

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2006-309640
(P2006-309640A)

(43) 公開日 平成18年11月9日(2006.11.9)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 13/10 (2006.01)	G06F 13/10 340A	5B005
G06F 3/06 (2006.01)	G06F 3/06 302A	5B014
G06F 12/08 (2006.01)	G06F 3/06 304B	5B065
G06F 13/12 (2006.01)	G06F 12/08 501E	
	G06F 12/08 507Z	
審査請求 未請求 請求項の数 20 O L (全 25 頁) 最終頁に続く		

(21) 出願番号	特願2005-133978 (P2005-133978)	(71) 出願人	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17年5月2日(2005.5.2)	(74) 代理人	100079108 弁理士 稲葉 良幸
		(74) 代理人	100093861 弁理士 大賀 真司
		(72) 発明者	渡辺 治明 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	山神 憲司 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
		Fターム(参考)	5B005 JJ01 MM11 PP21 5B014 EB04 GD13 GD23
		最終頁に続く	

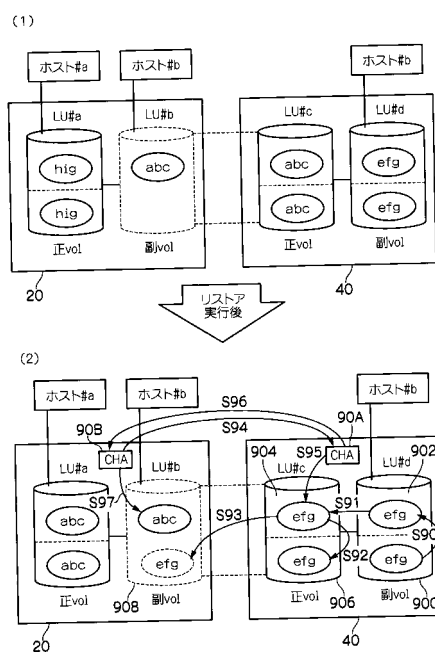
(54) 【発明の名称】 記憶制御システム及び記憶制御方法

(57) 【要約】

【課題】 複数の記憶制御装置を接続し、一方の記憶制御装置から他方の記憶制御装置にリストアデータを転送する記憶制御システムにおいて、上位装置がリストアデータの転送先の論理ボリュームにアクセスしても、リストアデータに確実にアクセスできるようにした記憶制御システムを提供する。

【解決手段】 第1の記憶制御装置と第2の記憶制御装置とが互いに通信可能に接続して構成され、上位装置からのリクエストに応じたデータ処理を行う記憶制御システムであって、前記第1の記憶制御装置は、前記第2の記憶制御装置の論理ボリュームに対応付けられた仮想ボリュームと、当該仮想ボリュームに対応付けられたキャッシュメモリと、当該キャッシュメモリと仮想ボリュームと間のデータ処理を制御する制御部とを備え、前記制御部は、前記論理ボリュームのリストアデータを前記キャッシュメモリに記憶する際に、前記上位装置がアクセスする前記仮想ボリュームに対応する、前記キャッシュメモリの記憶領域をページする。

【選択図】 図9



【特許請求の範囲】**【請求項 1】**

第 1 の記憶制御装置と第 2 の記憶制御装置とが互いに通信可能に接続して構成され、上位装置からのリクエストに応じたデータ処理を行う記憶制御システムであって、

前記第 1 の記憶制御装置は、

前記第 2 の記憶制御装置の論理ボリュームに対応付けられた仮想ボリュームと、

当該仮想ボリュームに対応付けられたキャッシュメモリと、

当該キャッシュメモリと仮想ボリュームと間のデータ処理を制御する制御部とを備え、

前記制御部は、前記論理ボリュームのリストアデータを前記キャッシュメモリに記憶する際に、前記上位装置がアクセスする前記仮想ボリュームに対応する、前記キャッシュメモリの記憶領域をパージするように構成された記憶制御システム。

10

【請求項 2】

第 1 の記憶制御装置と第 2 の記憶制御装置とが互いに通信可能に接続して構成され、上位装置からのリクエストに応じたデータ処理を行う記憶制御システムであって、

前記第 1 及び第 2 の記憶制御装置はそれぞれ、

記憶デバイスと、

当該記憶デバイスの論理的な記憶構造であり、前記上位装置からアクセス可能な論理ボリュームと、

前記論理ボリュームのデータを記憶するキャッシュメモリと、

前記記憶デバイス、前記キャッシュメモリ、そして前記論理ボリュームとの間でのデータ処理を制御する制御部と、

20

を備え、

さらに、前記第 1 の記憶制御装置は、前記第 1 の記憶制御装置の記憶デバイスに対応付けられることなく、前記第 2 の記憶制御装置の論理ボリュームに対応付けられ、かつ前記上位装置がアクセス可能な仮想ボリュームを備え、

前記第 1 の記憶制御装置の制御部は、前記第 2 の記憶制御装置の論理ボリュームのリストアデータを前記第 1 の記憶制御装置のキャッシュメモリに記憶する際に、前記上位装置がアクセスする前記仮想ボリュームの管理領域に対応する当該キャッシュメモリの記憶領域をパージするように構成されてなる、記憶制御システム。

【請求項 3】

30

前記第 2 の記憶制御装置の論理ボリュームは、コピーペアの関係を形成可能な第 1 の論理ボリュームと第 2 の論理ボリュームとを備えてなり、当該第 1 の論理ボリュームが前記仮想ボリュームに対応付けられている、請求項 2 記載の記憶制御システム。

【請求項 4】

前記第 2 の記憶制御装置の制御部は、前記リストアの命令を受領すると、前記第 2 の論理ボリュームの記憶データを前記第 1 の論理ボリュームにコピーし、さらに、当該第 1 の論理ボリュームのコピーされたデータを前記仮想ボリュームに転送する、請求項 3 記載の記憶制御システム。

【請求項 5】

前記第 1 の記憶制御装置の論理ボリュームは前記仮想ボリュームとコピーペアの関係になっている請求項 2 記載の記憶制御システム。

40

【請求項 6】

前記第 1 の記憶制御装置の制御部が前記第 2 の記憶制御装置の制御部に、前記リストアの命令を発行する請求項 4 記載の記憶制御システム。

【請求項 7】

前記上位装置が前記第 2 の記憶制御装置の制御部に、前記リストア命令を発行する請求項 4 記載の記憶制御システム。

【請求項 8】

前記第 1 の記憶制御装置の制御部による前記パージが完了後、前記第 2 の記憶制御装置の制御部が前記リストアを開始する請求項 1 又は 2 記載の記憶制御システム。

50

【請求項 9】

前記仮想ボリュームは、前記第 1 の記憶制御装置のキャッシュメモリの記憶領域との対応関係を制御する情報が設定されたディレクトリ構造を備えており、このディレクトリ構造は、複数の管理領域を備え、各管理領域には、当該管理領域に対応する前記第 1 の記憶制御装置のキャッシュメモリの記憶領域を特定するための情報が設定されている、請求項 1 又は 2 記載の記憶制御システム。

【請求項 10】

前記各管理領域は、前記キャッシュメモリ記憶領域の記憶データの有効又は無効を制御する情報を備えてなる請求項 9 記載の記憶制御システム。

【請求項 11】

前記第 1 の記憶制御装置の制御部は、
前記仮想ボリュームの全管理領域の前記特定情報をパージし、
次いで、前記第 1 の記憶制御装置のキャッシュメモリに前記第 2 の記憶制御装置の論理ボリュームに保存されているデータを割り当て、
前記上位装置が当該データにアクセスできるように、
構成されてなる請求項 1 又は 2 記載の記憶制御システム。

10

【請求項 12】

前記第 1 の記憶制御装置の制御部は、
前記仮想ボリュームの一部の管理領域の前記特定情報をパージし、
次いで、前記第 1 の記憶制御装置のキャッシュメモリに前記第 2 の記憶制御装置の論理ボリュームに保存されているデータを割り当て、
前記上位装置が当該データにアクセスできるように、
構成されてなる請求項 1 又は 2 記載の記憶制御システム。

20

【請求項 13】

前記第 1 の記憶制御装置の制御部は、前記仮想ボリュームに全体世代番号を、当該仮想ボリュームの各管理領域に個別の世代番号を設定し、前記上位装置がアクセスした前記仮想ボリュームの管理領域の世代番号と前記全体番号を比較し、両者が一致している場合には、当該管理領域に対応するキャッシュメモリ記憶領域のデータを前記リストアデータと判定し、両者が一致していない場合には、前記管理領域に対応するキャッシュメモリの記憶領域のデータをパージし、次いで、当該記憶領域に前記リストアデータを記憶するよう

30

【請求項 14】

前記第 1 の記憶制御装置の制御部は、前記仮想ボリュームの代表世代番号を、前記上位装置からの前記パージのコマンドを受ける都度、更新するようにした請求項 13 記載の記憶制御システム。

【請求項 15】

前記第 2 記憶制御装置の論理ボリュームに他の上位装置が接続され、当該上位装置からのリストアデータが当該論理ボリュームに記憶されている請求項 2 記載の記憶制御システム。

【請求項 16】

前記仮想ボリュームとコピーペアの関係を持つ他の論理ボリュームが前記第 1 の記憶制御装置に備わっている、請求項 3 記載の記憶制御システム。

40

【請求項 17】

前記パージ処理は、前記第 1 の記憶制御装置の制御部が、前記仮想ボリュームの管理領域に前記キャッシュメモリの記憶領域との対応関係が無いことを示す制御コードを設定するものである請求項 9 記載の記憶制御システム。

【請求項 18】

前記パージ処理は、前記第 1 の記憶制御装置の制御部が、前記管理領域に、前記キャッシュメモリ記憶領域の記憶データが無効であることを示す制御コードを設定するものである請求項 2 記載の記憶制御システム。

50

【請求項 19】

前記第1の記憶制御装置の制御部は、前記ページ処理された、前記キャッシュメモリの記憶領域に、前記リストアデータを記憶する請求項1又は2記載の記憶制御システム。

【請求項 20】

第1の記憶制御装置と第2の記憶制御装置とが互いに通信可能に接続された記憶制御システムに対して適用され、上位装置からのリクエストに応じたデータ処理を行う記憶制御方法であって、

前記第1の記憶制御装置は、

前記第2の記憶制御装置の論理ボリュームに対応付けられた仮想ボリュームと、

当該仮想ボリュームに対応付けられたキャッシュメモリと、

当該キャッシュメモリと仮想ボリュームと間のデータ処理を制御する制御部とを備え、

前記制御部は、前記論理ボリュームのリストアデータを前記キャッシュメモリに記憶する際に、前記上位装置がアクセスする前記仮想論理ボリュームに対応する前記キャッシュメモリの記憶領域をページするステップを実行する記憶制御方法。

10

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、ホストコンピュータが接続されて、当該ホストコンピュータとの間で記憶データの通信を行う記憶制御装置を備えてなる記憶制御システムに関するものである。

