

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4338068号
(P4338068)

(45) 発行日 平成21年9月30日(2009.9.30)

(24) 登録日 平成21年7月10日(2009.7.10)

(51) Int.Cl.

F I

G 0 6 F 3/06 (2006.01)

G 0 6 F 3/06 3 0 4 B

G 0 6 F 3/06 3 0 2 A

G 0 6 F 3/06 3 0 2 B

請求項の数 12 (全 30 頁)

(21) 出願番号 特願2002-77499 (P2002-77499)
 (22) 出願日 平成14年3月20日(2002.3.20)
 (65) 公開番号 特開2003-280824 (P2003-280824A)
 (43) 公開日 平成15年10月2日(2003.10.2)
 審査請求日 平成16年9月24日(2004.9.24)

前置審査

(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
 (74) 代理人 100093861
 弁理士 大賀 真司
 (72) 発明者 藤本 和久
 東京都国分寺市東恋ヶ窪一丁目280番地
 株式会社日立製作所中央研究所内

審査官 梅景 篤

最終頁に続く

(54) 【発明の名称】 ストレージシステム

(57) 【特許請求の範囲】

【請求項1】

ディスク装置と、

個々が、ホストコンピュータとのインターフェースを有するチャンネルインターフェース部及び前記ディスク装置とのインターフェースを有するディスクインターフェース部とを有する複数のディスク制御クラスタと、

前記ディスク装置に対しリード/ライトされるデータと、前記データの転送に関する制御情報と、前記ディスク装置の管理情報と、前記複数のディスク制御クラスタの負荷情報若しくは障害情報を含む管理情報を格納するグローバル情報制御部と、

前記複数のディスク制御クラスタを相互接続する相互結合網と、

前記複数のディスク制御クラスタの個々が有する前記チャンネルインターフェース部の相互を接続するスイッチと、を有し、

前記複数のディスク制御クラスタにおいて、前記ホストコンピュータからのデータのリード/ライト要求に対し、前記チャンネルインターフェース部は前記ホストコンピュータとのインターフェースと前記グローバル情報制御部との間のデータ転送を実行し、前記ディスクインターフェース部は前記ディスク装置と前記グローバル情報制御部との間のデータ転送を実行し、

前記スイッチは、定期的に若しくは必要に応じて、前記グローバル情報制御部に格納された前記負荷情報若しくは障害情報を含む管理情報の複製を格納するメモリを有し、

前記グローバル情報制御部は前記相互結合網と前記スイッチに接続され、

10

20

前記スイッチは、自身が有する複数のポート間の接続切換え情報を前記メモリの複製された管理情報に基づいて更新し、当該接続切換え情報に基づいて前記複数のポート間の接続切換えを行うことを特徴とするストレージシステム。

【請求項 2】

個々が、ホストコンピュータとのインターフェースを有するチャンネルインターフェース部と、ディスク装置とのインターフェースを有するディスクインターフェース部と、前記ディスク装置に対しリード/ライトされるデータ、前記データの転送に関する制御情報及び前記ディスク装置の管理情報を格納するローカル共有メモリ部とを有し、前記ホストコンピュータからのデータのリード/ライト要求に対し、前記チャンネルインターフェース部は前記ホストコンピュータとのインターフェースと前記ローカル共有メモリ部との間のデータ転送を実行し、前記ディスクインターフェース部は前記ディスク装置と前記ローカル共有メモリ部との間のデータ転送を実行することにより、データのリード/ライトを行う複数のディスク制御クラスタと、

前記ディスクインターフェース部と接続されるディスク装置と、

前記複数のディスク制御クラスタ負荷情報若しくは障害情報を含む管理情報を格納するグローバル情報制御部と、

前記複数のディスク制御クラスタを相互接続する相互結合網と、

前記複数のディスク制御クラスタの個々が有するチャンネルインターフェース部を相互に接続するスイッチと、を有し、

前記スイッチは、定期的に若しくは必要に応じて、負荷情報若しくは障害情報を含む管理情報の複製を格納するメモリを有し、自身が有する複数のポート間の接続切換え情報を前記メモリに格納された前記管理情報に基づいて更新し、当該更新後の接続切換え情報に基づいて前記複数のポート間の接続切換えを行うことを特徴とするストレージシステム。

