博司

(56)参考文献 特開2007-066162(JP,A)

特開2008-065425(JP,A)

特開2007-249447(JP,A)

特開2007-272510(JP,A)

特開2006-053912(JP,A)

特開2004-319013(JP,A)

特開2004-227528(JP,A)

特表2006-527875(JP,A)

特開2008-134986(JP,A)

米国特許出願公開第2004/0260899(US,A1)

米国特許出願公開第2008/0104443(US,A1)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06