20

【背景技術】**【0002】**

例えば、データセンタ等のような大規模なデータを取り扱うデータベースシステムでは、ホストコンピュータとは別に構成された記憶制御システムを用いてデータを管理する。この記憶制御システムは、例えば、ディスクアレイ装置を備えて構成される。ディスクアレイ装置は、多数の記憶デバイスをアレイ状に配設して構成されるもので、例えば、RAID (Redundant Array of Independent Inexpensive Disks) に基づいて構築されている。記憶デバイス群が提供する物理的な記憶領域上には少なくとも1つ以上の論理的な記憶領域である論理ボリュームが形成され、この論理ボリュームが、上位装置としてのホストコンピュータ（より詳しくは、ホストコンピュータ上で稼働するデータベースなどのプログラム）に提供される。ホストコンピュータは、所定のコマンドをディスクアレイ装置に送信することにより、論理ボリュームに対してデータの書込み、読み出しができる。

30

【0003】

この種の記憶制御システムにおいて、処理の実行中に発生した障害などによって破壊されたデータをバックアップデータから高速に回復することが行われる。例えば、特開2001-216185号に記載されたディスクアレイ装置では、データを復旧させるために次のように構成されている。ホストは、ペア形成要求およびペア分割要求をディスクアレイ装置に発行し、ボリュームのスナップショットを作成した後、処理を開始する。ディスクサブシステムは、スナップショットの作成以降、更新が行われたデータの位置を差分情報として記憶する。ホストからのデータ復旧要求を受領すると、ディスクアレイ装置は、差分情報に基づいて、更新が行われた箇所のデータをスナップショットからリストアする。

40

【特許文献1】特開2001-216185

【発明の開示】**【発明が解決しようとする課題】****【0004】**

前記従来のディスクサブシステムは、一台の記憶制御装置内でのリストアに関するものであり、複数の記憶制御装置を接続し、一方の記憶制御装置から他方の記憶制御装置にデータをリストアする際のデータ復旧のことまでは配慮していない。複数の記憶制御装置間でのデータのリストアは、ホストのパフォーマンスを落とさないようにするために、ホス

50

トを經由することなく記憶制御装置間で行われることが望まれるが、リストアデータの送達を受けた記憶制御装置のキャッシュメモリ上にリストア前のデータが残っていると、リストア前のデータがキャッシュヒットしてホストに転送されてしまう。そのため、ホストが、他の記憶制御装置の記憶資源に割り当てられた論理ボリュームにアクセスすると、リストアされるべきデータではなく、リストア前のデータがホストに見えてしまう問題がある。

【0005】

そこで、本発明は、複数の記憶制御装置を接続し、一方の記憶制御装置から他方の記憶制御装置にリストアデータを転送する記憶制御システムにおいて、上位装置がリストアデータの転送先の論理ボリュームにアクセスしても、リストアデータに確実にアクセスできるようにした記憶制御システム及び記憶制御方法を提供することを目的とする。本発明の他の目的は、キャッシュメモリのページに要する時間を短縮し、ホストが早くリストアデータにアクセスできるようにした記憶制御システムを提供することにある。

10

【課題を解決するための手段】

【0006】

前記目的を達成するために、本発明は、記憶制御装置のキャッシュメモリに残っているリストア前のデータをページした後、キャッシュメモリの記憶領域にリストアデータを記憶させるように構成されている。

【0007】

すなわち、本発明は、第1の記憶制御装置と第2の記憶制御装置とが互いに通信可能に接続して構成され、上位装置からのリクエストに応じたデータ処理を行う記憶制御システムであって、前記第1の記憶制御装置は、前記第2の記憶制御装置の論理ボリュームに対応付けられた仮想ボリュームと、当該仮想ボリュームに対応付けられたキャッシュメモリと、当該キャッシュメモリと仮想ボリュームと間のデータ処理を制御する制御部とを備え、前記制御部は、前記論理ボリュームのリストアデータを前記キャッシュメモリに記憶する際に、前記上位装置がアクセスする前記仮想論理ボリュームに対応する、前記キャッシュメモリの記憶領域をページするように構成されたことを特徴とするものである。

20

【0008】

さらに本発明は、第1の記憶制御装置と第2の記憶制御装置とが互いに通信可能に接続して構成され、上位装置からのリクエストに応じたデータ処理を行う記憶制御システムであって、前記第1及び第2の記憶制御装置はそれぞれ、記憶デバイスと、当該記憶デバイスの論理的な記憶構造であり、前記上位装置からアクセス可能な論理ボリュームと、前記論理ボリュームのデータを記憶するキャッシュメモリと、前記記憶デバイス、前記キャッシュメモリ、そして前記論理ボリュームとの間でのデータ処理を制御する制御部と、を備え、さらに、前記第1の記憶制御装置は、前記第1の記憶制御装置の記憶デバイスに対応付けされることなく、前記第2の記憶制御装置の論理ボリュームに対応付けられ、かつ前記上位装置がアクセス可能な仮想ボリュームを備え、前記第1の記憶制御装置の制御部は、前記第2の記憶制御装置の論理ボリュームのリストアデータを前記第1の記憶制御装置のキャッシュメモリに記憶する際に、前記上位装置がアクセスする前記仮想ボリュームの管理領域に対応する当該キャッシュメモリの記憶領域をページするように構成されてなる、ことを特徴とするものである。

30

40

【0009】

本発明の好適な形態では、前記第2の記憶制御装置の論理ボリュームは、コピーペアの関係形成可能な第1の論理ボリュームと第2の論理ボリュームとを備えてなり、当該第1の論理ボリュームが前記仮想ボリュームに対応付けられている。第2の論理ボリュームから第1の論理ボリュームへデータのリストアが実施される。

【0010】

前記第2の記憶制御装置の制御部は、前記リストアの命令を受領すると、前記第2の論理ボリュームの記憶データを前記第1の論理ボリュームにコピーし、さらに、当該第1の論理ボリュームのコピーされたデータを前記仮想ボリュームに転送する。前記第1の記憶

50

制御装置の論理ボリュームは前記仮想ボリュームとコピーペアの関係になっている。前記第1の記憶制御装置の制御部が前記第2の記憶制御装置の制御部に、前記リストアの命令を発行する。前記上位装置が前記第2の記憶制御装置の制御部に、前記リストア命令を発行する。

【0011】

前記第1の記憶制御装置の制御部による前記ページが完了後、前記第2の記憶制御装置の制御部が前記リストアを開始する。前記仮想ボリュームは、前記第1の記憶制御装置のキャッシュメモリの記憶領域との対応関係を制御する情報が設定されたディレクトリ構造を備えており、このディレクトリ構造は、複数の管理領域を備え、各管理領域には、当該管理領域に対応する前記第1の記憶制御装置のキャッシュメモリの記憶領域を特定するための情報が設定されている。前記各管理領域は、前記キャッシュメモリ記憶領域の記憶データの有効又は無効を制御する情報を備えてなる。前記第1の記憶制御装置の制御部は、前記仮想ボリュームの全管理領域の前記特定情報をページし、次いで、前記第1の記憶制御装置のキャッシュメモリに前記第2の記憶制御装置の論理ボリュームに保存されているデータを割り当て、前記上位装置が当該データにアクセスできるように、構成されてなる。前記第1の記憶制御装置の制御部は、前記仮想ボリュームの一部の管理領域の前記特定情報をページし、次いで、前記第1の記憶制御装置のキャッシュメモリに前記第2の記憶制御装置の論理ボリュームに保存されているデータを割り当て、前記上位装置が当該データにアクセスできるように、構成されてなる。

10

【0012】

前記第1の記憶制御装置の制御部は、前記仮想ボリュームに全体世代番号を、当該仮想ボリュームの各管理領域に個別の世代番号を設定し、前記上位装置がアクセスした前記仮想ボリュームの管理領域の世代番号と前記全体番号を比較し、両者が一致している場合には、当該管理領域に対応するキャッシュメモリ記憶領域のデータを前記リストアデータと判定し、両者が一致していない場合には、前記管理領域に対応するキャッシュメモリの記憶領域のデータをページし、次いで、当該記憶領域に前記リストアデータを記憶するようにした。前記第1の記憶制御装置の制御部は、前記仮想ボリュームの代表世代番号を、前記上位装置からの前記ページのコマンドを受ける都度、更新するようにした。前記第2記憶制御装置の論理ボリュームに他の上位装置が接続され、当該上位装置からのリストアデータが当該論理ボリュームに記憶されている。前記仮想ボリュームとコピーペアの関係を

20

30

【0013】

前記ページ処理は、前記第1の記憶制御装置の制御部が、前記仮想ボリュームの管理領域に前記キャッシュメモリの記憶領域との対応関係が無いことを示す制御コードを設定するものである。前記ページ処理は、前記第1の記憶制御装置の制御部が、前記管理領域に、前記キャッシュメモリ記憶領域の記憶データが無効であることを示す制御コードを設定するものである。前記第1の記憶制御装置の制御部は、前記ページ処理された、前記キャッシュメモリの記憶領域に、前記リストアデータを記憶する。

【発明の効果】

【0014】

以上説明したように、本発明によれば、複数の記憶制御装置を接続し、一方の記憶制御装置から他方の記憶制御装置にリストアデータを転送する記憶制御システムにおいて、上位装置がリストアデータの転送先の論理ボリュームにアクセスしても、リストアデータに確実にアクセスできるようになった。

40

【発明を実施するための最良の形態】

【0015】

本発明の実施形態を以下に説明する。以下に説明する代表的な記憶制御システムは、第1の記憶制御装置に仮想的に論理記憶領域を設定し、この仮想領域に第1の記憶制御装置の外部に存在する第2の記憶制御装置の論理記憶領域を対応付けて(マッピングして)、第1の記憶制御装置が第2の記憶制御装置の記憶領域をあたかも自分の記憶領域であるかの

50

ようにしてこれを上位装置に提供している。

【0016】

図1は、記憶制御システムの要部の構成を示すブロック図である。上位装置としてのホスト装置10は、例えば、CPU(Central Processing Unit)やメモリ等の情報処理資源を備えたコンピュータ装置であり、例えば、パーソナルコンピュータ、ワークステーション、メインフレーム等として構成される。ホスト装置10は、例えば、キーボード、スイッチやポインティングデバイス、マイクロフォン等の情報入力装置(図示せず)と、例えば、モニタディスプレイやスピーカー等の情報出力装置(図示せず)とを備えている。さらに、ホスト装置10には、例えば、第1の記憶制御装置20が提供する記憶領域を使用するデータベースソフトウェア等のアプリケーションプログラム11と、通信ネットワークCN1を介して第1の記憶制御装置20にアクセスするためのアダプタ12とが設けられている。

