

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5057366号
(P5057366)

(45) 発行日 平成24年10月24日(2012.10.24)

(24) 登録日 平成24年8月10日(2012.8.10)

(51) Int.Cl. F I
G 0 6 F 3/06 (2006.01) G O 6 F 3/06 3 O 5 F
 G O 6 F 3/06 3 O 1 Z

請求項の数 4 (全 72 頁)

(21) 出願番号	特願2007-85680 (P2007-85680)	(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19年3月28日(2007.3.28)	(74) 代理人	100093861 弁理士 大賀 真司
(65) 公開番号	特開2008-134987 (P2008-134987A)	(74) 代理人	100129218 弁理士 百本 宏之
(43) 公開日	平成20年6月12日(2008.6.12)	(72) 発明者	岩村 卓成 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21年6月26日(2009.6.26)	(72) 発明者	二瀬 健太 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
(31) 優先権主張番号	特願2006-293485 (P2006-293485)		
(32) 優先日	平成18年10月30日(2006.10.30)		
(33) 優先権主張国	日本国(JP)		

最終頁に続く

(54) 【発明の名称】 情報システム及び情報システムのデータ転送方法

(57) 【特許請求の範囲】

【請求項1】

ファイルシステム、HBAデバイスドライバを備えるオペレーティングシステム、及び前記ファイルシステムを通じてライトリクエストを送信するアプリケーションを備えるホストコンピュータと、前記ホストコンピュータに接続される第1のストレージ装置と、前記第1のストレージ装置及び前記ホストコンピュータに接続される第2のストレージ装置と、前記第1のストレージ装置及び前記第2のストレージ装置に接続され、複数の第1のHDDとコントローラとキャッシュメモリから構成され、前記複数の第1のHDDの一部の領域から構成される第1のボリュームを有する第3のストレージ装置と、を有する情報システムのデータ転送方法であって、

前記第1のストレージ装置は、前記第3のストレージ装置の前記第1のボリュームに実体が存在する第1の仮想ボリュームを定義し、

前記第2のストレージ装置は、前記第3のストレージ装置の前記第1のボリュームに実体が存在し、当該第1のボリュームを前記第1の仮想ボリュームと共有する第2の仮想ボリュームを定義し、

前記第1のストレージ装置及び前記第2のストレージ装置は、それぞれ複数の第2のHDDとコントローラとキャッシュメモリとを備え、それぞれ前記複数の第2のHDDの一部の領域から構成される第2のボリュームも有し、

前記第1のストレージ装置及び前記第2のストレージ装置は、前記第1の仮想ボリュームのデータを前記第2の仮想ボリュームへコピーするリモートコピーを設定し、

前記ホストコンピューターは、前記ファイルシステムがライトリクエスト処理を行う際、前記リモートコピーに関する情報を管理するI/Oパスマネージャーが前記ファイルシステムを通じた前記アプリケーションからのライトリクエストを、前記リモートコピーに関する情報に基づいて、前記HBAデバイスドライバーを通じて前記第1の仮想ボリュームが存在する前記第1のストレージ装置へ送信し、

前記第1のストレージ装置は、受信した前記ライトリクエストが前記第1の仮想ボリューム宛なのか前記第1のストレージ装置の前記第2のボリューム宛なのかを判断し、

前記第1の仮想ボリューム宛の場合、前記ライトリクエストのデータを前記第1のストレージ装置及び前記第2のストレージ装置の前記キャッシュメモリに記憶し、

前記第1のストレージ装置の前記キャッシュメモリからデステージするデータを特定し、当該特定したデータをデステージしてから前記第2のストレージ装置の前記キャッシュメモリに記憶された当該データの破棄を指示する

10

ことを特徴とする情報システムのデータ転送方法。

【請求項2】

請求項1記載の方法であって、

前記第1のストレージ装置と、前記第2のストレージ装置との間の通信に障害が発生した場合、前記第1のストレージ装置は、前記通信障害を前記ホストコンピューターに通知し、前記ホストコンピューターは、前記第1のストレージ装置と、前記第2のストレージ装置とにライトリクエストを発行する

ことを特徴とする情報システムのデータ転送方法。

20

【請求項3】

請求項1記載の方法であって、

前記第1のストレージ装置の障害にかかわらずアプリケーションが継続処理可能であることを特徴とする情報システムのデータ転送方法。

【請求項4】

ファイルシステム、HBAデバイスドライバーを備えるオペレーティングシステム、及び前記ファイルシステムを通じてライトリクエストを送信するアプリケーションを備えるホストコンピューターと、前記ホストコンピューターに接続される第1のストレージ装置と、前記第1のストレージ装置及び前記ホストコンピューターに接続される第2のストレージ装置と、前記第1のストレージ装置及び前記第2のストレージ装置に接続され、複数の第1のHDDとコントローラーとキャッシュメモリから構成され、前記複数の第1のHDDの一部の領域から構成される第1のボリュームを有する第3のストレージ装置と、を有する情報システムのデータ転送方法であって、

30

前記第1のストレージ装置は、前記第3のストレージ装置の前記第1のボリュームに実体が存在する第1の仮想ボリュームを定義し、

前記第2のストレージ装置は、前記第3のストレージ装置の前記第1のボリュームに実体が存在し、当該第1のボリュームを前記第1の仮想ボリュームと共有する第2の仮想ボリュームを定義し、

前記第1のストレージ装置及び前記第2のストレージ装置は、それぞれ複数の第2のHDDとコントローラーとキャッシュメモリとを備え、それぞれ前記複数の第2のHDDの一部の領域から構成される第2のボリュームも有し、

40

前記第1のストレージ装置及び前記第2のストレージ装置は、前記第1の仮想ボリュームのデータを前記第2の仮想ボリュームへコピーするリモートコピーを設定し、

前記ホストコンピューターは、前記ファイルシステムがライトリクエスト処理を行う際、前記リモートコピーに関する情報を管理するI/Oパスマネージャーが前記ファイルシステムを通じた前記アプリケーションからのライトリクエストを、前記リモートコピーに関する情報に基づいて、前記HBAデバイスドライバーを通じて前記第1の仮想ボリュームが存在する前記第1のストレージ装置へ送信し、

前記第1のストレージ装置は、受信した前記ライトリクエストが前記第1の仮想ボリューム宛なのか前記第1のストレージ装置の前記第2のボリューム宛なのかを判断し、

50

前記第1の仮想ボリューム宛の場合、前記ライトリクエストのデータを前記第1のストレージ装置及び前記第2のストレージ装置の前記キャッシュメモリに記憶し、

前記第1のストレージ装置の前記キャッシュメモリからデステージするデータを特定し、当該特定したデータをデステージする一方、

前記第1のストレージ装置と、前記第2のストレージ装置との間の通信に障害が発生した場合、前記第1のストレージ装置は、前記通信障害を前記ホストコンピュータに通知し、前記ホストコンピュータは、前記第1のストレージ装置と、前記第2のストレージ装置とにライトリクエストを発行し、

前記第1のストレージ装置及び前記第2のストレージ装置は、自らを正系として前記ライトリクエストを処理する場合には、前記第3のストレージ装置の前記第1のボリュームに対して排他制御を行った後に、当該第1のストレージ装置又は当該第2のストレージ装置の前記キャッシュメモリに記憶した前記データを前記第1のボリュームにデステージする

10

ことを特徴とする情報システムのデータ転送方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数の記憶領域を有するストレージシステムと、ストレージシステムに接続されたホストコンピュータと、を備えた情報システムのデータ転送方法に関する。

【背景技術】

20

【0002】

一般に、情報システムでは、記憶デバイスとしてHDD（ハードディスクドライブ）を用いるストレージ装置が備えられ、そのストレージ装置を含むストレージシステムが、ストレージエリアネットワーク（SAN：Storage Area Network）経由で、複数の上位装置（例えばホスト）からアクセスされる。一般的にストレージ装置では、RAID（Redundant Array of Independent (or Inexpensive) Disks）技術に従う高信頼化方法が採用されることでHDD単体の信頼性を超えた信頼性をストレージ装置として提供している。しかし、近年の情報化社会の進化によって上記RAIDによる信頼性が貢献する情報システムの可用性（サービス継続性）では不足してきた。

【0003】

30

このような状況に対応する高可用化技術として、特許文献1に開示された技術がある。当該技術では、ホストコンピュータ（以後ホストと省略する）とストレージ装置をそれぞれ有するプロダクションサイトとバックアップサイトを用意し、プロダクションサイトのストレージ装置が保存するデータをバックアップサイトのストレージ装置にミラーリングする。もし、プロダクションサイトのストレージ装置が障害停止した場合は、バックアップサイトのストレージ装置とホストを用いて装置停止の結果停止していたアプリケーション処理を再開させる。本技術は一般的にリモートコピー又はリモートミラーリングと呼ばれる。

【0004】

【特許文献1】特開平7 244597号

40

【特許文献2】米国特許 7,080,197号

【発明の開示】

【発明が解決しようとする課題】

【0005】

特許文献1の技術ではストレージ装置の障害停止の結果、異なるホストでアプリケーションを再開させるため、アプリケーションの再起動処理が必要になる。当然ながらアプリケーション停止から再起動完了まではアプリケーションは通常動作ができないため、可用性については問題がある。本発明は、2台以上のストレージ装置間でリモートコピーを行うストレージシステムと当該ストレージシステムを利用するホストを含めた情報システムの可用性を向上させることを目的とする。

50

【課題を解決するための手段】

【0006】

かかる課題を解決するため本発明においては、ファイルシステム、HBAデバイスドライバを備えるオペレーティングシステム、及び前記ファイルシステムを通じてライトリクエストを送信するアプリケーションを備えるホストコンピュータと、前記ホストコンピュータに接続される第1のストレージ装置と、前記第1のストレージ装置及び前記ホストコンピュータに接続される第2のストレージ装置と、前記第1のストレージ装置及び前記第2のストレージ装置に接続され、複数の第1のHDDとコントローラとキャッシュメモリから構成され、前記複数の第1のHDDの一部の領域から構成される第1のボリュームを有する第3のストレージ装置と、を有する情報システムのデータ転送方法であって、前記第1のストレージ装置は、前記第3のストレージ装置の前記第1のボリュームに実体が存在する第1の仮想ボリュームを定義し、前記第2のストレージ装置は、前記第3のストレージ装置の前記第1のボリュームに実体が存在し、当該第1のボリュームを前記第1の仮想ボリュームと共有する第2の仮想ボリュームを定義し、前記第1のストレージ装置及び前記第2のストレージ装置は、それぞれ複数の第2のHDDとコントローラとキャッシュメモリとを備え、それぞれ前記複数の第2のHDDの一部の領域から構成される第2のボリュームも有し、前記第1のストレージ装置及び前記第2のストレージ装置は、前記第1の仮想ボリュームのデータを前記第2の仮想ボリュームへコピーするリモートコピーを設定し、前記ホストコンピュータは、前記ファイルシステムがライトリクエスト処理を行う際、前記リモートコピーに関する情報を管理するI/Oパスマネージャが前記ファイルシステムを通じた前記アプリケーションからのライトリクエストを、前記リモートコピーに関する情報に基づいて、前記HBAデバイスドライバを通じて前記第1の仮想ボリュームが存在する前記第1のストレージ装置へ送信し、前記第1のストレージ装置は、受信した前記ライトリクエストが前記第1の仮想ボリューム宛なのか前記第1のストレージ装置の前記第2のボリューム宛なのかを判断し、前記第1の仮想ボリューム宛の場合、前記ライトリクエストのデータを前記第1のストレージ装置及び前記第2のストレージ装置の前記キャッシュメモリに記憶し、前記第1のストレージ装置の前記キャッシュメモリからデステージするデータを特定し、当該特定したデータをデステージしてから前記第2のストレージ装置の前記キャッシュメモリに記憶された当該データの破棄を指示することを特徴とする。

10

20

30

【0007】

また本発明においては、ファイルシステム、HBAデバイスドライバを備えるオペレーティングシステム、及び前記ファイルシステムを通じてライトリクエストを送信するアプリケーションを備えるホストコンピュータと、前記ホストコンピュータに接続される第1のストレージ装置と、前記第1のストレージ装置及び前記ホストコンピュータに接続される第2のストレージ装置と、前記第1のストレージ装置及び前記第2のストレージ装置に接続され、複数の第1のHDDとコントローラとキャッシュメモリから構成され、前記複数の第1のHDDの一部の領域から構成される第1のボリュームを有する第3のストレージ装置と、を有する情報システムのデータ転送方法であって、前記第1のストレージ装置は、前記第3のストレージ装置の前記第1のボリュームに実体が存在する第1の仮想ボリュームを定義し、前記第2のストレージ装置は、前記第3のストレージ装置の前記第1のボリュームに実体が存在し、前記第1のボリュームを前記第1の仮想ボリュームと共有する第2の仮想ボリュームを定義し、前記第1のストレージ装置及び前記第2のストレージ装置は、それぞれ複数の第2のHDDとコントローラとキャッシュメモリとを備え、それぞれ前記複数の第2のHDDの一部の領域から構成される第2のボリュームも有し、前記第1のストレージ装置及び前記第2のストレージ装置は、前記第1の仮想ボリュームのデータを前記第2の仮想ボリュームへコピーするリモートコピーを設定し、前記ホストコンピュータは、前記ファイルシステムがライトリクエスト処理を行う際、前記リモートコピーに関する情報を管理するI/Oパスマネージャが前記ファイルシステムを通じた前記アプリケーションからのライトリクエストを、前記リモートコピーに関す

40

50

る情報に基づいて、前記HBAデバイスドライバーを通じて前記第1の仮想ボリュームが存在する前記第1のストレージ装置へ送信し、前記第1のストレージ装置は、受信した前記ライトリクエストが前記第1の仮想ボリューム宛なのか前記第1のストレージ装置の前記第2のボリューム宛なのかを判断し、前記第1の仮想ボリューム宛の場合、前記ライトリクエストのデータを前記第1のストレージ装置及び前記第2のストレージ装置の前記キャッシュメモリに記憶し、前記第1のストレージ装置の前記キャッシュメモリからデステージするデータを特定し、当該特定したデータをデステージする一方、前記第1のストレージ装置と、前記第2のストレージ装置との間の通信に障害が発生した場合、前記第1のストレージ装置は、前記通信障害を前記ホストコンピュータに通知し、前記ホストコンピュータは、前記第1のストレージ装置と、前記第2のストレージ装置とにライトリクエストを発行し、前記第1のストレージ装置及び前記第2のストレージ装置は、自らを正系として前記ライトリクエストを処理する場合には、前記第3のストレージ装置の前記第1のボリュームに対して排他制御を行った後に、当該第1のストレージ装置又は当該第2のストレージ装置の前記キャッシュメモリに記憶した前記データを前記第1のボリュームにデステージすることを特徴とする。

10

【0008】

さらに、本発明においては、実施の形態として、前記第1のストレージ装置と、前記第2のストレージ装置との間の通信に障害が発生した場合、前記第1のストレージ装置は、前記通信障害を前記ホストコンピュータに通知し、前記ホストコンピュータは、前記第1のストレージ装置と、前記第2のストレージ装置とに書き込み要求を発行する。また、前記第1のストレージ装置の障害にかかわらずアプリケーションが継続処理可能である。

20

【発明の効果】

【0009】

本発明によれば、2台以上のストレージ装置間でリモートコピーを行う情報システムと当該情報システムを利用するデータ転送の可用性を向上させることができる。

【発明を実施するための最良の形態】

【0010】

以下、図面を参照して、本発明の実施の形態を説明する。

【0011】

(1) 第1の実施の形態

30

< 1. 情報システムの構成 >

図1は、本発明の一実施の形態に係る情報システムのハードウェア構成の一例を示す図である。

【0012】

情報システムは、例えば、ストレージ装置1500、ホストコンピュータ（以後ホストと省略する）1100、管理ホスト1200と、2台以上の仮想化ストレージ装置1000とから構成される。ストレージ装置1500、ホストコンピュータ（以後ホストと省略する）1100、管理ホスト1200の数は、それぞれ、1以上とすることができる。仮想化ストレージ装置1000とホスト1100は、I/Oネットワーク1300を介して、相互に接続される。仮想化ストレージ装置1000とストレージ装置1500と管理ホスト1200は、管理ネットワーク（図示せず）又はI/Oネットワーク1300を介して相互に接続される。

40

【0013】

ホスト1100には、ホスト内部ネットワーク1104があり、そのネットワーク1104に、プロセッサ（図中ではProcと略記）1101と、メモリ（図中ではMemと略記）1102と、I/Oポート（図中ではI/O Pと略記）1103とが接続されている。管理ホスト1200も、ホスト1100と同じハードウェア構成を有することができる。なお、I/Oポートをホスト1100に追加する拡張カードをHBA（Host Bas Adapter）と呼ぶことがある。

【0014】

50

管理ホスト1200は、表示装置を有し、その表示装置に、仮想化ストレージ装置1000とストレージ装置1500の管理用の画面を表示することができる。また、管理ホスト1200は、管理操作リクエストを、ユーザー（例えば管理ホスト1200のオペレーター）から受け付け、その受け付けた管理操作リクエストを、仮想化ストレージ装置1000やストレージ装置1500に送信することができる。管理操作リクエストは、仮想化ストレージ装置1000やストレージ装置1500の操作のためのリクエストであり、例えば、パリティグループ作成リクエスト、内部LU（Logical Unit）作成リクエスト、バス定義リクエスト、及び仮想化機能に関する操作がある。

【0015】

I/Oネットワーク1300は、ファイバーチャネルによる接続が第一に考えられるが、それ以外でも、FICON（Fibre CONnection：登録商標）やEthernet（登録商標）とTCP/IP（Transmission Control Protocol/Internet Protocol）とiSCSI（internet SCSI（Small Computer System Interface））の組み合わせや、Ethernet（登録商標）とNFS（Network File System）やCIFS（Common Internet File System）等のネットワークファイルシステムの組み合わせ等が考えられる。さらに、I/Oネットワーク1300は、I/Oリクエストを転送可能な通信装置であればこれ以外でもよい。また、仮想化ストレージ装置1000とストレージ装置500を接続するネットワークについてもI/Oネットワーク1300と同様である。

【0016】

仮想化ストレージ装置1000は、コントローラー（図中はCTLと表記）1010と、キャッシュメモリ（図中はCMと表記）1020と、複数のHDD1030とを備える。好ましい形態としては、コントローラー1010及びキャッシュメモリ1020は、それぞれ複数のコンポーネントから構成することが考えられる。なぜなら、コンポーネント単体に障害が発生して閉塞した場合でも、残りのコンポーネントを用いてリードやライトに代表されるI/Oリクエストを引き続き受け取ることができるためである。

【0017】

コントローラー1010は、仮想化ストレージ装置1000の動作を制御する装置（例えば回路基盤）である。コントローラー1010には、内部ネットワーク1017があり、その内部ネットワーク1017に、I/Oポート1013、キャッシュポート（図中ではCPと表記）1015、管理ポート（図中ではMPと表記）1016、バックエンドポート（図中ではB/E Pと表記）1014、プロセッサ（例えばCPU（Central Processing Unit））1011及びメモリ1012が接続されている。コントローラー1010同士とキャッシュメモリ1020は、ストレージ内部ネットワーク1050にて相互に接続される。また、コントローラー1010と各HDD1030は、複数のバックエンドネットワーク1040にて相互接続される。

【0018】

ストレージ装置1500のハードウェア構成は仮想化ストレージ装置1000と同種の部品から構成される。なお、仮想化ストレージ装置1000がHDDを持たない仮想化専用装置またスイッチの場合は、ストレージ装置1500は仮想化ストレージ装置1000と同種の部品から構成されなくてもいい。さらに、ホスト1100及び仮想化ストレージ装置1000の内部のネットワークは、好ましくは、I/Oポート1013の有する転送帯域より広帯域であり、また、バスやスイッチ型のネットワークによって全てまた一部が代替されてもよい。また、図1では、I/Oポート1013は、コントローラー1010に一つ存在することになっているが、実際には、複数のI/Oポート1013がコントローラー1010に存在してもよい。

【0019】

以上のハードウェア構成によって、仮想化ストレージ装置1000やストレージ装置1500のHDDに保存された全て又は一部のデータを、ホスト1100が読出したり書き込んだりすることができるようになる。なお、以後の説明では、データ保存を担当するシステムをストレージクラスタと呼ぶ。また、ストレージクラスタ内部に当該システムを2

10

20

30

40

50

系統含むことで高可用化を実現するサブシステムで、仮想化ストレージ装置 1000 とストレージ装置 1500 の片方又は両方を含むサブシステムをストレージサブシステムと呼ぶ。

【0020】

< 2 . 本実施の形態の概要 >

本実施の形態では、他のストレージ装置内のボリューム等の記憶領域を仮想化する仮想化機能を有する仮想化ストレージ装置 1000 を含むストレージシステムの可用性を向上させるため、もう一台の仮想化ストレージ装置 1000 を用いた二重化構成を採用する。図 2 はその概要を示した図である。

【0021】

本概要では、ストレージシステムに仮想化ストレージ装置 1000 L、仮想化ストレージ装置 1000 R、ストレージ装置 1500 L、ストレージ装置 1500 R が含まれる。なお、以下においては、説明を容易にするため、仮想化ストレージ装置 1000 L 及びストレージ装置 1500 L を正系（プロダクション系）、仮想化ストレージ装置 1000 R 及びストレージ装置 1500 R を副系（バックアップ系）の役割をもっているものとする。しかし、それぞれの仮想化ストレージ装置 1000 L、1000 R がホスト 1100 へ提供するボリュームが二つ以上の場合は、仮想化ストレージ装置単位で正系・副系を担当する代わりにボリューム単位で正系を担当する仮想化ストレージ装置 1000 L、1000 R が定まっていればよい。

【0022】

それぞれの仮想化ストレージ装置 1000 L、1000 R は自身が有する HDD 1030 を構成要素とするパリティグループ（RAID 技術によって構成される）の一部又は全ての領域をボリューム 3000 LA やボリューム 3000 RA としてホスト 1100 に提供する（図中の円柱内に 'A' と記された部分に対応）。また、仮想化ストレージ装置 1000 はオプションとして仮想化機能による仮想ボリューム 3000 LB、3000 RB（対応する HDD 等の不揮発記憶領域が仮想化ストレージ装置 1000 L、1000 R の外部に存在するボリュームのこと）を提供することができる。本概要ではストレージ装置 1500 L、1500 R が提供するボリューム 3500 LB、3500 RB の一部又は全てを対応する不揮発記憶領域として用いている。なお、以後の説明では「ボリュームのデータ」と書いた場合は、HDD 1030 に保存されたデータに加えてキャッシュメモリ 1020 に一時保存されたデータも含む。また、後ほど述べる「仮想ボリュームのデータ」に関してはストレージ装置 1500 L、1500 R のボリューム 3500 LB、3500 RB に保存されたデータに加えて仮想化ストレージ装置 1000 L、1000 R のキャッシュメモリ 1020 に一時保存されたデータを含む。

【0023】

一方、ホスト 1100 上ではアプリケーションプログラム（以後、アプリケーションと略すことがある）2010 と、OS と、OS の設定・処理を補佐するデーモンや管理プログラムに代表されるシステムプログラムとが動作している。OS はアプリケーション 2010 に対して仮想化ストレージ装置 1000 L、1000 R が提供するボリューム 3000 LA、3000 LB、3000 RA、3000 RB 内に存在するデータに対する I/O リクエスト用インターフェースを提供し、アプリケーション 2010 からの要求に応じて適切な仮想化ストレージ装置 1000 L、1000 R 及びボリューム 3000 LA、3000 LB、3000 RA、3000 RB に対する I/O リクエストを送信する。通常状態ではホスト 1100 は仮想化ストレージ装置 1000 L のボリューム 3000 LA、3000 LB に対してリードやライトに代表される I/O リクエストを発行し、データの送受信を行う。つまり、リードリクエストを受け取った場合、仮想化ストレージ装置 1000 L は、リクエスト対象のボリューム 3000 LA、3000 LB、3500 LB が仮想化ストレージ装置 1000 L 内部の HDD 1030 に対応している場合は当該 HDD 1030 からデータを読み上げてこれをホスト 1100 に返したり、ストレージ装置 1500 L に対してリードリクエストを発行することで必要なデータを取得し、そのデータ（の全て

10

20

30

40

50

又は一部)をホスト1100に返す。

【0024】

ライトリクエストの場合は、データの冗長化のために、ライトデータを受け取った仮想化ストレージ装置1000Lは副系である仮想化ストレージ装置1000Rへライトデータを送信し、仮想化ストレージ装置1000Lがライトデータの受け取り完了メッセージを仮想化ストレージ装置1000Rから受け取った後にホスト1100に対してライト完了メッセージを返す。なお、仮想化ストレージ装置1000Lに対するライトデータも仮想化ストレージ装置1000Rが仮想化ストレージ装置1000Lを経由して受け取ったライトデータも、各仮想化ストレージ装置1000L, 1000R内のキャッシュメモリ1020L, 1020Rに一時保持されてもよい。なお、本実施の形態の一つとして、このライトデータの転送はストレージリモートコピーによって行われる。

10

【0025】

図3は通常状態で仮想化ストレージ装置1000Lに障害が発生した後の情報システムの処理概要を示している。

【0026】

正系の仮想化ストレージ装置1000Lが障害によって停止した場合、ホスト1100上のシステムプログラムはその障害を検知し、I/Oリクエストの発行先を正系の仮想化ストレージ装置1000Lから副系の仮想化ストレージ装置1000Rへ切り替える。ただし、その場合もアプリケーション2010はI/Oリクエストの発行先が切り替わったことを認識せずにI/Oを継続することができる。そのために、通常時からシステムプログラムはOSレイヤ(より具体的にはファイルシステムより下位のレイヤ)にて、アプリケーション2010やファイルシステムからI/Oリクエスト時に指定されるボリューム識別子として仮想的なボリューム識別子(又はデバイスファイル)を指定させるようにしておき、OSの下位レイヤは当該識別子と実際のボリュームに対して割り当てられた識別子(又はデバイスファイル)の対応を管理しておく。I/Oリクエストの発行先を切り替える場合は、その対応関係をこれまでの仮想化ストレージ装置1000Lのボリューム3000LA、ボリューム3000LB宛であったものを仮想化ストレージ装置1000Rのボリューム3000RAとボリューム3000RB宛に切り替えることでアプリケーション2010に対して透過に切り替えを実現する。

20

【0027】

さらに、仮想化ストレージ装置1000Rもホスト1100からの当該ボリューム3000RA, 3000RBに対するライトリクエストの到着やその他明示的なフェイルオーバー要求に応じて、ライトリクエストを処理できるようにする。この変更処理の一例としては、仮想化ストレージ装置1000Lから仮想化ストレージ装置1000Rに対するデータコピーに伴い、仮想化ストレージ装置1000Rのボリューム3000RA, 3000RBに対するホスト1100からのライトリクエストの拒否が設定されている場合はそれを解除する。また、リモートコピーを用いてライトデータの転送を行っている場合はリモートコピーのコピー状態の変更を行うことも考えられる。

30

【0028】

図4は仮想化ストレージ装置1000L, 1000R間のネットワークに障害が発生した後の情報システムの処理概要を示している。

40

【0029】

ネットワーク障害を検知した仮想化ストレージ装置1000Lはホスト1100に当該障害を通知する。障害通知を受けたホスト1100は、副系の仮想化ストレージ装置1000Rに対してライトリクエストを処理できるように要求し、以後のライトリクエストは正系の仮想化ストレージ装置1000L及び副系の仮想化ストレージ装置1000Rの両方に発行することで、正系と副系のデータを同一にする。

【0030】

< 3. ホスト1100で実行されるプログラム及び情報 >

図5はホスト1100上で実行されるソフトウェアプログラムと、当該ソフトウェアプ

50

プログラムが用いる情報とに加えて、各ソフトウェアプログラムが提供する概念について記した図である。なお、当該ソフトウェアプログラムはメモリ 1102 (図1) とプロセッサ 1101 (図1) とによって保持と実行がされるが、その一部をハードウェア化して実行してもよい。

【0031】

ホスト 1100 上ではアプリケーション 2010、リモートコピーマネージャ 5030 に加えて、OS 又は Kernel 内部のプログラムモジュールとしてファイルシステム 5020、I/O パスマネージャ 5000 及び HBA デバイスドライバ 5010 が実行される (ファイルシステム 5020、I/O パスマネージャ 5000 又は HBA デバイスドライバ 5010 は、全ての処理が Kernel 内部で実行される必要はない)。

10

【0032】

HBA デバイスドライバ 5010 は HBA に搭載された I/O ポート 1103 (図1) を通じて I/O リクエストやそれに伴うデータを送受信したり、その他の仮想化ストレージ装置 1000L, 1000R やストレージ装置 1500L, 1500R 等との通信を制御するプログラムである。HBA デバイスドライバ 5010 は、また、上位レイヤに対して仮想化ストレージ装置 1000L, 1000R が提供するボリューム 3000LA, 3000LB, 3000RA, 3000RB に対応する識別子を提供し、その識別子を伴った I/O リクエストを受け付けることができる。ボリューム 5040 はその概念を示したもので、仮想ストレージ装置 1000L, 1000R が提供するボリューム 3000LA, 3000LB, 3000RA, 3000RB にそれぞれ対応している。

20

【0033】

I/O パスマネージャ 5000 は、アプリケーション 2010 の I/O リクエスト発信先を切り替えるためのモジュールである。当該モジュールは HBA デバイスドライバ 5010 が提供するボリューム 5040 に対応する識別子と同種のホスト 1100 内での仮想的なボリュームに対応する識別子及び I/O リクエスト用インターフェースをファイルシステム 5020 に対して提供する。このホスト 1100 内での仮想的なボリュームに対応する識別子は当該モジュール内で HBA デバイスドライバ 5010 が提供するボリューム 5040 に対応する識別子と対応しており、デバイス関係テーブル 5001 がその対応関係を保持している。ボリューム 5050 はこのホスト 1100 内での仮想的なボリュームの概念を示したもので、本図ではその対応関係の一例として仮想化ストレージ装置 1000L のボリューム 3000LA, 3000LB に対応する識別子と対応している (他の言い方をすると、ホスト 1100 内での仮想的なボリューム 5050 の実体は仮想化ストレージ装置 1000L のボリューム 3000LA, 3000LB であるともいえる)。

30

【0034】

ここまでのレイヤでの I/O リクエストは通常固定長ブロックアクセス形式で指定する。ただし、ホスト 1100 がメインフレームの場合はこれに限定されず、CKD (Count Key Data) 形式で指定してもよい。