【請求項 3】

前記チャンネルインターフェース部、前記ディスクインターフェース部及び前記ローカル共有メモリ部とを相互に接続する接続部を前記複数のディスク制御クラスタの個々が有し、

前記複数のディスク制御クラスタの個々が有する前記制御部は、他のディスク制御クラスタが有する接続部と前記相互結合網を介して接続され、

前記グローバル情報制御部は前記相互結合網及び前記スイッチに接続されることを特徴とする請求項 2 記載のストレージシステム。

【請求項 4】

前記複数のディスク制御クラスタの個々が有する前記ローカル共有メモリ部が前記相互結合網を介して相互に接続され、

前記グローバル情報制御部は前記相互結合網と前記スイッチに接続されることを特徴とする請求項 2 記載のストレージシステム。

【請求項 5】

ディスク装置と、

個々が、ホストコンピュータとのインターフェースを有するチャンネルインターフェース部と、ディスク装置とのインターフェースを有するディスクインターフェース部と、前記ディスク装置に対しリード/ライトされるデータを格納する第 1 のメモリと前記チャンネルインターフェース部及び前記ディスクインターフェース部と前記第 1 のメモリとの間のデータ転送に関する制御情報及び前記ディスク装置の管理情報を格納する第 2 のメモリとを有するローカル共有メモリ部とを有し、前記ホストコンピュータからのデータのリード/ライト要求に対し、前記チャンネルインターフェース部は前記ホストコンピュータとのインターフェースと前記ローカル共有メモリ部内の前記第 1 のメモリとの間のデータ転送を実行し、前記ディスクインターフェース部は前記ディスク装置と前記ローカル共有メモリ部内の前記第 1 のメモリとの間のデータ転送を実行することにより、データのリード/ライトを行う複数のディスク制御クラスタと、

前記複数のディスク制御クラスタの負荷情報若しくは障害情報を含む管理情報を格納す

るグローバル情報制御部と、

前記複数のディスク制御クラスタを相互に接続する第 1、第 2 の相互結合網と、

前記複数のディスク制御クラスタの個々が有する前記チャンネルインターフェース部を相互に接続するスイッチと、を有し、

前記スイッチは、定期的に若しくは必要に応じて、前記グローバル情報制御部に格納された前記負荷情報若しくは障害情報を含む管理情報の複製を格納するメモリを有し、自身が有する複数のポート間の接続切換え情報を前記メモリに格納された前記管理情報に基づいて更新し、当該更新後の接続切換え情報に基づいて前記複数のポート間の接続切換えを行うことを特徴とするストレージシステム。

【請求項 6】

前記複数のディスク制御クラスタの個々が有する前記チャンネルインターフェース部と前記ディスクインターフェース部が、他のディスク制御クラスタが有する前記チャンネルインターフェース及び前記ディスクインターフェース部と前記第 1 の相互結合網を介して接続され、

前記グローバル情報制御部は前記第 1 の相互結合網と前記スイッチに接続されることを特徴とする請求項 5 記載のストレージシステム。

【請求項 7】

ディスク装置と、

個々が、ホストコンピュータとのインターフェースを有するチャンネルインターフェース部と、前記ディスク装置とのインターフェースを有するディスクインターフェース部と、前記ディスク装置に対しリード/ライトされるデータと前記データの転送に関する制御情報と前記ディスク装置の管理情報を格納するローカル共有メモリ部と、前記チャンネルインターフェース部、前記ディスクインターフェース部及び前記ローカル共有メモリ部とを相互に接続する第 1 の接続部とを有し、前記ホストコンピュータからのデータのリード/ライト要求に対し、前記チャンネルインターフェース部は前記ホストコンピュータとのインターフェースと前記ローカル共有メモリ部との間のデータ転送を実行し、前記ディスクインターフェース部は前記ディスク装置と前記ローカル共有メモリ部との間のデータ転送を実行することにより、データのリード/ライトを行う複数のディスク制御クラスタと、