10

【0017】

ホスト装置10は、スイッチ回路SWを備えて構成される通信ネットワークCN1を介して第1の記憶制御装置20に接続されている。通信ネットワークCN1としては、例えば、LAN、SAN、インターネット、専用回線、公衆回線等の場合に応じて適宜用いることができる。LANを介するデータ通信は、例えば、TCP/IP(Transmission Control Protocol/Internet Protocol)プロトコルに従って行われる。ホスト装置10がLANを介して第1の記憶制御装置20に接続される場合、ホスト装置10は、ファイル名を指定してファイル単位でのデータ入出力を要求する。一方、ホスト装置10がSANを介して第1の記憶制御装置20等に接続される場合、ホスト装置10は、ファイバチャネルプロトコルに従って、複数のディスク記憶装置(ディスクドライブ)により提供される記憶領域のデータ管理単位であるブロックを単位としてデータ入出力を要求する。通信ネットワークCN1がLANである場合、アダプタ12は、例えばLAN対応のネットワークカードである。通信ネットワークCN1がSANの場合、アダプタ12は、例えばホストバスアダプタである。

20

【0018】

スイッチ回路SWは通信ネットワークに接続されたルータや交換機を備えて構成されている。スイッチ回路SWは、第1の記憶制御装置20のエクスターナルポート21Aに、第2の記憶制御装置40のターゲットポート41及び第1の記憶制御装置のターゲットポート21Bを切り換えて接続できるように構成されている。なお、第1の記憶制御装置20は、エクスターナルポート21Aとスイッチ回路SWを介して自筐体のターゲットポート21Bに接続できるように構成されているために、所謂自己ループ式に分類されるものである。各ポート及びスイッチ回路にはネットワーク上のアドレスが設定されている。

30

【0019】

第1の記憶制御装置20は、例えば、ディスクアレイサブシステムとして構成されるものである。但し、これに限らず、第1の記憶制御装置20を、高機能化されたインテリジェント型のファイバチャネルスイッチとして構成することもできる。第1の記憶制御装置20は、後述のように、第2の記憶制御装置40の有する記憶資源を自己の論理ボリューム(Logical Unit)としてホスト装置10に提供する。

40

【0020】

第1の記憶制御装置20は、コントローラ部と記憶装置部とに大別することができ、コントローラ部は、例えば、複数のチャネルアダプタ(CHA)21と、複数のディスクアダプタ(DKA)22と、サービスプロセッサ(SVP)23と、キャッシュメモリ24と、共有メモリ25と、接続部26とを備えている。このコントローラ部が請求項記載の制御部に対応する。

【0021】

チャネルアダプタ(CHA)21は、ホスト装置10との間のデータ通信を行うものである。他のチャネルアダプタ21は、第1の記憶制御装置のエクスターナルポート21A、第2の記憶制御装置のターゲットポート41を介して、第2の記憶制御装置40の内部

50

論理ボリュームとデータ通信を行うものである。各チャンネルアダプタ21は、それぞれマイクロプロセッサやメモリ等を備えたマイクロコンピュータシステムとして構成されており、ホスト装置10から受信した各種コマンドを解釈して実行する。各チャンネルアダプタ21には、それぞれを識別するためのネットワークアドレス(例えば、IPアドレスやWWN)が割り当てられているので、各チャンネルアダプタ21は、それぞれが個別にNAS(Network Attached Storage)として振る舞うことができるようになっている。複数のホスト装置10が存在する場合、各チャンネルアダプタ21はホスト装置毎に設けられ、各ホスト装置10からの要求をそれぞれ個別に受け付けることができるように構成されている。

【0022】

各ディスクアダプタ(DKA)22は、記憶装置30の記憶デバイス31,32との間のデータ授受を行うものである。各ディスクアダプタ22は、記憶デバイス31,32に接続するための通信ポート22Aを備えている。また、各ディスクアダプタ22は、マイクロプロセッサやメモリ等を備えたマイクロコンピュータシステムとして構成されている。各ディスクアダプタ22は、チャンネルアダプタ21がホスト装置10から受信したデータを、ホスト装置10からのリクエストに基づいて(書込み命令)、所定の記憶デバイス31,32の所定のアドレスに書込み、また、ホスト装置10からのリクエストに基づいて(読み出し命令)、所定の記憶デバイス31,32の所定のアドレスからデータを読み出し、ホスト装置10に送信させる。記憶デバイス31,32との間でデータ入出力を行う場合、各ディスクアダプタ22は、論理的なアドレスを物理的なアドレスに変換する。各ディスクアダプタ22は、記憶デバイス31,32がRAIDに従って管理されている場合は、RAID構成に応じたデータアクセスを行う。

【0023】

サービスプロセッサ(SVP)23は、装置全体の作動を制御するものである。SVP23には、例えば、管理用クライアント(図示せず)が接続される。SVP23は、装置内の障害発生を監視して管理用クライアントに表示させたり、管理用クライアントからの指令に基づいて記憶ディスクの閉塞処理等を指示するようになっている。さらに、SVP23の管理用クライアントは後述のように仮想ボリュームを定義する処理を実行する。管理用クライアントは、例えばJAV A(登録商標)アプレット上に管理用プログラムが実装して構成されている。

【0024】

キャッシュメモリ24は、ホスト装置10から受信したデータや、記憶デバイス31,32から読み出したデータを一時的に記憶するものである。共有メモリ25には、第1の記憶制御装置の作動に使用するための制御情報等が格納される。また、共有メモリ25には、ワーク領域が設定されるほか、後述するマッピングテーブルTm等の各種テーブル類も格納される。また、キャッシュメモリ130と共有メモリ140とは、それぞれ別々のメモリとして構成することもできるし、同一のメモリの一部の記憶領域をキャッシュ領域として使用し、他の記憶領域を制御領域として使用することもできる。なお、記憶デバイス31,32のいずれか1つあるいは複数、キャッシュ用のディスクとして使用してもよい。

【0025】

接続部26は、各チャンネルアダプタ21,各ディスクアダプタ22,SVP23,キャッシュメモリ24,共有メモリ25を相互に接続させる。接続部26は、例えば、高速スイッチング動作によってデータ伝送を行う超高速クロスバスイッチ等のような高速バスとして構成することができる。

【0026】

記憶装置30は、複数の記憶デバイス31を備えている。記憶デバイス31としては、例えば、ハードディスク、フレキシブルディスク、磁気テープ、半導体メモリ、光ディスク等のようなデバイスを用いることができる。また、例えば、FC(Fibre Channel)ディスクやSATA(Serial AT Attachment)ディスク等のように、異種類のディスクを記憶装置

10

20

30

40

50

30内に混在させることもできる。記憶装置30内に点線で示される記憶デバイス32は、第2の記憶制御装置40の有する記憶デバイス42を第1の記憶制御装置20側に取り込んだ状態を示すものである。即ち、第1の記憶制御装置20から見て外部に存在する記憶デバイス42を、第1の記憶制御装置20の内部記憶デバイスかのように認識できるようにし、ホスト装置10に外部記憶デバイス42の記憶資源を提供する。このことは後述のように、第1の記憶制御装置20内の中間論理記憶領域である仮想ボリュームに第2の記憶制御装置40の論理ボリュームをマッピングすることにより可能になる。キャッシュメモリ24の記憶空間を利用して、仮想ボリュームが構築される。第1の記憶制御装置20に形成された仮想ボリュームは、第1の記憶制御装置20内の実論理ボリュームとともにホスト装置10に認識される。なお、図1に点線で示される外部の記憶デバイス32とのデータ交換は第2の記憶制御装置40のディスクアダプタによって実行される。 10

【0027】

第2の記憶制御装置40は、通信ポート(ターゲットポート)41と記憶デバイス42とを備えている。このほか、チャンネルアダプタやディスクアダプタ等を備えることもできる。第2の記憶制御装置40は、スイッチ回路SWを介して第1の記憶制御装置20に接続されており、第2の記憶制御装置40の記憶デバイス42は、第1の記憶制御装置20の内部記憶デバイスとして扱われるようになっている。スイッチ回路SWには複数の外部記憶制御装置40Aを接続可能である。

【0028】

続いて図2について説明する。図2は、第1の記憶制御装置20及び第2の記憶制御装置40の1つの記憶構造の概略を示すブロック図である。まず、第1の記憶制御装置の構成から先に説明する。 20

【0029】

第1の記憶制御装置の記憶構造は、例えば、物理的記憶階層と論理的記憶階層とに大別することができる。物理的記憶階層は、物理的なディスクであるPDEV(Physical Device)161により構成される。PDEVは、ディスクドライブに該当する。

【0030】

論理的記憶階層は、複数の(例えば2種類の)階層から構成することができる。一つの論理的階層は、VDEV(Virtual Device)162と、VDEV162のように扱われる仮想的なVDEV(以下、「V-VOL」とも呼ぶ)163とから構成可能である。他の一つの論理的階層は、LDEV(Logical Device)164から構成することができる。 30

【0031】

VDEV162は、例えば、4個1組(3D+1P)、8個1組(7D+1P)等のような所定数のPDEV161をグループ化して構成される。グループに属する各PDEV161がそれぞれ提供する記憶領域が集合して一つのRAID記憶領域が形成されている。このRAID記憶領域がVDEV162となる。

【0032】

VDEV162が物理的な記憶領域上に構築されるのと対照的に、V-VOL163は、物理的な記憶領域を必要としない仮想的な中間記憶デバイスである。V-VOL163は、物理的な記憶領域に直接関係づけられるものではなく、第2の記憶制御装置41のLU(Logical Unit)250をマッピングするための受け皿となる仮想的な存在である。 40

【0033】

LDEV164は、VDEV162またはV-VOL163上に、それぞれ少なくとも一つ以上設けることができる。LDEV164は、例えば、VDEV162を固定長で分割することにより構成することができる。ホスト10がオープン系ホストの場合、LDEV164がLU165にマッピングされることにより、ホスト10は、LDEV164を一つの物理的なディスクとして認識する。オープン系のホストは、LUN(Logical Unit Number)や論理ブロックアドレスを指定することにより、所望のLDEV164にアクセスする。なお、メインフレーム系ホストの場合は、LDEV164を直接認識する。

【0034】

LU165は、SCSIの論理ユニットとして認識可能なデバイスである。各LU165は、ターゲットポート21Bを介してホスト10に接続される。各LU165には、少なくとも一つ以上のLDEV164をそれぞれ関連付けることができる。一つのLU165に複数のLDEV164を関連付けることにより、LUサイズを仮想的に拡張することもできる。

【0035】

CMD (Command Device) 166は、ホスト10上で稼働するI/O制御プログラムと記憶制御装置20のコントローラ (CHA21, DKA22) との間で (図1参照)、コマンドやステータスを受け渡すために使用される専用のLUである。ホスト10からのコマンドは、CMD166に書き込まれる。記憶制御装置20のコントローラは、CMD166に書き込まれたコマンドに応じた処理を実行し、その実行結果をステータスとしてCMD166に書き込む。ホスト装置10は、CMD166に書き込まれたステータスを讀出して確認し、次に実行すべき処理内容をCMD166に書き込む。このようにして、ホスト装置10は、CMD166を介して、記憶制御装置20に各種の指示を与えることができる。

10

【0036】

なお、ホスト装置10から受信したコマンドを、CMD166に格納することなく、処理することもできる。また、実体のデバイス (LU) を定義せずに、CMDを仮想的なデバイスとして生成し、ホスト装置10からのコマンドを受け付けて処理するように構成してもよい。即ち、例えば、CHA21は、ホスト装置10から受信したコマンドを共有メモリ25に書き込み、この共有メモリ25に記憶されたコマンドを、CHA21又はDKA22が処理する。その処理結果は共有メモリ25に書き込まれ、CHA21からホスト装置10に送信される。