【0035】

ファイルシステム 5020 は、HBA デバイスドライバ 5010 が提供するボリューム 5040 に対応する識別子及び I/O インターフェースと、I/O パスマネージャ 5000 が提供するホスト 1100 内での仮想的なボリューム 5050 に対応する識別子及び I/O インターフェースとを通じて、仮想化ストレージ装置 1000L, 1000R への I/O リクエストを送信したり、データの送受信を行うモジュールである。図5では例としてファイルシステム 5020 内部にディレクトリツリーの構造を示し、そのツリー構造の一部 5052 が、I/O パスマネージャ 5000 がホスト 1100 内での仮想化で提供したボリューム 5050 に保存されている状態を示している (これまで説明した通り、より正確には I/O パスマネージャ 5000 のホスト 1100 内での仮想的なボリューム 5050 の提供は識別子を通じたものであり、さらに、そのボリューム 5050 に保存されていると書いたデータは実際にはデバイス関係テーブル 5001 にて示される仮想

40

50

化ストレージ装置 1000L, 1000R が提供するボリューム 3000LA, 3000LB, 3000RA, 3000PB に保存されている)。ファイルシステム 5020 はアプリケーション 2010 に対してファイル I/O のインターフェースを提供する。ファイル I/O インターフェースを通じてアプリケーション 2010 から呼び出されたファイルシステム 5020 は、ファイル名とファイル内でのデータオフセットを伴ったリード又はライトリクエストをディレクトリファイルや inode といったファイルシステム 5020 内の構造化情報を参照しつつ、ブロック形式のリード又はライトリクエストに変換し、I/O パスマネージャ 5000 又は HBA デバイスドライバ 5010 ヘリッド又はライトリクエストを渡す。

【0036】

なお、Unit 系や Windows (登録商標) 系の OS ではファイル I/O のインターフェースを用いて直接ボリュームのデータを操作するためのインターフェースとしてデバイスファイルシステムと呼ばれる機能を提供している。通常、デバイスファイルシステムはファイル空間の '/dev' ディレクトリ配下に展開されており、当該ディレクトリ以下のファイル(図中の例では、rsda 等)のファイル名はファイルシステム 5020 の下位レイヤ(HBA デバイスドライバ 5010 や I/O パスマネージャ 5000)が提供するボリューム 5040, 5050 に対応する。そして、当該ボリューム 5040, 5050 に保存されたデータはデバイスファイル 5070, 5080 に保存されたデータであるかのようにファイル I/O 用インターフェースで読み書き可能となる。なお、図 5 では例としてデバイスファイル 5070 (rsda, rsdb, rsdc, rsdd) は HBA デバイスドライバ 5010 が認識し、提供しているボリューム 5040 に対応し、デバイスファイル 5080 (vsda, vsdb) は I/O パスマネージャ 5000 が提供しているボリューム 5050 に対応している。このデバイスファイル 5070, 5080 は、アプリケーション 2010 がデータベースである場合に、独自のデータ編成やバッファ管理を実現する目的で使われることがある。

【0037】

リモートコピーマネージャ 5030 は仮想化ストレージ装置 1000L, 1000R との間のデータ転送を実現するリモートコピーの状態を取得したり、ホスト 1100 や I/O パスマネージャ 5000 がリモートコピーの操作を行うためのプログラムで、当該プログラムを使用するプログラム、ユーザー又は I/O パスマネージャ 5000 の要求に応じて仮想化ストレージ装置 1000L, 1000R と通信を行う。

【0038】

なお、これまで説明した通り HBA デバイスドライバ 5010 や I/O パスマネージャ 5000 は一部又は全ての機能が Kernel 内部のモジュールとしてインストールやアンインストールすることができることが望ましい。なぜならば、HBA デバイスドライバ 5020 は HBA を制御するプログラムであるが故、HBA の製造会社が提供することが多い。同様に I/O パスマネージャ 5000 は仮想化ストレージ装置 1000L, 1000R の処理を前提として処理が決定されるため、一部又は全てのモジュールが仮想化ストレージ装置 1000L, 1000R の製造会社が提供することが考えられる。したがって、当該プログラムがインストール・アンインストールできることによって幅広い HBA と仮想化ストレージ装置 1000L, 1000R の組み合わせによる情報システムを構築することができる。また、本発明ではアプリケーション 2010 に対して透過に正系と副系の切り替えを行うために Kernel 内部で処理を実行することでアプリケーション 2010 の再コンパイル等が不要な透過的な切り替えが可能である。さらに、I/O パスマネージャ 5000 がファイルシステム 5020 と HBA デバイスドライバ 5010 の中間レイヤに存在することで、ファイルシステム 5020 に対する再コンパイル等を不要とし、さらにファイルシステム透過性も確保している。そして、I/O パスマネージャ 5000 が HBA デバイスドライバ 5010 の機能を利用することができるようになっている。

【0039】

また、Kernel内部にいるI/Oパスマネージャ5000がリモートコピーマネージャ5030を呼び出す場合やその逆の通信方法として以下の二通りが考えられる。

【0040】

(A) I/Oパスマネージャ5000は通信用の仮想的なボリュームを作成し、ファイルシステム5020はこの通信用ボリュームをデバイスファイルとしてファイル空間に作成する。リモートコピーマネージャ5030は定期的にデバイスファイルに対してリードシステムコールを実行した状態で待つ。I/Oパスマネージャ5000はリモートコピーマネージャ5030からのI/Oリクエストを受信するが、内部で保留する。そして、当該モジュールがリモートコピーマネージャ5030に対するメッセージ送信をする必要が出てきたらI/Oリクエストの返り値として定められたメッセージを含むデータをファイルシステム5020を通じてリモートコピーマネージャ5030に返す。なおこの際リモートコピーマネージャが発行するリードシステムコールは長時間Kernel内部で待たされることになる。それが好ましくない場合は、I/Oパスマネージャ5000が、一定時間経過後に何もメッセージがない旨のデータをファイルシステム5020を通じてリモートコピーマネージャ5030へ返し、それを受信したリモートコピーマネージャ5030が再度リードシステムコールを実行すればよい。

10

【0041】

(B) Unix(登録商標)ドメインソケットを用いて仮想的なネットワーク通信として扱う。具体的には、ソケットの一方のエンドをリモートコピーマネージャ5030が操作し、残りのエンドをI/Oパスマネージャ5000が操作する。

20

【0042】

なお、以後の説明ではI/Oパスマネージャ5000がリモートコピーの操作や状態参照を行う場合はこのような通信によってリモートコピーマネージャ5030を呼び出すことで操作を行っているものとする。

【0043】

< 4. 仮想ストレージ装置1000で実行されるプログラム及び情報 >

図6は、仮想化ストレージ装置1000(1000L, 1000R)とストレージ装置1500(1500L, 1500R)とで実行されるプログラムと、当該プログラムにより管理される情報とについて示した図である。なお、当該プログラムはメモリ1012(図1)と、プロセッサ1011(図1)と、キャッシュメモリ1020とによって保持と実行がされるが、その一部をハードウェア化して実行してもよい。

30

【0044】

< 4.1. I/O処理プログラム6020、パリティグループ情報6060及びボリューム情報6050 >

パリティグループ情報6060には、パリティグループ毎の以下の構成に関連する情報が含まれる。

(1) パリティグループを構成するHDD1030の識別子。パリティグループには複数のHDD1030が参加しているため、当該情報はパリティグループ毎に複数存在する。

(2) RAIDレベル

【0045】

また、ボリューム情報6050には、ボリューム毎の以下の構成に関連する情報が含まれる。

(1) ボリューム容量

(2) ボリュームに対応するデータが保存されるパリティグループの識別子とパリティグループ内の領域(開始アドレスと終了アドレスの片方又は両方)。

40

【0046】

I/O処理プログラム6020は、ボリューム情報6050やパリティグループ情報6060を参照してホスト1100から受信したI/Oリクエストに関する以下の処理を実行する。

【0047】

50

(A) ステージング：HDD 1030に保存されたデータをキャッシュメモリ1020上にコピーする。

(B) デステージング：キャッシュメモリ1020に保存されたデータをHDD 1030へコピーする。なお、その前の処理としてRAID技術による冗長データを作成してもよい。

【0048】

(C) リード処理：ホスト1100から受信したリードリクエストに対して、当該リクエストに対応するデータがキャッシュメモリ1020上に存在するかどうか判定する。そして、当該リクエストに対応するデータがキャッシュメモリ1020上に存在しない場合は、ステージング処理を実行して当該データをキャッシュメモリ1020上にコピーした後、そのデータをホスト1100に対して送信する。なお、キャッシュメモリ1020上にかかるデータが存在する場合は、当該データをホスト1100に対して送信する。

10

【0049】

(D) ライト処理：ホスト1100から受信したライトデータをキャッシュメモリ1020上に保存する。なお、当該処理時にキャッシュメモリ1020上に十分な空き領域が無い場合はデステージング処理を実行して適切なデータをHDD 1030上にコピーした後にキャッシュメモリ1020上の当該領域を流用する。また既にキャッシュメモリ1020上に保存された領域がライトリクエストに含まれる場合は、そのまま既存のキャッシュメモリ1020上の領域へ上書きすることもある。

【0050】

20

(E) キャッシュアルゴリズム：キャッシュメモリ1020上のデータの参照頻度や参照時期等を元にLRU等のアルゴリズムによってステージングすべきHDD 1030上のデータやデステージングすべきキャッシュメモリ1020上のデータを決定する。

【0051】

< 4.2. 仮想化プログラム6030と仮想化情報6070 >

仮想化情報6070には、仮想化ボリューム毎の以下の構成に関連する情報が含まれる。

【0052】

(1) ストレージ装置1500内のボリューム内の領域とその領域が仮想ボリューム上のアドレス空間のどの領域としてホスト1100に提供するかに関する以下の情報。仮想ボリュームが複数で構成される場合は下記情報も複数存在する。

30

(1-1) 仮想ボリュームを構成する、ストレージ装置1500の識別子(又はポートの識別子)と、ボリュームの識別子と、ボリューム内の領域(開始アドレスと終了アドレス)

(1-2) 仮想ボリュームにおける領域(開始アドレスと終了アドレス)

【0053】

(2) 仮想ボリュームの容量

仮想化プログラム6030は、仮想化ストレージ装置1000が、ストレージ装置1500が提供するボリュームを用いてホスト1100にボリュームを提供するためのプログラムである。なお、仮想化プログラム6030が提供する仮想ボリュームと、それに対応するストレージ装置1500上のボリュームとの対応関係として、以下のパターンがある。

40

【0054】

(A) ストレージ装置1500上のボリューム全体を仮想ボリュームの記憶領域として用いる場合。この場合、仮想ボリュームの容量は選択したボリュームとおおよそ同容量となる(制御情報や冗長情報をストレージ装置1500上のボリュームに保存する場合。当該情報等がない場合は同一容量)。

(B) ストレージ装置1500上のボリュームの一部の領域を仮想化ボリュームに対応する保存領域として用いる場合。この場合、仮想ボリュームの容量は当該利用対象の領域容量と大体同じとなる。

50

【 0 0 5 5 】

(C) 複数のストレージ装置 1 5 0 0 上の複数のボリュームを仮想ボリュームの記憶領域として結合して用いる場合。この場合、仮想ボリュームの容量は各ボリューム容量の合計値とおおよそ同容量となる。なお、この結合方式としてはストライピングやConcatenate (複数ボリュームを連結して一つのボリュームとして扱う方法) 等がある。

(D) (C) のパターンに付随してパリティ情報やミラーデータを保存する場合。この場合、仮想ボリュームの容量はミラーデータを保存する場合は (C) の半分で、パリティを保存する場合はパリティ計算方式に依存する。ストレージ装置 1 5 0 0 内部でRAIDによる高信頼化と組み合わせることによって仮想ボリュームに保存されたデータについての信頼性がより向上する。

10

【 0 0 5 6 】

なお、いずれのパターンについても、I/Oリクエストで指定するストレージ装置識別子 (又はポート識別子) とボリューム識別子 (I/Oリクエストで用いる、仮想化ストレージ装置内又はポート配下のボリュームを識別する情報で、LUN (Logical Unit Number) や、CKD形式のCU番号とLDEV (Logical DEvice) 番号等がある) が元々のボリュームと異なる。

【 0 0 5 7 】

仮想化プログラム 6 0 3 0 は、ステージングやデステージング対象となるデータが仮想ボリュームに対応する場合にI/O処理プログラム 6 0 2 0 により呼び出され、仮想化情報 6 0 7 0 を用いて以下の処理を実行する。

20

【 0 0 5 8 】

(A) ステージング : 仮想化ボリュームとストレージ装置 1 5 0 0 のボリュームの対応関係を元に、どのストレージ装置 1 5 0 0 のボリュームに保存されたデータをキャッシュメモリ 1 0 2 0 上にコピーすべきかを決定した後に、キャッシュメモリ 1 0 2 0 上へデータコピーする。

【 0 0 5 9 】

(B) デステージング : 仮想化ボリュームとストレージ装置 1 5 0 0 のボリュームの対応関係を元に、どのストレージ装置 1 5 0 0 のボリュームへキャッシュメモリ 1 0 2 0 上のデータをコピーすべきかを決定した後に、ストレージ装置 1 5 0 0 のボリュームへデータコピーする。なお、その前の処理としてRAID技術による冗長データを作成してもよい。

30

【 0 0 6 0 】

< 4 . 3 . リモートコピープログラム 6 0 1 0 とコピーペア情報 6 0 4 0 >

コピーペア情報 6 0 4 0 はリモートコピーのコピー元ボリュームとコピー先ボリュームのコピーペア (ペアと省略することがある) 毎に以下の情報を持つ。なお、本実施の形態では、コピー元ボリューム及びコピー先ボリュームは高可用性を実現する対象ボリュームが指定されることになる :

【 0 0 6 1 】

(1) コピー元ボリュームを持つ仮想化ストレージ装置 1 0 0 0 の識別子及びボリュームの識別子

40

(2) コピー先ボリュームを持つ仮想化ストレージ装置 1 0 0 0 の識別子とボリュームの識別子

(3) コピーペアの状態 (詳細は後ほど述べる)

【 0 0 6 2 】

リモートコピープログラム 6 0 1 0 は、コピー元ボリュームに保存されたデータをコピー先ボリュームにミラーリングするプログラムであり、コピーペア情報 6 0 4 0 を参照して処理を行う。以下にリモートコピー (特に同期リモートコピー) の処理概要とペア状態について説明する。

【 0 0 6 3 】

< 4 . 3 . 1 . 同期リモートコピーのコピー処理動作 >

50

同期リモートコピーとは、前述の様に、コピー元の仮想化ストレージ装置1000がホスト1100からコピー元ボリュームに対するライトリクエストを受け付けた場合、ライトデータをコピー先の仮想化ストレージ装置1000に送信した後に、ホスト1100に対してライトリクエスト完了を返すリモートコピー方法である。

【0064】

同期リモートコピーが実行される際、コピー元ボリュームとコピー先ボリュームとのペア間におけるリモートコピーの状況を管理1200に表示したり、リモートコピーの状態を操作するために、仮想化ストレージ装置1000のコントローラ1010は、コピーペア状態(Simplex、Initial Copying、Duplex、Suspend及びDuplex Pending)と呼ばれる情報を管理する。図7に同期リモートコピーのペア状態に関する状態遷移図を示す。以下、各ペア状態について説明する。

10

【0065】

<4.3.1.1. Simplex状態>

Simplex状態は、ペアを構成するコピー元ボリュームとコピー先ボリュームとの間でコピーが開始されていない状態である。

【0066】

<4.3.1.2. 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ペアを構成するコピー元ボリューム及びコピー先ボリュームのデータ内容が同一となった状態である。本状態では、書き込み途中の領域を除けば、コピー元ボリュームのデータ及びコピー先ボリュームのデータの内容は同じとなる。なお、Duplex中及びDuplex Pending及びInitial Copying状態ではホスト1100からコピー先ボリュームへのライトリクエストは拒否される。

20

【0067】

<4.3.1.3. 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、必要ならばコピー元ボリュームからコピー先ボリュームへの初期化コピー(コピー元ボリュームに既に格納されていたデータのコピー先ボリュームへのコピー)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。

30

【0068】

<4.3.1.4. Suspend状態>

Suspend状態は、コピー元ボリュームに対する書き込みの内容をコピー先ボリュームに反映させない状態である。この状態では、ペアを構成しているコピー元ボリューム及びコピー先ボリュームのデータの内容は同じでない。ユーザーやホスト1100からの指示を契機に、ペア状態は他の状態からSuspend状態へ遷移する。それ以外に、仮想化ストレージ装置1000間のネットワーク障害等が原因で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Suspend状態に遷移することが考えられる。

【0069】

以後の説明では、後者の場合、即ち障害により生じたSuspend状態を障害Suspend状態と呼ぶことにする。障害Suspend状態となる代表的な原因としては、ネットワーク障害のほかに、コピー元ボリュームやコピー先ボリュームの障害、コントローラ1010の障害が考えられる。

40

【0070】

Suspend状態となった場合、コピー元ストレージ1000は、Suspend状態となった時点以降にコピー元ボリュームに対するライトリクエストがあると、ライトリクエストに従ってライトデータを受信し、コピー元ボリュームに保存するが、コピー先の仮想化ストレージ装置1000にはライトデータを送信しない。またコピー元の仮想化ストレージ装置1000は、書き込まれたライトデータのコピー元ボリューム上での書き込

50

み位置を差分ビットマップ等として記憶する。

【0071】

なお Suspend 状態となった時点以降にコピー先ボリュームに対してライトリクエストがあった場合には、コピー先の仮想化ストレージ装置 1000 も上記の動作を行う。また、ペアが障害 Suspend 状態となるより前に、当該ペアに対してフェンスと呼ばれる設定を行った場合、ペア状態が障害 Suspend に遷移するとコピー元ボリュームに対するライトを拒否する。なお、コピー先の仮想化ストレージ装置 1000 は障害 Suspend 状態中のコピー先ボリュームに対するライトリクエストを拒否してもよい。

【0072】

< 4.3.1.5. 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、コピー元ボリューム及びコピー先ボリュームのデータの内容を一致させるために、コピー元ボリュームからコピー先ボリュームへのデータのコピーが実行される。コピー元ボリューム及びコピー先ボリュームのデータの内容が同一になった後、ペア状態は Duplex となる。

【0073】

なお、Duplex Pending 状態におけるデータのコピーは、Suspend 状態の間、コピー元の仮想化ストレージ装置 1000 又はコピー先の仮想化ストレージ装置 1000 が記録した書き込み位置（例えば上述の差分ビットマップ等）を利用して、更新が必要な部分（即ちコピー元ボリュームとコピー先ボリュームとのデータの不一致部分）だけをコピーする差分コピーによって実行される。

【0074】

また、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ホスト 1200 の画面に表示したり、状態を遷移させても良い。

【0075】

< 4.3.1.6. ペア操作指示 >

ペア状態はホスト 1100 や管理ホスト 1200 からの以下の指示によって他の状態へ遷移する。

【0076】

(A) 初期化指示：Simplex 状態にて本指示を受信すると Initial Copying 状態へ遷移する。

(B) 再同期指示：Suspend 状態又は障害 Suspend 状態にて本指示を受信すると Duplex Pending 状態へ遷移する。

(C) 分割指示：Duplex 状態にて本指示を受信すると Suspend 状態へ遷移する。

(D) コピー方向反転指示：Duplex 状態、Suspend 状態又は障害 Suspend 状態にて本指示を受信すると、コピー元とコピー先との関係が反転する。Duplex 状態の場合は、本指示を受信することでコピー方向も反転する。

【0077】

なお、初期化指示はコピー元の仮想化ストレージ装置 1000 及びコピー元ボリュームと、コピー先の仮想化ストレージ装置 1000 及びコピー先ボリュームとを指定することが考えられ、その他の指示については既にペア関係が出来上がっているため当該関係を示す識別子（コピー元の仮想化ストレージ装置 1000 及びコピー元ボリュームと、コピー先の仮想化ストレージ装置 1000 及びコピー先ボリュームとの組み合わせもその識別子の一つである）を指示すればよい。

【0078】

< 5. ストレージ装置 1500 で実行されるプログラム及び情報 >

図 6 にはストレージ装置 1500 にて実行されるプログラム及び情報について記されているが、それぞれのプログラム及び情報は仮想化ストレージ装置 1000 と同様の動作を

10

20

30

40

50

行う。

【 0 0 7 9 】

< 6 . デバイス関係テーブル 5 0 0 1 >

図 8 はデバイス関係テーブル 5 0 0 1 が有する情報を示した図である。デバイス関係テーブル 5 0 0 1 は、I / O パスマネージャ 5 0 0 0 が提供するホスト 1 1 0 0 内で仮想的なボリューム（より正確には当該ボリュームに対応する識別子）毎に以下の情報を管理する。

【 0 0 8 0 】

(A) ホスト 1 1 0 0 内で仮想的なボリュームの識別子

(B) 関係ボリューム識別子リスト：上記ホスト 1 1 0 0 で仮想的なボリュームの実体となりうるストレージ装置 1 5 0 0 のボリュームの識別子が入る。なお、個々の識別子は I / O パスマネージャ 5 0 0 0 の下位レイヤである H B A デバイスドライバー 5 0 1 0 が割り当てた識別子を用いる。本実施の形態においては、正系の仮想化ストレージ装置 1 0 0 0 (1 0 0 0 L) が有するボリュームと副系の仮想化ストレージ装置 1 0 0 0 (1 0 0 0 R) が有するボリュームの識別子がリストアップされる（通常状態ならば）。

10

【 0 0 8 1 】

(C) 正系ボリューム：(B) でリストアップしたどちらのボリュームが正系かを示す。

(D) 障害状態

(E) ペア状態

【 0 0 8 2 】

なお、ファイルシステム 5 0 2 0 の視点からは (A) の識別子も (B) の識別子も同様の扱いとするため、(A) や (B) の識別子はそれぞれ重複が許されない。また (A) と (B) をあわせた場合にも重複が許されないため、I / O パスマネージャ 5 0 0 0 はその点を考慮して (A) の識別子を生成する必要がある。

20

【 0 0 8 3 】

< 7 . 初期化处理 >

図 9 は、I / O パスマネージャ 5 0 0 0 の初期化处理について記したフローチャートである。以下、このフローチャートを参照して、かかる初期化处理について説明する。なお、以下においては各種処理の処理主体を「I / O パスマネージャ 5 0 0 0」として説明する場合があるが、実際上は、ホスト 1 1 0 0 のプロセッサ 1 1 0 1 (図 1) が「I / O パスマネージャ 5 0 0 0」というプログラムに基づいて対応する処理を実行することは言うまでもない。

30

【 0 0 8 4 】

(S 9 0 0 1) I / O パスマネージャ 5 0 0 0 は、管理ホスト 1 2 0 0 やホスト 1 1 0 0 のユーザーからの以下の情報を含んだ初期化指示を受信する。尚、二重化システムの初期化处理として、H A (ハイ アベイラビリティ) 初期化指示ともいう。

(A) 正系の仮想化ストレージ装置 1 0 0 0 とその中のボリューム

(B) 副系の仮想化ストレージ装置 1 0 0 0 とその中のボリューム

【 0 0 8 5 】

(S 9 0 0 2) I / O パスマネージャ 5 0 0 0 は、S 9 0 0 1 で指示された仮想化ストレージ装置 1 0 0 0 の両方と通信をしてボリュームの存在の有無及び容量を取得する。

(S 9 0 0 3) I / O パスマネージャ 5 0 0 0 は、S 9 0 0 1 で指定されたボリュームが存在し、同容量であることを確認する。確認できない場合は、I / O パスマネージャ 5 0 0 0 は指示発信元へエラーを返す。

40

【 0 0 8 6 】

(S 9 0 0 4) I / O パスマネージャ 5 0 0 0 は、仮想化ストレージ装置 1 0 0 0 の一つ又は両方に対して、リモートコピー初期化指示を送信する。この初期化指示には正系のボリュームをコピー元ボリューム、副系のボリュームをコピー先ボリュームとして指示を出す。本指示によって仮想化ストレージ装置 1 0 0 0 はリモートコピーを開始する。

【 0 0 8 7 】

50

(S 9 0 0 5) I / O パスマネージャ 5 0 0 0 は、デバイス関係テーブル 5 0 0 1 に以下の情報を登録し、その後初期化指示の発信元へ初期化開始応答を返す。

(A) ホスト 1 1 0 0 内で仮想的なボリュームの識別子 (= I / O パスマネージャ 5 0 0 0 が作成した値)

(B) 関係ボリューム識別子リスト (= S 9 0 0 1 で指定された仮想化ストレージ装置 1 0 0 0 とボリュームに対応する識別子が二つ (正系及び副系の両方)) 。

(C) 正系ボリューム (= S 9 0 0 1 で指定された正系ボリューム) の識別子

(D) 障害状態 (= 副系準備中)

(E) ペア状態 (= I n i t i a l - C o p y i n g)

【 0 0 8 8 】

(S 9 0 0 6) I / O パスマネージャ 5 0 0 0 は、リモートコピーのペア状態を監視し、Duplex 状態に遷移したらデバイス関係テーブル 5 0 0 1 を以下の情報に更新する。

(D) 障害状態 (= 通常状態)

(E) ペア状態 (= Duplex)

【 0 0 8 9 】

以上の処理によって、I / O パスマネージャ 5 0 0 0 は、ユーザー指示に応じてリモートコピーの設定を含めた高可用化のための準備を開始することができる。なお、実際には S 9 0 0 5 の直後に I / O パスマネージャ 5 0 0 0 がホスト 1 1 0 0 内で仮想的なボリュームを提供できるため、ファイル形式でアクセスしたいユーザーは当該ボリュームに対するマウント指示等を出して、ファイル I / O を開始することができる。また、別な方法として I / O パスマネージャ 5 0 0 0 はリモートコピー設定前に既に高可用化すべきボリュームに対応するホスト 1 1 0 0 内で仮想的なボリュームを定義し、ファイルシステム 5 0 2 0 も当該ボリュームをマウントした状態から、ユーザーが副系となるボリュームを指定することによって上記の処理を開始してもよい。

【 0 0 9 0 】

< 8 . ライトリクエスト処理フロー >

図 1 0 は、I / O パスマネージャ 5 0 0 0 がファイルシステム 5 0 2 0 からライトリクエストを受信した時の処理フローを示した図である。

【 0 0 9 1 】

(S 1 0 0 0 1) I / O パスマネージャ 5 0 0 0 は、ファイルシステム 5 0 2 0 より、ライト先となるホスト 1 1 0 0 内の仮想的なボリュームの識別子と、当該ボリュームのライト位置と、ライト長とを含むライトリクエスト関数を呼び出される (又はメッセージを受信する) 。

(S 1 0 0 0 2) I / O パスマネージャ 5 0 0 0 は、当該仮想的なボリュームの障害状態を確認し、リモートコピー失敗状態ならば S 1 0 0 0 2 の両書き処理に制御を移し、それ以外ならば S 1 0 0 0 3 を実行する。

【 0 0 9 2 】

(S 1 0 0 0 3) I / O パスマネージャ 5 0 0 0 は、正系ボリュームに対してライトリクエストを発行する。なお、当該ライトリクエストの発行は実際は下位レイヤの H B A デバイスドライバ 5 0 1 0 を呼び出すことで実現する。

(S 1 0 0 0 4) I / O パスマネージャ 5 0 0 0 は、ライトリクエストの応答を確認し、正常終了ならばファイルシステム 5 0 2 0 に対して完了応答を返し、リモートコピー失敗なら S 1 0 0 0 2 の両書き処理に制御を移し、無応答など、これ以外の場合は S 1 0 0 1 0 の切り替え処理に制御を移す。

【 0 0 9 3 】

なお、S 1 0 0 0 2 の両書き処理は以下のステップで実行される。

(S 1 0 0 2 1) リモートコピーの設定によって、正系又は副系のボリュームに対するライトが拒否されている場合は、I / O パスマネージャ 5 0 0 0 はこの設定を解除する。

(S 1 0 0 2 2) I / O パスマネージャ 5 0 0 0 は、正系ボリュームに対してライトリ

10

20

30

40

50

クエストを発行する。

(S10023) I/Oパスマネージャ5000は、副系ボリュームに対してライトリクエストを発行する。I/Oパスマネージャ5000は、正系と副系の両方からのライトリクエスト応答の到着を待って、ファイルシステム5020に対して完了応答を返す。

【0094】

< 8.1. 切り替え処理のフロー >

以下、引き続き切り替え処理にて実現される処理を説明する。

【0095】

(S10011) I/Oパスマネージャ5000は、まず、デバイス関係テーブル5001の障害状態を参照することで副系ボリュームが使用可能であるか確認し、使用不可能だと判断した場合はファイルシステム5020に対してエラー応答を返し、利用可能であればS10012を実行する。なお、使用不可能と判断できる状態としては、副系なし(障害によって副系の仮想化ストレージ装置1000が機能していない場合や、初めから副系の仮想化ストレージ装置1000を設定していないボリュームの場合)の状態や、前述の初期化準備中の状態がある。

10

【0096】

(S10012) I/Oパスマネージャ5000は、副系の仮想化ストレージ装置1000に対してリモートコピーの停止指示を発行し、コピー状態がSuspend状態となったことを確認後、コピー方向反転指示を指示する。

(S10013) I/Oパスマネージャ5000は、副系の仮想化ストレージ装置1000に対してリモートコピーの再同期指示を発行する。なお、実際に再同期が完了してペア状態がDuplex状態に遷移するまで待つ必要はない。

20

【0097】

(S10014) I/Oパスマネージャ5000は、デバイス関係テーブル5001の正系ボリューム識別子をこれまで副系であったボリューム識別子に更新し、正系と副系を入れ替える。そして新たに正系となったボリュームに対してライトリクエストを、HBAデバイスドライバ5010を通じて送信する。

(S10015) I/Oパスマネージャ5000は、ライトリクエストの応答を確認し、正常終了ならばファイルシステム5020に対して完了応答を返し、エラーならばエラー応答を返して終了する。

30

【0098】

< 8.1.1. 両書き処理中のライトリクエスト失敗への対策 >

S10020の両書き処理中にS10022の正系ボリュームに対するライトリクエストが失敗に終わった場合は、S10010の切り替え処理に制御を移すことが考えられる。また、S10023の副系ボリュームに対するライトリクエストが失敗に終わった場合は、デバイス関係テーブル5001の障害状態を'副系なし'に変更し、ライト完了とする。

【0099】

また、両書き処理中はペア状態が障害Suspend状態であるため、仮想化ストレージ装置1000のボリュームにはリモートコピーの差分ビットマップによってライト位置が記される。しかし、両書き処理によって両ボリュームに書かれるライトデータは同一であるため、両書き処理が正常に行われている間はこの差分ビットマップへの記録を回避し、通信障害回復後の再同期処理では差分データだけコピーできるようにすることが望ましい。その解決策として、両書き処理が正常に行われている間は正系と副系両方の仮想化ストレージ装置1000の当該ボリュームの差分ビットマップを一定時間ごとに繰り返しクリアすることが考えられる。この方式ではクリア指示をライトリクエスト毎に発行する必要がなく、かつリモートコピーの再同期では対象ボリュームの全領域コピーは回避できる。なぜならば、直近に実施したクリア以後に行われた両書きのライトリクエストは両書きが失敗したライトリクエストと共にライト位置が差分ビットマップに記録されるが、両書きにて記録されたデータ領域が再同期でコピーされた場合もコピー先のデータ内容が変わ