前記複数のディスク制御クラスタの負荷情報若しくは障害情報を含む管理情報を格納するグローバル情報制御部と、

前記複数のディスク制御クラスタの個々が有する前記チャンネルインターフェース部を相互に接続するスイッチと、

前記複数のディスク制御クラスタの相互を接続する第 2 の接続部を有し、

前記グローバル情報制御部は、前記第 2 の接続部と前記スイッチに接続パスで接続され、

前記複数のディスク制御クラスタの個々が有する前記第 1 の接続部は前記第 2 の接続部に接続パスで接続され、

前記スイッチは、定期的に若しくは必要に応じて、前記グローバル情報制御部に格納された前記負荷情報若しくは障害情報を含む管理情報の複製を格納するメモリを有し、自身が有する複数のポート間の接続切換え情報を前記メモリに格納された前記管理情報に基づいて更新し、当該更新後の接続切換え情報に基づいて前記複数のポート間の接続切換えを行うことを特徴とするストレージシステム。

【請求項 8】

前記複数のディスク制御クラスタの個々が有するローカル共有メモリ部は、自身が属するディスク制御クラスタ内の各部位の負荷情報及び障害情報、並びに前記自身が属するディスク制御クラスタが管理する記憶領域情報を格納しており、

前記グローバル情報制御部は第 2 のメモリを有し、前記第 2 のメモリは前記複数のディスク制御クラスタの前記負荷情報、障害情報、記憶領域情報を格納しており、

前記グローバル情報制御部は、前記ローカル共有メモリ部内の前記負荷情報、障害情報、記憶領域情報のある時間間隔で参照し、前記第 2 のメモリ内の該負荷情報、障害情報、

10

20

30

40

50

記憶領域情報を更新することを特徴とする請求項 2 記載のストレージシステム。

【請求項 9】

前記グローバル情報制御部は第 2 のメモリを有し、前記第 2 のメモリは前記複数のディスク制御クラスタの各部位の障害情報、負荷情報を格納しており、

前記複数のディスク制御クラスタの各々は、自身内のある部位で障害が発生した時点、あるいはある部位の負荷が予め設定した値より高くなった時点で、前記グローバル情報制御部の前記第 2 のメモリ内の前記障害情報あるいは負荷情報を更新することを特徴とする請求項 2 記載のストレージシステム。

【請求項 10】

前記グローバル情報制御部は第 2 のメモリを有し、前記第 2 のメモリは前記複数のディスク制御クラスタが管理する記憶領域情報を格納しており、

前記複数のディスク制御クラスタの各々は、自身内のある記憶領域のデータを、他のディスク制御クラスタ内のある記憶領域にコピーあるいは移動した時点で、前記グローバル情報制御部の前記第 2 のメモリ内の前記記憶領域情報を更新することを特徴とする請求項 2 記載のストレージシステム。

【請求項 11】

前記グローバル情報制御部は、該メモリ内の前記負荷情報、障害情報、あるいは記憶領域情報が更新された時点で、前記負荷情報、障害情報、あるいは記憶領域情報を前記スイッチ内のメモリにコピーすることを特徴とする請求項 8、9 及び 10 に記載のうちのいずれか一つのストレージシステム。

【請求項 12】

前記スイッチは、前記接続切換え情報としての接続切換えテーブルを有し、前記管理情報、あるいは、記憶領域情報が更新された時点で、前記接続切換えテーブルに登録された情報を更新することを特徴とする請求項 11 記載のストレージシステム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、データを複数のディスク装置に格納するストレージシステムとそれを構成するディスク制御クラスタに関する。