20

【0037】

さて、第1の記憶制御装置20の外部接続用のエクスターナルポート (External Port) 21Aには、既述の通りスイッチ回路SWを介して第2ストレージ装置40のターゲットポート41又は第1の記憶制御装置20のターゲットポート21Bが接続可能である。

【0038】

第2の記憶制御装置40は、複数のPDEV220と、PDEV220の提供する記憶領域上に設定されたVDEV230と、VDEV230上に少なくとも一つ以上設定可能なLDEV240とを備えている。各LDEV240は、LU250にそれぞれ関連付けられている。

【0039】

そして、本実施形態では、第2の記憶制御装置40のLU250 (即ち、LDEV240) は、仮想的な中間記憶デバイスであるV-VOL163にマッピングされており、第1の記憶制御装置20の内部論理ボリュームとして取り扱われる。

30

【0040】

例えば、第2の記憶制御装置40の「LDEV1」、「LDEV2」は、「LU1」、「LU2」を介して、第1の記憶制御装置20の「V-VOL1」、「V-VOL2」にそれぞれマッピングされている。そして、「V-VOL1」、「V-VOL2」は、それぞれ「LDEV3」、「LDEV4」にマッピングされ、「LU3」、「LU4」を介して、ホスト装置10が第2の記憶制御装置の記憶領域にアクセス可能になっている。

【0041】

なお、VDEV162、V-VOL163は、RAID構成を適用することができる。即ち、一つのディスクドライブ161を複数のVDEV162、V-VOL163に割り当てることもできるし (スライシング)、複数のディスクドライブ161から一つのVDEV162、V-VOL163を形成することもできる (ストライピング)。

40

【0042】

そして、第1の記憶制御装置20の「LDEV1」または「LDEV2」が第1の記憶制御装置20に記憶領域を持つ、実論理ボリューム190に該当する。第1の記憶制御装置20の「LDEV3」または「LDEV4」が、第1の記憶制御装置に記憶領域を持たず、第2の記憶制御装置に記憶領域を持つ仮想ボリューム191に該当する。第2の記憶制御装置40の「LDEV1」または「LDEV2」が、仮想ボリュームに対応する (仮想ボリュームがマッピング

50

された)論理ボリューム260に該当する。このように、実論理ボリューム190は、第1の記憶制御装置1内に設けられた物理的な記憶デバイス(例えば、ディスクドライブ等)に基づいて形成されるものである。仮想ボリューム191は、仮想的な存在であり、データを記憶する実体は、第2の記憶制御装置40内に存在する。即ち、第1の記憶制御装置20が有する記憶階層の所定の層に、第2の記憶制御装置40の有する論理ボリューム260がマッピングされるように、仮想ボリューム191が構築されている。

【0043】

図2によれば、第1の記憶制御装置20のCHA21は、エクスターナルポート21Aから参照できるボリュームを外部デバイスとして認識する。この外部記憶デバイスにマッピングされる既述の仮想ボリュームが第1の記憶制御装置20内に定義される。この定義動作は、既述のSVP23の管理クライアントである、例えばJAVA(登録商標)アプレットによって実行される。この定義情報はマッピングテーブルTmとして共有メモリ25上に置かれる。但し、図2の符号300に示すように、第1の記憶制御装置の論理ボリューム190に仮想ボリュームをマッピングすることはできない。

10

【0044】

次に仮想ボリュームを第1の記憶制御装置20に設定するための動作について説明する。図3はこの設定動作のフローを示す。管理クライアント上のJAVA(登録商標)アプレットは、SVPにエクスターナルポート21A(図1)から参照できる論理ボリューム情報取得(ディスカバリー)をリクエストする(S1)。このリクエストを受けたSVPは、第1の記憶制御装置のCHA21(図1)に、ディスカバリーの実行を指令する(S2)。CHAは共有メモリ25を参照してスイッチ回路SWのIPアドレスを入手してエクスターナルポート21Aを介してスイッチ回路SWにアクセスし、このスイッチ回路から参照できる全てのボリュームの情報を取得する。すなわち、CHAは自己の筐体の内部論理ボリュームか他の記憶制御装置の内部論理ボリュームかを区別することなく、エクスターナルポート21Aに接続可能な内部論理ボリュームのディスカバリーを実行する(S3)。図1に示すように、CHA21がエクスターナルポート21Aに接続されるスイッチ回路SWを介して参照できるボリュームは、第1の記憶制御装置20の論理ボリュームと第2の記憶制御装置40の論理ボリュームである。

20

【0045】

次いで、第1の記憶制御装置20のCHA21及び第2の記憶制御装置40のCHA(図示せず)は、それぞれの共有メモリを参照して、それぞれの内部論理ボリュームの特性データを取得する(S4)。この特性データを受けたCHA21はSVP32を介して管理クライアントのJAVA(登録商標)アプレットに特性データを送信する(S5、S6)。JAVA(登録商標)アプレット上には、特性データ中の「ストレージのベンダー名、製品型名、製造番号」等の情報からディスカバリーによって参照された複数の内部論理ボリュームのそれぞれについて、抽出された内部論理ボリュームが属する筐体が、CHA21が属する自己の記憶制御装置か或いは外部記憶制御装置かを判断するストレージ管理用アプリケーションが実装されている。

30

【0046】

この管理用アプリケーションは、ディスカバリーによって抽出された複数の内部論理ボリュームについて、既述の仮想ボリュームに対する候補として良いか否かのフィルタリングを行うモジュールを備えている。管理用アプリケーションは、このフィルタリングモジュールによって複数の論理ボリュームについての各抽出情報をフィルタリング(S7)し、合致した条件の論理ボリュームについて仮想ボリュームのマッピング先の候補群として画面表示処理が行われる。

40

【0047】

本実施形態では、既述のとおり、第1の記憶制御装置20の内部論理ボリュームが、仮想ボリュームがマッピングされるボリュームとしての候補群として表示されないように管理用アプリケーションがプログラムされている。次いで、管理用クライアントはSVPに仮想ボリュームの定義処理の実行をリクエストする(S8、S9)。SVPはCHAに仮

50

想ボリュームの定義処理、すなわち、既述のマッピングテーブル作成の処理動作を実行させ（S 1 0）、その処理結果を受けてこれを管理用クライアントに送る（S 1 1, S 1 2）。管理用クライアントのユーザは自筐体（第 1 の記憶制御装置 2 0）の内部論理ボリュームが自筐体内の仮想ボリュームのマッピング先として表示されないの、自筐体の内部論理ボリュームに対して仮想ボリュームを設定することができない。このように、S V P 及び管理用クライアントによって仮想ボリューム設定制御が実現される。この設定制御処理は C H A 或いは D K A のプロセッサによっても実現可能である。

【 0 0 4 8 】

次に、仮想ボリュームの定義処理について詳細を説明する。図 4 は、既述のマッピングテーブル T m の構造の一例である。マッピングテーブル T m は、例えば、V D E V をそれぞれ識別するための V D E V 番号と外部の記憶デバイスの情報とをそれぞれ対応付けることにより、構成ができる。外部デバイス情報としては、例えば、デバイス識別情報と、記憶デバイスの記憶容量と、デバイスの種別を示す情報と（例えば、テープ系デバイスかディスク系デバイスか等）、記憶デバイスへのパス情報とを含んで構成することができる。また、パス情報は、各通信ポート（図 1 の 2 1 A, 2 1 B）に固有の識別情報（W W N）と、図 2 の L U を識別するための L U 番号とを含んで構成できる。

【 0 0 4 9 】

なお、図 4 中に示すデバイス識別情報や W W N 等は、説明の便宜上の値であって特に意味はない。また、図 4 中の下側に示す V D E V 番号「3」の V D E V は、3 個のパス情報が対応付けられている。即ち、この V D E V（# 3）にマッピングされる外部記憶デバイス（図 1 の 4 2）は、その内部に 3 つの経路を有する交代パス構造を備えているが、V D E V には、この交代パス構造を認識してマッピングされている。これら 3 つの経路のいずれを通っても同一の記憶領域にアクセスできることが判明しているため、いずれか 1 つまたは 2 つの経路に障害等が発生した場合でも、残りの正常な経路を介して所望のデータにアクセスできる。図 4 に示すようなマッピングテーブル T m を採用することにより、第 1 の記憶制御装置 2 0 内の 1 つ以上の V D E V に対し、1 つまたは複数の外部の記憶デバイス 4 2 をマッピングすることができる。

【 0 0 5 0 】

次に、図 5 を参照して、外部の記憶デバイス 4 2 を V D E V 1 0 1 にマッピングする方法の一例を説明する。図 5 は、マッピング時に、第 1 の記憶制御装置 2 0 と第 2 の記憶制御装置 4 0 との間で行われる処理の要部を示すタイムチャートである。まず、第 1 の記憶制御装置 2 0 は、チャンネルアダプタ 2 1 のエクスターナルポート（2 1 A）からスイッチ回路 S W を介して、第 2 の記憶制御装置 4 0 にログインする（S 1）。第 2 の記憶制御装置 4 0 が、第 1 の記憶制御装置 2 0 のログインに対して応答を返すことにより、ログインが完了する（S 2）。次に、第 1 の記憶制御装置 2 0 は、例えば、SCSI（Small Computer System Interface）規格で定められている照会コマンド（inquiryコマンド）を、第 2 の記憶制御装置 4 0 に送信し、第 2 の記憶制御装置 4 0 の有する記憶デバイス 4 2 の詳細について応答を求める（S 3）。これは既述のディスカバーの動作と同じである。

【 0 0 5 1 】

照会コマンドは、照会先の装置の種類及び構成を明らかにするために用いられるもので、照会先装置の有する階層を透過してその物理的構造を把握することができる。照会コマンドを使用することにより、第 1 の記憶制御装置 2 0 は、例えば、装置名、デバイスタイプ、製造番号（プロダクト I D）、L D E V 番号、各種バージョン情報、ベンダ I D 等の情報を第 2 の記憶制御装置 4 0 から取得できる（S 4）。第 2 の記憶制御装置 4 0 は、問合わせられた情報を第 1 の記憶制御装置 2 0 に送信し、応答する（S 5）。ディスカバリーで抽出された情報は共有メモリに保存されている。C H A はこの保存された情報を利用することによって S 3 ~ S 5 のステップを省略しても良い。

【 0 0 5 2 】

第 1 の記憶制御装置 2 0 は、第 2 の記憶制御装置 4 0 から取得した情報を、マッピングテーブル T m の所定箇所に登録する（S 6）。次に、第 1 の記憶制御装置 2 0 は、第 2 の

10

20

30

40

50

記憶制御装置 40 から記憶デバイス 42 の記憶容量を読み出す (S7)。第 2 の記憶制御装置 40 は、第 1 の記憶制御装置 20 からの問合せに対して、記憶デバイス 42 の記憶容量を返信し (S8)、応答を返す (S9)。第 1 の記憶制御装置 20 は、記憶デバイス 42 の記憶容量をマッピングテーブル Tm の所定箇所に登録する (S10)。