40

50

らないため、データ不整合やコピー漏れ領域が発生しないからである。

【0100】

なお、上記解決策では正系と副系両方の差分ビットマップをクリアするために一時的にライトリクエストの処理を停止してもよい。その停止方法としてはI/Oパスマネージャー5000がファイルシステム5020から受け取ったライトリクエストを、両方の差分ビットマップのクリアが完了するまで、仮想化ストレージ装置1000へ転送しない方法が考えられるし、正系の仮想化ストレージ装置1000にて、両方の差分ビットマップのクリアが完了するまでライトリクエストの処理を保留する方法も考えられる。

【0101】

第2の回避策としては、正系と副系のボリュームに対してそれぞれ2面の差分ビットマップを割り当てる方式がある。以下にその処理内容を示す。

10

【0102】

(初期状態)正系と副系の仮想化ストレージ装置1000は、それぞれ2面の差分ビットマップの片面に対してライトリクエストの位置を記録する。そのために、両仮想化ストレージ装置1000は、アクティブ面(ライトリクエスト到着時にライト位置を記録する面を指し、もう一面の差分ビットマップは非アクティブ面と呼ぶ)に関する情報を保持・管理する。また、非アクティブ面の差分ビットマップは何も記録されていない状態が望ましい。

【0103】

(Step1)正系の仮想化ストレージ装置1000は、アクティブ面の管理情報を非アクティブ面になっていたもう一つの差分ビットマップへ更新することで、ライトリクエストの位置の記録先となる差分ビットマップを切り替え、以後のライトリクエストは切り替え後の差分ビットマップへ記録する。副系の仮想化ストレージ装置1000も同様に切り替える。なお、当該切り替え処理開始の契機はI/Oパスマネージャー5000が両仮想化ストレージ装置1000へ与える。なお、正系と副系の切り替え処理はどちらが先に実行してもよく、並列に実行してもよい。

20

【0104】

(Step2)I/Oパスマネージャー5000は、両仮想化ストレージ装置1000からの切り替え完了の応答を待ってから、両仮想化ストレージ装置1000に対して差分ビットマップのクリア指示を出す。クリア指示を受信した仮想化ストレージ装置1000は、非アクティブ面となっている差分ビットマップのライト位置をクリアし、I/Oパスマネージャー5000へ応答を返す。切り替え処理と同様に、正系と副系のクリア処理はどちらが先に実行してもよく、並列に実行してもよい。

30

【0105】

(Step3)I/Oパスマネージャー5000は、両仮想化ストレージ装置1000からのクリア完了の応答を待ち、時間経過後にStep1から再度実行する。

【0106】

本解決策の場合、通信障害回復後の再同期処理では、正系と副系のビットマップ4面の論理和を計算することで、Duplex Pending状態中に差分コピーを行う領域を決定することができる。また本方式ではビットマップの面数が多いものの、ライトリクエストの保留は必要ない。

40

【0107】

第3の解決策としては、上記第2の解決策の変形の以下の方式がある。

(初期状態)正系及び副系の仮想化ストレージ装置1000は、それぞれ2面の差分ビットマップの両面に対してライトリクエストの位置を記録する。また、両仮想化ストレージ装置1000は前回クリアを行った差分ビットマップ面に関する情報を保持・管理しておく。

【0108】

(Step1)I/Oパスマネージャー5000は、両仮想化ストレージ装置1000に対して差分ビットマップのクリア指示を出す。クリア指示を受信した仮想化ストレージ装

50

置 1 0 0 0 は、前回クリアした差分ビットマップでないもう一つの差分ビットマップのライト位置をクリアし、I/Oパスマネージャへ応答を返す。

(Step 3) I/Oパスマネージャ 5 0 0 0 は、両仮想化ストレージ装置 1 0 0 0 からのクリア完了の応答を待ち、時間経過後に Step 1 から再度実行する。

【0 1 0 9】

< 9 . リードリクエスト処理フロー >

図 1 1 は I/Oパスマネージャ 5 0 0 0 がファイルシステム 5 0 2 0 からリードリクエストを受信したときの処理内容を示すフローチャートである。

【0 1 1 0】

(S 1 1 0 0 1) I/Oパスマネージャ 5 0 0 0 は、ファイルシステム 5 0 2 0 より、リード先となるホスト内の仮想的なボリュームの識別子と、当該ボリュームのライト位置と、ライト長とを含むライトリードリクエスト関数を呼び出される(又はメッセージを受信する)。

10

【0 1 1 1】

(S 1 1 0 0 2) I/Oパスマネージャ 5 0 0 0 は、当該仮想的なボリュームの障害状態を確認し、通常状態でかつ正系ボリュームに対する I/O 負荷が高い場合(たとえば、一定 IOPS を超える場合や一定帯域を超える場合等)と判断したときには S 1 1 0 2 1 を実行し、それ以外の状態(副系なし、副系準備中、通常状態等)のときには S 1 1 0 0 3 を実行する。

【0 1 1 2】

(S 1 1 0 0 3) I/Oパスマネージャ 5 0 0 0 は、正系ボリュームに対してリードリクエストを発行する。

20

(S 1 1 0 0 4) I/Oパスマネージャ 5 0 0 0 は、リードリクエストの応答を確認し、正常終了ならばファイルシステム 5 0 2 0 に対して完了応答を返し、それ以外ならば S 1 1 0 1 0 の切り替え処理に制御を移す。

【0 1 1 3】

(S 1 1 0 2 1) I/Oパスマネージャ 5 0 0 0 は、副系ボリュームに対してリードリクエストを発行する。

(S 1 1 0 2 2) I/Oパスマネージャ 5 0 0 0 は、リードリクエストの応答を確認し、正常終了ならばファイルシステム 5 0 2 0 に対して完了応答を返し、それ以外ならば S 1 1 0 2 3 を実行する。

30

(S 1 1 0 2 3) I/Oパスマネージャ 5 0 0 0 は、デバイス関係テーブル 5 0 0 1 の障害状態を'副系なし'に更新し、S 1 1 0 0 3 を実行する。

【0 1 1 4】

< 9 . 1 . 切り替え処理のフロー >

以下、引き続き切り替え処理にて実現される処理を説明する。

(S 1 1 0 1 1) I/Oパスマネージャ 5 0 0 0 は、まず、デバイス関係テーブル 5 0 0 1 の障害状態を参照することで副系ボリュームが使用可能であるか確認し、使用不可能だと判断した場合はファイルシステム 5 0 2 0 に対してエラー応答を返し、利用可能だと判断した場合は S 1 1 0 1 2 を実行する。なお、使用不可能と判断できる状態としては、副系なし(障害によって副系の仮想化ストレージ装置 1 0 0 0 が機能してない場合や、初めから副系の仮想化ストレージ装置 1 0 0 0 を設定していないボリュームの場合)の状態や、前述の初期化準備中の状態がある。

40

【0 1 1 5】

(S 1 0 0 1 2) I/Oパスマネージャ 5 0 0 0 は、副系の仮想化ストレージ装置 1 0 0 0 に対してリモートコピーの停止指示を発行し、コピー状態が S u s p e n d 状態となったことを確認後、コピー方向反転指示を指示する。

(S 1 0 0 1 3) I/Oパスマネージャ 5 0 0 0 は、副系の仮想化ストレージ装置 1 0 0 0 に対してリモートコピーの再同期指示を発行する。なお、実際に再同期が完了してペア状態が D u p l e x 状態に遷移するまで待つ必要はない。

50

【 0 1 1 6 】

(S 1 0 0 1 4) I / O パスマネージャ 5 0 0 0 は、デバイス関係テーブル 5 0 0 1 の正系ボリューム識別子をこれまで副系であったボリュームの識別子に更新し、正系と副系を入れ替える。そして新たに正系となったボリュームに対してリードリクエストを、H B A デバイスドライバー 5 0 1 0 を通じて送信する。

(S 1 0 0 1 5) I / O パスマネージャ 5 0 0 0 は、リードリクエストの応答を確認し、正常終了ならばファイルシステム 5 0 2 0 に対して完了応答を返し、エラーならばエラー応答を返して終了する。

【 0 1 1 7 】

< 1 0 . 障害対策処理フロー >

本章では、I / O パスマネージャ 5 0 0 0 が障害を検知してから回復を完了するまでの処理の流れを説明する。なお、本処理は定期的にバックグラウンドで実行される。

【 0 1 1 8 】

< 1 0 . 1 . 仮想化ストレージ装置 1 0 0 0 間のネットワーク障害 >

(S t e p 1) I / O パスマネージャ 5 0 0 0 は、リモートコピーのペア状態を監視し、障害 S u s p e n d 状態を発見することで何らかの障害発生を検知する。

【 0 1 1 9 】

(S t e p 2) I / O パスマネージャ 5 0 0 0 は、副系の仮想化ストレージ装置 1 0 0 0 に対してリモートコピーの停止指示を発行し、コピー状態が S u s p e n d 状態となった事を確認後、コピー方向を反転し、各仮想化ストレージ装置 1 0 0 0 に対して状態問い合わせを行い、仮想化ストレージ装置 1 0 0 0 自体に障害が発生しておらず、ネットワーク障害が原因であることを確認したら、デバイス関係テーブル 5 0 0 1 の障害状態を ' リモートコピー失敗 ' に更新する。なお、本処理はストレージ管理者が行った作業結果を利用してよい。

【 0 1 2 0 】

(S t e p 3) 当該ネットワークが回復するまで待つ。

(S t e p 4) I / O パスマネージャ 5 0 0 0 は、正系の仮想化ストレージ装置 1 0 0 0 に対してペアの再同期指示を発行する。

(S t e p 5) I / O パスマネージャ 5 0 0 0 は、デバイス関係テーブル 5 0 0 1 の障害状態を ' 副系準備中 ' に更新する。

(S t e p 6) I / O パスマネージャ 5 0 0 0 は、ペア状態が D u p l e x になるまで待った後に、デバイス関係テーブル 5 0 0 1 の障害状態を ' 通常状態 ' に更新する。

【 0 1 2 1 】

< 1 0 . 2 . 正系仮想化ストレージ装置 1 0 0 0 の障害停止 >

(S t e p 1) I / O パスマネージャ 5 0 0 0 は、正系の仮想化ストレージ装置 1 0 0 0 の状態を監視することで障害発生を検知する。

(S t e p 2) I / O パスマネージャ 5 0 0 0 は、デバイス関係テーブル 5 0 0 1 の正系ボリュームの識別子を副系ボリュームの識別子に変更することで以後の I / O リクエスト先を副系の仮想化ストレージ装置 1 0 0 0 に切り替え、さらに障害状態を ' 副系なし ' に更新する。

(S t e p 3) I / O パスマネージャ 5 0 0 0 は、旧正系 (S t e p 2 にて切り替えたので現副系) の仮想化ストレージ装置 1 0 0 0 が回復するまで待つ。

【 0 1 2 2 】

(S t e p 4) I / O パスマネージャ 5 0 0 0 は、正系の仮想化ストレージ装置 1 0 0 0 に対してペアの再同期指示又は初期化指示を発行する。

(S t e p 5) I / O パスマネージャ 5 0 0 0 は、デバイス関係テーブル 5 0 0 1 の障害状態を ' 副系準備中 ' に更新する。

(S t e p 6) I / O パスマネージャ 5 0 0 0 は、ペア状態が D u p l e x になるまで待った後に、デバイス関係テーブル 5 0 0 1 の障害状態を ' 通常状態 ' に更新する。

【 0 1 2 3 】

10

20

30

40

50

< 10.3. 副系仮想化ストレージ装置1000の障害停止 >

(Step 1) I/Oパスマネージャ5000は、副系の仮想化ストレージ装置1000の状態を監視することで障害発生を検知する。

(Step 2) I/Oパスマネージャ5000は、デバイス関係テーブル5001の障害状態を'副系なし'に更新する。

(Step 3) I/Oパスマネージャ5000は、副系の仮想化ストレージ装置1000が回復するまで待つ。

【0124】

(Step 4) I/Oパスマネージャ5000は、正系の仮想化ストレージ装置1000に対してペアの再同期指示又は初期化指示を発行する。

(Step 5) I/Oパスマネージャ5000は、デバイス関係テーブル5001の障害状態を'副系準備中'に更新する。

(Step 6) I/Oパスマネージャ5000は、ペア状態がDuplexになるまで待った後に、デバイス関係テーブル5001の障害状態を'通常状態'に更新する。

【0125】

< 11. もう一つの初期化方法 >

これまでの説明では、I/Oパスマネージャ5000に出された初期化要求に応じて仮想化ストレージ装置1000にリモートコピーの設定を行ったが、以下に示す逆の方法も考えられる。

【0126】

(Step 1) 管理ホスト1200は、仮想化ストレージ装置1000に対してリモートコピーのペア初期化指示を出すことで、リモートコピーを開始する。

(Step 2) I/Oパスマネージャ5000は、スキヤニング要求を受信する。

(Step 3) I/Oパスマネージャ5000は、HBAデバイスドライバ5010を通じて各ボリュームに対するリモートコピーの設定(リモートコピー設定の有無やコピー元かコピー先か、ペアの相手となる仮想化ストレージ装置1000とボリューム)を取得する。なお、この取得方法としてI/Oネットワーク上でSCSIコマンドを使うことも考えられるし、それ以外の通信ネットワークを用いて情報を取得してもよい。

【0127】

(Step 4) I/Oパスマネージャ5000は、前ステップで取得した情報を元に、デバイス関係テーブル5001を作成し、これまで説明してきた処理を開始する。なお、当該デバイス関係テーブル5001の作成例としては以下がある。

(A) ホスト1100内で仮想的なボリュームの識別子 = I/Oパスマネージャ5000が作成した値

(B) 関係ボリューム識別子リスト = リモートコピーのコピー元ボリュームとコピー先ボリュームの識別子

(C) 正系ボリューム = リモートコピーのコピー元ボリューム

(D) 障害状態 = 仮想化ストレージ装置1000から取得したペア状態がDuplex状態ならば'通常状態'、Initial Copying又はDuplex Pending状態ならば'副系準備中'、Suspend又は障害Suspend状態ならば'リモートコピー失敗'

(E) ペア状態 = 仮想化ストレージ装置1000から取得したペア状態

【0128】

以上、これまで説明したハードウェア及びプログラムの動作によって本実施の形態では高可用性を実現する。なお、図10と図11等に記した切り替え処理に長時間要する場合の対策として、I/Oパスマネージャ5000がI/Oリクエストを再送信する必要が出てきた場合に、予備処理として前記切り替え処理の一部を実行してもよい。この場合、再送信したI/Oリクエストが正常応答で返ってきた場合は先行して行った切り替え処理を元に戻せば良く、一方で再送信したI/Oリクエストがエラー応答で返ってきたり、まったく応答がなければ前記切り替え処理の残り部分を実行すればよい。また、本実施の形

10

20

30

40

50

態は全てのボリュームが仮想化ストレージ装置 1000 によって仮想化され、実体がストレージ装置 1500 にある仮想ボリュームで、仮想化ストレージ装置 1000 は仮想化専用のアプライアンスであってもよく、またその逆に全てのボリュームの実体が仮想化ストレージ装置 1000 の内部にある構成であってもよい。また、仮想化ストレージ装置 1000 が提供するボリュームには容量以外にもさまざまな属性が設定されることがある（たとえば、エミュレーションタイプや SCSI 規格で定められた Inquiry コマンドで取得可能なボリューム識別番号がある）。

【0129】

こうした属性情報や属性変更リモートコピーによって正系の仮想化ストレージ装置から副系の仮想化ストレージ装置へ転送し、両方の仮想化ストレージ装置にて管理することも考えられる。

10

【0130】

< 12 . もう一つのリード・ライト処理 >

図 10 や図 11 に記したライト・リード処理では、I/O パスマネージャ 5000 が明示的にリモートコピーの操作を仮想化ストレージ装置 1000 へ転送する。しかし、当該リモートコピーの操作が仮想化ストレージ装置 1000 のベンダー毎に異なる場合があるため、I/O パスマネージャ 5000 のライト処理やリード処理に含めないほうが好ましい場合がある。図 25 ~ 図 27 にこうした形態での処理内容を示す。なお、以下においては各種処理の処理主体を「仮想化ストレージ装置 1000」として説明する場合があるが、実際上は、その仮想化ストレージ装置 1000 内のプロセッサ 1011 (図 1) がメモリ 1012 (図 1) に格納されたプログラムに基づいて対応する処理を実行することは言うまでもない。

20

【0131】

< 12 . 1 . I/O パスマネージャのライト処理 >

図 25 は、I/O パスマネージャ 5000 で実行される図 10 の大体処理を示したフローチャートである。以下の点が図 10 と異なる。

(相違点 1) リモートコピーの操作 S10012、S10013、S10021 がスキップされる。

(相違点 2) リモートコピー失敗時のフロー S10020 に到達しない。ただし、本相違点は通常のリード/ライト処理ではリモートコピー失敗を意味するエラーメッセージを識別できない場合に限った話である。

30

【0132】

< 12 . 2 . ストレージ装置 1000 の処理 >

図 27 は、仮想化ストレージ装置 1000 がライトリクエストを受信した時に行うリモートコピーの操作について示した図である。

【0133】

(S27001) 仮想化ストレージ装置 1000 は、ライトリクエストを受信する。

(S27002) 仮想化ストレージ装置 1000 は、ライトリクエストが対象とするボリュームがリモートコピーに関係するかどうか判断し、無関係の場合は S27003 を実行し、関係する場合は S27004 を実行する。

40

(S27003) 仮想化ストレージ装置 1000 は、通常のライト処理を行い、ホスト 1100 へ応答を返して終了する。

【0134】

(S27004) 仮想化ストレージ装置 1000 は、ライトリクエストが対象とするボリュームのリモートコピーの属性を判断し、コピー元属性の場合は S27005 を実行し、コピー先属性の場合は S27011 を実行する。

(S27005) 仮想化ストレージ装置 1000 は、同期リモートコピー処理を実行し、副系ストレージヘライトデータを転送し、応答を待つ。

(S27006) 仮想化ストレージ装置 1000 は、コピーが成功したかどうか判断し、成功ならば S27008 を実行し、失敗ならば S27007 を実行する。

50

(S 2 7 0 0 7) 仮想化ストレージ装置 1 0 0 0 は、対象ボリュームがコピー元となるリモートコピーペアの状態を障害 S u s p e n d 状態に遷移する。ただし、当該ボリュームに対するライトは禁止しない。

【 0 1 3 5 】

(S 2 7 0 0 8) 仮想化ストレージ装置 1 0 0 0 は、通常のライト処理を行い、ホスト 1 1 0 0 へ応答を返して終了する。

(S 2 7 0 1 1) 仮想化ストレージ装置 1 0 0 0 は、リモートコピーを停止し、コピー元とコピー先の関係を反転する。

(S 2 7 0 1 2) 仮想化ストレージ装置 1 0 0 0 は、再同期処理を開始する。

(S 2 7 0 1 3) 仮想化ストレージ装置 1 0 0 0 は、通常のライト処理を行い、ホスト 1 1 0 0 へ応答を返して終了する。

10

【 0 1 3 6 】

なお、S 2 7 0 1 2 の再同期処理は完了まで待たなくても良い。なぜならば、S 2 7 0 1 2 を実行する仮想化ストレージ装置 1 0 0 0 は副系であり、正系の仮想化ストレージ装置 1 0 0 0 が正常動作しているとは限らないこと、及び再同期処理が完了するまでの時間が長いことが考えられるからである。なお、こうしたケースは < 1 0 . 障害対策処理フロー > で述べた処理によって回復される点はこれまでと同じである。

【 0 1 3 7 】

< 1 2 . 3 . I / O パスマネージャのリード処理 >

図 2 6 は、I / O パスマネージャ 5 0 0 0 で実行される図 1 1 の大体処理を示したフローチャートである。以下の点が図 1 1 と異なる。

20

(相違点 1) リモートコピーの操作 S 1 1 0 1 2 、 S 1 1 0 1 3 がスキップされる。

【 0 1 3 8 】

なお、図 1 1 ではリード処理に応じてリモートコピーの向きが反転したが、本処理では反転させない。なぜならば、副系の仮想化ストレージ装置 1 0 0 0 に対するリードリクエストは正系の仮想化ストレージ装置 1 0 0 0 が (ホスト = 仮想化ストレージ装置間の通信障害による原因を含めて) 応答を返さない場合に加えて、正系の仮想化ストレージ装置 1 0 0 0 の過負荷が原因の場合もあるからである。そのため、副系の仮想化ストレージ装置 1 0 0 0 がコピー先ボリュームに対するリードリクエストを契機としてリモートコピーのペア反転を行うと、たまたま副系の仮想化ストレージ装置 1 0 0 0 に出されたリードリクエストでペアが反転し、その次のリードリクエストで再びペアが反転してしまうため、リード性能が悪化する結果となるからである。

30

【 0 1 3 9 】

ただし、S 1 1 0 2 1 の実行が抑制される場合は、仮想化ストレージ装置 1 0 0 0 はリード処理に際して以下の処理を行うことでリモートコピーのペア反転を行っても良い。

(S t e p 1) 仮想化ストレージ装置 1 0 0 0 は、リードリクエストを受信する。

(S t e p 2) 仮想化ストレージ装置 1 0 0 0 は、通常のリード処理を行う。

(S t e p 3) 仮想化ストレージ装置 1 0 0 0 は、リード対象のボリュームがリモートコピーのコピー先ボリュームであるかどうかを判断し、該当する場合は次の S t e p 4 を実行し、そうでない場合は終了する。

40

(S t e p 4) 仮想化ストレージ装置 1 0 0 0 は、リモートコピーを停止し、コピー元とコピー先の関係を反転する。

【 0 1 4 0 】

(2) 第 2 の実施の形態

次に第 2 の実施の形態について図 1 2 を用いて説明する。第 1 の実施の形態と異なる点は、ストレージ装置 1 5 0 0 L が複数の仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R に接続され、これら仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R がストレージ装置 1 5 0 0 L 内のボリュームを共有することによって、仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R の片方が停止した場合でも第 1 の実施の形態よりも低コストでサービスが継続できるようになる点である。

50

【 0 1 4 1 】

ただし、仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R はキャッシュメモリ 1 0 2 0 L , 1 0 2 0 R を有するため、仮想化ボリュームに対してライトデータを書き込んだ直後に正系の仮想化ストレージ装置 1 0 0 0 L が災害停止した場合に備えて、ライトデータを副系の仮想化ストレージ装置 1 0 0 0 R のキャッシュメモリ 1 0 2 0 R にも保存する必要があり、また両方の仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R のデステージングやステージングに対して工夫が必要となる。

【 0 1 4 2 】

通常状態におけるライトリクエストは以下のステップにて処理される。

【 0 1 4 3 】

(Step 1) ホスト 1 1 0 0 からライトリクエストを受信した正系の仮想化ストレージ装置 1 0 0 0 L は当該ライトリクエストが当該仮想化ストレージ装置 1 0 0 0 L 内部の HDD 1 0 3 0 に対応するボリューム 3 0 0 0 L A 宛なのか、両方の仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R がストレージ装置 1 5 0 0 L のボリューム 3 5 0 0 L を共有して提供する仮想化ボリューム (以後、共有仮想化ボリュームと呼ぶ) 3 0 0 0 L B 宛なのか、通常の仮想化ボリューム宛なのかを判断する。なお、共有仮想化ボリューム 3 0 0 0 L B 以外の処理については第 1 の実施の形態と同様の処理を行う。

【 0 1 4 4 】

(Step 2) 正系の仮想化ストレージ装置 1 0 0 0 L は自身のキャッシュメモリ 1 0 2 0 L に当該ライトデータを保存すると共に、当該ライトデータをリモートコピープログラムによって副系の仮想化ストレージ装置 1 0 0 0 R のキャッシュメモリ 1 0 2 0 R に保存した後に、ホスト 1 1 0 0 に対して正常応答を返す。

【 0 1 4 5 】

(Step 3) 正系の仮想化ストレージ装置 1 0 0 0 L のキャッシングアルゴリズムがデステージすべきキャッシュメモリ 1 0 2 0 L 上のデータを決定し、当該データをストレージ装置 1 5 0 0 L のボリュームにデステージする。

【 0 1 4 6 】

(Step 4) デステージ完了後、正系の仮想化ストレージ装置 1 0 0 0 L はデステージしたキャッシュメモリ 1 0 2 0 L 上のデータのアドレスを破棄するように副系の仮想化ストレージ装置 1 0 0 0 R に指示する。なお、指示を受けた副系の仮想化ストレージ装置 1 0 0 0 R は指示を受けたデータをキャッシュメモリ 1 0 2 0 R から破棄する。

【 0 1 4 7 】

なお、本構成では仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R 間のネットワークが切断された状態で副系の仮想化ストレージ装置 1 0 0 0 R に I / O リクエストの切り替えを行った場合、仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R の両方が正系として自立的にデステージングを行う場合がある。そういった状況を回避するため、両仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R は自らを正系として処理する場合は先にストレージ装置 1 5 0 0 L 内のかかる共有化されたボリューム 3 5 0 0 L に対して SCSI Reserve 等の機能を用いて排他制御を行ってもよい。また、これ以外の方式として共有仮想化ボリューム 3 0 0 0 L B については仮想化ストレージ装置 1 0 0 0 L のキャッシングを無効化してもよく、この場合は当該共有仮想ボリューム 3 0 0 0 L B のアクセス権限がリードオンリーのアクセス権限へ変更された場合は当該変更に応じてキャッシングを有効にすることが考えられる。

【 0 1 4 8 】

(3) 第 3 の実施の形態

次に第 3 の実施の形態について図 1 3 を用いて説明する。本実施の形態はこれまでの実施の形態に記した情報システムをこれまでのプロダクションサイトと異なる遠隔地 (バックアップサイト) に別途用意し、リモートコピーを行うもので、これによりプロダクションサイト被災時にバックアップサイトでサービスを再開することができる。

【 0 1 4 9 】

なお、これ以後の説明では、上述の「仮想化ストレージ装置」をストレージ装置と、「コピー元ボリューム」を正ボリュームと、「コピー先ボリューム」を副ボリュームと、「正系」をアクティブ側と、「副系」をスタンバイ側と呼ぶことがある。また、プロダクションサイトとバックアップサイトの情報システムをあわせてリモートコピーシステムと呼ぶことがある。

【0150】

< 1 . リモートコピーシステムの構成 >

本実施の形態では、各サイトはホスト13010, 13020と複数のストレージサブシステム13001, 13002, 13003, 13004とから構成されている。そしてプロダクションサイトでは、ストレージサブシステム13001, 13002同士でこれまで説明してきた高可用化構成を採用している。またバックアップサイトでも同様に、ストレージサブシステム13003, 13004同士でかかる高可用化構成を採用している。

10

【0151】

さらに本実施の形態では、プロダクションサイトのアクティブ側のストレージサブシステム(コピー元ボリュームを持つ)13001からバックアップサイトのアクティブ側のストレージサブシステム(コピー先ボリュームを持つ)13003に対して同期又は非同期リモートコピーを行う。そしてプロダクションサイト被災時にはバックアップサイトのホスト13010が高可用構成のストレージサブシステム13003, 13004のいずれがアクティブな側に対してI/Oリクエストを発行することで、再起動したアプリケーション2010が処理を再開する。

20

【0152】

なお、前述の通り、ストレージサブシステムとは仮想化ストレージ装置1000(図1)の仮想化機能を用いない設定の構成や、仮想化ストレージ装置1000とストレージ装置1500(図1)の組み合わせで仮想化ストレージ装置1000が仮想化機能を用いて仮想化ボリュームを提供している構成のどちらの概念も含んだものとして呼んでいる。また、本実施の形態では個々のストレージサブシステム13001, 13002, 13003, 13004が別々な内部構成(例えば、ストレージサブシステム13001だけ仮想化ストレージ装置1000のみで構成し、仮想化機能を用いない場合や、バックアップサイトのストレージサブシステム13003と13004でストレージ装置1500(図1)を共有し、プロダクションサイト側では共有しない場合)を採用してもよい。

30

【0153】

なお、以下においては各種処理の処理主体を「ストレージサブシステム」として説明する場合があるが、実際には、そのストレージサブシステム内のプロセッサが当該ストレージサブシステム内のメモリに格納されたプログラムに基づいて対応する処理を実行することは言うまでもない。

【0154】

< 2 . 処理 >

プロダクションサイトのホスト13010のアプリケーション2010がライトリクエストを発行すると、OSによってプロダクションサイト内のアクティブ側のストレージサブシステムを判断し、そちらにライトリクエストを転送する。なお、本図ではストレージサブシステム13001がこれに対応する。

40

【0155】

プロダクションサイトのアクティブ側のストレージサブシステム13001は同期リモートコピーによってライトデータをプロダクションサイト内のスタンバイ側のストレージサブシステム(本図では13002が対応する)へ転送する。また、アクティブ側のストレージサブシステム13001はバックアップサイトのアクティブ側のストレージサブシステム(本図では13003が対応する)へ向けて同期又は非同期のリモートコピーとしてライトデータを転送する(本実施の形態による高可用構成ではアクティブ側にのみライトリクエストを処理するようにしているため、リモートコピーであっても同様にアクティ

50

ブ側にて処理を行う)。ライトデータを受信したバックアップサイト内のアクティブ側のストレージサブシステム13003は受け取ったライトデータをサイト内のスタンバイ側のストレージサブシステム13004へ同期リモートコピーによって転送する。

【0156】

そのため、プロダクションサイトのストレージサブシステム13001, 13002はバックアップサイトのアクティブ側のストレージサブシステムを把握しており、バックアップサイトのストレージサブシステム13003, 13004も想定外のストレージサブシステムからのリモートコピーを受け付けないために、プロダクションサイトのアクティブなストレージサブシステム(ストレージサブシステム1301)を把握している。

【0157】

以上の処理によってプロダクションサイト、バックアップサイト共にサイト内の高い可用性を実現している。ただしバックアップサイト側では、コスト削減のために高可用構成をとらない構成であってもよい。

【0158】

<3. 非同期リモートコピー>

これまで説明してきた同期リモートコピーとは異なり、非同期リモートコピーはホスト13010からのライトリクエストが到着した時点でライトデータを転送するのではなく、当該リクエスト完了応答後に転送する(言い方を変えると、非同期リモートコピーはホスト13010へのリクエスト応答とは独立なタイミングでライトデータを転送する)。そのため、非同期リモートコピーはサイト間の距離が長く通信遅延が大きな場合でもライトリクエストの応答時間を低下させずにリモートコピーを行うことができる。しかし、非同期リモートコピーではプロダクションサイト側のストレージサブシステム13001にてライトデータをバッファリングする必要がある。このライトデータのバッファリング方式としては以下が考えられる。