【0002】

【従来の技術】

半導体記憶装置を記憶媒体とするコンピュータの主記憶の I/O 性能に比べて、磁気ディスクを記憶媒体とするディスクサブシステム（以下「サブシステム」という。）の I/O 性能は 3～4 桁程度小さく、従来からこの差を縮めること、すなわちサブシステムの I/O 性能を向上させる努力がなされている。サブシステムの I/O 性能を向上させるための 1 つの方法として、複数のディスク装置でサブシステムを構成し、データを複数のディスク装置に格納する、いわゆるディスクと呼ばれるシステムが知られている。

【0003】

例えば、従来技術では、図 2 に示すようにホストコンピュータ 3 とディスク制御装置 4 との間のデータ転送を実行する複数のチャンネル I/F 部 11 と、ディスク装置 2 とディスク制御装置 4 間のデータ転送を実行する複数のディスク I/F 部 16 と、ディスク装置 2 のデータとディスク制御装置 4 に関する制御情報（例えば、ディスク制御装置 4 内のデータ転送制御に関する情報、ディスク装置 2 に格納するデータの管理情報）を格納する共有メモリ部 20 とを備え、1 つのディスク制御装置 4 内において、共有メモリ部 20 は全てのチャンネル I/F 部 11 及びディスク I/F 部 16 からアクセス可能な構成となっている。このディスク制御装置 4 では、チャンネル I/F 部 11 及びディスク I/F 部 16 と共有メモリ部 20 との間は相互結合網 30 で接続される。

【0004】

チャンネル I/F 部 11 は、ホストコンピュータ 3 と接続するためのインターフェース及びホストコンピュータ 3 に対する入出力を制御するマイクロプロセッサ（図示せず）を有して

10

20

30

40

50

いる。また、ディスクＩＦ部１６は、ディスク装置２と接続するためのインターフェース及びディスク装置２に対する入出力を制御するマイクロプロセッサ（図示せず）を有している。また、ディスクＩＦ部１６は、ＲＡＩＤ機能の実行も行う。

インターネットの普及等により企業で扱うデータは爆発的に増大しており、データセンタ等では一台のディスク制御装置で扱えるデータ量以上のデータを記憶する必要がある。このため、図４に示すようにディスク制御装置４を複数台設置し、それらのホストコンピュータ３とのインターフェースをＳＡＮスイッチ５を介して、ホストコンピュータ３に接続していた。

また、データ量の増大に伴いＳＡＮスイッチ５に接続するディスク制御装置４の台数が増え、ホストコンピュータ３とＳＡＮスイッチ５を含めたシステム全体（このシステムをストレージ・エリア・ネットワーク（ＳＡＮ）と呼ぶ）の管理が複雑化する。それに対処するため、ＳＡＮスイッチ５にＳＡＮアプライアンス６を接続し、ＳＡＮアプライアンス６においてＳＡＮスイッチ５に繋がる全てのディスク制御装置４が管理するデータのディレクトリサービスを行い、ホストコンピュータ３に対して複数のディスク制御装置４を１つのストレージシステムに見せる処理、言い換えると、個々のディスク制御装置４が提供する記憶領域を１つの大きな記憶領域の固まりに見せ、その中から必要な量の記憶領域をホストコンピュータ３に割当てるという処理を行っていた。

また、特開２００１－２５６００３に開示されている他の従来技術では、図２３に示すように、１つのストレージシステム１は複数のディスク制御クラスタ１－１乃至１－ｎから構成される。各ディスク制御クラスタは、ホストコンピュータ３と該ディスク制御クラスタの間のデータ転送を実行する複数のチャンネルＩＦ部１１と、ディスク装置２と該ディスク制御クラスタの間のデータ転送を実行するディスクＩＦ部１６と、共有メモリ部２５を有し、チャンネルＩＦ部１１及びディスクＩＦ部１６と共有メモリ部２５の間は複数のディスク制御クラスタ１－１乃至１－ｎに跨る相互結合網３１を介して接続されている。共有メモリ部２５にストレージシステムの制御情報が格納されている。共有メモリ部２５は、相互結合網３１を介して、全てのチャンネルＩＦ部１１及びディスクＩＦ部１２からアクセス可能な構成となっており、共有メモリ部２５を介して制御情報のやりとりをすることにより、複数のディスク制御クラスタが１つのストレージシステムとして動作していた。