【0053】

以上の処理を行うことにより、マッピングテーブル Tm を構築できる。第 1 の記憶制御装置 20 の VDEV にマッピングされた外部の記憶デバイス 42 (外部 LUN、即ち外部の LDEV) との間でデータの入出力を行う場合は、後述する他のテーブルを参照してアドレス変換等を行う。

【0054】

図 6 ~ 図 8 を参照して、第 1 の記憶制御装置 20 と第 2 の記憶制御装置 40 との間のデータ入出力について説明する。まず最初に、データを書き込む場合について、図 6 及び図 7 に基づいて説明する。図 6 は、データ書き込み時の処理を示す模式図である。図 7 は、図 6 中の処理の流れを各種テーブルとの関係で示す説明図である。

【0055】

ホスト装置 10 は、第 1 の記憶制御装置 20 が提供する論理ボリューム (LDEV102) にデータを書き込むことができる。例えば、SAN の中に仮想的な SAN サブネットを設定するゾーニングや、アクセス可能な LUN のリストをホスト装置 10 が保持する LUN マスキングという手法により、ホスト装置 10 を特定の LDEV102 に対してのみアクセスさせるように設定できる。

【0056】

ホスト装置 10 がデータを書き込もうとする LDEV102 が、VDEV101 を介して内部の記憶デバイスである記憶デバイス 31 に接続されている場合、通常の処理によってデータが書き込まれる。即ち、ホスト装置 10 からのデータは、いったんキャッシュメモリ 24 に格納され、キャッシュメモリ 24 からディスクアダプタ 22 を介して、所定の記憶デバイス 31 の所定アドレスに格納される。この際、ディスクアダプタ 22 は、論理的なアドレスを物理的なアドレスに変換する。また、RAID 構成の場合、同一のデータが複数の記憶デバイス 31 に記憶等される。

【0057】

これに対し、ホスト装置 10 が書き込もうとする LDEV102 が、VDEV101 を介して外部の記憶デバイス 42 に接続されている場合、図 6 に示すような流れでデータが書き込まれる。図 6 (a) は記憶階層を中心に示す流れ図であり、図 6 (b) はキャッシュメモリ 24 の使い方を中心に示す流れ図である。

【0058】

ホスト装置 10 は、書き込み先の LDEV102 を特定する LDEV 番号とこの LDEV102 にアクセスするための通信ポート 21A を特定する WWN とを明示して、書き込みコマンド (Write) を発行する (S21)。第 1 の記憶制御装置 20 は、ホスト装置 10 からの書き込みコマンドを受信すると、第 2 の記憶制御装置 40 に送信するための書き込みコマンドを生成し、第 2 の記憶制御装置 40 に送信する (S22)。第 1 の記憶制御装置 20 は、ホスト装置 10 から受信した書き込みコマンド中の書き込み先アドレス情報等を、外部 LDEV43 に合わせて変更することにより、新たな書き込みコマンドを生成する。

【0059】

次に、ホスト装置 10 は、書き込むべきデータを第 1 の記憶制御装置 40 に送信する (S23)。第 1 の記憶制御装置 20 に受信されたデータは、LDEV102 から VDEV101 を介して (S24)、外部の LDEV43 に転送される (S26)。ここで、第 1 の記憶制御装置 20 は、ホスト装置 10 からのデータをキャッシュメモリ 24 に格納した時点で、ホスト装置 10 に対し書き込み完了の応答 (Good) を返す (S25)。第 2 の記憶制御装置 40 は、第 1 の記憶制御装置 20 からデータを受信した時点で (あるいは記憶デバイス 42 に書き込みを終えた時点で)、書き込み完了報告を第 1 の記憶制御装置 20 に送信する (S26)。即ち、第 1 の記憶制御装置 20 がホスト装置 10 に対して書き込み完了を

10

20

30

40

50

報告する時期（S25）と、実際にデータが記憶デバイス42に記憶される時期とは相違する（非同期方式）。従って、ホスト装置10は、実際にデータが記憶デバイス42に格納される前にデータ書込み処理から解放され、別の処理を行うことができる。

【0060】

図6（b）のように、キャッシュメモリ24には、多数のサブブロック24Aが設けられている。第1の記憶制御装置20は、ホスト装置10から指定された論理ブロックアドレスをサブブロックのアドレスに変換し、キャッシュメモリ24の所定箇所にデータを格納する（S24）。

【0061】

図7を参照して、各種テーブルを利用してデータが変換される様子を説明する。図7の上部に示すように、ホスト装置10は、所定の通信ポート21Aに対し、LUN番号（LUN#）及び論理ブロックアドレス（LBA）を指定してデータを送信する。第1の記憶制御装置20は、LDEV102用に入力されたデータ（LUN#+LBA）を、図7（a）に示す第1の変換テーブルT1に基づいて、VDEV101用のデータに変換する。第1の変換テーブルT1は、内部のLUN103を指定するデータをVDEV101用データに変換するための、LUN-LDEV-VDEV変換テーブルである。このテーブルT1は、例えば、LUN番号（LUN#）と、そのLUN103に対応するLDEV102の番号（LDEV#）及び最大スロット数と、LDEV102に対応するVDEV101の番号（VDEV#）及び最大スロット数等に対応付けることにより構成される。このテーブルT1を参照することにより、ホスト装置10からのデータ（LUN#+LBA）は、VDEV101用のデータ（VDEV#+SLOT#+SUBBLOCK#）に変換される。

【0062】

次に、第1の記憶制御装置20は、図7（b）に示す第2の変換テーブルT2を参照して、VDEV101用のデータを、第2の記憶制御装置40の外部LUN（LDEV）用に送信して記憶させるためのデータに変換する。第2の変換テーブルT2には、例えば、VDEV101の番号（VDEV#）と、そのVDEV101からのデータを第2の記憶制御装置40に送信するためのイニシエータポートの番号と、データ転送先の通信ポート41を特定するためのWWNと、その通信ポートを介してアクセス可能なLUN番号とが対応付けられている。この第2の変換テーブルT2に基づいて、第1の記憶制御装置20は、記憶させるべきデータの宛先情報を、イニシエータポート（ターゲットポート）番号#+WWN+LUN#+LBAの形式に変換する。このように宛先情報が変更されたデータは、指定されたイニシエータポートから通信ネットワークCN1を介して、指定された通信ポート41に到達する。そして、データは、指定されたLUN43でアクセス可能なLDEVの所定の場所に格納される。LDEVは、複数の記憶デバイス42上に仮想的に構築されているので、データのアドレスは物理アドレスに変換されて、所定のディスクの所定アドレスに格納される。

【0063】

図7（c）は、別の第2の変換テーブルT2aを示す。この変換テーブルT2aは、外部記憶デバイス42に由来するVDEV101に、ストライプやRAIDを適用する場合に使用される。変換テーブルT2aは、VDEV番号（VDEV#）と、ストライプサイズと、RAIDレベルと、第2の記憶制御装置40を識別するための番号（SS#（ストレージシステム番号））と、イニシエータポート番号と、通信ポート41のWWN及びLUN43の番号とを対応付けることにより構成されている。図7（c）に示す例では、1つのVDEV101は、SS#（1,4,6,7）で特定される合計4つの外部記憶制御装置を利用してRAID1を構成する。また、SS#1に割り当てられている3個のLUN（#0,#0,#4）は、同一デバイス（LDEV#）に設定されている。なお、LUN#0のボリュームは、2個のアクセスデータパスを有する交代パス構造を備えている。このように、本実施例では、外部に存在する複数の論理ボリューム（LDEV）からVDEV101を構成することにより、ストライピングやRAID等の機能を追加した上でホスト装置10に提

供することができる。

【0064】

図8を参照して、第2の記憶制御装置40のLDEVからデータを読み出す場合の流れを説明する。まず、ホスト装置10は、通信ポート21Aを指定して第1の記憶制御装置20にデータの読み出しコマンドを送信する(S31)。第1の記憶制御装置20は、読み出しコマンドを受信すると、要求されたデータを第2の記憶制御装置40から読み出すべく、読み出しコマンドを生成する。第1の記憶制御装置20は、生成した読み出しコマンドを第2の記憶制御装置40に送信する(S32)。第2の記憶制御装置40は、第1の記憶制御装置20から受信した読み出しコマンドに応じて、要求されたデータを記憶デバイス42から読み出して、第1の記憶制御装置20に送信し(S33)、正常に読み出しが完了した旨を報告する(S35)。第1の記憶制御装置20は、図8(b)に示すように、第2の記憶制御装置40から受信したデータを、キャッシュメモリ24の所定の場所に格納させる(S34)。

10

【0065】

第1の記憶制御装置20は、キャッシュメモリ24に格納されたデータを読み出し、アドレス変換を行った後、LUN103等を介してホスト装置10にデータを送信し(S36)、読み出し完了報告を行う(S37)。これらデータ読み出し時の一連の処理では、図7と共に述べた変換操作が逆向きで行われる。

【0066】

図8では、ホスト装置10からの要求に応じて、第1の制御装置のチャンネルアダプタ21が第2の記憶制御装置40からデータを読み出し、仮想ボリュームを構成するキャッシュメモリ24に保存することを示している。符号103は、仮想ボリュームを構成するLUNを示し、キャッシュメモリ24は仮想ボリュームの実記憶領域に相当する。外部LUNは後述のリストアデータが記憶された論理ボリュームである。

20

【0067】

次に、本願発明者が今回新たに認知した、本願発明を成立させる契機となった課題について説明する。図9は、第2の記憶制御装置40が第1の記憶制御装置20に接続した記憶制御システムのブロックである。図9(1)はリストア前のデータの記憶状況を示したブロック図であり、(2)はリストア後のデータの記憶状況を示すブロック図である。第1の記憶制御システム20内には、実論理ボリュームLUaと仮想ボリュームLUbとが、特開2001-216185号公報に記載のコピーペアの関係を形成可能に構成されている。

30

【0068】

LUaが正ボリューム番号によって識別される正ボリュームであり、LUbが副ボリューム番号で識別される副ボリュームである。第2の記憶制御装置40も同様に、互いにコピーペアの関係になる実論理ボリュームLUc(正ボリューム)と実論理ボリュームd(副ボリューム)とを備えている。これら正ボリューム及び副ボリュームとがコピーペア形成中であれば、正ボリュームのデータと副ボリュームの一方がコピー元、他方がコピー先となり、両者間で全コピーが行われる。

【0069】

第1の記憶制御装置の仮想ボリュームLUbは図2で説明したように、論理ボリュームLUcにマッピングされている。論理ボリュームLUはLDEVであっても良い。実論理ボリュームLUa, LUc, LUdのそれぞれは、ディスクドライブの物理的な記憶領域とキャッシュメモリに対応しており、データのREAD命令及びWRITE命令の実行に伴うデータの交換はキャッシュメモリを介してHDDなどの記憶デバイスとの間で行われる。