【0159】

(1) プロダクションサイトのストレージサブシステム13001は、コピー元ボリュームへのライトデータとライトデータの順序情報を含むジャーナルを作成し、これを自身のキャッシュメモリ又は専用ボリュームに保存すると共に、このジャーナルをバックアップサイトのストレージサブシステム13003へ転送し、バックアップサイトのストレージサブシステム13003はジャーナルの順序情報を参考にコピー先ボリュームへライトデータを保存する。これにより、プロダクションサイト災害時にはライト順序が守られた(より正確には依存関係のあるライトデータ)データをバックアップサイト側で提供できる。

【0160】

(2) プロダクションサイトのストレージサブシステム13001は、ある期間毎のコピー元ボリュームへライトされたデータをグループ化して自身のキャッシュメモリ又は専用ボリュームへ保存し、非同期にバックアップサイトのストレージサブシステム13003へ転送し、当該グループ単位でバックアップサイトのストレージサブシステム13003が有するコピー先ボリュームへデータを保存する。

【0161】

そのため、これら非同期リモートコピーのためにバッファリングされるライトデータもスタンバイ側のストレージサブシステム13002で保持しなければ、アクティブ側ストレージサブシステム13001が停止したときに非同期リモートコピーを引き継ぐことができない。よって、プロダクションサイトのアクティブ側のストレージサブシステム13001はライトデータだけではなく、コピー先ボリュームの情報や、前述の順序情報や、グループ化するタイミング等をスタンバイ側のストレージサブシステム13002へ伝え、スタンバイ側のストレージサブシステム13002はそれに従ってアクティブ側と同じ非同期リモートコピーのためのバッファリングデータを作成する。

【0162】

なお、バックアップサイトのストレージサブシステム13003もプロダクションサイ

10

20

30

40

50

トから受け取ったライトデータを直ぐにコピー先ボリュームへ保存せずに、バッファリングを行うため、プロダクションサイト側と同様にアクティブ側の指示に従ってスタンバイ側も同様のバッファリングデータを作成し、また同様のタイミングでコピー先ボリュームにライトデータを保存する必要がある。

【 0 1 6 3 】

(4) 第 4 の実施の形態

次に第 4 の実施の形態について図 1 4 を用いて説明する。本実施の形態では、2 台のストレージ装置により先に説明した同期リモートコピーを用いて冗長構成された情報システムにおいて、ストレージ装置が提供する機能を制御するインターフェース (機能 I / F) の構成について述べる。

10

【 0 1 6 4 】

なお、本実施の形態から第 1 4 の実施の形態までは、これまで仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R、ストレージ装置 1 5 0 0 L , 1 5 0 0 R と呼んでいたコンポーネントを、それぞれストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B 及び外部ストレージ装置 1 6 0 0 0 A , 1 6 0 0 0 B と呼ぶ。また、以下においては各種処理の処理主体を「ストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B 」や「外部ストレージ装置 1 6 0 0 0 A , 1 6 0 0 0 B 」として説明する場合があるが、実際上は、そのストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B 内の図示しないプロセッサやその外部ストレージ装置 1 6 0 0 0 A , 1 6 0 0 0 B 内のプロセッサが当該ストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B 又は外部ストレージ装置 1 6 0 0 0 A , 1 6 0 0 0 B 内のメモリに格納されたプログラムに基づいて対応する

20

処理を実行することは言うまでもない。

【 0 1 6 5 】

本実施の形態は、ホスト 1 4 0 0 0 からの機能制御要求が、ストレージ装置 1 5 0 0 0 A に送信された後、ストレージ装置 1 5 0 0 0 A が機能制御要求をストレージ装置 1 5 0 0 0 B に転送し、ストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B の双方が当該機能制御要求を解釈し実行する例を示している。

【 0 1 6 6 】

コマンドデバイス 1 5 0 0 2 A , コマンドデバイス 1 5 0 0 2 B はそれぞれストレージ装置 1 5 0 0 0 A、ストレージ装置 1 5 0 0 0 B が提供する論理ボリュームであり、機能を制御するホスト 1 4 0 0 0 とのインターフェースとなる。なお、本実施の形態ではコマンドデバイス 1 5 0 0 2 A がアクティブ側と仮定している。

30

【 0 1 6 7 】

また、同期リモートコピーにより、コマンドデバイス 1 5 0 0 2 A の内容はコマンドデバイス 1 5 0 0 2 B の内容と常に一致している。コマンドデバイス 1 5 0 0 2 A、コマンドデバイス 1 5 0 0 2 B はオペレーティングシステム 1 4 0 0 1 が提供するパス管理機能 (I / O パスマネージャ 5 0 0 0 (図 1) が提供する機能に相当する) によりひとつのボリューム 1 4 0 0 4 として機能管理プログラム 1 4 0 0 3 に提供される。

【 0 1 6 8 】

論理ボリューム 1 5 0 0 1 A、論理ボリューム 1 5 0 0 1 B はそれぞれストレージ装置 1 5 0 0 0 A、ストレージ装置 1 5 0 0 0 B が提供する論理ボリュームであり、機能制御対象の論理ボリュームである。なお、本実施の形態では論理ボリューム 1 5 0 0 1 A がアクティブ側と仮定している。

40

【 0 1 6 9 】

また、同期リモートコピーにより、論理ボリューム 1 5 0 0 1 A の内容は、論理ボリューム 1 5 0 0 1 B の内容と常に一致している。論理ボリューム 1 5 0 0 1 A、論理ボリューム 1 5 0 0 1 B はオペレーティングシステム 1 4 0 0 1 が提供するパス管理機能によりひとつのボリューム 1 4 0 0 5 としてアプリケーションプログラム 1 4 0 0 2 に提供される。

【 0 1 7 0 】

なお、ここで説明した機能制御対象の論理ボリュームは複数あってもよい。

50

【 0 1 7 1 】

機能管理プログラム 1 4 0 0 3 の機能制御要求処理部 1 4 0 0 5 は、ユーザーあるいはホスト 1 4 0 0 0 内の他のプログラムあるいはホスト 1 4 0 0 0 とは別のホスト（管理ホストなど）内のプログラムから、機能制御要求を受け付ける。機能制御要求を受け付けた機能制御要求処理部 1 4 0 0 5 はボリューム 1 4 0 0 4 に対する制御要求の内容をボリューム 1 4 0 0 4 に対してライト/リードする。本実施の形態ではコマンドデバイス 1 5 0 0 2 A がアクティブ側であるため、ライト/リードはコマンドデバイス 1 5 0 0 2 A に対して発行される。

【 0 1 7 2 】

コマンドデバイス 1 5 0 0 2 A に対するライトは機能制御を起動するときに用いられ、コマンドデバイス 1 5 0 0 2 A に対するリードは機能制御の結果の出力値を得るために用いられる。

10

【 0 1 7 3 】

機能制御要求処理部 1 4 0 0 5 が受け付ける制御要求には制御対象のストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B を一意に識別する情報（装置情報とも呼ぶ）と、制御対象の論理ボリューム 1 5 0 0 1 A , 1 5 0 0 0 1 B を一意に識別する情報（ボリューム情報とも呼ぶ）と、機能制御に付随する情報とが含まれる。

【 0 1 7 4 】

ストレージ装置 1 5 0 0 0 A の制御 I / F 処理部 1 5 0 0 3 A はコマンドデバイス 1 5 0 0 2 A に制御要求がライトされたことを検出する。制御 I / F 処理部 1 5 0 0 3 A は制御要求の装置情報が自ストレージ装置（ストレージ装置 1 5 0 0 0 A ）に一致するか判定する（判定 1 0 0 ）。本実施の形態ではコマンドデバイス 1 5 0 0 2 A がアクティブ側なので、判定の結果は「一致する」となる。一致した場合、制御 I / F 処理部 1 5 0 0 3 A はボリューム情報に対応する論理ボリューム 1 5 0 0 1 A に対して所定の機能制御を実行するよう機能処理部 1 5 0 0 4 A を呼び出す。具体的な例としては、ストレージ装置 1 5 0 0 0 A が提供する機能のひとつであるローカルコピー機能（後で説明）のペア状態の参照操作がある。当該操作が論理ボリューム 1 5 0 0 1 A に対して呼び出された場合、機能処理部 1 5 0 0 4 A は、ローカルコピー機能の管理情報を参照し、ペア状態を取得した後、制御 I / F 処理部 1 5 0 0 3 A 、コマンドデバイス 1 5 0 0 2 A 及びボリューム 1 4 0 0 4 を介して、機能制御要求処理部 1 4 0 0 5 に対して、ペア状態を送信する。

20

30

【 0 1 7 5 】

一方、ストレージ装置 1 5 0 0 0 B の制御 I / F 処理部 1 5 0 0 3 B も同様の処理を行うが、本実施の形態では、コマンドデバイス 1 5 0 0 2 B はスタンバイ側なので、判定 1 0 0 の結果は「一致しない」となる。この場合、制御 I / F 処理部 1 5 0 0 3 B は同期リモートコピーのペアの管理情報を参照し、ボリューム情報（論理ボリューム 1 5 0 0 1 A に対応）に対応する自ストレージ装置（ストレージ装置 1 5 0 0 0 B ）内の論理ボリューム（論理ボリューム 1 5 0 0 1 B に対応）を特定する。そして、制御 I / F 処理部 1 5 0 0 3 B は論理ボリューム 1 5 0 0 1 B に対して所定の機能制御を実行するよう機能処理部 1 5 0 0 4 B を呼び出す。

【 0 1 7 6 】

以上により、ストレージ装置 1 5 0 0 0 A の論理ボリューム 1 5 0 0 1 A 、ストレージ装置 1 5 0 0 0 B の論理ボリューム 1 5 0 0 1 B に対して、所定の機能の制御が実行される。

40

【 0 1 7 7 】

本実施の形態では、ストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B が提供するローカルコピー機能のペア状態の参照操作を例にとって説明したが、（ 1 ）ローカルコピー機能のその他のペア操作（ペアの作成、ペアの分割等）、（ 2 ）ストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B が提供するローカルコピー機能の各種ペア操作、（ 3 ）ストレージ装置 1 5 0 0 0 A , 1 5 0 0 0 B が提供する論理ボリューム 1 5 0 0 1 A , 1 5 0 0 1 B に対するセキュリティ機能（後で説明する L D E V ガード機能）の操作、（ 4 ）ストレージ装置 1 5

50

000A, 15000Bが提供する論理スナップショット機能(後で説明)の操作、等、ストレージ装置15000A, 15000Bが提供する各種機能の操作について適用できる。

【0178】

なお、別な実行形態としては、アクティブ側とスタンバイ側両方のストレージ装置15000A, 15000Bに発行すべきコマンドを受けた場合は、アクティブ側のストレージ装置15000Aは受取ったコマンドを処理すると共に、スタンバイ側のストレージ装置15000Bへ転送してコマンド処理をしてもらうことで、ホスト14000からは1回のコマンドで両方のストレージ処理を開始することも考えられる。また、プログラムの状態取得に関するコマンドの場合は、コマンドを受取ったアクティブ側のストレージ装置15000Aがスタンバイ側のストレージ装置15000Bに同じコマンドを転送して状態を取得し、アクティブ側のストレージ装置15000Aが両方の状態を比較した後にコマンド発信元へ状態を返すことも考えられる。

10

【0179】

(5) 第5の実施の形態

本実施の形態では機能I/Fの別の構成について述べる。図15を用いて本実施の形態の構成を説明する。

【0180】

本実施の形態の構成は図14とほぼ同様である。図14との違いは、

(1) コマンドデバイス15002A、コマンドデバイス15002Bが同期リモートコピのペアでない。

20

(2) 機能管理プログラム14003からはコマンドデバイス15002A及びコマンドデバイス15002Bが別々のボリューム14004A、14004Bとして認識されている。

(3) 機能制御要求処理部14005は機能制御要求をコマンドデバイス15002A及びコマンドデバイス15002Bに送信する。

という3点である。

【0181】

本実施の形態では、第4の実施の形態と同様に、機能制御要求処理部14005が受け付ける制御要求には制御対象のストレージ装置15000A, 15000Bを一意に識別する情報(装置情報とも呼ぶ)と、制御対象の論理ボリューム15001A, 15001Bを一意に識別する情報(ボリューム情報とも呼ぶ)と、機能制御に付随する情報とが含まれる。

30

【0182】

本実施の形態では、第4の実施の形態と異なり、前述のように、ユーザーあるいはホスト14000内の他のプログラムあるいはホスト14000とは別のホスト内のプログラムから機能制御要求を受け付けた機能制御要求処理部14005は、両方のコマンドデバイス15002A、15002Bに制御要求を送信する。

【0183】

なお、機能制御要求処理部14005が装置情報を判定し、コマンドデバイス15002Aに対しては、ボリューム情報として論理ボリューム15001Aを指定し、コマンドデバイス15002Bに対しては、ボリューム情報として論理ボリューム15001Bを指定するように制御要求を書き換えてもよい。

40

【0184】

さらにまた、ユーザーあるいはホスト14000内の他のプログラムあるいはホスト14000とは別のホスト内のプログラムがストレージ装置15000A, 15000Bを識別し、ストレージ装置15000A、15000Bに対して二重に異なる制御要求を出してもよい。即ち、コマンドデバイス15002Aに対して、論理ボリューム15001Aの制御要求を出し、コマンドデバイス15002Bに対して、論理ボリューム15001Bの制御要求を出す。

50

【 0 1 8 5 】

(6) 第 6 の実施の形態

本実施の形態では機能 I / F の更に別の構成について述べる。図 1 6 を用いて本実施の形態の構成を説明する。

【 0 1 8 6 】

第 6 の実施の形態は第 4 の実施の形態とほぼ同様である。第 4 の実施の形態との違いは以下の点である。

(1) ホスト 1 4 0 0 0、ストレージ装置 1 5 0 0 0 A、ストレージ装置 1 5 0 0 0 B は互いに LAN (Local Area Network) のような相互結合網により接続されている。なお、これらは LAN により直結されていてもよいし、スイッチを経由して接続されていてもよい。

10

【 0 1 8 7 】

(2) コマンドデバイスがない構成であり、3者(ホスト 1 4 0 0 0、ストレージ装置 1 5 0 0 0 A、ストレージ装置 1 5 0 0 0 B)間の通信は LAN を介して行なわれる。

(3) 機能制御要求処理部 1 4 0 0 5 は LAN を介して、制御要求を制御 I / F 処理部 1 5 0 0 3 A に送信する。

(4) 制御要求を受け取った制御 I / F 処理部 1 5 0 0 3 A は LAN を介して、制御要求を制御 I / F 処理部 1 5 0 0 3 B に送信する。

【 0 1 8 8 】

制御 I / F 処理部 1 5 0 0 3 A、1 5 0 0 3 B が受け取った制御要求を処理する点は第 4 の実施の形態と同様であり、第 6 の実施の形態は第 4 の実施の形態と同等の機能 I / F を提供することができる。

20

【 0 1 8 9 】

(7) 第 7 の実施の形態

本実施の形態では機能 I / F の更に別の構成について述べる。図 1 7 を用いて本実施の形態の構成を説明する。

【 0 1 9 0 】

第 7 の実施の形態は第 6 の実施の形態とほぼ同様である。第 6 の実施の形態との違いは以下の点である。

(1) 機能制御要求処理部 1 4 0 0 5 は LAN を介して、制御要求を両方の制御 I / F 処理部 1 5 0 0 3 A、1 5 0 0 3 B に送信する。

30

(2) 制御 I / F 処理部 1 5 0 0 3 A は制御 I / F 処理部 1 5 0 0 3 B に対して、制御要求を送信しない。

【 0 1 9 1 】

制御 I / F 処理部 1 5 0 0 3 A、1 5 0 0 3 B が受け取った制御要求を処理する点は第 6 の実施の形態と同様であり、第 7 の実施の形態は第 6 の実施の形態と同等の機能 I / F を提供することができる。

【 0 1 9 2 】

(8) 第 8 の実施の形態

本実施の形態では、ストレージ装置内の論理ボリュームに対してセキュリティ機能 (L D E V セキュリティ機能) を適用する場合の例を説明する。

40

【 0 1 9 3 】

図 1 8 は L D E V セキュリティ機能の一実施の形態を示したものである。本実施の形態の構成は第 4 の実施の形態の図 1 4 とほぼ同一である。図 1 4 と異なる点は論理ボリュームセキュリティ情報 1 5 0 0 5 A、1 5 0 0 5 B が追加された点である。論理ボリュームセキュリティ情報 1 5 0 0 5 A、1 5 0 0 5 B は、ホスト 1 4 0 0 0 からストレージ装置 1 5 0 0 0 A、1 5 0 0 0 B 内の論理ボリューム 1 5 0 0 1 A、1 5 0 0 1 B に対するアクセス制御を行うために用いられる。アクセス制御の例としては、論理ボリューム 1 5 0 0 1 A、1 5 0 0 1 B 内のデータの改ざんを抑制するために論理ボリューム 1 5 0 0 1 A、1 5 0 0 1 B に対するライトアクセスを一切禁止する制御がある。また、別の例として

50

は、法令等により一定期間の保存を義務付けられたデータに対して、所定の期間ライトを禁止する機能がある。さらに、別の例としては、機密情報の保護の観点等から特定のホストからのリード/ライトアクセスを禁止する機能がある。

【0194】

図18のように2台のストレージ構成15000A, 15000Bを用いて同期リモートコピーにより冗長化を図った構成においてもLDEVセキュリティ機能を適用したい場合が考えられる。この場合においても第4の実施の形態で説明した機能I/Fを用いてLDEVセキュリティ機能を制御することができる。具体的には、機能処理部15004において、対象ボリュームに対するセキュリティ情報を格納する論理ボリュームセキュリティ情報15005A、15005Bに、LDEVセキュリティに関するパラメータを設定したり、参照したりすればよい。

10

【0195】

(9)第9の実施の形態

本実施の形態では、ストレージ装置内の論理ボリュームにローカルコピー機能を適用した場合の例を説明する。

【0196】

ローカルコピー機能とは、ユーザーから指定されたボリュームの複製を、コピー元ボリュームと同じストレージ装置内において作成する機能である。本機能を用いて作成されたボリュームの複製はデータマイニングやテープバックアップのためにホストがアクセスを行ったり、あるいはバックアップデータとして長時間保存される。ローカルコピー機能は複製を作成したいボリュームと複製先ボリュームをペア(コピーペア)として指定し、そのペアに対してユーザーが操作を行うことでユーザーは複製を作成する。以後の説明では複製対象ボリュームを正ボリュームと呼び、複製先ボリュームを副ボリュームと呼ぶことがある。本実施の形態ではこのローカルコピー機能についてもアクティブ側のストレージとスタンバイ側のストレージで連携することで可用性を向上させる。

20

【0197】

図19はローカルコピー機能の一実施の形態を示したものである。図19においては、ホスト14000はストレージ装置15000Aとストレージ装置15000Bに接続されている。また、ストレージ装置15000Aは外部ストレージ装置16000Aと接続され、ストレージ装置15000Bは外部ストレージ装置16000Bと接続されている。また、ローカルコピー機能及び差分ビットマップ(正ボリューム15006A, 15006Bと副ボリューム15007A, 15007Bの間の差分の有無を示す情報)がストレージ装置15000Aとストレージ装置15000Bにて実行及び管理される。

30

【0198】

本実施の形態は正ボリューム15006A, 15006Bがストレージ装置15000A, 15000B内にあり、副ボリューム15007A, 15007Bが外部ストレージ装置16000A, 16000B内にある構成例を示している。正ボリューム15006Aと副ボリューム15007Aはペアであり、副ボリューム15007Aの実体は外部ボリューム16001A内にある。同様に、正ボリューム15006Bと副ボリューム15007Bはペアであり、副ボリューム15007Bの実体は外部ボリューム16001B内にある。

40

【0199】

< Duplex 状態における動作 >

Duplex 状態とはペア状態のひとつで正ボリューム15006A, 15006Bから副ボリューム15007A, 15007Bへ後述するバックグラウンドコピーが行われている状態である。

【0200】

以下ではDuplex 状態におけるリード/ライト処理について述べる。なお、以下のリード/ライト処理の説明は、リード/ライト処理の対象ボリューム(正ボリューム15006A, 15006B)のアクティブ側がストレージ装置15000Aであるという前

50

提である。

【0201】

まずリード処理について説明する。アプリケーションプログラム14002からリード要求を受け付けたオペレーティングシステム14001はパス管理機能により、(リード対象の正ボリュームに関して)アクティブ側のストレージがストレージ装置15000Aとストレージ装置15000Bのどちらかを判断し、アクティブ側のストレージ装置15000Aにリード要求を発行する。リード要求を受信したストレージ装置15000Aはリード対象データをホスト14000に送信する。アプリケーション14002はオペレーティングシステム14001を介してリード対象データを受信する。以上によりリード処理は完了する。

10

【0202】

次にライト処理について説明する。アプリケーションプログラム14002からライトリクエストを受け付けたオペレーティングシステム14001はパス管理機能により、(リード対象の正ボリュームに関して)アクティブ側のストレージ装置がストレージ装置15000Aとストレージ装置15000Bのどちらかを判断し、アクティブ側のストレージ装置15000Aにライトリクエストを発行する。ライトリクエストを受信したストレージ装置15000Aは、ライトデータを受信し、図示しないキャッシュメモリにライトデータを格納すると共にライトデータに対応する差分ビットマップのビットを1(オン)に設定する。

20

【0203】

その後、当該ライトデータはリモートコピー機能により、ストレージ装置15000A内のキャッシュメモリからストレージ装置15000B内の正ボリューム15006Bにコピー(同期リモートコピー)する。なお、同期リモートコピーの方法はこれまで説明した通りである。同期リモートコピーによりストレージ装置15000Aからライトデータを受信したストレージ装置15000Bは、図示しないキャッシュメモリにライトデータを格納すると共にライトデータに対応する差分ビットマップのビットを1(オン)に設定する。その後、ストレージ装置15000Bはストレージ装置15000Aに対してライト完了報告に送信し、ライト完了報告を受信したストレージ装置15000Aはホスト14000に対してライト完了報告を送信する。

30

【0204】

なお、ストレージ装置15000Aの正ボリューム15006A、ストレージ装置15000Bの正ボリューム15006Bにライトされたデータは、正ボリューム15006A, 15006Bへのライトとは非同期に副ボリューム15007A, 15007Bへコピーされる(以後、本処理をバックグラウンドコピー処理と呼ぶ)。バックグラウンドコピーは、差分ビットマップ定期的に監視し、差分あり(すなわちビットがオン)と記録された領域のデータを正ボリューム15006A, 15006Bから副ボリューム15007A, 15007Bへコピーし、コピーが終了したらビットをクリア(オフ又は0に)することにより行なわれる。

40

【0205】

一方、スタンバイ側のストレージ装置15000Bも同期リモートコピーによってライトデータが到着した時点を契機として同様の処理を行う。

【0206】

なお、本発明は上記例以外の構成、たとえば正ボリューム15006A, 15006Bは外部ストレージ装置16000A内にあってもよいし、ストレージ装置15000A, 15000B内にあってもよい。副ボリューム15007A, 15007Bもまた、外部ストレージ装置16000A内にあってもよいし、ストレージ装置15000A, 15000B内にあってもよい。

【0207】

何らかの障害が発生し、アクティブ側の正ボリューム15006Aに対するI/Oリクエストが処理できなくなった場合には、すでに説明した通り、オペレーティングシステム

50

14001は、I/Oリクエストの対象を正ボリューム15006Bに切り替えてアクセスを継続する。この場合でも、ストレージ装置15000B内にはローカルコピー機能のペアが存在するため、副ボリューム15007Bを用いて先に述べたバックアップ等の処理を行なうことができる。

【0208】

<ペアSplitとSplit状態の動作>

Split状態とはペア状態のひとつで、副ボリュームのイメージが確定した状態のことを指す。この状態では、正ボリュームと副ボリュームの内容が一致しておらず、正ボリュームと副ボリュームの間の差分が差分ビットマップで管理されている。また、この状態においては、副ボリュームが静止した状態になるため、ユーザーは先に述べたバックアップ等の処理を行なうことができる。

10

【0209】

ホスト14000はローカルコピーのDuplex状態のペアをSplit状態にする場合、これまで説明してきたバックグラウンドコピーの動作を停止させる(これをペアSplitと呼ぶ)。ペアSplitは第4～第7の実施の形態で説明した機能I/Fを介して実施する。

【0210】

(1)ホスト14000は機能I/Fを介してストレージ装置15000A,15000Bにローカルコピーの停止命令を出す。通常、ホスト側ではこの停止命令直前にI/Oリクエストの発行を停止する。

20

(2)アクティブ側及びスタンバイ側のストレージ装置15000A,15000Bはそれぞれ差分ビットマップ上でオンとなった領域のバックグラウンドコピーを完了させる。ホスト14000は両ストレージ装置15000A,15000Bにおけるバックグラウンドコピーが完了ことを認識するメッセージをアクティブ側のストレージ装置15000A、もしくは両ストレージ装置15000A,15000Bから受領する。

(3)ホスト14000は当該メッセージを受領した後、I/O発行を再開する。

【0211】

(2)までの処理により、アクティブ側及びスタンバイ側のストレージ装置15000A,15000B内のペアはSplit状態になったことが確定する。この時点で両ストレージ装置15000A,15000B内のペア状態はSplit状態となっている。なお、Split中に行われた正ボリュームまたは副ボリュームへ行われたライトリクエストのライト位置は、後ほど説明するペア再同期のために差分ビットマップに記録される。

30

【0212】

その後のI/Oリクエストの処理はDuplex状態とほぼ同様である。Duplex状態との違いは、バックグラウンドコピー処理が動作しない点である。

【0213】

<ペア作成>

正ボリュームと副ボリュームがペア関係にない状態をSimplex状態と呼ぶ。Simplex状態からDuplex状態に遷移させるための処理をペア作成と呼ぶ。ペア状態がSimplex状態からDuplex状態に遷移している過渡状態をInitial Copying状態と呼ぶ。

40

【0214】

ペア作成の指示は、第4～第7の実施の形態で説明した機能I/Fを介して実施する。

(1)ホスト14000は機能I/Fを介して、ストレージ装置15000Aに対してペア作成指示を出す。この結果アクティブ側及びスタンバイ側の両ストレージ装置15000A,15000Bでペア作成処理が開始される。

(2)両ストレージ装置15000A,15000Bはペア状態をInitial-Copying状態に設定し、差分ビットマップ上のビットを全てオンにし、バックグラウンドコピーを開始する。

(3)バックグラウンドコピーが差分ビットマップの最後まで完了したら、ストレージ装

50

置 15000A, 15000B はペア状態を Duplex 状態に設定する。

【0215】

Initial Copying 状態におけるリード/ライト処理は Duplex 状態におけるリード/ライト処理と同様である。

【0216】

<ペア再同期>

ペア状態を Suspend 状態から Duplex 状態に遷移させる操作をペア再同期と呼ぶ。ペア状態が Suspend 状態から Duplex 状態に遷移している過渡状態を Duplex Pending 状態と呼ぶ。

【0217】

ペア再同期の指示は、第4～第7の実施の形態で説明した機能 I/F を介して実施する。

(1) ホスト 14000 は機能 I/F を介して、ストレージ装置 15000A に対してペア再同期指示を出す。この結果アクティブ側及びスタンバイ側の両ストレージ装置 15000A, 15000B でペア再同期処理が開始される。

(2) 両ストレージ装置 15000A, 15000B はペア状態を Duplex Pending に設定し、バックグラウンドコピーを開始する。

(3) バックグラウンドコピーが差分ビットマップの最後まで完了したら、ストレージ装置 15000A, 15000B はペア状態を Duplex 状態に設定する。

【0218】

Duplex Pending 状態におけるリード/ライト処理は Duplex 状態におけるリード/ライト処理と同様である。

【0219】

(10) 第10の実施の形態

本実施の形態では第9の実施の形態とは異なるローカルコピー機能の実施の形態を説明する。本実施の形態の一構成例を図20に示す。

【0220】

まず、本実施の形態と第9の実施の形態との構成の違いは、外部ストレージ装置 16000B が存在せず、副ボリューム 15007A, 15007B の実体がいずれも外部ストレージ装置 16000A 内の外部ボリューム 16001A となるようにマッピングされている点である。その他の構成は第9の実施の形態と同様である。

【0221】

このように構成することにより、副ボリューム 15007A, 15007B に必要とされる物理的な記憶装置を削減することができる。

【0222】

本実施の形態と第9の実施の形態の処理動作との大きな違いはスタンバイ側のストレージ装置 15000B が外部ボリューム 16001A に対するバックグラウンドコピーを行わず、ストレージ装置 15000A との通信により、ペアに関する制御情報であるペア状態と差分ビットマップ 15010B のみを操作する点である。

【0223】

以下では処理動作を詳細に説明する。

【0224】

< Duplex 状態における動作 >

以下では Duplex 状態におけるリード/ライト処理について述べる。

【0225】

まず、リード処理は第9の実施の形態でのリード処理と同様である。

【0226】

次にライト処理について説明する。アプリケーションプログラム 14002 からライトリクエストを受け付けたオペレーティングシステム 14001 はパス管理機能により、(リード対象の正ボリューム 15006A に関して) アクティブ側のストレージ装置がスト

10

20

30

40

50

レージ装置 15000A 及びストレージ装置 15000B のどちらかを判断し、アクティブ側のストレージ装置 15000A にライトリクエストを発行する。ライトリクエストを受信したストレージ装置 15000A は、ライトデータを受信し、図示しないキャッシュメモリにライトデータを格納すると共にライトデータに対応する差分ビットマップ 15010A のビットを 1 (オン) に設定する。

【0227】

その後、当該ライトデータは同期リモートコピー機能により、ストレージ装置 15000A 内の正ボリューム 15006A からストレージ装置 15000B 内の正ボリューム 15006B にコピーされる。なお、同期リモートコピーの方法はこれまで説明した通りである。同期リモートコピー機能によりストレージ装置 15000A からライトデータを受信したストレージ装置 15000B は、図示しないキャッシュメモリにライトデータを格納すると共にライトデータに対応する差分ビットマップ 15010B のビットを 1 (オン) に設定する。その後、ストレージ装置 15000B はストレージ装置 15000A に対してライト完了報告に送信し、ライト完了報告を受信したストレージ装置 15000A はホスト 14000 に対してライト完了報告を送信する。