【０００５】

【発明が解決しようとする課題】

銀行、証券、電話会社等に代表される大企業では、従来各所に分散していたコンピュータ及びストレージを、データセンターの中に集中化してコンピュータシステム及びストレージシステム構成することにより、コンピュータシステム及びストレージシステムの運用、保守、管理に要する費用を削減する傾向にある。

【０００６】

このような傾向の中で、大型／ハイエンドのディスク制御装置には、数百台以上のホストコンピュータへ接続するためのチャンネルインターフェースのサポート（コネクティビティ）、数百テラバイト以上の記憶容量のサポートが要求されている。

【０００７】

一方、近年のオープン市場の拡大、ストレージ・エリア・ネットワーク（ＳＡＮ）の普及により、大型／ハイエンドのディスク制御装置と同様の高機能・高信頼性を備えた小規模構成（小型筐体）のディスク制御装置への要求が高まっている。

【０００８】

前者の要求に対しては、従来の大型／ハイエンドのディスク制御装置を複数接続して超大規模なストレージシステムを構成する方法が考えられる。

【０００９】

また後者の要求に対しては、従来の大型／ハイエンドのディスク制御装置の最小構成のモデルにおいて筐体を小型化した装置を構成する方法が考えられる。また、この小型化した装置を複数台接続することにより、従来のディスク制御装置がサポートしている中規模から大規模の構成をサポートするストレージシステムを構成する方法が考えられる。

【 0 0 1 0 】

ストレージシステムでは、上記のように、小規模な構成から超大規模な構成まで、同一の高機能・高信頼なアーキテクチャで対応可能な、スケーラビリティのある構成のシステムが必要となっており、そのためには、複数のディスク制御装置をクラスタリングし、1つのシステムとして運用できるストレージシステムが必要となる。

図2に示す従来技術では、複数のディスク制御装置4をS A Nスイッチ5を介してホストコンピュータ3に接続し、S A Nアプライアンス6によりホストコンピュータ3に対して複数のディスク制御装置4を1つのストレージシステムに見せていた。

しかし、S A Nアプライアンス6上で動作するソフトウェアで複数のディスク制御装置4を1つのシステムとして運用するため、従来の単体の大型のディスク制御装置に比べて信頼性、可用性が低いという問題があった。また、S A Nアプライアンス6上でホストコンピュータ3から要求されたデータが存在するディスク制御装置4を検索するため、性能が低下するという問題があった。

また、あるディスク制御装置4のチャンネルI F部11に障害が発生し、このディスク制御装置に繋がるディスク装置のデータにアクセスできなくなった場合、ホストコンピュータ3からそのディスク制御装置4へのアクセスをいったん停止してから、チャンネルI F部11を交換する必要があるため、ホストコンピュータ3上で動作しているアプリケーションプログラムに影響を及ぼすという問題があった。

図23に示す従来技術では、複数のディスク制御クラスタが共有メモリ部25を介して制御情報のやり取りをすることにより、1つのストレージシステムとして動作し、スケーラビリティの高いストレージシステムを提供していた。

しかし、以下のように、ユーザにとって使い勝手が悪い点があった。すなわち、性能を上げるためには、ホストコンピュータ3に割当てて記憶領域は、そのホストコンピュータ3が繋がるディスク制御クラスタに接続されたディスク装置2上の記憶領域にする必要があった。また、ホストコンピュータ3とディスク制御クラスタの間のインターフェースの障害のためにディスク制御クラスタへアクセスすることが不可能になることを防ぐため、1つのホストコンピュータから複数のディスク制御クラスタに接続パスを繋ぐ必要があった。また、複数の接続パスを接続した場合にもホストコンピュータに通知せずに接続パスの切換えを行うことは不可能であった。