40

【0070】

仮想ボリュームLUbは対応する記憶デバイスを持たず、キャッシュメモリのみから構成されている。論理ボリュームLUcのデータに仮想ボリュームLUbを介してホストbがアクセスができる。各論理ボリュームの上段がキャッシュメモリに記憶された

50

データに対応し、下段がディスクドライブ内のデータに対応している。

【0071】

図9に示すシステムによれば、第1の記憶制御装置に第2の記憶制御装置が接続している関係において、正ボリュームから副ボリュームにデータをコピーするコピーペアの機能を利用して、データの世代管理が可能となる。すなわち、ホスト a がアクセス可能な実論理ボリューム LU a には最新世代のデータを記憶し、ホスト b がアクセス可能な仮想ボリューム b 及び実論理ボリューム c には前の世代のデータを記憶し、ホスト b がアクセス可能な実論理ボリューム LU d にはさらに前の世代のデータを記憶することができる。例えば、最新世代のデータをオンラインデータとすると、一つ前の世代のデータを過去データ1（例えば一日前のデータ）とし、さらに前の世代のデータを過去データ2（例えば二日前のデータ）として、図9の記憶制御システムでデータの世代管理をすることができる。

10

【0072】

本発明に於けるリストアは、第2の記憶制御装置の論理ボリューム LU d のデータを論理ボリューム LU c に記憶させることをいう。ホスト b は仮想ボリューム LU b にアクセスすることによって、論理ボリューム LU c に記憶されたリストアデータにアクセスすることができる。リストアによって、記憶制御装置20のデータを記憶制御装置40の記憶資源に記憶された過去データに基づいて復旧することができる。

【0073】

リストアの処理について説明する。既述のとおり、図9の(1)はリストア前の記憶制御システムのデータの記憶状態を示し、(2)はリストア実行後の状態を示している。第2の記憶制御装置40のDKAが、ホスト b からリストア命令を受けると、論理ボリューム LU d の記憶デバイス900からリストアの対象となるデータ(リストアデータ)「def」を読み込み、これを論理ボリューム LU d のキャッシュメモリ902の記憶領域に保存する(S90)。CHA90Aはこのキャッシュメモリ902の記憶領域のデータを論理ボリューム LU c のキャッシュメモリ904の記憶領域に記憶する(S91)。そして、この記憶領域のデータ「def」をDKAが論理ボリューム LU c の記憶デバイス906の記憶領域に記憶する(S92)。

20

【0074】

一方、第1の記憶制御装置20のCHA90Bは、ホスト b からのリストア命令を実行し、リストアデータを保存するキャッシュメモリ908のブロック領域を設定し、この設定情報を第2の記憶制御装置40のCHA90Aに転送する(S94)。CHA90Aはリストアデータが記録されているキャッシュメモリ904の記憶領域のデータ「def」を読み込み(S95)、第1の記憶制御装置40のCHA90Bに渡す(S96)。CHA90Bは予め決められたキャッシュメモリ908のブロック記憶領域にこのデータ「def」を保存する(S93)。

30

【0075】

なお、図9(2)には示されていないが、論理ボリューム LU a と仮想ボリューム b とがコピーペアを形成すれば、仮想ボリューム LU b に対応するキャッシュメモリ908の記憶領域のデータが論理ボリューム LU a の記憶領域(キャッシュメモリ及びHDD)にコピーされる。

40

【0076】

いま、第2の記憶制御装置40がリストア命令を実行すると、図9(2)に示すように、論理ボリューム LU c のデータは「abc」から「def」になる。この論理ボリューム LU c にマッピングされている仮想ボリューム LU b には「def」が転送される。しかしながら、仮想ボリュームに対応するキャッシュメモリの記憶領域にはリストア前の「abc」のデータが残っているので、ホスト b にはリストアデータである「def」ではなく、リストア前のデータである「abc」が見てしまう。

【0077】

さらに、仮想ボリューム LU b とコピーペアの関係にある正論理ボリューム LU a

50

にも「d e f」ではなく「a b c」をコピーしてしまうことになる。これでは、第2の記憶制御装置40のリストアデータに、ホスト bが正しくアクセスできない問題がある。

【0078】

そこで、第1の記憶制御装置10は、ホスト bのページコマンドを受けて、リストア命令の実行の際に仮想ボリュームLU bに対応するキャッシュメモリのデータ「a b c」をページ（「無効化」或いは「開放」）することとした（S97）。ここで、ページとは、仮想ボリュームとキャッシュメモリの記憶領域との対応関係をクリアするための処理に相当する。この処理が終了した後、或いはこの処理に同期して、第1の記憶制御装置10は、キャッシュメモリの記憶領域にリストアデータ「d e f」を記憶することとした。

【0079】

図10は、仮想ボリュームのディレクトリ構造とキャッシュメモリとの対応関係を示すブロック図である。仮想ボリュームLU bは複数の管理領域100から構成されており、100個のLBA毎に1個の管理領域100が割り当てられている。一つの管理領域はキャッシュメモリの記憶領域を指す番号（キャッシュ領域）100Aと、キャッシュ領域の中のデータが有効かどうかを示す有効ビットマップ100Bとから構成される。

【0080】

図10は、LBA#100～LBA#199にキャッシュメモリの記憶領域（キャッシュ領域）#1の各ブロックが対応していることを示している。キャッシュ領域は、各LBAに対応したブロックに分割されており、各ブロックに各LBAに対応するデータが置かれる。ビットマップの各ビットは、一つのLBA及びこれに対応するキャッシュ領域のブロック単位に対応しており、有効ビットが「0」の場合には、キャッシュ領域の該当するブロック単位のデータが無効であることを示し、有効ビットが「1」の場合には、対応ブロック単位にあるデータが有効であることを示している。有効ビットが「0」の場合には、ホストから対応するLBAにリード要求が来た場合には、第2の記憶制御装置の論理ボリュームのデータをリードする必要がある。

【0081】

キャッシュ領域#1においては、LBA#100とLBA#101に対応するブロック単位のデータが有効であり、LBA 102, 103に対応するブロック単位のデータが無効である事を示している。管理領域100の[NULL]は、キャッシュ領域に対応する部分がないことを示している。仮想ボリュームの管理領域に割り当てていないキャッシュ領域は、フリーキューとして共有メモリ（図1の25）の特定記憶領域に記憶されている。

【0082】

ホストからのリード命令が、キャッシュ領域の割り当てがない、仮想ボリュームの管理領域に来た場合には、第1の制御装置のCHAは共有メモリのフリーキューの中からキャッシュ領域の情報を取得し、この情報を仮想ボリュームの管理領域に割り当てる。図10に示す仮想ボリュームのディレクトリ構造を規定する制御情報は共有メモリに記憶される。なお、図10は仮想ボリュームについて説明したが、実論理ボリュームでも同様である。

【0083】

図11は、この実施形態におけるページ処理を説明するフローチャートである。図9（2）のホスト bがページ処理の対象となる仮想ボリュームを指定し、この仮想ボリュームの第1の管理領域に対して（S1100）、ページ処理を実行する（S1102 - S1106）。

【0084】

ホストが第1の記憶制御装置40にページコマンドを発行すると、CHA90Bは、共有メモリの仮想ボリュームのディレクトリ構造（図10）をリードし、指定された管理領域にキャッシュメモリの記憶領域の情報が設定されているかを判定する（S1102）。CHA90Bは、これが肯定された場合には、割り当てられたキャッシュ領域をフリーキューに戻して、このキャッシュ領域に「N U L L」を設定する（S1104）。S1102の

10

20

30

40

50

判定が否定された場合及びS 1 1 0 4の処理に続いて、CHAは全管理領域についてS 1 1 0 2及びS 1 1 0 4の処理が完了したか否かを判定する。それが肯定された場合は、ルーチンを終了する。否定判定の場合には、次の管理領域について、S 1 1 0 2、S 1 1 0 4の処理を継続する(S 1 1 0 8)。これにより、仮想ボリュームの全管理領域に対してキャッシュメモリの記憶領域のデータがページされたことになり、リストアされる前の、キャッシュメモリのデータが仮想ボリュームのディレクトリ構造を介してホストにアクセスされることが無いようにしている。

【0085】

このページが完了されると、第1の記憶制御装置のCHA90B(図9)はページ完了をホストbに通知する。ホストは第2の記憶制御装置のCHA90Aにリストア要求を発行する。リストアの命令は、コマンドデバイス(図2の166)を経由して第1の制御装置のCHA90Bが第2の制御装置のCHA90Aにホストbを介することなく送ることもできる。リストア命令を受けた第2の記憶装置は、図9に示すように論理ボリュームLU dとLU cとのコピーペアを形成して、論理ボリュームLU dのデータを論理ボリュームLU cにコピーする(リストア)。ホストからアクセスのあった仮想ボリュームの管理領域には、フリーキューの中から適宜キャッシュメモリの領域が選択されて割り当てられ、このキャッシュメモリの領域にリストアデータが記憶される。論理ボリュームLU cは仮想ボリュームLU bにマッピングされているので、論理ボリュームLU cのキャッシュメモリのリストアデータが、割り当てられたキャッシュ領域に転送される。したがって、ホストbがCHA90Bを介して、仮想ボリュームの管理領域にアクセスすることにより、ホストbは、この管理領域に対応するキャッシュメモリの記憶領域にあるリストアデータにアクセスすることができる。

10

20

【0086】

仮想ボリュームの管理領域に[NULL]を設定することにより、キャッシュメモリのリストア前のデータをページすることができる。なお、ページ処理の他の例として、有効ビット(図10の100B)に「0」を設定するようにしてもよい。さらに、キャッシュのページが終了する前に、第2の記憶制御装置はリストアを開始してもよい。ただし、仮想ボリュームの全ての管理領域についてキャッシュメモリのページが終了しないと、ホストbは仮想ボリュームLU bのリストアデータにアクセスすることができない。今後論理ボリュームの容量、キャッシュ容量がますます大きくなっていくため、ページに要する時間も長くなり、その結果リストアしたデータをホストが利用できるようになるまでに時間がかかってしまう。

30

【0087】

そこで、第2の実施形態では、ホストがアクセスした、仮想ボリュームに対応するキャッシュメモリの記憶領域のデータがリストアデータを反映したものであるか否かを第1の記憶制御装置が判断し、リストアデータを反映していない場合には、このキャッシュ領域のデータをページし、ついでこのキャッシュ領域にリストアデータを記憶するようにして、全ての管理領域についてページを待つことなく、ホスト装置がリストアデータにアクセスできるようにした。

【0088】

このことを可能にするための、仮想ボリュームのディレクトリ構造が図12に示されている。仮想ボリュームに相当するLU毎に全体世代番号120が設定され、各管理領域には図10に示す、キャッシュ領域100A、有効ビット100Bの他、個別の世代番号1120Aも設定される。図12では、仮想ボリュームのキャッシュ領域1がキャッシュメモリ24のキャッシュ領域1に、仮想ボリュームのキャッシュ領域5がキャッシュメモリ24のキャッシュ領域5に対応している。