10

【0228】

なお、ストレージ装置 15000A の正ボリューム 15006A にライトされたデータは、正ボリューム 15006A へのライトとは非同期に副ボリューム 15007A へバックグラウンドコピーされる。第 9 の実施の形態でのライト処理と異なり、ストレージ装置 15000B の正ボリューム 15006B にライトされたデータはバックグラウンドコピーされない。

20

【0229】

ストレージ装置 15000A におけるバックグラウンドコピーは、差分ビットマップ 15010A を定期的に監視し、差分あり (すなわちビットがオン) と記録された領域のデータを正ボリューム 15006A から副ボリューム 15007A へコピーし、コピーが終了したらビットをクリア (オフ又は 0 に) することにより行なわれる。なお、本実施の形態では、第 9 の実施の形態でのライト処理と異なり、ストレージ装置 15000B においてはバックグラウンドコピーが行なわれない。

【0230】

その後、第 9 の実施の形態でのライト処理と異なり、ストレージ装置 15000A はクリアした差分ビットマップ 15010A 上のビットの位置情報をストレージ装置 15000B に通知する。通知を受信したストレージ装置 15000B は当該ビットに対応するストレージ装置 15000B 内の差分ビットマップ 15010B 上のビット (差分ビット) をクリアする。なお、差分ビットの位置情報の通知はストレージ装置 15000B 内のコマンドデバイスを介して行なわれる。また、本実施の形態における構成では、コマンドデバイスを介して通知を行なったが、ストレージ装置 15000A、15000B 間が LAN で接続された構成である場合は、LAN を介した通信により通知を行なってもよい。以後、ストレージ装置 15000A とストレージ装置 15000B との間における、差分ビットやペア状態等といった機能の制御情報に関する通信はコマンドデバイスや LAN を介して行うものとする。

30

40

【0231】

何らかの障害が発生し、アクティブ側の正ボリューム 15006A に対する I/O リクエストが処理できなくなった場合、オペレーティングシステム 14001 は、第 9 の実施の形態と同様に、I/O リクエストの対象を正ボリューム 15006B に切り替えてアクセスを継続する。

【0232】

< ペア Split と Split 状態の動作 >

ホスト 14000 はローカルコピーの Duplex 状態のペアを Split 状態にする場合、第 9 の実施の形態と同様にペア Split を行なう。なお、ペア Split においては、バックグラウンドコピーの終了処理が行なわれるが、本実施の形態ではストレージ

50

装置 15000B においては、バックグラウンドコピーは動作していないため、実際には終了処理は行なわれない。

【0233】

その後の I/O リクエストの処理は Duplex 状態とほぼ同様である。Duplex 状態との違いは、ストレージ装置 15000B においてバックグラウンドコピー処理が動作しない点である。

【0234】

<ペア作成>

ペア作成の指示は、第 4～第 7 の実施の形態で説明した機能 I/F を介して実施されるのは第 9 の実施の形態と同様である。

10

【0235】

(1) ホスト 14000 は機能 I/F を介して、ストレージ装置 15000A に対してペア作成指示を出す。この結果アクティブ側及びスタンバイ側の両ストレージ装置 15000A, 15000B でペア作成処理が開始される。

(2) 両ストレージ装置 15000A, 15000B はペア状態を Initial-Copying 状態に設定する。ストレージ装置 15000A は差分ビットマップ 15010A 上のビットを全てオンにし、バックグラウンドコピーを開始する。第 9 の実施の形態と異なり、ストレージ装置 15000B は差分ビットマップ 15010B 上のビットを全てオンにするが、バックグラウンドコピーを行なわない。

【0236】

(3) ストレージ装置 15000A はバックグラウンドコピーが完了した領域に対応する差分ビットをクリアする処理とそれに付随する動作(差分ビットの位置情報の通知と差分ビットのクリア)は Duplex 状態における動作と同様である。

(4) 第 9 の実施の形態と異なり、ストレージ装置 15000A は、バックグラウンドコピーが差分ビットマップ 15010A の最後まで完了したら、ペア状態を Duplex 状態に設定し、ペア状態が Duplex 状態に変わったことをストレージ装置 15000B に通知する。通知を受信したストレージ装置 15000B はペア状態を Duplex 状態に設定する。

20

【0237】

Initial Copying 状態におけるリード/ライト処理は Duplex 状態におけるリード/ライト処理と同様である。

30

【0238】

<ペア再同期>

ペア再同期の指示は、第 4～第 7 の実施の形態で説明した機能 I/F を介して実施されるのは第 9 の実施の形態と同様である。

【0239】

(1) ホスト 14000 は機能 I/F を介して、ストレージ装置 15000A に対してペア再同期指示を出す。この結果アクティブ側及びスタンバイ側の両ストレージ装置 15000A, 15000B でペア再同期処理が開始される。

(2) ストレージ装置 15000A はペア状態を Duplex Pending に設定し、バックグラウンドコピーを開始する。第 9 の実施の形態と異なり、ストレージ装置 15000B においては、バックグラウンドコピーは行なわない。

40

【0240】

(3) ストレージ装置 15000A は、バックグラウンドコピーが差分ビットマップ 15010A の最後まで完了したら、ペア状態を Duplex 状態に設定する。ただし、第 9 の実施の形態と異なり、この処理を行なうのはストレージ装置 15000A のみである。その後、ストレージ装置 15000A は、ペア状態が Duplex 状態に変わったことをストレージ装置 15000B に通知する。通知を受信したストレージ装置 15000B はペア状態を Duplex 状態に設定する。

【0241】

50

Duplex Pending状態におけるリード/ライト処理はDuplex状態におけるリード/ライト処理と同様である。

【0242】

(11)第11の実施の形態

AOU (Allocation On Use) 機能の構成について述べる。AOU機能はホストから使用された(ライトされた)領域に関してのみ実記憶領域を割り当てる機能である。

【0243】

AOU機能はデータが実際に格納される実ボリュームの集合体であるプールと、ホストに見せるボリュームである仮想ボリュームから構成される。本実施の形態における仮想ボリュームはライトが行われた部分のみ実データが割り当てられるという意味で仮想的である。ホストに見せているボリュームの全アドレス空間に実データが割り当てられている訳ではない。なお、実ボリュームは外部ストレージ装置内にあってもよいし、仮想ボリュームと同じストレージ装置内にあってもよい。

10

【0244】

図21はAOU機能の一実施の形態を示したものである。図21においては、ホスト14000はストレージ装置15000Aとストレージ装置15000Bに接続されている。また、ストレージ装置15000Aは外部ストレージ装置16000Aと接続され、ストレージ装置15000Bは外部ストレージ装置16000Bと接続されている。

【0245】

本実施の形態は実ボリューム16002Aが外部ストレージ装置16000A, 16000B内にある構成例を示している。仮想ボリューム15008A内のデータはプール16003Aの実ボリューム16002A内のデータと対応付けられる。同様に仮想ボリューム15008B内のデータはプール16003Bの実ボリューム16002B内のデータと対応付けられる。また、仮想ボリューム15008Aと仮想ボリューム15008Bは同期リモートコピー機能により内容が一致するように構成される。同期リモートコピーの方法はこれまで説明したとおりである。

20

【0246】

次に本構成におけるリード/ライト処理について述べる。なお、以下のリード/ライト処理の説明は、リード/ライト処理の対象ボリュームのアクティブ側がストレージ装置15000Aであるという前提である。

30

【0247】

まずリード処理について説明する。アプリケーションプログラム14002からリードリクエストを受け付けたオペレーティングシステム14001はパス管理機能によりアクティブ側のストレージがストレージ装置15000A及びストレージ装置15000Bのどちらかを判断し、アクティブ側のストレージ装置15000Aにリードリクエストを発行する。リードリクエストを受け付けたストレージ装置15000Aは、仮想アドレス実アドレス変換テーブル15009Aを参照し、リードデータにプール16003A内の実領域が割り当てられているか判定する。

【0248】

前述の判定で実領域が割り当てられている場合、ストレージ装置15000Aは、当該実領域からリードデータを読み出してホスト14000に送信する。アプリケーション14002はオペレーティングシステム14001を介してリードデータを受信する。以上によりリード処理は完了する。

40

【0249】

次にライト処理について説明する。アプリケーションプログラム14002からライトリクエストを受け付けたオペレーティングシステム14001はパス管理機能によりアクティブ側のストレージ装置がストレージ装置15000Aとストレージ装置15000Bのどちらかを判断し、アクティブ側のストレージ装置15000Aにライトリクエストを発行する。ライトリクエストを受け取ったストレージ装置15000Aは仮想アドレス実アドレス変換テーブル15009Aを参照し、ライト対象データにプール16003A内

50

の実領域が割り当てられているか判定する（判定200）。

【0250】

前述の判定で実領域が割り当てられている場合、ストレージ装置15000Aは、ホスト14000からライトデータを受信し、当該実領域に対応する図示しないキャッシュメモリ内の領域にライトデータを格納する。そして、同期リモートコピー機能によりライトデータをストレージ装置15000Bにライトリクエストを送信する。ストレージ装置15000Aからライトリクエストを受信したストレージ装置15000Bは、ライトデータにプール16003A内の実領域が割り当てられているか判定する。ここで、仮想ボリューム15008Aの内容と仮想ボリューム15008Bの内容は同期リモートコピー機能により一致しているため、実領域は割り当てられていると判定される。その後、ストレージ装置15000Bは、ストレージ装置15000Aからライトデータを受信し、当該実領域に対応する図示しないキャッシュメモリ内の領域にライトデータを格納し、ストレージ装置15000Aにライト完了報告を行う。

10

【0251】

前述の判定（判定200）で実領域が割り当てられていない場合、ストレージ装置15000Aは、仮想アドレス実アドレス変換テーブル15009Aにライトデータのアドレスを登録し、実領域を確保する。その後、ストレージ装置15000Aは、ホスト14000からライトデータを受信し、当該実領域に対応する図示しないキャッシュメモリ内の領域にライトデータを格納する。そして、同期リモートコピー機能によりライトデータをストレージ装置15000Bにライトリクエストを送信する。

20

【0252】

ストレージ装置15000Aからライトリクエストを受信したストレージ装置15000Bは、ライトデータにプール16003B内の実領域が割り当てられているか判定する。ここで、仮想ボリューム15008Aの内容と仮想ボリューム15008Bの内容は同期リモートコピー機能により一致しているため、実領域は割り当てられていないと判定される。その後、ストレージ装置15000Bは、仮想アドレス実アドレス変換テーブル15009Bにライトデータのアドレスを登録し、実領域を確保する。そして、ストレージ装置15000Bは、ストレージ装置15000Bからライトデータを受信し、当該実領域に対応する図示しないキャッシュメモリ内の領域にライトデータを格納した後、ストレージ装置15000Aにライト完了報告を行う。ライト完了報告を受信したストレージ装置15000Aはホスト14000にライト完了報告を行う。ホスト14000がライト完了報告を受信し、ライト処理は完了する。

30

【0253】

なお、キャッシュメモリに格納されたデータはキャッシュメモリへの格納とは非同期に実ボリューム16002A、16002Bへライトされる。

【0254】

何らかの障害により、アプリケーション14002がストレージ装置15000A内の仮想ボリューム15008A経由でのリード/ライト処理が不可能になった場合、オペレーティングシステム14001の提供するパス管理機能は障害を検出し、リード/ライト処理のアクセス経路をストレージ装置15000B内の仮想ボリューム15008B経由に切り替える。仮想ボリューム15008Aの内容と仮想ボリューム15008Bの内容は同期リモート機能により一致しているため、アクセス経路が切り替わっても、継続して正常にリード/ライト処理を行うことができる。

40

【0255】

(12)第12の実施の形態

本実施の形態ではAOU機能の第11の実施の形態とは異なる実施の形態について述べる。本実施の形態の一構成例を図22に示す。

【0256】

まず、本実施の形態と第11の実施の形態との構成の違いは、外部ストレージ装置16000Bが存在せず、仮想ボリューム15008A、15008Bの実領域がいずれも外

50

部ストレージ装置 16000A 内のプール 16003A 内の領域に割り当てられている点である。その他の構成は第 11 の実施の形態と同様である。

【0257】

なお、本実施の形態はストレージ装置 15000A 及びストレージ装置 15000B が共通のプールとして、共通の外部ストレージ装置 16000A 内の実ボリューム 16002A を用いるため、第 11 の実施の形態と異なり、実ボリューム 16002A が外部ストレージ装置 16000A 内にある構成に限定される。

【0258】

このように構成することにより、プールに必要とされる物理的な記憶装置（HDD など）の容量を削減することができる。

10

【0259】

本実施の形態と第 11 の実施の形態の処理動作の大きな違いは、スタンバイ側のストレージ装置 15000B がキャッシュメモリから外部ストレージ装置 16000A 内の実ボリューム 16002A に対してライトを行わない点と、ストレージ装置 15000A が仮想アドレス実アドレス変換テーブル 15009A への更新をストレージ装置 15000B に通知し、通知を受けたストレージ装置 15000B が仮想アドレス実アドレス変換テーブル 15009B を更新する点である。

【0260】

以下では処理動作を詳細に説明する。

【0261】

まずリード処理は第 11 の実施の形態におけるリード処理と同様である。

20

【0262】

次にライト処理について説明する。アプリケーションプログラム 14002 からライトリクエストを受け付けたオペレーティングシステム 14001 はバス管理機能によりアクティブ側のストレージがストレージ装置 15000A 及びストレージ装置 15000B のどちらかを判断し、アクティブ側のストレージ装置 15000A にライトリクエストを発行する。ライトリクエストを受け取ったストレージ装置 15000A は仮想アドレス実アドレス変換テーブル 15009A を参照し、ライトデータにプール 16003A 内の実領域が割り当てられているか判定する（判定 300）。

【0263】

前述の判定で実領域が割り当てられている場合、ストレージ装置 15000A は、ホスト 14000 からライトデータを受信し、当該実領域に対応するキャッシュメモリ内の領域にライトデータを格納する。そして、同期リモートコピー機能によりライトデータをストレージ装置 15000B にライトリクエストを送信する。次に、本実施の形態では、第 11 の実施の形態と異なり、ストレージ装置 15000A からライトリクエストを受信したストレージ装置 15000B は、即座にストレージ装置 15000A からライトデータを受信し、キャッシュメモリに当該データを格納した後、ストレージ装置 15000A にライト完了報告を行なう。ストレージ装置 15000B からライト完了報告を受信したストレージ装置 15000A はホスト 14000 に対してライト完了報告を送信する。

30

【0264】

前述の判定（判定 300）で実領域が割り当てられていない場合、ストレージ装置 15000A は、仮想アドレス実アドレス変換テーブル 15009A にライトデータのアドレスを登録し、実領域を確保する。その後、ストレージ装置 15000A は、ホスト 14000 からライトデータを受信し、当該実領域に対応するキャッシュメモリ内の領域にライトデータを格納する。そして、ストレージ装置 15000A は、同期リモートコピー機能によりライトデータをストレージ装置 15000B にライトリクエストを送信する。

40

【0265】

次に、本実施の形態では、第 11 の実施の形態と異なり、ストレージ装置 15000A からライトリクエストを受信したストレージ装置 15000B は、即座にストレージ装置 15000A からライト対象データを受信し、キャッシュメモリに当該データを格納した

50

後、ストレージ装置15000Aにライト完了報告を行なう。ストレージ装置15000Aは、ストレージ装置15000Bからライト完了報告を受信した後、仮想アドレス実アドレス変換テーブル15009Aへの変更内容をストレージ装置15000Bに送信する。

【0266】

仮想アドレス実アドレス変換テーブル15009Aへの変更内容を受信したストレージ装置15000Bは、同様の変更を仮想アドレス実アドレス変換テーブル15009Bに対して行なう。これによりストレージ装置15000B内の仮想ボリューム15008B内の当該ライト領域の実領域が共通の外部ストレージ装置16000Aの実ボリューム16002A内の(ストレージ装置15000Aにより割り当てられた)実領域にマッピングされることになる。ストレージ装置15000Bは仮想アドレス実アドレス変換テーブル15009Bを更新した旨をストレージ装置15000Aに通知する。その後、通知を受信したストレージ装置15000Aはホスト14000に対してライト完了報告を行なう。なお、ストレージ装置15000Aは(1)同期リモートコピーのデータ送信と、(2)仮想アドレス実アドレス変換テーブル15009Aへの変更内容の送信を同時に行い、(1)及び(2)の処理の完了報告を受信した後ホスト14000に対してライト完了報告を行なってもよい。その後、ホスト14000がライト完了報告を受信し、ライト処理は完了する。

10

【0267】

なお、ストレージ装置15000A内のキャッシュメモリに格納されたデータはキャッシュメモリへの格納とは非同期に、ストレージ装置15000Aにより実ボリューム16002Aへライト(デステージ)される。デステージが完了した後、ストレージ装置15000Aはストレージ装置15000Bにその旨を通知する。通知を受けたストレージ装置15000Bは当該ライトに対応するキャッシュメモリの領域を破棄する。なお、破棄せずに当該ライトに対応するキャッシュメモリの領域の属性をクリーン(キャッシュメモリの内容と記憶装置(HDDなど)内のデータの内容が一致している状態)としてもよい。

20

【0268】

何らかの障害により、アプリケーション14002がストレージ装置15000A内の仮想ボリューム15008A経由でのリード/ライト処理が不可能になった場合、オペレーティングシステム14001の提供するパス管理機能は障害を検出し、リード/ライト処理のアクセス経路をストレージ装置15000B内の仮想ボリューム15008B経由に切り替える。仮想ボリューム15008Aの内容と仮想ボリューム15008Bの内容は同期リモート機能により一致しているため、アクセス経路が切り替わっても、継続して正常にリード/ライト処理を行うことができる。

30

【0269】

(13)第13の実施の形態

本実施の形態ではストレージ装置内のボリュームに論理スナップショット機能を適用した場合の例を説明する。

【0270】

論理スナップショット機能とは、ローカルレプリケーションと類似した機能であり、ユーザーの指示時点の複製データをホストに提供する機能である。しかし、複製データを有する副ボリュームは、プールに属する実ボリュームの領域に保存された複製作成指示以後のライトデータと、正ボリュームのデータを用いて提供される仮想的な存在である。仮想的な副ボリュームの実体は実ボリュームの集合体であるプールに保持される。正ボリュームと副ボリュームの関係をスナップショットペアもしくは単にペアと呼ぶこともある。論理スナップショット機能においては、静止化ポイントにおける正ボリュームの内容と同一の内容の論理ボリュームが実際に作成される訳ではないという意味で、副ボリュームは仮想的である。論理スナップショット機能は先に説明したローカルコピー機能とは異なり、正ボリュームのサイズと同一のサイズの副ボリュームが不要である。これにより、副ボリ

40

50

ュームの内容を保持するために必要な記憶装置（HDDなど）の容量を削減することが可能である。

【0271】

本実施の形態ではこの論理スナップショット機能についてもアクティブ側のストレージとスタンバイ側のストレージで連携することで可用性を向上させることができる。

【0272】

図23はスナップショット機能の一実施の形態を示したものである。図23においては、ホスト14000はストレージ装置15000Aとストレージ装置15000Bに接続されている。また、ストレージ装置15000Aは外部ストレージ装置16000Aと接続され、ストレージ装置15000Bは外部ストレージ装置16000Bと接続されている。また、スナップショット機能及び差分ビットマップ（静止化ポイントにおける正ボリューム15006A, 15006Bと現時点における正ボリューム15006A, 15006Bの間の差分の有無を示す情報）15010A, 15010Bと仮想アドレス実アドレス変換テーブル（仮想的な副ボリューム15007A, 15007Bの実体の位置を管理するテーブル）15009A, 15009Bがストレージ装置15000Aとストレージ装置15000Bにて実行及び管理される。更に、ストレージ装置15000A内の正ボリューム15006Aとストレージ装置15000B内の正ボリューム15006Bはリモートコピーのペアとなるように構成される。

10

【0273】

本実施の形態は正ボリューム15006A, 15006Bがストレージ装置15000A, 15000B内にあり、プール16003A, 16003Bが外部ストレージ装置16000A, 16000B内にある構成例を示している。なお、プール16003A, 16003Bはストレージ装置15000A, 15000B内にあってもよい。

20

【0274】

<論理スナップショット作成指示>

ホスト14000を利用するユーザーが論理スナップショット作成を指示すると、前記実施の形態に記載の方式によって、アクティブ側のストレージ装置15000Aとスタンバイ側のストレージ装置15000Bに作成指示を発行する。作成指示を受信したストレージ装置15000A, 15000Bは当該指示を受けて、仮想的な副ボリューム15007A, 15007Bを準備し、この副ボリューム15007A, 15007Bに全て0（差分なしの意味）の差分ビットマップ15010A, 15010Bと仮想アドレス実アドレス変換テーブル15009A, 15009Bとを割り当てる。

30

【0275】

<正ボリュームに対するリード処理>

これまで述べた実施の形態と同じである。

【0276】

<正ボリュームに対するライト処理>

アプリケーションプログラム14002からライトリクエストを受け付けたオペレーティングシステム14001はパス管理機能により、（ライト対象の正ボリュームに関して）アクティブ側のストレージがストレージ装置15000A及びストレージ装置15000Bのどちらかを判断し、アクティブ側のストレージ装置15000Aにライトリクエストを発行する。ライトリクエストを受信したストレージ装置15000Aは、ライト対象アドレスの差分ビットマップ15010Aをチェックする。結果、1であればキャッシュメモリに正ボリューム15006Aのライトデータとして、格納する。一方、0の場合は正ボリューム15006Aの更新前のデータを副ボリューム15007A用のデータとして用いるための以下に示すCopy On Write処理を行う。

40

【0277】

（Step1）プール16003Aに属する実ボリューム16002Aの記憶領域を確保する。

（Step2）正ボリューム15006Aから当該記憶領域へ更新前データをキャッシュ

50

メモリを利用しつつコピーする。

(Step 3) 退避する更新前データの保存先を管理するプール管理情報を更新し、当該データがプール16003A内の実ボリューム16002Aのどの領域に保存されたかわかるようにする。

(Step 4) 受信したライトデータをキャッシュメモリに正ボリューム15006Aの当該アドレス宛のデータとして保存し、ライト完了応答を返す。

【0278】

これと並行して、当該ライトデータはリモートコピー機能により、ストレージ装置15000A内の正ボリューム15006Aからストレージ装置15000B内の正ボリューム15006Bにコピーされ、同様の処理がなされる。そのため、各ストレージ装置15000A, 15000Bはそれぞれで仮想アドレス実アドレス変換テーブル15009A, 15009Bや差分ビットマップ15010A, 15010Bの管理を行う。

10

【0279】

<副ボリュームに対するリード処理>

アプリケーションプログラム14002からライトリクエストを受け付けたオペレーティングシステム14001はバス管理機能により、(リード対象の副ボリュームに関して)アクティブ側のストレージがストレージ装置15000A及びストレージ装置15000Bのどちらかを判断し、アクティブ側のストレージ装置15000Aにリードリクエストを発行する。リードリクエストを受信したストレージ装置15000Aは、正ボリューム15006Aに対して記録していた差分ビットマップ15010Aをチェックする。結果、リード対象アドレスのビットが0であれば正ボリューム15006Aの同じアドレスに保存されたデータをホスト14000へ返し、オペレーティングシステム14001は当該データをアプリケーション14002へ返す。一方、リード対象アドレスのビットが1の場合は仮想アドレス実アドレス変換テーブル15009Aを参照して、正ボリューム15006Aのリード対象アドレスに関する更新前のデータの場所を決定し、プール16003Aに属する実ボリューム16002Aからデータをホスト14000(アプリケーション14002)へ返す。

20

【0280】

<副ボリュームに対するライト処理>

アプリケーションプログラム14002からライトリクエストを受け付けたオペレーティングシステム14001はバス管理機能により、(ライト対象の副ボリュームに関して)アクティブ側のストレージがストレージ装置15000Aとストレージ装置15000Bのどちらかを判断し、アクティブ側のストレージ装置15000Aにライトリクエストを発行する。ライトリクエストを受信したストレージ装置15000Aは、正ボリューム15006Aに割り当てられたライト対象アドレスの差分ビットマップ15010Aをチェックする。結果、1であれば仮想アドレス実アドレス変換テーブル15009Aを参照することで、正ボリューム15006Aの当該アドレスの更新前データが保存されたプール16003A内の実ボリューム16002Aの記憶領域を探し、当該領域へライトデータを保存する。一方、0の場合は以下の処理を行う。

30

【0281】

- (A) プール16003Aに属する実ボリューム16002Aの領域を確保する。
- (B) 確保した領域にライトデータを保存し、仮想アドレス実アドレス変換テーブル15009Aを更新することで当該ライトデータがプール16003A内の実ボリューム16002Aのどの領域に保存されたかわかるようにする。
- (C) 差分ビットマップ15010Aの当該アドレスに対応するビットを1に更新する。

40

【0282】

これと並行して、当該ライトデータはリモートコピー機能により、ストレージ装置15000A内の正ボリューム15006Aからストレージ装置15000B内の正ボリューム15006Bにコピーされ、同様の処理がなされる。そのため、各ストレージ装置15000A, 15000Bはそれぞれで仮想アドレス実アドレス変換テーブル15009A

50

、15009Bや差分ビットマップ15010A、15010Bの管理を行う。

【0283】

<Copy After Write処理>

ストレージ装置15000A、15000Bは、正ボリューム15006A、15006Bに対するライト時に実行されるCopy On Write処理の代わりとして以下に示すCopy After Write処理を実行してもよい。

【0284】

(Step1)受信したライトデータをキャッシュメモリに正ボリューム15006A、15006Bの当該アドレス宛のデータとして保存し、ライト完了応答を返す。ただし、当該ライトデータのデステージングは抑制する。

(Step2)プール16003A、16003Bに属する実ボリューム16002A、16002Bの記憶領域を確保する。

(Step3)正ボリューム15006A、15006Bから当該記憶領域に更新前データをキャッシュメモリを利用しつつコピーする。

【0285】

(Step4)退避した更新前データの保存先を管理するプール管理情報を更新し、当該データがプール16003A、16003B内の実ボリューム16002A、16002Bのどの領域に保存されたかわかるようにする。

(Step5)抑制していたライトデータのデステージを許可する。

【0286】

<障害>

何らかの障害が発生し、アクティブ側の正ボリューム15006A及び副ボリューム15007Aに対するI/Oリクエストが処理できなくなった場合には、すでに説明した通り、オペレーティングシステム14001は、I/Oリクエストの対象を正ボリューム15006B及び副ボリューム15007Bに切り替えてアクセスを継続することができる。なお、前述の通り、好ましくはスナップショット機能の正ボリューム15006A、15006B及び副ボリューム15007A、15007Bは同一のストレージ装置15000A、15000Bに対してライトリクエストを発行したいため、正ボリューム15006A、15006Bに対する切り替えが必要な場合は副ボリューム15007A、15007Bも同時に行い、逆に副ボリューム15007A、15007Bに対する切り替えが必要な場合は正ボリューム15006A、15006Bに対しても切り替えを行う連携を行うことがこのましい。

【0287】

(14)第14の実施の形態

本実施の形態では、第13の実施の形態とは異なる論理スナップショット機能の実施の形態を説明する。図24に本実施の形態の一構成例を示す。

【0288】

まず、本実施の形態と第13の実施の形態との構成の違いは、外部ストレージ装置16000Bが存在せず、仮想的な副ボリューム15007A、15007Bの実領域がいずれも外部ストレージ装置16000A内のプール16003A内の領域に割り当てられている点である。その他の構成は第13の実施の形態と同様である。

【0289】

なお、本実施の形態はストレージ装置15000Aとストレージ装置15000Bが共通のプール16003Aとして、共通の外部ストレージ装置16000A内の実ボリューム16002Aを用いるため、第13の実施の形態と異なり、実ボリューム16002Aが外部ストレージ装置16000A内にある構成に限定される。

【0290】

このように構成することにより、プール16003Aに必要とされる物理的な記憶装置(HDDなど)の容量を削減することができる。

【0291】

10

20

30

40

50

本実施の形態と第13の実施の形態の処理動作の大きな違いは以下の通りである。

(A) 通常時はスタンバイ側のストレージ装置15000Bがキャッシュメモリから外部ストレージ装置16000A内の実ボリューム16002Aに対してライトを行なわない代わりに、アクティブ側のストレージ装置15000Aが正ボリューム15006A、副ボリューム15007A、プール16003A内の実ボリューム16002Aに対応するデータをデステージングする時にスタンバイ側のストレージ装置15000Bにこれを伝え、スタンバイ側のストレージ装置15000Bではこれによってキャッシュメモリ上のデータを破棄する。

【0292】

(B) ストレージ装置15000Aが仮想アドレス実アドレス変換テーブル15009Aへの更新をストレージ装置15000Bに通知し、通知を受けたストレージ装置15000Bが仮想アドレス実アドレス変換テーブル15009Bを更新する。

【0293】

また、(A)の処理に変えて、副ボリューム15007A、15007B又はプール16003A内の実ボリューム16002Aに対応するデータについてはキャッシングを無効化してもよい。この場合、前述のCopy On Write処理による更新前データ退避では正ボリューム15006A、15006Bに対するライト完了までにプール16003A内の実ボリューム16002Aへの退避データ保存が含まれるため、性能が悪化するが、Copy After Write方式ではそれが無いため、好適である。

【0294】

以上、本発明の幾つかの実施態様を説明したが、これらの実施の形態は本発明の説明のための例示にすぎず、本発明の範囲をそれらの実施の形態にのみ限定する趣旨ではない。本発明は、その要旨を逸脱することなく、その他の様々な態様でも実施することができる。例えば、HDD1030やキャッシュメモリ1020の代わりに不揮発性メモリを使用することができる。不揮発性メモリとしては、例えば、フラッシュメモリ(具体的には、例えば、NAND型フラッシュメモリ)、MRAM(Magnetoresistive Random Access Memory)、PRAM(Parameter Random Access Memory)など、種々の不揮発性メモリを採用することができる。

【0295】

(15) 第15の実施の形態

図1との対応部分に同一符号を付して示す図28は、仮想化機能を持つネットワークスイッチ(仮想化スイッチ)28000L、28000Rを適用した場合の実施の形態を示した概要図である。

【0296】

<1. 本実施の形態のハードウェア構成>

仮想化スイッチ28000L、28000Rは、複数のネットワークポートを持ち、ネットワークポート制御用のプロセッサが各ポートの転送制御や障害の検知や後術する仮想化を行う。なお、本概要図には図示されていないが、図1について上述した第1の実施の形態と同様に仮想化スイッチ28000L、28000Rには管理ホストが接続され、この管理ホストを介して仮想化スイッチ28000L、28000Rに対する設定を行ったり、仮想化スイッチ28000L、28000R間の設定コピーを行うことができる。なお、その他コンポーネントについては第1~第14の実施の形態と同じであるため、説明を省略する。