【 0 0 1 1 】

本発明の目的は、小規模な構成から超大規模な構成まで、同一の高信頼・高性能なアーキテクチャで対応可能な、スケーラビリティのある構成のストレージシステムを提供することにある。

【 0 0 1 2 】

より具体的には、本発明の目的は、複数台のディスク制御装置をまとめて1つのシステムとしたストレージシステムにおいて高信頼・高性能で使い勝手の良いシステムを提供することにある。

【 0 0 1 3 】

【課題を解決するための手段】

上記目的は、ホストコンピュータとのインターフェースを有する1または複数のチャンネルインターフェース部と、ディスク装置とのインターフェースを有する1または複数のディスクインターフェース部と、前記ディスク装置に対しリード/ライトされるデータと前記データの転送に関する制御情報と前記ディスク装置の管理情報を格納するローカル共有メモリ部とを有し、前記ホストコンピュータからのデータのリード/ライト要求に対し、前記チャンネルインターフェース部は前記ホストコンピュータとのインターフェースと前記ローカル共有メモリ部との間のデータ転送を実行し、前記ディスクインターフェース部は前記ディスク装置と前記ローカル共有メモリ部との間のデータ転送を実行することにより、データのリード/ライトを行う複数のディスク制御クラスタと、前記ディスク装置に対しリード/ライトされるデータと前記各ディスク制御クラスタの管理情報を格納するグローバル情報制御部と、前記複数のディスク制御クラスタを相互接続する相互結合網と、前記

10

20

30

40

50

複数のディスク制御クラスタ内のチャンネルインターフェース部を接続するスイッチを有するストレージシステムであり、

該スイッチは前記グローバル情報制御部に格納された管理情報が複写されたメモリを有するストレージシステムによって達成される。

【 0 0 1 4 】

また、前記各ディスク制御クラスタ内の前記チャンネルインターフェース部と前記ディスクインターフェース部と前記ローカル共有メモリ部とが接続された接続部と他の各ディスク制御クラスタ内の該接続部が前記相互結合網を介して接続され、前記グローバル情報制御部は前記相互結合網と前記スイッチに接続されたストレージシステムによって達成される。

10

【 0 0 1 5 】

また、前記各ディスク制御クラスタ内の前記チャンネルインターフェース部と前記ディスクインターフェース部が前記ローカル共有メモリ部に前記ディスク制御クラスタ内で直接接続され、該各ディスク制御クラスタ内の該ローカル共有メモリ部と他の各ディスク制御クラスタ内の該ローカル共有メモリ部が前記相互結合網を介して接続され、前記グローバル情報制御部は前記相互結合網と前記スイッチに接続されたストレージシステムによって達成される。

また、前記各ディスク制御クラスタ内の前記チャンネルインターフェース部と前記ディスクインターフェース部が前記ローカル共有メモリ部に前記ディスク制御クラスタ内で直接接続され、該各ディスク制御クラスタ内の該チャンネルインターフェース部と該ディスクインターフェース部との接続部と他の各ディスク制御クラスタ内の該接続部が前記相互結合網を介して接続され、前記グローバル情報制御部は前記相互結合網と前記スイッチに接続されたストレージシステムによって達成される。

20

【 0 0 1 6 】

また、ホストコンピュータとのインターフェースを有する 1 または複数のチャンネルインターフェース部と、ディスク装置とのインターフェースを有する 1 または複数のディスクインターフェース部と、前記ディスク装置に対しリード/ライトされるデータと前記データの転送に関する制御情報と前記ディスク装置の管理情報と前記ディスク制御クラスタの管理情報を格納するグローバル情報制御部とを有し、前記ホストコンピュータからのデータのリード/ライト要求に対し、前記チャンネルインターフェース部は前記ホストコンピュータとのインターフェースと前記グローバル情報制御部との間のデータ転送を実行し、前記ディスクインターフェース部は前記ディスク装置と前記グローバル情報制御部との間のデータ転送を実行することにより、データのリード/ライトを行う複数のディスク制御クラスタと、前記複数のディスク制御クラスタを相互接続する相互結合網と、前記複数のディスク制御クラスタ内のチャンネルインターフェース部を接続するスイッチを有するストレージシステムであり、