40

【0089】

図13はこの実施形態のページ処理のためのフローチャートを示すものであり、図9(2)に示す、ホストbから記憶制御装置20にページ要求が出力されると、CHA90Bは共有メモリの制御情報にアクセスして仮想ボリュームを指定し、仮想ボリュームのデ

50

ィレクトリ構造の全体世代番号120の値を1インクリメントする(S1300)。ついでCHA90Bはホストbに正常終了を通知する。ホストbの命令を受けて、記憶制御装置40は二つの内部論理ボリューム間でシャドウイメージを利用したリストア処理を開始する。

【0090】

ホストbがこの正常終了を受けると、図14に示す、リストアデータへのアクセス処理(I/O処理)がスタートする。ホストbが仮想ボリュームLUb(図9(2)参照)へI/Oを発行すると、第1の記憶制御装置のCHA90BはホストbからI/Oの要求を受けた仮想ボリュームのディレクトリ構造を共有メモリから参照し、I/O要求を受けたLBAのキャッシュ領域(図12)を取得する(S1400)。CHAはキャッシュ領域の制御データがNULLであるか否かをチェックする(S1401)。これが肯定されると、CHAはLBAが属する管理領域にキャッシュメモリのフリーキューのキャッシュ領域を割り当てる。さらに有効ビット(図12の100B)をオフし、世代番号に代表世代番号を設定する(S1402)。

10

【0091】

ステップS1401において、キャッシュ領域を示す情報がNULLでない場合には、CHAは、I/O要求を受けたLBAが属する管理領域の世代番号と全体世代番号とを比較する(S1404)。この比較結果が否の場合には管理領域の有効ビットがオフされる(S1406)。これが肯定されている場合には、S1406をスキップしてS1408にジャンプする。

20

【0092】

S1408では、CHAが、ホストからの命令がWRITE命令又はREAD命令かを判定し、後者の場合はS1410に移行する。CHAはLUcからリストアデータを読み出して、LBAで指定される、仮想ボリュームLUbにおける管理領域上のキャッシュ領域に該当するキャッシュメモリの記憶領域にリストアデータを書き込む。次いで、CHAは管理領域の有効ビットをオンし(S1412)、ホストにリストアデータを転送する(S1414)。

【0093】

S1408で、ホストからの命令がWRITE命令の場合には、CHAはホストからライトデータを受領して、I/O要求を受けたLBAのキャッシュ領域に該当するキャッシュメモリ上の記憶領域にライトデータを書き込んで(S1416)、有効ビットをオンする(S1418)。

30

【0094】

この一連の処理において、第1の記憶制御装置のCHAは、ホストからのI/O要求を受けたLBAに対応する、仮想ボリュームの管理領域の世代番号と全体世代番号とを比較し(S1404)、これが一致する場合には、キャッシュメモリの記憶領域に正しいリストアデータをアップロードできる(S1410)。一方、この世代番号が一致していない場合には、キャッシュメモリの記憶領域のデータはリストアデータに対応しない不適切なデータであるので、第1の記憶制御装置のCHAは有効ビットをオフしてなる、キャッシュメモリのデータのページを行い(S1406)、次いで、仮想ボリュームの管理領域にリストアデータを設定する。この実施形態では、第1の制御装置の制御部がキャッシュメモリ内のデータがページを必要とするリストア前のデータか、リストアデータであってホストをアクセスさせても良いデータかを識別し、ページが必要なデータである場合には、外部論理ボリュームからリストアデータをリードして、キャッシュメモリの記憶領域のリストア前のデータをリストアデータでリフレッシュする。この実施形態では、仮想ボリュームの全管理領域に対してキャッシュメモリのページを行わず、ホストがアクセスした管理領域について、キャッシュメモリのデータのページを行えばすむ為、ホストがリストアデータにアクセスするために要する時間を短縮化できるという効果がある。また、この実施形態はページ処理とリストアデータの転送とを並列処理するために、ホストがリストアデータを利用できるまでに要する時間を短縮化できる。

40

50

【 0 0 9 5 】

なお、既述の実施形態では、第2の記憶制御装置内の副ボリュームから正ボリュームに全コピーしてデータのリストアを行っていたが、論理ボリュームLU c (図9)にホストcが接続され、このホストに接続したテープデバイスなどの記憶資源からデータのリストアを行う場合にも本発明が適用される。また、仮想ボリュームに対するページ処理は仮想ボリューム全体、管理領域単位、或いはLBA(ロジカルブロックアドレス)の単位で行うことができる。

【 図面の簡単な説明 】

【 0 0 9 6 】

【 図 1 】 本発明の実施例に係わる記憶システムの全体構成を示すブロック図である。 10

【 図 2 】 記憶システムの論理的構成の概要を示す模式図である。

【 図 3 】 仮想論理ボリュームの設定動作の流れを示すフローである。

【 図 4 】 マッピングテーブルの概要を示す説明図である。

【 図 5 】 マッピングテーブルを構築するための処理の流れを示すフローである。

【 図 6 】 内部ボリュームとして仮想化される外部の記憶デバイスにデータを書き込む場合の概念図である。

【 図 7 】 書き込みデータのアドレス変換の様子を模式的に示す説明図である。

【 図 8 】 内部ボリュームとして仮想化される外部の記憶デバイスからデータを読み出す場合の概念図である。

【 図 9 】 キャッシュメモリのページ処理を説明するための、記憶制御システムのブロック図であり、(1)はリストア前のデータの記憶状態を、(2)はリストア後のデータの記憶状態を説明するブロック図である。 20

【 図 1 0 】 第1の実施形態における、仮想ボリュームのディレクトリ構造と、このディレクトリ構造に対応するキャッシュメモリの記憶構造とのブロック図である。

【 図 1 1 】 第1の実施形態のページ処理を説明するフローチャートである。

【 図 1 2 】 第2の実施形態における、仮想ボリュームのディレクトリ構造と、このディレクトリ構造に対応するキャッシュメモリの記憶構造とのブロック図である。

【 図 1 3 】 第2の実施形態における、ページ処理のフローチャートである。

【 図 1 4 】 ホストがリストアデータにアクセスする事示すフローチャートである。。

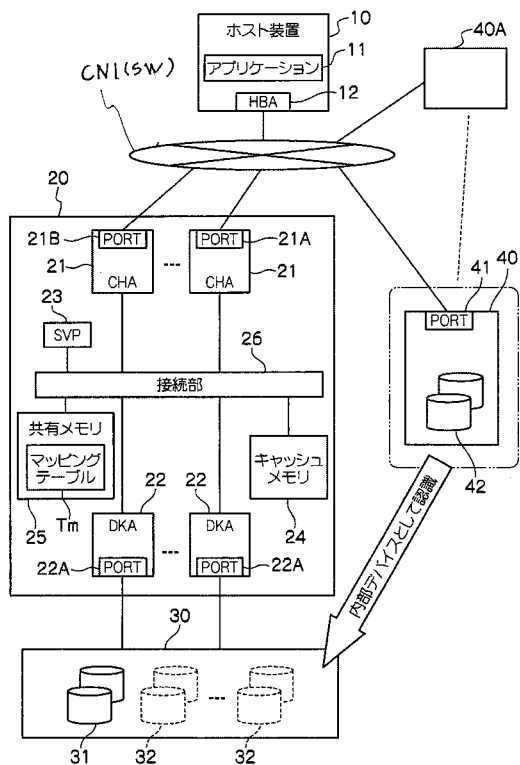
【 符号の説明 】

30

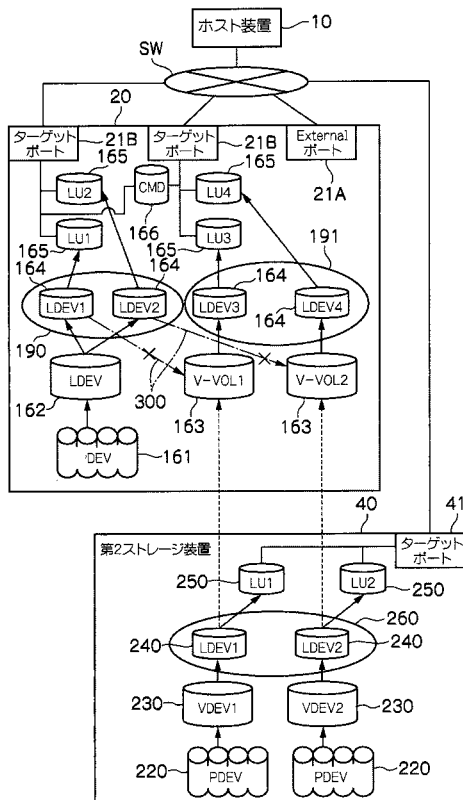
【 0 0 9 7 】

1 0 ... ホスト装置、 1 1 ... アプリケーションプログラム、 1 2 ... アダプタ、 2 0 ... 第1の記憶制御装置、 2 1 , 9 0 A , 9 0 B ... チャンネルアダプタ、 2 1 A ... 通信ポート、 2 2 ... ディスクアダプタ、 2 2 A ... 通信ポート、 2 3 ... コントロールユニット、 2 4 ... キャッシュメモリ、 2 4 A ... サブブロック、 2 5 ... 共有メモリ、 2 6 ... 接続部、 3 0 ... 記憶装置、 3 1 ... 記憶デバイス、 3 2 , LU b ... 記憶デバイス(仮想化された内部記憶デバイス/仮想ボリューム)、 4 0 ... 第2の記憶制御装置、 4 1 ... 各通信ポート、 4 2 ... 各記憶デバイス、 LU a , LU c , LU d ... 実論理ボリューム、 a , b ... ホスト