【0297】

<2. 仮想化スイッチを用いた本実施の形態の特徴>

仮想化スイッチ28000L、28000Rが提供する仮想化は第1~第14の実施の形態と異なる以下の特徴を持つ。

【0298】

(特徴1) 仮想的なWWN(又はポート名前)を提供できる。ファイバーチャネルスイッチのポートはFポート又はEポートと呼ばれ、通常のホストやストレージが持つNポー

10

20

30

40

50

ト（通信の始点又は終点となることを意味する）とは異なる属性を持つ。そのため、仮想化スイッチ 28000L, 28000R において仮想化を行う場合、内部で実際に接続されていない仮想的な WWN を仮想化スイッチ 28000L, 28000R の両方で作成・提供すれば、ホスト 1100 上のソフトウェアが明示的に I/O パスを切り替える必要がなくなる。なお、より具体的にはファイバーチャネルの通信は前述のポートネームによって行われるが、これはファイバーチャネルスイッチが割り当てる識別子であり、識別子内部にはルーティング用に用いられるスイッチを識別するための情報が含まれている。そのため、両仮想化スイッチ 28000L, 28000R は、ホスト 1100 に対して、仮想的な WWN を持つ N ポートが仮想的なスイッチを経由して両仮想化スイッチ 28000L, 28000R に接続されているかのごとく模擬できるように、ポートネームを割り当てて、ルーティングを行う。

10

【0299】

（特徴2）スイッチでキャッシングを行わない。ファイバーチャネルスイッチは、通常、制御ヘッダだけ参照して転送先を決定し、データのバッファリングを行わない、いわゆるカットスルー方式で転送制御を行うため、仮想化機能を提供する場合もキャッシングを行わないことが多い。なお、キャッシングを行う場合は、本特徴が関係する処理についてはこれまで説明してきた実施の形態と同様の処理によって実現される。また、キャッシングを行わない場合の仮想化スイッチ 28000L, 28000R のリード/ライト処理は、I/O リクエストを受け付けてからストレージ装置 1500L へのリクエスト処理転送及び処理完了を待ってホスト 1100 に処理完了が返る、ライトスルー型の制御に類似したものと考えることができる。

20

【0300】

（特徴3）本実施の形態での高可用性は両方の仮想化スイッチ 28000L, 28000R に同じ仮想化の設定を行うだけでよい。これは仮想化スイッチ 28000L, 28000R にてキャッシングを行っていないために可能となる。なお、仮想化スイッチ 28000L, 28000R がリモートコピーやローカルコピーを行う場合、差分ビットマップ等スイッチ内部に存在する情報がある場合は、これまでの実施の形態と同じく、正系と副系の両方で内部情報を保持する必要がある。

【0301】

なお、ここまで仮想化スイッチ 28000L, 28000R はファイバーチャネルスイッチであるものとして説明を行っているが、仮想化スイッチ 28000L, 28000R が Ethernet（登録商標）や、iSCSI 又は TCP/IP を用いたものでもよい。この場合、WWN が MAC アドレスで、ポートネームが IP アドレスに対応するものとするができるが、Ethernet（登録商標）や TCP/IP の場合は仮想的なスイッチを提供せずに、直接仮想的なポートとそれに割り当てた IP アドレスを外部へ提供し、当該 IP アドレスに対するルーティングを行えばよい。

30

【0302】

（16）第16の実施の形態

次に、第16の実施の形態について説明する。本実施の形態は第11の実施の形態及び第12の実施の形態にて説明した AOU 機能を高可用性構成の仮想化ストレージ装置が提供することに関する発明である。なお、以下において説明していない機能等については、第1～第15の実施の形態による情報システムと同様の構成を有する。

40

【0303】

前述の通り、AOU 機能とは仮想化ストレージ装置が当該機能によってホスト 1100 へ提供するボリューム（以後 AOU ボリュームと呼ぶ）について、使用開始時から AOU ボリュームの全アドレスに対して HDD の記憶領域を割り当ててのではなく、ホスト 1100 がライトを行ったアドレスに対して HDD の記憶領域（より正確には HDD から構成されるボリューム（プールボリュームと呼ぶ）の記憶領域の一部又は全て）を割り当てる機能である。AOU 機能は HDD を有効利用することができる。なぜならば、ホスト 1100 で動作するファイルシステムの一部の種類ではアクセス継続中に動的なデータ容量拡

50

張ができないため、ホスト 1100 の管理者は将来利用する可能性のあるデータ容量も含めてボリュームの容量設定を行う。そのため、従来技術ではボリュームの容量設定を行った時点では使用せず、また将来も確実に使うとは限らない HDD を搭載していなければならないからである。

【0304】

なお、HDD 容量の有効利用の観点から考えた場合、AOU ボリュームにライトが発生する前の領域に対してプールボリュームの領域が未割り当ての状態であったほうがよいが、他の目的（高性能化等）がある場合はこれに限られない。

【0305】

< 1 . 本実施の形態の概要 >

図 1 との対応部分に同一符号を付した図 29 に本実施の形態の概要を示す。本実施の形態による情報システムは、2 台の仮想化ストレージ装置 1000L, 1000R が共通のストレージ装置 1500L と接続されている。そして、高可用化された 2 台の仮想化ストレージ装置 1000L, 1000R が AOU 機能を有することで、情報システムのサービス停止時間を短縮している。なお、特に記載しない限り、ストレージ装置 1500L は両方の仮想化ストレージ装置 1000L, 1000R からアクセス可能な状態、言い換えれば共有された状態にあるものとするが、共有されていないストレージ装置が存在して当該ストレージ装置内のボリュームを AOU の記憶領域として用いてもよい。また図 29 には図示していないが、本実施の形態の場合、第 1 の実施の形態と同様に、仮想化ストレージ装置 1000L, 1000R には管理ホスト 1200 (図 1) が接続されている。

【0306】

ここでは、これまでに説明してきた実施の形態と異なる部分を中心に述べていく。2 台の仮想化ストレージ装置 1000L, 1000R は、AOU アドレス変換情報 31030L, 31030R を用いて AOU ボリューム 29010L, 29010R を生成し、これをホスト 1100 に提供する。AOU アドレス変換情報 31030L, 31030R には、AOU ボリューム 29010L, 29010R のアドレス空間と仮想化ストレージ装置 1000L, 1000R 内のプールボリュームの領域又はストレージ装置 1500L 内のプールボリュームの領域との対応関係が含まれる。

【0307】

ホスト 1100 から AOU ボリューム 29010L に対してライトリクエストが発行されると、正系の仮想化ストレージ装置 1000L は、リクエスト対象のアドレス範囲にプールボリュームの領域が割り当てられているか判断し、割り当てられていない場合は仮想化ストレージ装置 1000L 又はストレージ装置 1500L が有するプールボリュームの領域を割り当てる。そしてライトリクエストが処理されることで、正系の仮想化ストレージ装置 1000L のキャッシュ領域にライトデータが保存される。また、AOU ボリューム 29010L に対するライトデータは同期リモートコピーによって副系の仮想化ストレージ装置 1000R へ転送され、正系と同様にキャッシュ領域にライトデータが保存される。

【0308】

その後、両方の仮想化ストレージ装置 1000L, 1000R はデステージング処理を行うが、ストレージ装置 1500L に対応したライトデータに対しては仮想化ストレージ装置 1000L, 1000R の片方だけがデステージングを行う。なぜならば、両方の仮想化ストレージ装置 L, 1000R が独立にライトデータのデステージングを行うと、ストレージ装置 1500L に保存されるデータが不整合な状態（例えば、最後にライトしたデータが消えて、前回のライトデータに戻ってしまう等のデータ消失やライト順序の不整合がある）になってしまうからである。そのため、デステージングが必要になる前に予めどちらの仮想化ストレージ装置 1000L, 1000R がデステージングを行うか決めておく必要がある。本実施の形態ではその一例として正系の仮想化ストレージ装置 1000L がデステージングを行う場合について説明を行うが、副系の仮想化ストレージ装置 1000R が行ってもよく、又はデステージング対象のアドレス空間を元にどちらの仮想化ス

10

20

30

40

50

トレージ装置 1 0 0 0 L , 1 0 0 0 R が当該役割を担うか決定してもよい。

【 0 3 0 9 】

リードリクエストの場合も、正系の仮想化ストレージ装置 1 0 0 0 L は、まずはリクエスト対象のアドレス範囲にプールボリュームの領域が割り当てられているかどうかを判断する。判断の結果、割り当てられている領域に対しては、仮想化ストレージ装置 1 0 0 0 L は、該当するプールボリュームの領域（図示しないキャッシュメモリ上のデータを含む）からデータを読み出してホスト 1 1 0 0 へ転送し、割り当てられていない場合は予め定められた値（例えばゼロ）を返す。

【 0 3 1 0 】

図 3 0 は、正系の仮想化ストレージ装置 1 0 0 0 L の機能停止後に副系の仮想化ストレージ装置 1 0 0 0 R へ I / O リクエスト処理を切り替えた後の概要図である。本図にあるとおり、副系の仮想化ストレージ装置 1 0 0 0 R は A O U アドレス変換情報 3 1 0 3 0 R を元にストレージ装置 1 5 0 0 L や仮想化ストレージ装置 1 0 0 0 L 内部の A O U アドレス変換情報 3 1 0 3 0 R を用いて I / O リクエストを処理する。そのために、正系と副系の仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R は、通常時から通信を行うことで、A O U アドレス変換情報 3 1 0 3 0 L , 3 1 0 3 0 R のストレージ装置 1 5 0 0 L が関係する部分を同じ内容に維持する。これにより副系の仮想化ストレージ装置 1 0 0 0 R は、ストレージ装置 1 5 0 0 L の割り当て状況を引き継ぐことができる。また、副系の仮想化ストレージ装置 1 0 0 0 R は、正系の仮想化ストレージ装置 1 0 0 0 L 内のキャッシュメモリからデステージングされたデータでない限り、キャッシュメモリに格納されているデータを当該キャッシュメモリから削除しないようにする。これにより、機能停止時に正系の仮想化ストレージ装置 1 0 0 0 L 内のキャッシュメモリからデータが揮発した場合においてもデータ消失が発生しないようにすることができる。

【 0 3 1 1 】

< 2 . 仮想化ストレージ装置で実行されるプログラム及び情報 >

図 6 との対応部分に同一符号を付した図 3 1 は、仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R 上で実行されるソフトウェアプログラムと当該プログラムが用いる情報について示している。

【 0 3 1 2 】

この図 3 1 において、A O U 向け I / O 処理プログラム 3 1 0 1 0 は、仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R が受信した I / O リクエストを処理するプログラムで、第 1 ~ 第 1 4 の実施の形態における I / O 処理プログラム 6 0 2 0 （図 6 ）の機能を一部に含む。

【 0 3 1 3 】

A O U 管理プログラム 3 1 0 2 0 は、A O U 機能に関する設定や後述する重複削除（Duplication）処理を実行するためのプログラムである。また A O U アドレス変換情報 3 1 0 3 0 は、A O U ボリューム 2 9 0 1 0 L , 2 9 0 1 0 R とプールボリュームの領域との対応関係に関する情報である。さらに A O U プール管理情報 3 1 0 4 0 は、A O U 機能が用いるプールボリュームの集合（プール）を管理するための情報である。

【 0 3 1 4 】

< 2 . 1 . A O U アドレス変換情報 >

図 3 5 は A O U アドレス変換情報 3 1 0 1 0 の具体的な内容を示している。仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R は、ホスト 1 1 0 0 へ提供するボリュームの識別子と当該ボリューム内のアドレス空間を先頭から決められた大きさ（セグメントサイズ）に区切った領域（セグメント）のアドレスとでデータの保存領域等を管理する。なお、このセグメントサイズはプール定義時に設定される値である。

【 0 3 1 5 】

図 3 5 において、「A O U ボリューム識別子」及び「アドレス空間」は、対応するセグメントを含む A O U ボリューム 2 9 0 1 0 L , 2 9 0 1 0 R の識別子と当該 A O U ボリューム 2 9 0 1 0 L , 2 9 0 1 0 R 内における当該セグメントのアドレスとをそれぞれ示す

10

20

30

40

50

。またプールIDは、そのAOUボリューム29010L, 29010Rに領域を割り当てるプールの識別子を示す。

【0316】

「COW (Copy On Write) フラグ」は、そのセグメントに対するライトリクエストが到着したときに、対応するライトデータを別途確保したプールボリューム領域に保存する必要があるかどうかを示すフラグである。本フラグは、異なるセグメントが同じプールボリュームの領域に対応付けられていた場合に、ライトデータを他のプールボリュームに保存する必要があることを意味する「ON」となることがある。

【0317】

「プールボリューム領域識別子」は、そのセグメントが保存すべきデータを実際に保存しているプールボリューム領域の識別子を示す情報（識別情報）である。この識別情報は例えば以下の情報から構成される。

【0318】

(1) 仮想化ストレージ装置1000L, 1000Rの内部のボリュームの領域を用いている場合は、内部ボリュームの識別子及びアドレス範囲。

(2) ストレージ装置1500L内のボリュームの領域をもちいている場合は、ポートネーム等の装置又は通信先を識別する情報、LUN等の装置内のボリュームを識別する情報、及びアドレス範囲。

(3) 未割り当て領域の場合はNULL

【0319】

「引継ぎ領域」は、対応する「プールボリューム領域識別子」の欄に識別子が記載されたプールボリュームが正系と副系の仮想化ストレージ装置1000L, 1000Rの両方から管理されるかどうかを示す情報（両方から管理される場合は「Yes」、両方から管理されない場合は「No」）である。

【0320】

「対AOUボリューム識別子」は、対応するAOUボリューム識別子で特定されるボリュームと対を形成するAOUボリューム29010L, 29010Rの識別子が保持される。この識別子としては、対応する仮想化ストレージ装置1000L, 1000Rの識別子と、対応するAOUボリューム29010L, 29010Rの識別子とを組み合わせたものが用いられる。

【0321】

なお、上述のようにAOUの領域管理をセグメントを単位として行うのは、ブロック単位で管理するとAOUアドレス変換情報31030等の管理情報が大きくなりすぎてしまうためにI/O性能が悪化してしまうことが理由の一つである。

【0322】

< 2.2. AOUプール管理情報 >

図36は、AOUプール管理情報31040の具体的な構成を示している。AOUプール管理情報31040はプール毎に以下の情報を保持する。

【0323】

(1) セグメントサイズ

(2) プールに割り当てられたボリューム（プールボリューム）のリスト

(3) プールボリュームの領域で割り当てられていない領域のリスト

(4) 空き容量

(5) 容量が不足してきたことを警告するアラートを出すスレッシュホールド値

(6) プール対の相手が設定された仮想化ストレージ装置の識別子と当該装置内のプールID。なお、「プール対」については後ほど説明する。

【0324】

< 3. 初期化 >

本実施の形態の初期化は以下の手順で行われる。

【0325】

10

20

30

40

50

1. プールの初期化
2. AOUボリュームの生成
3. AOUボリューム同士の関連付け
4. 同期リモートコピーの設定

【0326】

以下に、その詳細について説明する。なお、以下においては、一部処理の処理主体を「管理ホスト」や「プログラム」として説明しているが、「管理ホスト」の部分については、その管理ホスト内のプロセッサが当該管理ホスト内のメモリに格納された対応するプログラムに基づいてその処理を実行し、「プログラム」の部分については、対応する仮想化ストレージ装置1000L, 1000R内のプロセッサ1011がそのプログラムに基づいて処理を実行することは言うまでもない。

10

【0327】

< 3. 1. プールの初期化 >

以下の手順で初期化を行う。

【0328】

(Step 1) 管理ホスト1200からの指示によって、仮想化ストレージ装置1000L, 1000Rの片方で実行されるAOU管理プログラム31020は、プールを作成する。この時、当該指示にはセグメントサイズが含まれる。また、プール作成の過程で、AOU管理プログラム31020はプールIDを含むAOUプール管理情報31040の該当エントリを作成する。

20

(Step 2) Step 1と同様の処理によって、仮想化ストレージ装置1000R, 1000Lのもう片方にもプールを作成する。

【0329】

(Step 3) 管理ホスト1200は、仮想化ストレージ装置1000L, 1000Rの両方に対して、Step 1とStep 2で作成したプールをプール対に設定すべき指示を発行する。当該指示にはプール対となるプールのIDと、そのプールを提供する仮想化ストレージ装置1000L, 1000Rの識別子との組が含まれる。当該指示を受信したAOU管理プログラム31020は、プール対となる相手の仮想化ストレージ装置1000L, 1000RのAOU管理プログラム31020と通信を行い、両プールに設定されたセグメントサイズが等しく、両プールとも既にプール対になっていないことが確認できたときに、それらプールをプール対に設定する。なお、AOU管理プログラム31020は、かかるプールをプール対に設定するに際して、相手のプールIDの識別子をAOUプール管理情報31040に登録する。

30

【0330】

(Step 4) 管理ホスト1200は、プールボリューム作成の指示を仮想化ストレージ装置1000L, 1000Rの片方に発行する。なお、当該指示には仮想化ストレージ装置1000L, 1000R内部に定義されたボリュームの識別子が含まれている。当該指示を受信した仮想化ストレージ装置1000L, 1000RのAOU管理プログラム31020は、指定されたボリュームの属性をプールボリュームに変更し、指定されたボリュームの識別子をAOUプール管理情報31040のプールボリュームリストに追加する。

40

【0331】

(Step 5) 管理ホスト1200は、Step 3と同様の指示を仮想化ストレージ装置1000R, 1000Lのもう片方へ発行する。指示を受け取った仮想化ストレージ装置1000R, 1000Lのもう片方はStep 3と同様の処理を行う。

【0332】

なお、管理者が仮想化ストレージ装置1000内部のボリュームをAOUに用いないと判断した場合はStep 4及びStep 5は省略することができる。

【0333】

(Step 6) 管理ホスト1200は、ストレージ装置1500Lのボリュームをプールボリュームに設定する指示を仮想化ストレージ装置1000L, 1000Rのどちらかに

50

発行する。なお、理解を容易にするために、以後の説明では、指示発行先を仮想化ストレージ装置 1000L、支持発行先と対の仮想化ストレージ装置 1000を仮想化ストレージ装置 1000Rであるものとするが、逆の関係であってもよい。ここで、当該指示にはストレージ装置 1500Lと当該ボリュームを識別する情報のほかに、当該ボリュームがプール対になっている相手の仮想化ストレージ装置 1000Rが引き継ぐことを示す情報が含まれる。指示を受信した仮想化ストレージ装置 1000LのAOU管理プログラム 31020は、対の仮想化ストレージ装置 1000Rと以下に示す連携を行う。

【0334】

(A) 指示を受信した仮想化ストレージ装置 1000Lは、指示に含まれるストレージ装置 1500Lのボリュームに対してリードリクエスト(又はライトリクエスト)を発行することで、当該ストレージ装置 1500L及び当該ボリュームのいずれもが存在し、かつ当該ボリュームにアクセス可能であるかどうかを確認する。ストレージ装置 1500Lやかかるボリュームが存在しなかったり、当該ボリュームにアクセスができなかった場合は管理ホスト 1200にエラーを返し、アクセス可能であった場合は次に進む。なお、当該エラーにはストレージ装置 1500Lに対するアクセスができなかったことを示す情報を添付するものとし、この情報を管理ホスト 1200において表示するようにしてもよい。

10

【0335】

(B) 指示を受信した仮想化ストレージ装置 1000Lは、対の仮想化ストレージ装置 1000Rに対してプールボリューム作成指示を転送する。なお、当該指示には管理ホスト 1200からの指示に含まれていた対象ボリュームを識別する情報と、当該ボリュームがプール対に属する両方のプールで管理することを示す情報とが含まれている。なお、プールボリューム作成指示の転送先は、AOUプール管理情報 31040における「プール対の識別情報」を参照することで特定できる。

20

【0336】

(C) 仮想化ストレージ装置 1000Rは、(B)の指示を受信すると、(A)と同様の処理を行うことでストレージ装置 1500L内のかかるボリュームにアクセス可能であることを確認する。そして、かかるボリュームにアクセス可能ならば、AOUプール管理情報 31040のプールボリュームリストに当該ボリュームを共通管理であることを示す情報と共に追加し、前述の指示を受信した仮想化ストレージ装置 1000Lへ結果を返す。一方、確認の結果、かかるボリュームにアクセスが不可能だった場合は、対の仮想化ストレージ装置 1000Rからストレージ装置 1500Lへのアクセスができなかったことを示す情報を添付して失敗を意味する結果を返す。

30

【0337】

(D) 結果を受け取った前述の指示を受信した仮想化ストレージ装置 1000Lは、かかるボリュームへのアクセス結果が失敗だった場合は、理由と共にその結果を管理ホスト 1200へ転送し、一連の処理を終了する。一方、かかるボリュームへのアクセス結果が成功だった場合は、この結果をAOUプール管理情報 31040のプールボリュームリストに当該ボリュームが共通管理すべきボリュームであることを示す情報と共に追加し、管理ホスト 1200へ成功を意味する結果を転送し、終了する。

【0338】

なお、(C)及び(D)でボリュームをプールボリュームリストに追加した場合、AOU管理プログラム 21020は、対応する「空き容量」の欄に格納されていた空き容量を、追加したボリュームの容量を加算した値に更新し、そのボリュームの領域を空き領域リストに追加する。また、Step 5の処理は管理ホスト 1200から仮想化ストレージ装置 1000L, 1000Rの両方に対して別々に指示を発行することで行っても良い。

40

【0339】

(Step 7)管理ホスト 1200は、仮想化ストレージ装置 1000L, 1000Rの両方に対して容量警告の設定値の設定指示を転送し、当該指示を受信した仮想化ストレージ装置 1000L, 1000Rの各々は指示に含まれる値をAOUプール管理情報 31040に設定する。

50

【0340】

< 3.2. AOUボリュームの作成 >

AOUボリューム29010L, 29010Rの作成は、仮想化ストレージ装置1000L, 1000Rの各々に対して指示を出すことで行われる。以下にその手順を示す。

【0341】

(Step1)管理ホスト1200は、ボリュームの容量とプールIDを伴ったAOUボリューム作成指示を、仮想化ストレージ装置1000L, 1000Rの各々へ転送する。

(Step2)当該指示を受信した仮想化ストレージ装置1000Lは、新しいAOUボリューム29010Lに関するAOUアドレス変換情報31030Lを作成する。このとき全てのセグメントについて、対応する「COWフラグ」及び「引継ぎ領域」を「No」に設定し、「プールボリューム領域識別子」も「NULL」に設定する。そして仮想化ストレージ装置1000Lは作成完了応答を返す。

(Step3)同様に、当該指示を受信した仮想化ストレージ装置1000Rは、新しいAOUボリューム29010Rに関するAOUアドレス変換情報31030Rを作成する。作成の詳細はStep2と同様である。

【0342】

なお、AOUボリューム29010L, 29010Rの作成は、管理ホスト1200から仮想化ストレージ装置1000L, 1000Rのいずれか片方に指示を出し、指示を受けた仮想化ストレージ装置1000L, 1000Rが対の仮想化ストレージ装置R, 1000Lに指示を出し直してもよい。なお、ボリューム作成指示にポート名前やLUNを含めることでAOUボリューム29010L, 29010Rに管理者が指定したポート名前配下のLUNを割り当ててもよい。また、AOUボリューム29010L, 29010Rの作成指示後にポート名前とLUNを割り当ててもよい。

【0343】

< 3.3. AOUボリューム同士の関連付け >

それぞれの仮想化ストレージ装置1000L, 1000Rに作成したAOUボリューム29010L, 29010R同士を関連付ける。そのために、管理ホスト1200は、仮想化ストレージ装置1000L, 1000Rにそれら2つのAOUボリューム29010L, 29010Rの識別子を含む関連付け指示を転送する。当該指示を受けた仮想化ストレージ装置1000L, 1000Rは、AOUアドレス変換情報31030の該当する「AOUボリューム識別子」の欄に、対となるAOUボリューム29010L, 29010Rを登録する。本指示は、それぞれの仮想化ストレージ装置1000L, 1000Rに対して行うことでAOUボリューム29010L, 29010R同士の関連付けが行われるが、他の実施の形態に開示されている通り、片方の仮想化ストレージ装置1000L, 1000Rがもう片方の仮想化ストレージ装置1000R, 1000Lに本指示を転送する形で実現してもよい。

【0344】

なお、上記関連付けの指示の際、指示に含まれるAOUボリューム29010L, 29010Rの存在を確認すると共に、AOUボリューム29010L, 29010Rの片方がプール対の片方のプールから生成されており、もう片方のAOUボリューム29010R, 29010Lがプール対のもう片方のプールから生成されていることを確認することで、プール管理の実装を簡単にしてもよい。また、本関連付けはAOUボリューム29010L, 29010Rの作成や同期リモートコピーの設定に伴って行われても良い。

【0345】

< 3.4. 同期リモートコピーの設定 >

これまで説明した同期リモートコピーではInitial Copying状態でボリュームの全ての領域をコピーする必要があったが、本実施の形態では、形成コピーは以下に示す手順で行う。また、理解を簡単にするために、以後の説明では正系の仮想化ストレージ装置1000を仮想化ストレージ装置1000Lとし、副系の仮想化ストレージ装置1000を仮想化ストレージ装置1000Rであるものとして説明する。

【0346】

(Step 1) コピー元となる(すなわち当該ボリュームに対しては正系となる)仮想化ストレージ装置1000Lは、変数*i*にAOUボリューム29010Lの先頭セグメントを代入する。

【0347】

(Step 2) コピー元の仮想化ストレージ装置1000Lは、AOUアドレス変換情報31030におけるセグメント*i*の「引継ぎ領域」及び「プールボリューム領域識別子」をそれぞれ確認し、それぞれの条件下で以下の処理をおこなう。

(A) 「引継ぎ領域」が「No」の場合には、通常の形成コピーに従ってセグメント*i*のデータをコピーする。仮想化ストレージ装置1000L内部のプールボリュームの領域のため、冗長性確保のためにコピーしなければならないからである。

10

(B) 「引継ぎ領域」が「Yes」の場合には、セグメント*i*に関するその仮想化ストレージ装置1000L内の図示しないキャッシュメモリ上のダーティデータをデステージングするか、形成コピーでコピー先(すなわち当該ボリュームに対しては副系となる)仮想化ストレージ装置1000Rのキャッシュ領域へコピーする。キャッシュメモリ上のデータを除けばデータは正系の仮想化ストレージ装置1000Lの外部にあるため、キャッシュメモリ上のデータを当該仮想化ストレージ1000Lの外部へ追い出せば正系の仮想化ストレージ装置1000Lが機能停止しても失うデータはないからである。

(C) 「プールボリューム領域識別子」が「NULL」の場合には、セグメント*i*には正系・副系共に領域が割り当てられていないため、コピーは行わない。

20

【0348】

(Step 3) コピー元の仮想化ストレージ装置1000Lは、セグメント*i*が最後の場合は形成コピーを終えてペア状態をDuplex状態に遷移させ、そうでない場合は変数*i*に次のセグメントを設定してStep 1へ戻る。

【0349】

なお、上記処理は、仮想化ストレージ装置1000L, 1000R間の再同期処理で用いてもよく、片方の仮想化ストレージ装置1000L, 1000Rが機能停止し、復旧した後の処理で用いてもよい。

【0350】

< 4. I/Oリクエスト処理について >

30

ここから、本実施の形態のI/Oリクエスト処理について説明する。

【0351】

< 4.1. ライトリクエスト処理 >

図32は、AOU向けI/O処理プログラム31010がライトリクエストを受信したときに実行する処理内容を示すフローチャートである。なお、これまでの説明では、ライトリクエストを構成するコマンドとライトデータの個々についてフローチャートを用いた説明はしなかったが、本処理はライトリクエスト対象の一部の領域が割り当て済みで他の領域が未割り当てである場合もあるため、フローチャートを参照しながら詳細に説明する。

【0352】

40

(S32001) AOU向けI/O処理プログラム31010は、ライトリクエストを構成するライトコマンドを受信する。このライトコマンドにはアドレス(位置)及びデータ長が含まれる。

【0353】

(S32100) AOU向けI/O処理プログラム31010は、受信したライトコマンドを元に割り当て処理を実行する。本処理を実行することで、AOU向けI/O処理プログラム31010は、セグメント毎にプールボリュームの領域が割り当てられているか否かをチェックし、プールボリュームの領域が未割り当てのセグメントや、他のセグメントと共有の領域を割り当てられている場合で「COWフラグ」が「ON」のセグメント(ライトの際には共有領域以外に書き込む必要があるセグメント)に対しては、プールボリ

50

ームの領域を割り当てる。またA O U向けI / O処理プログラム3 1 0 1 0は、かかるプールボリュームの領域の割り当て結果をA O Uアドレス変換情報3 1 0 3 0に反映させる。

【0354】

(S32003) A O U向けI / O処理プログラム3 1 0 1 0は、A O Uボリューム29010R, 29010Lの属性を確認し、当該A O Uボリューム29010R, 29010Lがコピー元ボリュームの場合にはS32004を実行し、そうでなければS32005を実行する。

(S32004) A O U向けI / O処理プログラム3 1 0 1 0は、リモートコピープログラム6010を呼び出すことでコピー先ボリュームを有する仮想化ストレージ装置(副系の仮想化ストレージ装置)1000Rに同期リモートコピーのコマンドを転送する。

10

(S32005) A O U向けI / O処理プログラム3 1 0 1 0は、S32001に対応するライトリクエストを構成するライトデータ(の一部又は全て)を受信する。

【0355】

(S32006) A O U向けI / O処理プログラム3 1 0 1 0は、A O Uボリューム29010R, 29010Lの属性を確認し、当該A O Uボリューム29010R, 29010Lがコピー元ボリュームの場合はS32007を実行し、そうでなければS32008を実行する。

(S32007) A O U向けI / O処理プログラム3 1 0 1 0は、リモートコピープログラム6010を呼び出すことで、コピー先ボリュームを有する仮想化ストレージ装置(副系の仮想化ストレージ装置)1000Rにライトデータを転送する。

20

【0356】

(S32008) A O U向けI / O処理プログラム3 1 0 1 0は、A O Uアドレス変換情報3 1 0 3 0に基づいてA O Uボリューム29010R, 29010L上のアドレスから、実際にライトデータが保存されるプールボリュームの領域を求める。そして求めた領域に対するライトデータをキャッシュメモリ上で保存・管理する。

(S32009) A O U向けI / O処理プログラム3 1 0 1 0は、ライトデータ受信の続きの有無を判断し、続きがある場合はS32005を再び実行する。

(S32010) A O U向けI / O処理プログラム3 1 0 1 0は、ライト完了の応答を正系の仮想化ストレージ装置1000L又はホスト1100に転送し、このライトリクエスト処理を完了する。

30

【0357】

なお、副系の仮想化ストレージ装置1000Rは、同期リモートコピーのコマンドの受信をホスト1100からのライトコマンドの受信と同様に扱う。同様に、仮想化ストレージ装置1000Rは、同期リモートコピーのデータ転送によるデータ受信をホスト1100からのライトデータの受信と同様に扱う。これによって副系の仮想化ストレージ装置1000Rでのライトリクエスト処理が理解されるであろう。

【0358】

< 4.1.1. 割り当て処理 >

以下に図32の割り当て処理について説明する。

40

【0359】

(S32101) A O U向けI / O処理プログラム3 1 0 1 0は、ライトコマンドで指定されたライト範囲(即ちライトアドレスとデータ長)をセグメント毎に分割する。

(S32102) A O U向けI / O処理プログラム3 1 0 1 0は、分割で生成した複数のセグメントの最初のセグメントを変数iに代入する。

【0360】

(S32103) A O U向けI / O処理プログラム3 1 0 1 0は、セグメントiの割り当て状態やC O W (Copy On Write) が必要かどうかを判断する。なお、このときの判断にはA O Uアドレス変換情報3 1 0 3 0を用いる。かかる判断の結果、プールボリュームの領域の割り当てが不要な場合はS32105を実行し、プールボリュームの領域が未割り

50

当ての場合や、割り当て済みであってもCOWフラグが立っている場合は（例えば、他のAOUボリューム29010R, 29010L上のセグメントと割り当て領域を共有している場合）S32104を実行する。

【0361】

（S32104）AOU向けI/O処理プログラム31010は、セグメントiに割り当てるため、プールボリュームの領域から未使用の領域を探す。そして、探した領域をAOUアドレス変換情報31030の「プールボリューム領域識別子」へ登録する。なお、未使用領域が見つからない場合は、ライトコマンドが失敗したことを示す応答を転送して、この割り当て処理を終了する。

【0362】

なお、失敗応答を転送する際には、この失敗応答と共に何らかのエラーメッセージを返しても良く、当該失敗応答の原因としてプールの容量不足が原因であることを示す情報を含めても良い。さらに、「COWフラグ」が立っている場合の領域割り当ての場合、AOU向けI/O処理プログラム31010は、領域割り当てに際して旧領域（共有領域）から割り当て領域へデータコピーを行うようにしてもよい。ただし、セグメントi全体がライト対象の場合はこのデータコピーを省略することができる。また、領域の割り当てに伴って、AOU向けI/O処理プログラム31010は、AOUプール管理情報の空き領域リストを編集し、空き容量の削減を行うようにしてもよい。

【0363】

さらに、AOU向けI/O処理プログラム31010は、割り当てたプールボリューム上の領域と当該領域を割り当てたAOUボリューム29010R, 29010Lのセグメントの情報を副系の仮想化ストレージ装置1000Rへ転送する。なお、当該割り当て情報は同期リモートコピーのコマンドと共に転送してもよい。

【0364】

（S32105）AOU向けI/O処理プログラム31010は、次のセグメントが存在するかどうか確認し、存在する場合はS32106を実行し、存在しない場合は本処理を終了し、ライトリクエスト処理へ戻る。

（S32106）AOU向けI/O処理プログラム31010は、変数iに次のセグメントを代入する。

【0365】

以上の処理によって、仮想化ストレージ装置1000Lは、セグメント毎の割り当て状況を確認し、必要ならばセグメントにプールボリュームの領域を割り当てる。

【0366】

< 4.1.2. 副系のプールボリューム領域割り当て方法 >

副系の仮想化ストレージ装置1000Rのプールボリューム領域割り当てステップ（S32104）は、正系の仮想化ストレージ装置1000Lから受信した割り当て情報を元に以下の方法によって、セグメントに対して領域を割り当てる。

【0367】

（A）正系の仮想化ストレージ装置1000Lが共有のストレージ装置（すなわちストレージ装置1500L）のプールボリュームから領域を割り当てた場合には、副系の仮想化ストレージ装置1000Rは、AOUアドレス変換情報31030における対応するセグメントの「引継ぎ領域」を「Yes」に、「プールボリューム領域識別子」を受信した領域識別子に設定する。これによって、共有ストレージ装置1500Lに関するプールボリューム領域の割り当ては正系と副系で同じ対応になる。

【0368】

（B）正系の仮想化ストレージ装置1000Lが仮想化ストレージ装置1000R内部のボリュームから領域を割り当てた場合には、副系の仮想化ストレージ装置1000Rは、内部ボリュームの空き領域を探して該当するセグメントに割り当てる。その結果、AOUアドレス変換情報31030における当該セグメントの「引継ぎ領域」は「No」に、「プールボリューム領域識別子」は内部ボリュームの領域が設定される。これによって、正

10

20

30

40

50

系の仮想化ストレージ装置 1000L が内部ボリュームの領域を割り当てたセグメントは副系の仮想化ストレージ装置 1000R でも内部ボリュームを割り当てることができる。

【0369】

< 4.2. リードリクエスト処理 >

図 33 は、AOU 向け I/O 処理プログラム 31010 が、リードリクエストを受信したときに実行する処理内容を示すフローチャートである。以下に当該フローチャートを参照して、かかる処理内容について説明する。

【0370】

(S33001) AOU 向け I/O 処理プログラム 31010 は、リードリクエストを構成するリードコマンドを受信する。なお、受信したリードコマンドにはアドレス（位置）及びデータ長が含まれる。

10

(S33002) AOU 向け I/O 処理プログラム 31010 は、リードコマンドで指定されたリード範囲（即ちライトアドレスとデータ長）をセグメント毎に分割する。

(S33003) AOU 向け I/O 処理プログラム 31010 は、分割で生成した複数のセグメントの最初のものを変数 i に代入する。

【0371】

(S33004) AOU 向け I/O 処理プログラム 31010 は、セグメント i にプールボリュームの領域が割り当てられているかどうかを判断する。なお、判断には AOU アドレス変換情報 31030 を用いる。かかる判断の結果、プールボリュームの領域が割り当てられる場合には S33006 を実行し、プールボリュームの領域が未割り当ての場合には S33005 を実行する。

20

(S33005) AOU 向け I/O 処理プログラム 31010 は、その仮想化ストレージ装置 1000L, 1000R 内のキャッシュメモリ上に当該セグメント向けのキャッシュ領域を確保し、確保したキャッシュ領域をゼロで初期化し、ホスト 1100 へゼロデータを転送する。

【0372】

(S33006) AOU 向け I/O 処理プログラム 31010 は、割り当てられたプールボリュームの領域に保存されたデータを転送する。なお、当該プールボリュームの領域が既にキャッシュ領域に存在する場合（ステージング済みの場合）には、かかるデータをそのキャッシュ領域から転送し、キャッシュ領域に存在しない場合はステージング後に、当該データの転送を行う。

30

【0373】

(S33008) AOU 向け I/O 処理プログラム 31010 は、続きのセグメントがあるかどうかを判断し、ある場合は S33009 を実行し、ない場合は S33010 を実行する。

(S33009) AOU 向け I/O 処理プログラム 31010 は、変数 i に次のセグメントを代入し、再び S33004 を実行する。

(S33010) AOU 向け I/O 処理プログラム 31010 は、リード完了の応答をホスト 1100 に転送し、完了する。

【0374】

なお、処理の単純化のために、仮想化ストレージ装置 1000L はプールボリュームのある決められた領域に対して予め定められた値（ゼロ）を保存しておき、当該領域に保存されたデータを AOU ボリューム 29010R, 29010L の未割り当て領域に対するリードで転送してもよい。

40

【0375】

< 4.3. AOU 向けデステージング処理 >

図 34 は AOU 向け I/O 処理プログラム 31010 が実行するデステージング処理の処理内容を示すフローチャートである。以下、かかるデステージング処理について、当該フローチャートを参照しながら説明する。

【0376】

50

(S 3 4 0 0 1) A O U 向け I / O 処理プログラム 3 1 0 1 0 は、キャッシュアルゴリズムによってデステージ対象とするキャッシュメモリ上のデータを決定する。なお、キャッシュアルゴリズムは L R U (Less Recently Used) アルゴリズムを用いてダーティデータを対象として決定する方法が一般的であるが、これ以外のアルゴリズムを用いて決定してもよい。

【 0 3 7 7 】

(S 3 4 0 0 2) A O U 向け I / O 処理プログラム 3 1 0 1 0 は、デステージ対象のデータが共有ストレージ装置 (すなわち、ストレージ装置 1 5 0 0 L) が有するボリュームに対応するものかどうかを判断し、対応する場合は S 3 4 0 0 3 を実行し、対応しない場合は S 3 4 0 0 4 を実行する。

10

【 0 3 7 8 】

(S 3 4 0 0 3) A O U 向け I / O 処理プログラム 3 1 0 1 0 は、デステージング処理を実行し、その後、この一連の処理を終了する。なお、デステージング処理は他の実施の形態と同様に行われる。

(S 3 4 0 0 4) A O U 向け I / O 処理プログラム 3 1 0 1 0 は、デステージ対象のデータが格納されたボリュームのボリューム属性を判断し、当該ボリュームがコピー元ボリュームである場合には S 3 4 0 0 5 を実行し、当該ボリュームがコピー先ボリュームの場合には S 3 4 0 0 7 を実行し、それ以外の場合は S 3 4 0 0 3 を実行する。

【 0 3 7 9 】

(S 3 4 0 0 5) A O U 向け I / O 処理プログラム 3 1 0 1 0 は、デステージング処理を実行する。

20

(S 3 4 0 0 6) A O U 向け I / O 処理プログラム 3 1 0 1 0 は、デステージが終了したデータの R C デステージ許可指示を副系の仮想化ストレージ装置 1 0 0 0 R へ転送し、処理を終了する。

【 0 3 8 0 】

(S 3 4 0 0 7) A O U 向け I / O 処理プログラム 3 1 0 1 0 は、R C デステージ許可フラグが O N かどうかを確認し、O F F の場合は S 3 4 0 0 1 を再び実行し、別なデステージ対象のデータを選択し直す。なお、R C デステージ許可フラグは、同期リモートコピーによってキャッシュメモリ上にライトデータが保存又は更新された時点では O F F が設定され、S 3 4 0 0 6 で送信された指示を受信すると O N が設定される。

30

【 0 3 8 1 】

(S 3 4 0 0 8) A O U 向け I / O 処理プログラム 3 1 0 1 0 は、デステージング処理を実行し、処理を終了する。

【 0 3 8 2 】

本アルゴリズムによって、以下のキャッシュ制御が実現される。

【 0 3 8 3 】

(A) 共有ストレージ装置向けでなく、正系と副系の仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R でデステージを連携する必要のないキャッシュデータは両系独立にデステージを行う。

(B) 正系の仮想化ストレージ装置 1 0 0 0 L でのデステージ処理後に送信されるメッセージによって、副系の仮想化ストレージ装置 1 0 0 0 R のキャッシュデータのデステージが行われる。

40

【 0 3 8 4 】

なお、ステージング処理は、第 1 ~ 第 1 4 の実施の形態と同様に行なわれる。

【 0 3 8 5 】

< 4 . 3 . 1 . R C デステージ許可指示 >

R C デステージ許可指示の転送は、非同期に指示を送信してもよい。ただし、正系及び副系の仮想化ストレージ装置 1 0 0 0 L , 1 0 0 0 R はリモートコピーを契機として、R C デステージフラグに未反映の当該指示を無効化してもよい。

【 0 3 8 6 】

50

< 4.4. プールの空き領域監視 >

AOU管理プログラム31020は、定期的に各プールの空き領域を監視し、ユーザーが設定したスレッシュールド値を下回った場合は、管理ホスト1200へメッセージを送信する。これによって、容量不足に伴うホスト1100からのライトリクエストの失敗を回避することができる。さらに、AOU管理プログラム31020は、空き領域の監視を共有のストレージ装置1500Lと共有でないストレージ装置とで分けて管理し、容量不足の際に転送するメッセージを使い分けても良い。

【0387】

< 5. 正系の仮想化ストレージ装置障害時の切り替え >

正系の仮想化ストレージ装置1000Lが障害などで機能を停止した場合は、他の実施の形態と同様の処理を行うことでホスト1100は引きつづきアプリケーションを動作させることができる。

【0388】

一方で、ホスト1100は、コピー元ボリュームに対するライトリクエストが容量不足で失敗したことを契機として副系の仮想化ストレージ装置1000RにI/Oリクエスト先を切り替える場合もある。副系の仮想化ストレージ装置1000Rが有するプール容量が正系よりも多い場合は、当該切り替えによって、ホスト1100においてI/Oリクエストを発行しているアプリケーション2010(図30)の処理を継続することができるからである。

【0389】

なお、この場合はリクエスト先の切り替えによってリモートコピーの向きは反転するが、リモートコピーは停止する。なぜならば、旧正系の仮想化ストレージ装置1000Lはライトリクエスト時のプール容量不足で当該リクエストが失敗しているため、同期リモートコピーによって新正系(旧副系)の仮想化ストレージ装置1000Rに対してライトデータを書き込もうとしても失敗するからである。

【0390】

ただし、旧正系の仮想化ストレージ装置1000Lに対するリクエスト(特にリード)は継続可能であるため、本障害は仮想化ストレージ装置1000L, 1000R間の通信路障害と見分けがつかず、ホスト1100が旧正系の仮想化ストレージ1000Lの古いデータをリードする可能性がある。

【0391】

こうした状況を回避するため、リモートコピー失敗の理由がプール容量不足である場合には、ホスト1100からの旧正系の仮想化ストレージ装置1000Lへのリードリクエスト発行を抑制してもよい。あるいは、リモートコピー失敗の理由が絞れない間は、ホスト1100からの副系の仮想化ストレージ装置1000Rまたは1000Lに対するリードを抑制し、通信路障害であることが判明した時点でかかる抑制を解除するようにしてもよい。

【0392】

以上の処理によって、本実施の形態によるストレージシステムがサービス継続性の高いAOU機能を持ったストレージサービスを提供することができる。また、AOU機能は、I/Oリクエスト毎にAOUアドレス変換情報31030L, 31030Rを参照・変更する必要があり、通常のストレージI/Oよりもコントローラーの負荷が高い。したがって、ホスト1100が必要とするボリュームの一部(又は半分)については片方の仮想化ストレージ装置1000L, 1000Rが正系としてリードとライトを担当し、残りのボリュームについてはもう片方の仮想化ストレージ装置1000R, 1000Lが正系としてリードとライトを担当するようにしてもよい。このような構成を採用することで、ストレージシステムの可用性を維持しつつ、仮想化ストレージ装置1000L, 1000Rとの間でのAOU機能のコントローラー負荷の平準化を実現できる。

【0393】

< 6. プールボリューム領域の割り当てとデータ移行について >

10

20

30

40

50

これまでに述べたとおり、本実施の形態では仮想化ストレージ装置1000L, 1000R内部のボリュームとストレージ装置1500Lのボリュームの両方をプールボリュームとすることができる。そのため、アクセス頻度の高いデータが格納される又は格納されたセグメントに対して仮想化ストレージ装置1000L, 1000R内部のボリュームを割り当てることによって、アクセス性能の向上が図れるほかに、仮想化ストレージ装置1000L, 1000Rとストレージ装置1500Lとの間の通信ネットワークのボトルネック化を回避することもできる。

【0394】

しかし、AOUでは最初のライトリクエストによってセグメントにプールボリュームの領域を割り当てるため、仮想化ストレージ装置1000L, 1000R単体でアクセス頻度を考慮した割り当てをすることは難しい。こうした課題を解決する方法として以下の方法が考えられる。

10

【0395】

<6.1. AOUボリュームに属性を付加する方法>

AOUボリューム29010L, 29010Rを作成する時点でアクセス頻度に関する属性を与え、AOU向けI/O処理プログラム31010がセグメントにプールボリュームの領域の割り当てを行う際に、そのセグメントに書き込まれるデータのアクセス頻度がある程度分かっているときには、かかるアクセス頻度属性を参照して、アクセス頻度の高いデータが格納されるセグメントについてはその仮想化ストレージ装置1000L, 1000R内部のボリュームを割り当て、アクセス頻度の低いデータ(例えばバックアップデータ)が格納されるセグメントについてはストレージ装置1500L内のボリュームの領域を割り当てる。

20

【0396】

<6.2. プールボリューム領域のデータ移行>

AOUボリューム29010L, 29010Rに対するアクセス頻度をセグメント単位(又は複数セグメント単位)で測定し、アクセス頻度の高いセグメントに格納されているデータは仮想化ストレージ装置1000L, 1000R内部のプールボリュームの領域に移動させる。この場合、データの移行に伴って、AOUボリューム29010L, 29010Rにおける当該データの移行が行なわれたセグメントの対応先をストレージ装置1500L内のボリューム内のセグメントから、仮想化ストレージ装置1000L, 1000Rにおけるデータの移行先のセグメントに変更する必要があるが、AOU機能では元々仮想化ストレージ装置1000L, 1000R内においてアドレス変換を行っているため、ホスト1100に対して透過的にデータ移行を行うことができる。

30

【0397】

なお、本実施の形態でこのようなデータ移行を行う場合、対象となるセグメントのデータは、正系と副系の両方の仮想ストレージ装置1000L, 1000R内部のプールボリュームに保存されることが望ましい。しかし、他に効果がある場合(以下に列挙した)は片方のセグメントだけ仮想化ストレージ装置1000L, 1000R内部のプールボリュームの領域が割り当てられた形態にデータ移行を行うことも考えられる。

【0398】

(例1) どちらかの仮想化ストレージ装置1000L, 1000Rが先に内部のプールボリュームを使い果たし、共有のストレージ装置1500Lしかない場合。

(例2) コピー元のAOUボリューム29010Lに対するリードリクエストの負荷が大きく、正系の仮想化ストレージ装置1000Lとストレージ装置1500Lの間のネットワーク性能を圧迫する場合。

40

【0399】

こうした場合、正系の仮想化ストレージ装置1000Lは、ストレージ装置1500L内部のプールボリュームの領域から自身のプールボリュームの領域にセグメントのデータをコピーし、コピー先の領域を用いてAOUボリューム29010Lを提供する。一方の副系の仮想化ストレージ装置1000Rはコピー元のストレージ装置1500Lのプール

50

ボリュームの領域を用いてA O Uボリューム29010Rを提供することもできる。この場合、ストレージ装置1500Lのプールボリュームの領域に対するライトデータの反映は副系の仮想化ストレージ装置1000Rが行っても良い。

【0400】

さらに、リードもライトも含めたアクセス性能向上のためのセグメントのデータ移行の中間状態として、前述の正系の仮想化ストレージ装置1000Lだけ内部のプールボリュームの領域を用い、副系の仮想化ストレージ装置1000Rがストレージ装置1500Lのプールボリュームを用いる構成を採用してもよい。

【0401】

<7. 本実施の形態のバリエーション>

10

<7.1. ステージングやデステージング処理でアドレス変換を実施>

これまで述べてきた本実施の形態では、リードリクエスト処理やライトリクエスト処理でアドレス変換を行っている。本方法は、ライトリクエストの受付の時点で、プールボリュームの容量不足を契機とした失敗応答を返すことができる反面、リクエスト毎にアドレス変換を行うため、性能上の課題がある。こうした課題を解決する方法としてステージングやデステージング処理でアドレス変換を行う方法が考えられる。ただし、この方法ではデステージングの時点でセグメントに対してプールボリュームの領域の割り当てを行うため、HDD1030の二重閉塞等が原因のボリューム閉塞時と類似するデータ消失が発生する。そのため、後者の方式では空き容量に余裕が少なくなってきた時点からリクエストの処理を遅らせるか停止する等の処理を行っても良い。

20

【0402】

なお、具体的な処理はこれまで図32及び図33について説明してきた処理内容を以下に示すとおりに変更すればよい。

(ライトとデステージング) 図32のS32100の割り当て処理をデステージング処理のS34001の後に移動する。

(リードとステージング) 図33のS33004~S33006で行っているアドレス変換を伴った割り当て有無の判断と未割り当て時のゼロデータの転送とを、ステージングにて行う。

【0403】

さらに、両者の利点を併せ持つために、A O U向けI / O処理プログラム31010が、プールボリュームの空き容量がスレッシュホールド値以上の場合にはステージング/デステージング処理で変換を行い、かかる空き容量がスレッシュホールド値以下になった場合にはI / Oリクエスト処理で変換を行ってもよい。

30

【0404】

<7.2. De Duplication>

A O U管理プログラム31010は、I / Oリクエストとは独立にDe duplicationと呼ばれる以下の処理を行っても良い。

【0405】

(Step1) A O U管理プログラム31010は、各プールボリュームの領域のデータをスキャンし、重複するセグメントを探す。

40

(Step2) A O U管理プログラム31010は、プールボリューム領域同士で保存するデータが重複している事を知った場合、いずれか一つの領域だけを残し他の領域は空き領域として開放する。そして、A O Uアドレス変換情報31030における開放した領域に対応したセグメントの「プールボリューム領域識別子」は一つだけ残した領域に更新し、「COWフラグ」を「ON」にする。

【0406】

ここで、重複検知の方法としては、プールボリュームの領域毎のハッシュ値を計算後に、領域毎に、そのハッシュ値を他の領域のハッシュ値と順次比較し、同じ値の場合はさらに実際のデータを比較する2段階方式を採用してもよい。さらに、ハッシュ値の計算とデータの比較は負荷の高い処理であるため、副系の仮想化ストレージ装置1000Rにて処

50

理を行うことで負荷分散を行うことも考えられる。

【図面の簡単な説明】

【0407】

【図1】図1は、第1の実施の形態にかかる情報システムのハードウェア構成の一例を示すブロック図である。

【図2】図2は、第1の実施の形態の概要を示す第1の概念図である。

【図3】図3は、第1の実施の形態の概要を示す第2の概念図である。

【図4】図4は、第1の実施の形態の概要を示す第3の概念図である。

【図5】図5は、ホスト上のソフトウェア構成を表した概念図である。

【図6】図6は、仮想化ストレージ装置及びストレージ装置上のソフトウェア構成を表したブロック図である。

【図7】図7は、リモートコピーのペア状態とペア状態の遷移を表した概念図である。

【図8】図8は、I/Oパスマネージャが管理するデバイス関係テーブルを示す概念図である。

【図9】図9は、I/Oパスマネージャが初期化処理を行うときのフローを示したフローチャートである。

【図10】図10は、I/Oパスマネージャがライト処理を行うときのフローを示したフローチャートである。

【図11】図11は、I/Oパスマネージャがリード処理を行うときのフローを示したフローチャートである。

【図12】図12は、第2の実施の形態の概要を示す概念図である。

【図13】図13は、第3の実施の形態の概要を示す概念図である。

【図14】図14は、第4の実施の形態の概要を示す概念図である。

【図15】図15は、第5の実施の形態の概要を示す概念図である。

【図16】図16は、第6の実施の形態の概要を示す概念図である。

【図17】図17は、第7の実施の形態の概要を示す概念図である。

【図18】図18は、第8の実施の形態の概要を示す概念図である。

【図19】図19は、第9の実施の形態の概要を示す概念図である。

【図20】図20は、第10の実施の形態の概要を示す概念図である。

【図21】図21は、第11の実施の形態の概要を示す概念図である。

【図22】図22は、第12の実施の形態の概要を示す概念図である。

【図23】図23は、第13の実施の形態の概要を示す概念図である。

【図24】図24は、第14の実施の形態の概要を示す概念図である。

【図25】図25は、I/Oパスマネージャがライト処理を行うときの別なフローを示したフローチャートである。

【図26】図26は、I/Oパスマネージャがリード処理を行うときの別なフローを示したフローチャートである。

【図27】図27は、I/Oパスマネージャが図25に記したライト処理を行うときに、ストレージ装置にて行うライトリクエストに応じたペア操作を示したフローチャートである。

【図28】図28は、第15の実施の形態の概要を示す概念図である。

【図29】図29は、第16の実施の形態の概要を示す概念図である。

【図30】図30は、第16の実施の形態の概要を示す概念図である。

【図31】図31は、本実施の形態における仮想化ストレージ装置及びストレージ装置上のソフトウェア構成を表したブロック図である。

【図32】図32は、仮想化ストレージ装置がライト処理を行うときのフローを示したフローチャートである。

【図33】図33は、仮想化ストレージ装置がリード処理を行うときのフローを示したフローチャートである。

【図34】図34は、AOU向けでステージング処理のフローを示したフローチャートで

10

20

30

40

50

ある。

【図35】AOUアドレス変換情報の具体的内容の説明に供する概念図である。

【図36】AOUプール管理情報の具体的内容の説明に供する概念図である。

【符号の説明】

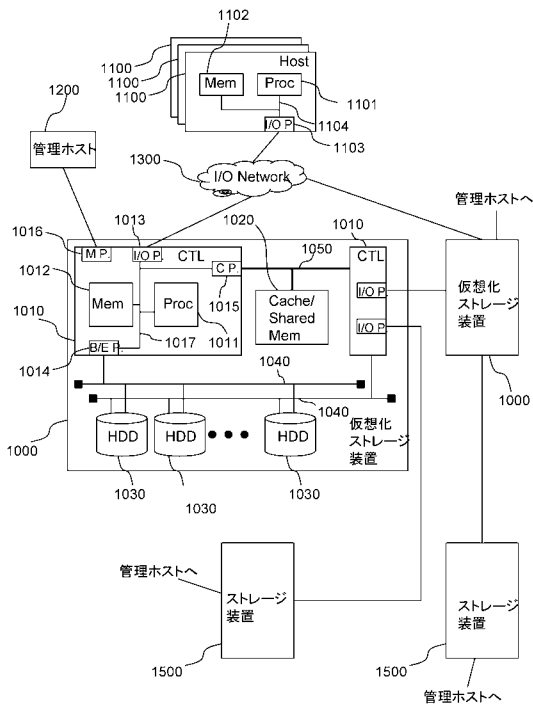
【0408】

1000, 1000L, 1000R.....仮想化ストレージ装置、1010.....コントローラ、1011, 1101.....プロセッサ、1020, 1020L, 1020R.....キャッシュメモリ、1030.....HDD、1100, 13010, 14000.....ホスト、1500, 1500L, 1500R, 15000, 15000L, 15000R.....ストレージ装置、2800L, 2800R.....仮想化スイッチ、3500LB, 3500RB, 5040, 5050.....ボリューム、2010, 14002.....アプリケーションプログラム、5000.....I/Oパスマネージャ、5010.....HBAデバイスドライバー、5020.....ファイルシステム、13001, 13002.....ストレージサブシステム、15002A, 15002B.....コマンドデバイス、15010A, 15010B.....差分ビットマップ、16000.....外部ストレージ装置。

10

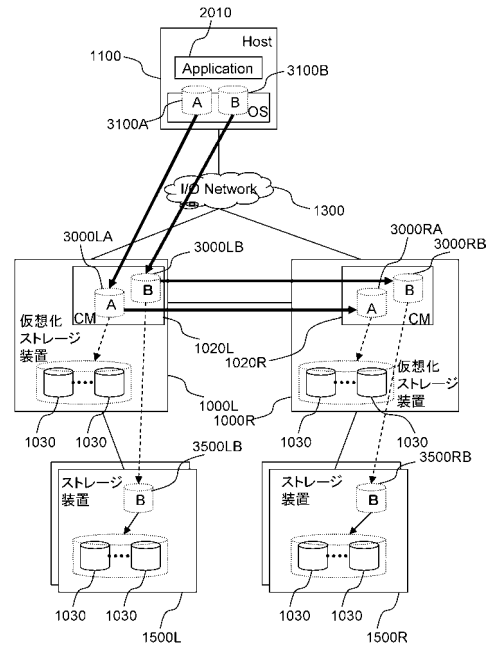
【図1】

図1



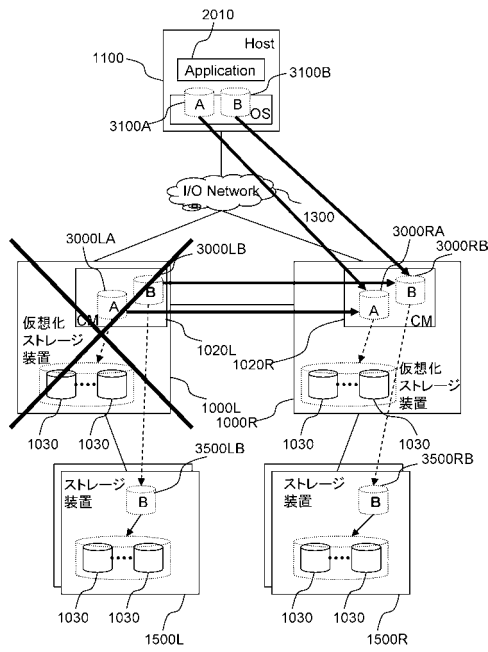
【図2】

図2



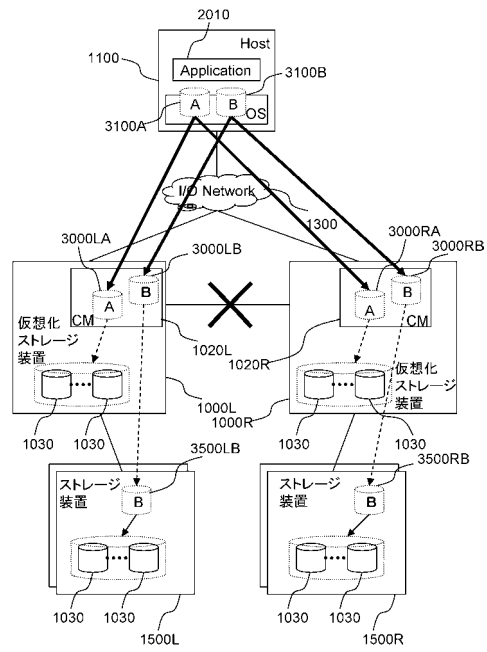
【図3】

図3



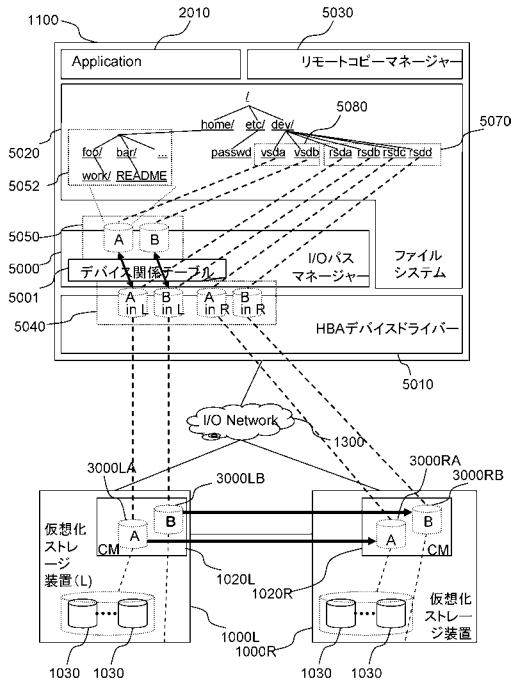
【図4】

図4



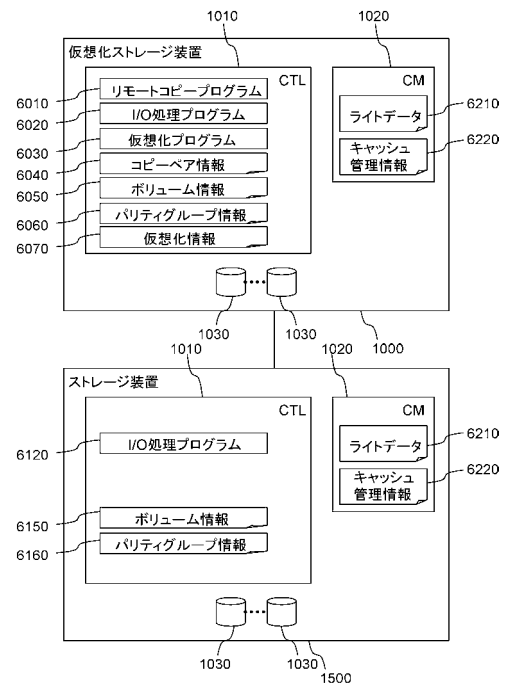
【図5】

図5



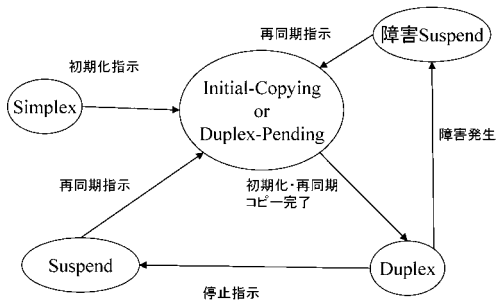
【図6】

図6



【図7】

図7



【図8】

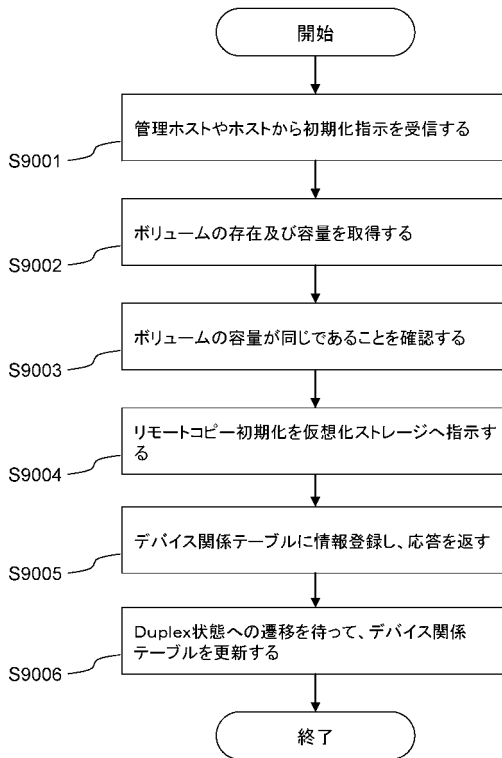
図8

5001

ホスト内で仮想的なボリュームの識別子	関係ボリューム識別子リスト	正系ボリューム識別子	障害状態	ペア状態
vsda	rsda, rsdb	rsda	通常状態	Duplex
vsdb	rsdc, rsdd	rsdc	副系準備中	Duplex-Pending
..