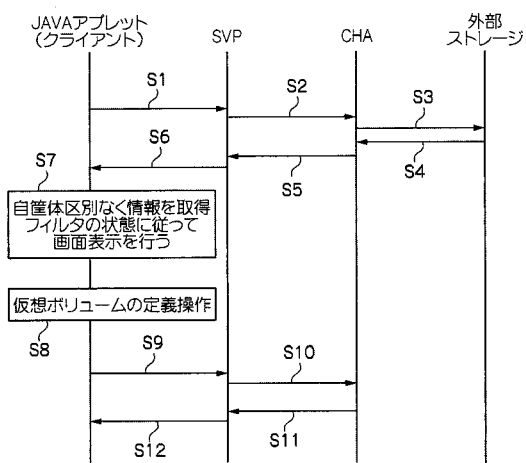
【 図 1 】



【 図 2 】



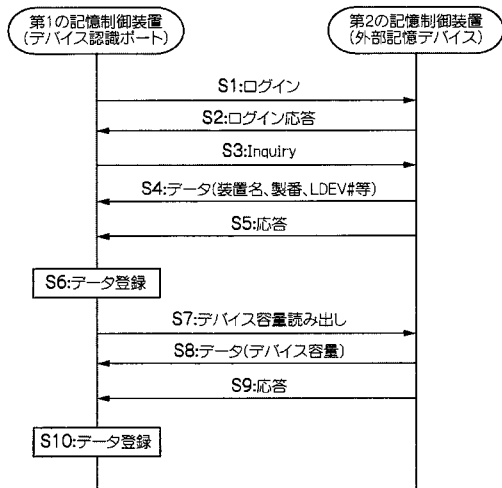
【 図 3 】



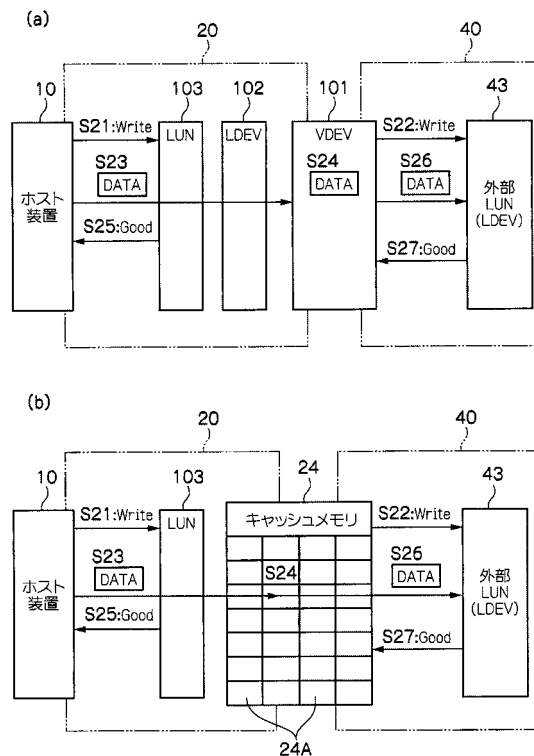
【 図 4 】

VDEV	外部デバイス情報			バス情報	
	デバイス識別情報	容量(KB)	デバイス種別	WWN	LUN
0	DRFGTFNEIEK	657,456	DISK	0xAABCCDD	0
1	ADRF-GTFNEIE	89,854	DISK	0xAABBEFF	3
2	GGRRFFDDERT	-	TAPE	0x445566AAB	5
3	AABCCDDEE	5,544,223	DISK	0x77DE12345	6
				0x77DE12345	3
				0x377DE7890	5

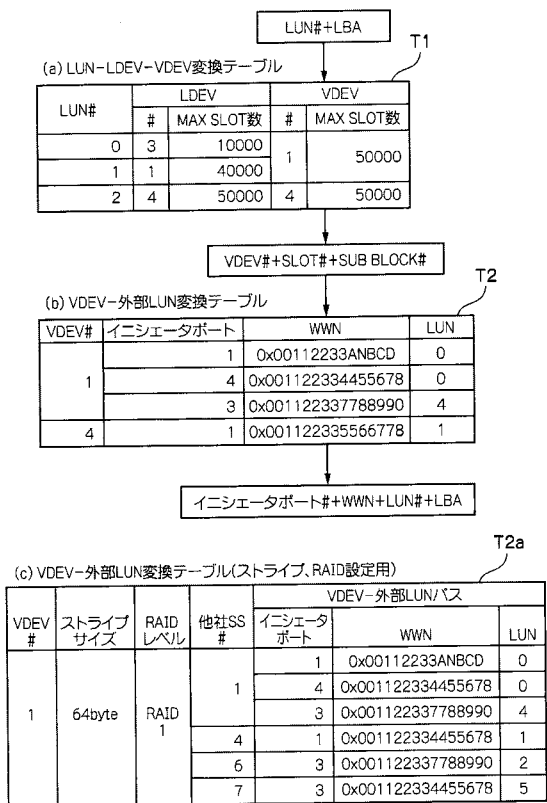
【 図 5 】



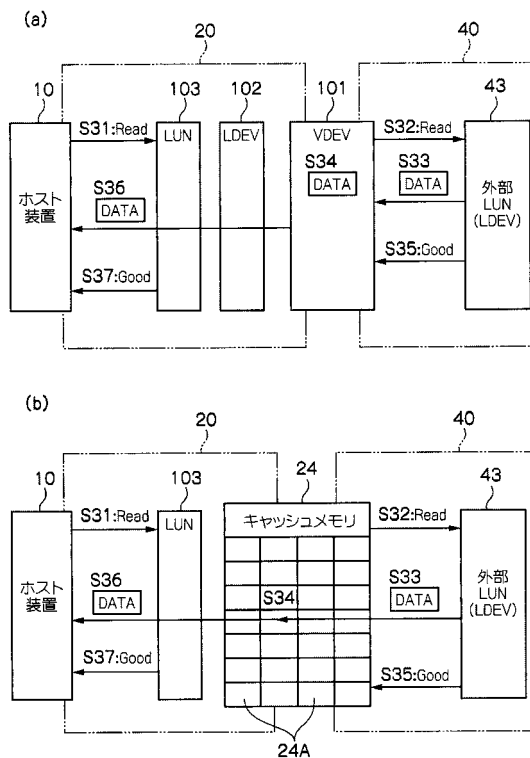
【 図 6 】



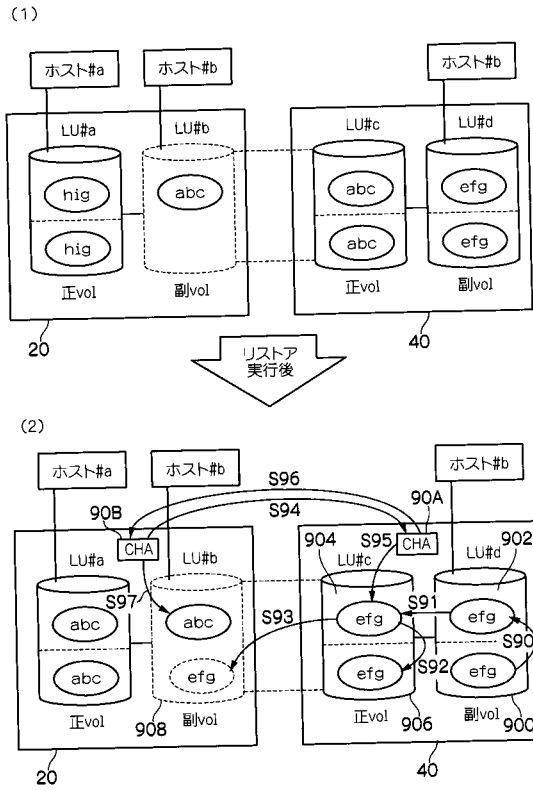
【 図 7 】



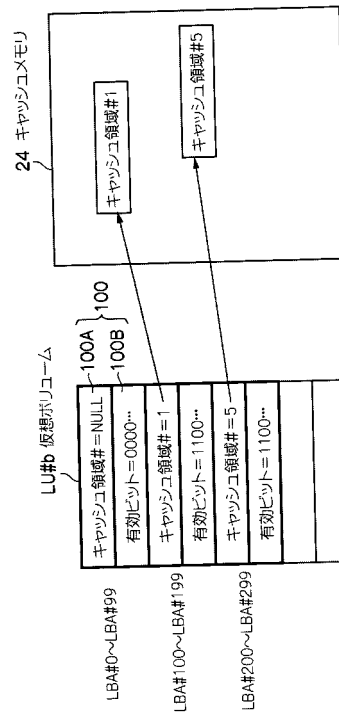
【 図 8 】



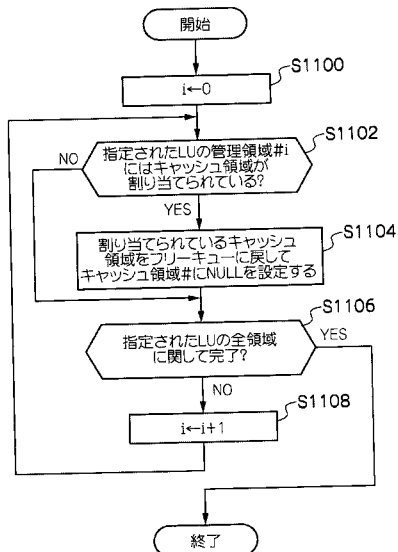
【 図 9 】



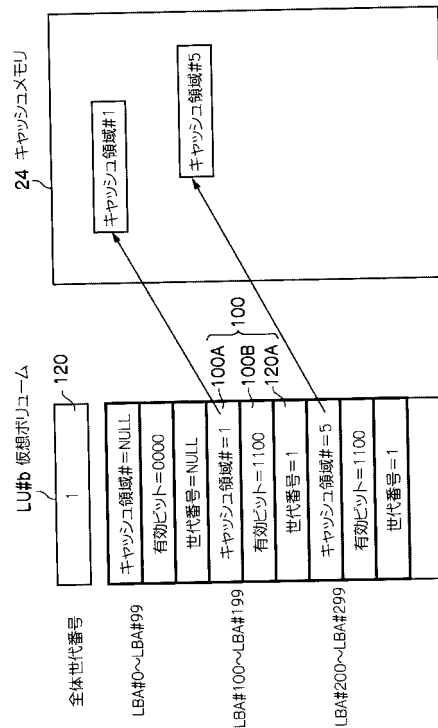
【 図 10 】



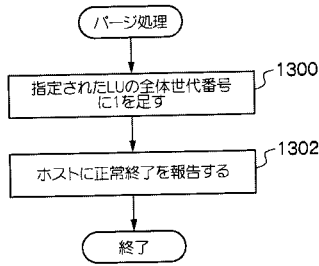
【 図 11 】



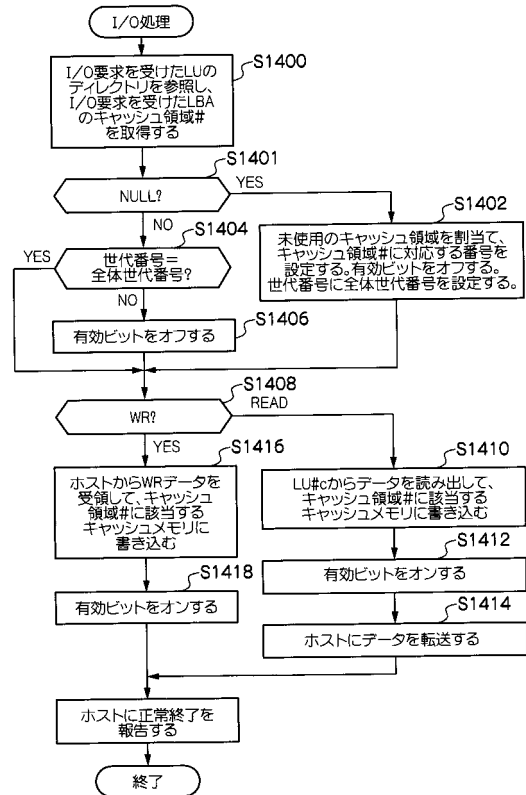
【 図 12 】



【 図 1 3 】



【 図 1 4 】



フロントページの続き

(51) Int.Cl.

F I

テーマコード(参考)

G 0 6 F	12/08	5 3 1 Z
G 0 6 F	12/08	5 5 7
G 0 6 F	12/08	5 5 9 Z
G 0 6 F	13/12	3 4 0 B

Fターム(参考) 5B065 BA01 CA12 CC03 CC08 CE12 CH01 EA02 EA25