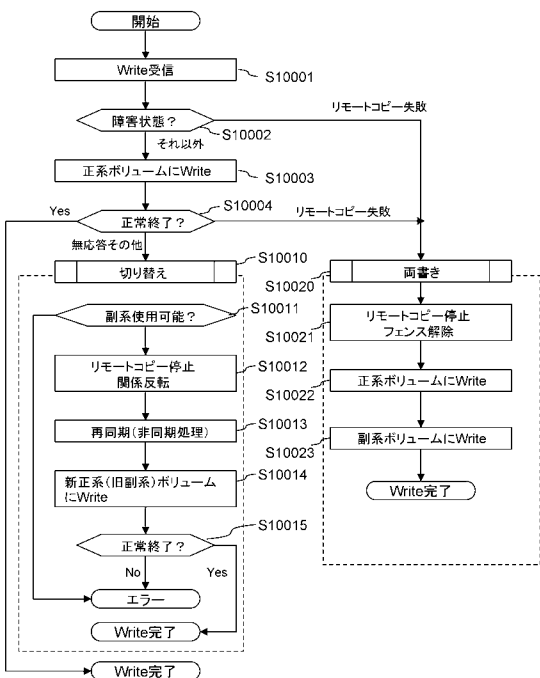
【図9】

図9



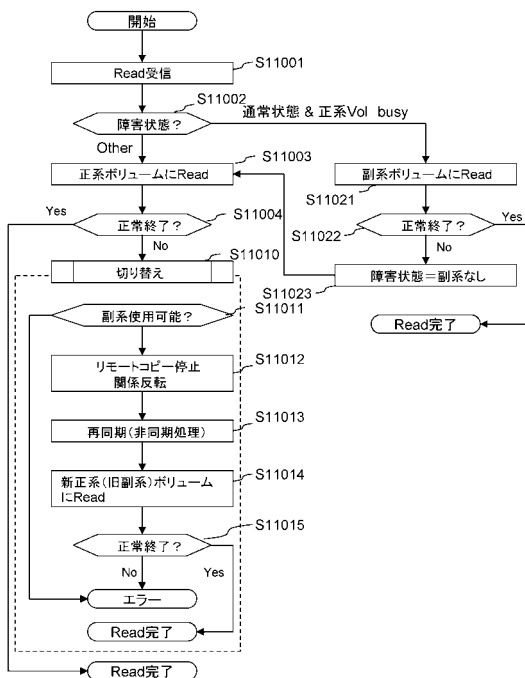
【図10】

図10



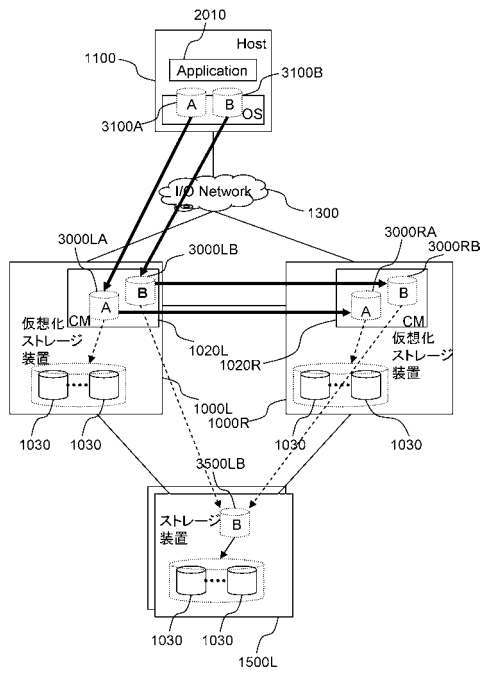
【図11】

図11



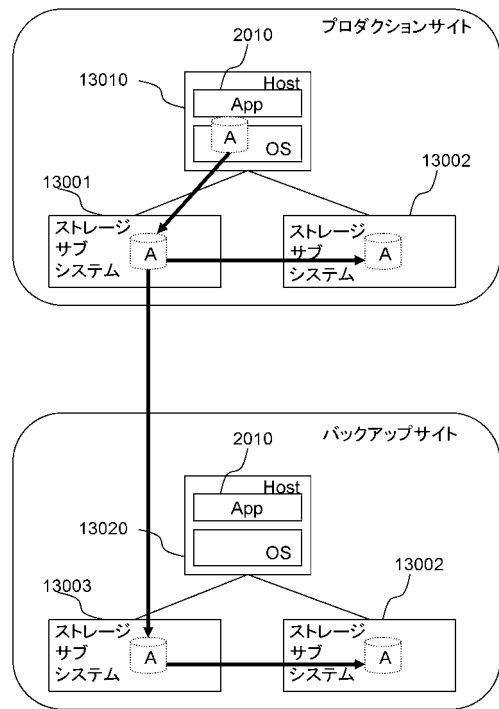
【図12】

図12



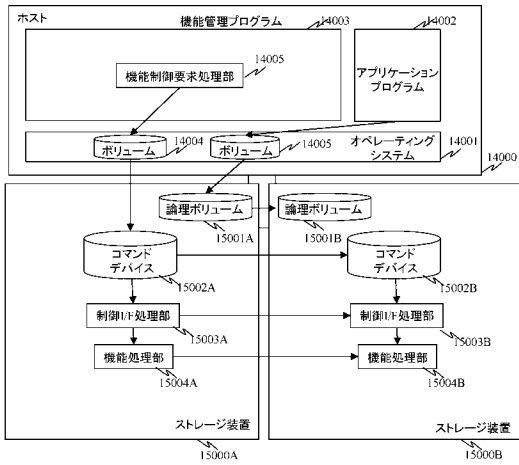
【図13】

図13



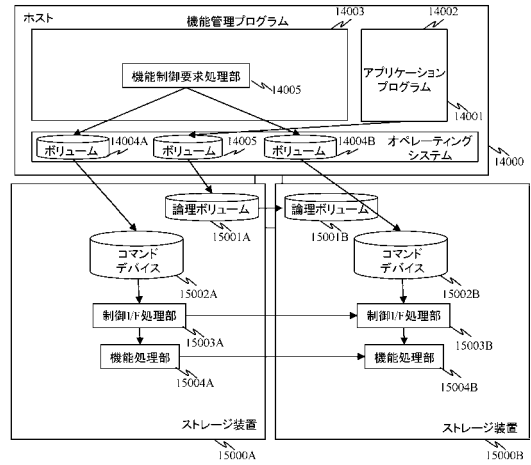
【図14】

図14



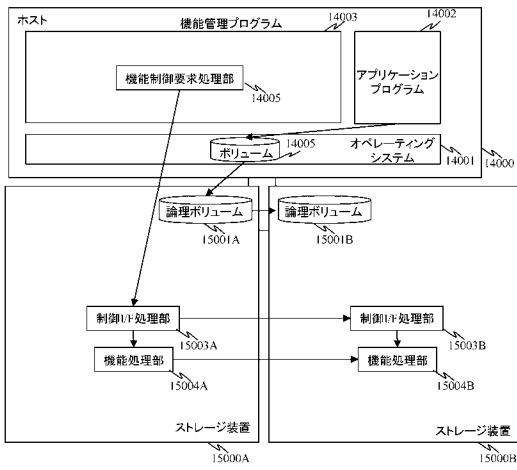
【図15】

図15



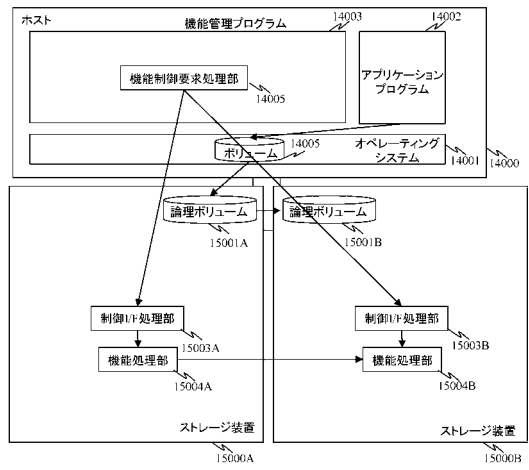
【図16】

図16



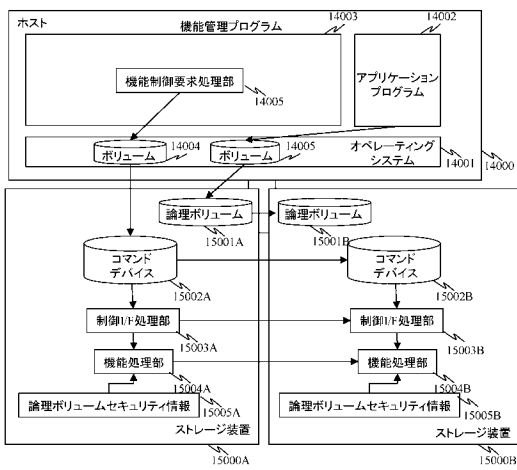
【図17】

図17



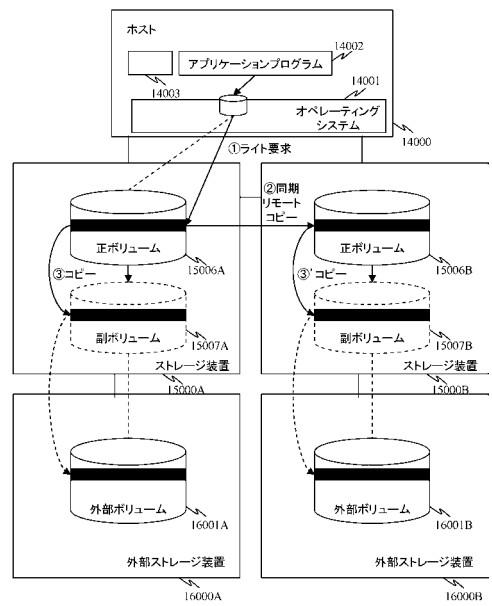
【図18】

図18



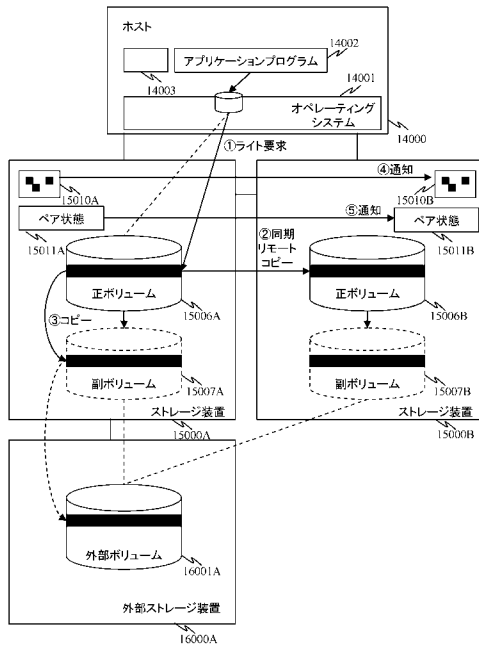
【図19】

図19



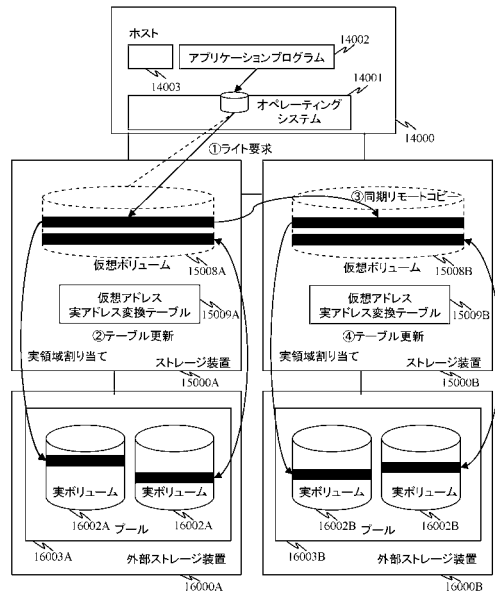
【図20】

図20



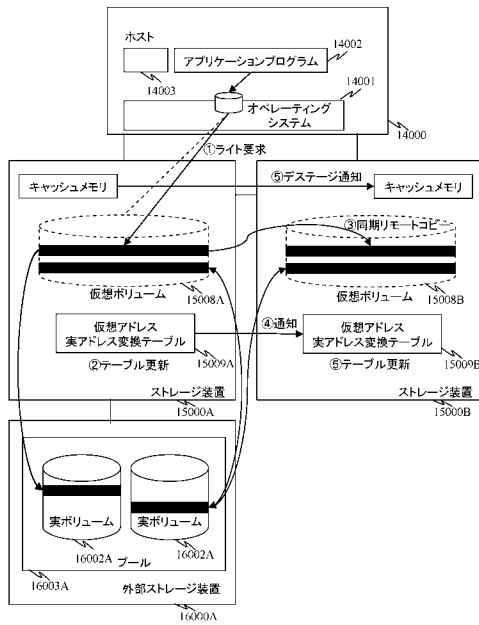
【図21】

図21



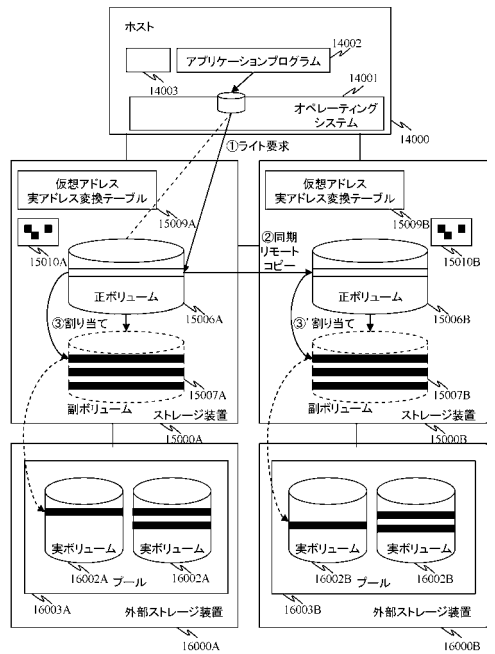
【図22】

図22



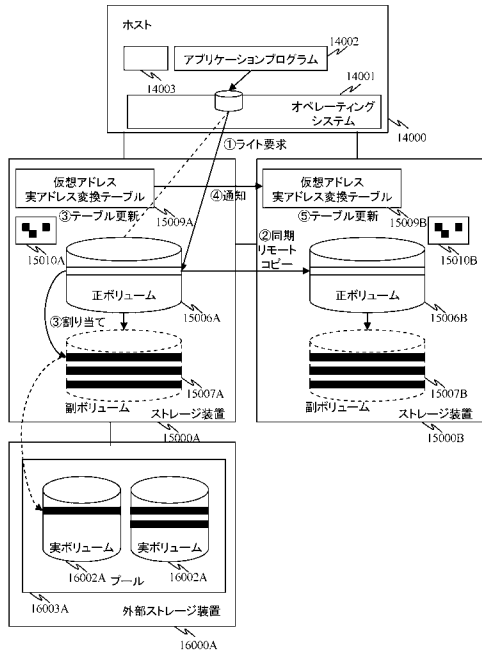
【図23】

図23



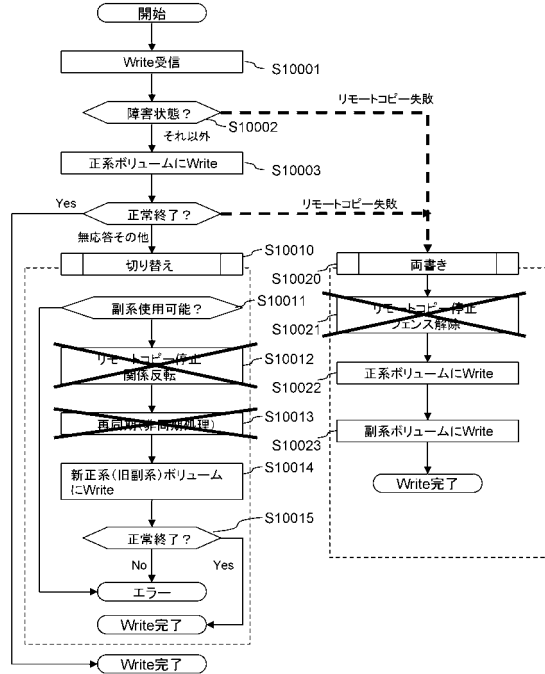
【図24】

図24



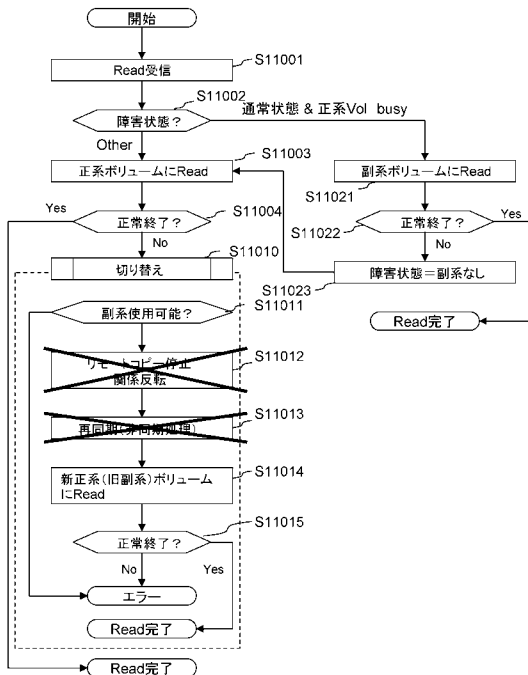
【図25】

図25



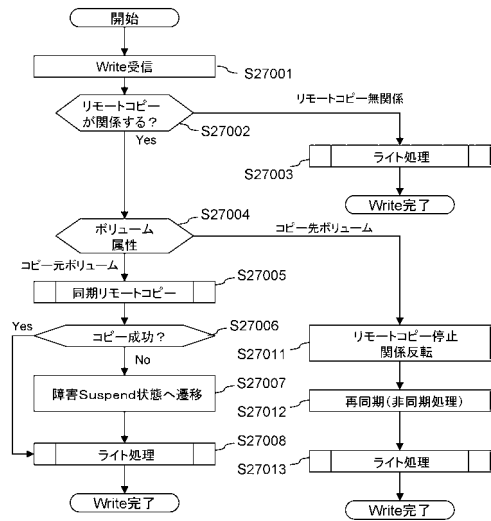
【図26】

図26



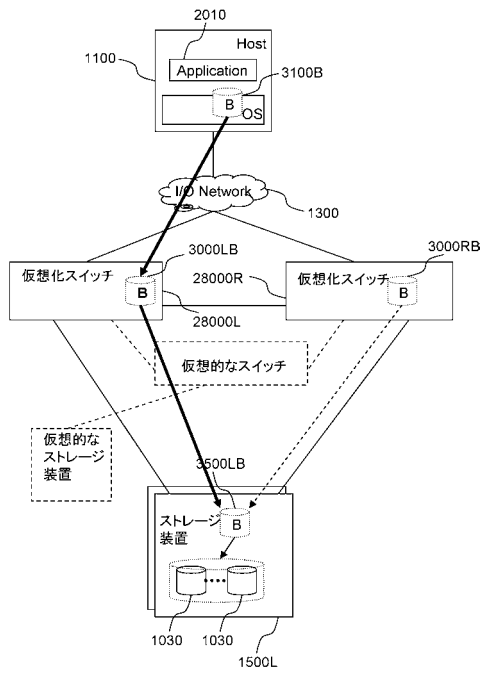
【図27】

図27



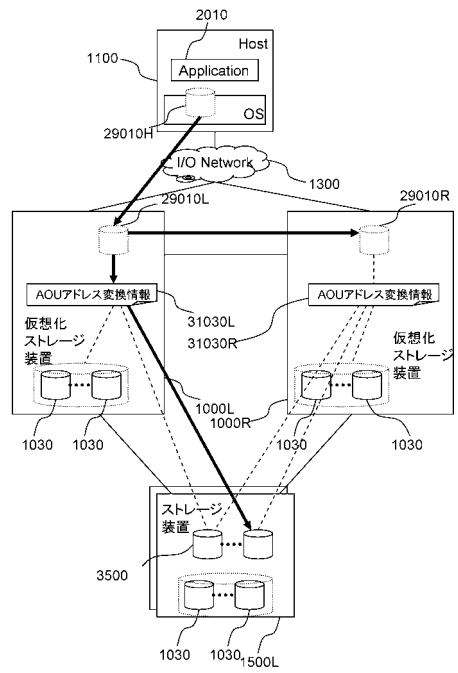
【図28】

図28



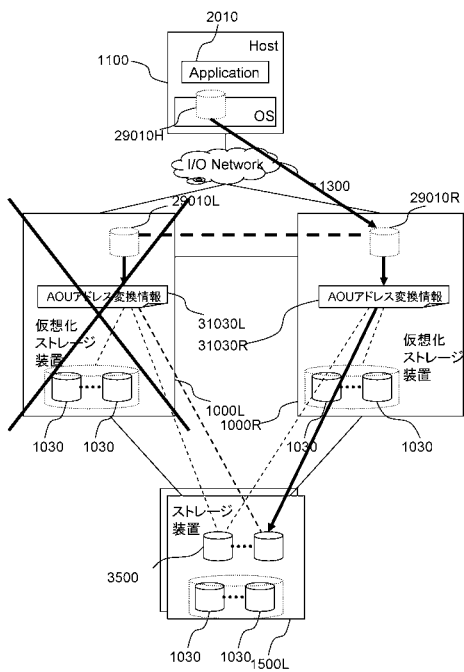
【図29】

図29



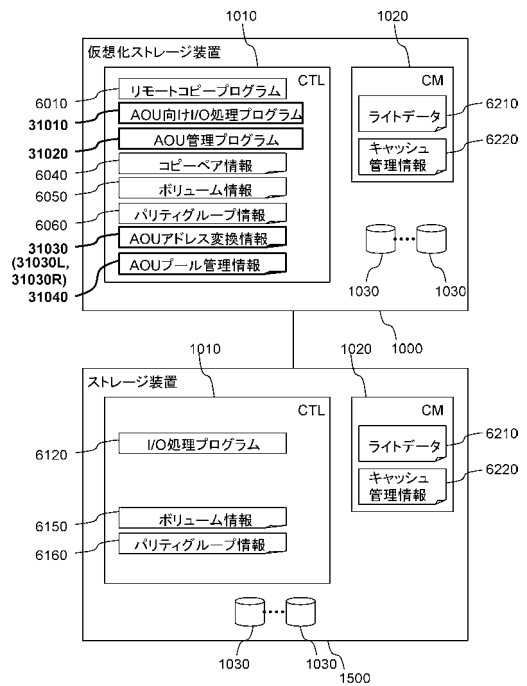
【図30】

図30



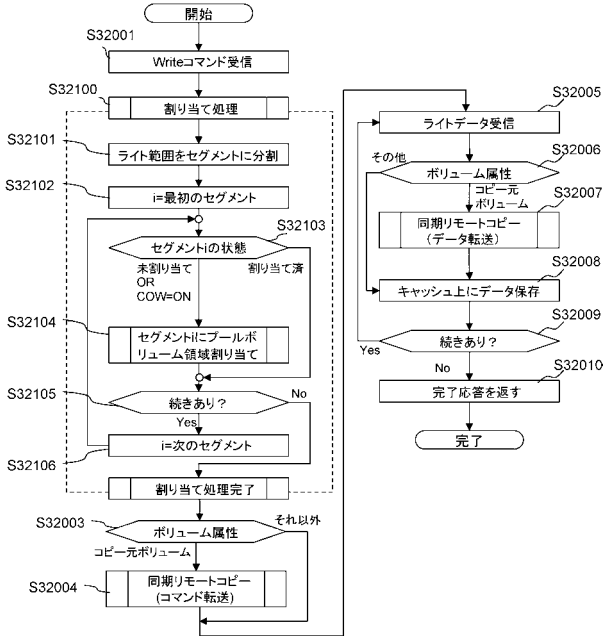
【図31】

図31



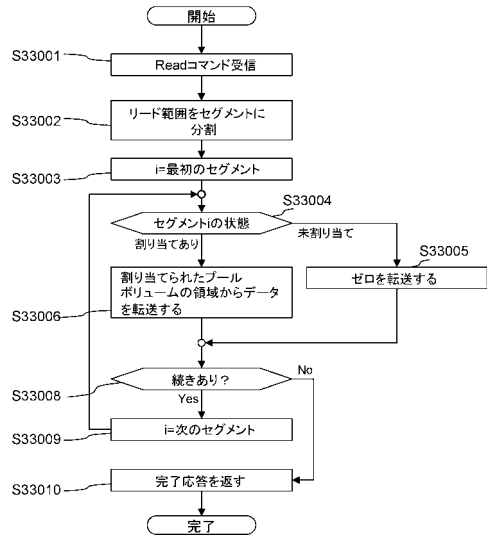
【図 3 2】

図32



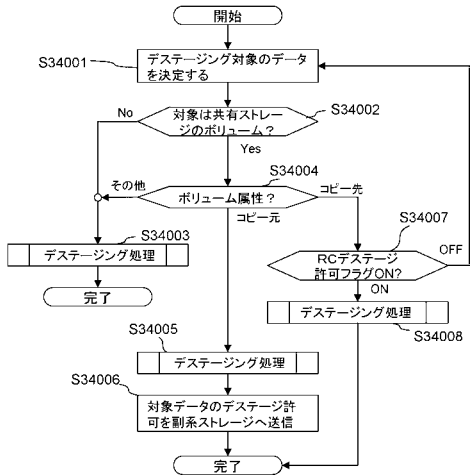
【図 3 3】

図33



【図 3 4】

図34



【図 3 5】

図35

プール ID	AOU ボリューム 識別子	対 AOU ボリューム 識別子	アドレス 空間	COW フラグ	引継ぎ 領域	プールボリューム 領域識別子
1	A-LDEV1	装置C A-LDEV5	0-1999	No	No	LDEV1(0-1999)
			2000-3999	No	Yes	Port=xxx. LUN=yyy (8000-9999)
		
1	A-LDEV2	装置C A-LDEV6	18000-19999	No	No	NULL
			0-1999	No	No	LDEV2(0-1999)
			2000-3999	Yes	Yes	Port=xxx. LUN=yyy (10000-11999)
...	
48000-49999	No	No	NULL			

【図36】

図36

31040

プールID	属性	値
1	セグメントサイズ	2000block
	プールボリュームリスト	LDEV1, LDEV2, Port=xxx&LUN=yyy, Port=zzz&LUN=aaa,....
	空き領域リスト	LDEV1(2000-40000), LDEV2(2000-10000), Port=xxx&LUN=yyy,....
	空き容量	3000000block
	容量警告	10000block
	プール対の識別子	装置L プールID=3
2	セグメントサイズ	1000block
	プールボリュームリスト	LDEV8, Port=aaa&LUN=bbb,....
	空き領域リスト	LDEV8(0000-40000), Port=aaa&LUN=bbb,....
	空き容量	1000000block
	容量警告	10000block
	プール対の識別子	装置L プールID=3

フロントページの続き

- (72)発明者 江口 賢哲
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
- (72)発明者 渡辺 恭男
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
- (72)発明者 本間 久雄
神奈川県小田原市中里322番2号 株式会社日立製作所RAIDシステム事業部内
- (72)発明者 山本 康友
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

審査官 古河 雅輝

- (56)参考文献 特開2006-048676(JP,A)
米国特許出願公開第2004/0260736(US,A1)
国際公開第2005/071544(WO,A1)
特開2005-267216(JP,A)
特開2006-024215(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08
G06F 12/00
G06F 13/10 - 13/14