30

該スイッチは前記グローバル情報制御部に格納された管理情報が複写されたメモリを有するストレージシステムによって達成される。

また、前記各ディスク制御クラスタ内の前記チャンネルインターフェース部と前記ディスクインターフェース部との接続部と他の各ディスク制御クラスタ内の該接続部が前記相互結合網を介して接続され、前記グローバル情報制御部は前記相互結合網と前記スイッチに接続されたストレージシステムによって達成される。

40

【 0 0 1 7 】

また、ホストコンピュータとのインターフェースを有する 1 または複数のチャンネルインターフェース部と、ディスク装置とのインターフェースを有する 1 または複数のディスクインターフェース部と、前記ディスク装置に対しリード/ライトされるデータを格納する第 1 のメモリと、前記チャンネルインターフェース部及び前記ディスクインターフェース部と前記第 1 のメモリとの間のデータ転送に関する制御情報及び前記ディスク装置の管理情報を格納する第 2 のメモリとを有するローカル共有メモリ部とを有し、前記ホストコンピュータからのデータのリード/ライト要求に対し、前記チャンネルインターフェース部は前記

50

ホストコンピュータとのインターフェースと前記ローカル共有メモリ部内の前記第 1 のメモリとの間のデータ転送を実行し、前記ディスクインターフェース部は前記ディスク装置と前記ローカル共有メモリ部内の前記第 1 のメモリとの間のデータ転送を実行することにより、データのリード/ライトを行う複数のディスク制御クラスタと、前記各ディスク制御クラスタの管理情報を格納するグローバル情報制御部と、前記複数のディスク制御クラスタを相互接続する 2 つの異なる第 1、第 2 の相互結合網と、前記複数のディスク制御クラスタ内のチャネルインターフェース部を接続するスイッチを有するストレージシステムであり、

該スイッチは前記グローバル情報制御部に格納された管理情報が複写されたメモリを有するストレージシステムによって達成される。

10

【0018】

また、前記各ディスク制御クラスタ内の前記チャネルインターフェース部と前記ディスクインターフェース部が前記ローカル共有メモリ部内の第 2 のメモリに前記ディスク制御クラスタ内で直接接続され、該各ディスク制御クラスタ内の該チャネルインターフェース部と該ディスクインターフェース部との第 1 の接続部と他の各ディスク制御クラスタ内の該第 1 の接続部が前記第 1 の相互結合網を介して接続され、前記各ディスク制御クラスタ内の前記チャネルインターフェース部と前記ディスクインターフェース部と前記ローカル共有メモリ部内の第 1 のメモリとが接続された第 2 の接続部と他の各ディスク制御クラスタ内の該第 2 の接続部が前記第 2 の相互結合網を介して接続され、前記グローバル情報制御部は前記第 1 の相互結合網と前記スイッチに接続されたストレージシステムによって達成

20

【0019】

その他、本願が開示する課題、及びその解決方法は、発明の実施形態の欄及び図面により明らかにされる。

【0020】

【発明の実施の形態】

以下、本発明の実施例を図面を用いて説明する。

[実施例 1]

図 1、図 3、図 12、及び図 13 に、本発明の一実施例を示す。

以下の実施例において、相互結合網はスイッチを利用したものを例にして説明してあるが、相互に接続され制御情報やデータが転送されれば良いのであり、例えばバスで構成されても良い。

30

図 1 に示すように、ストレージシステム 1 は複数のディスク制御クラスタ 1-1 乃至 1-n とフロントエンドスイッチ 7 から構成される。

ディスク制御クラスタ 1-1 は、ホストコンピュータ 3 とのインターフェース部（チャネル I/F 部）11 と、ディスク装置 2 とのインターフェース部（ディスク I/F 部）16 と、ローカル共有メモリ部 22 を有し、チャネル I/F 部 11 及びディスク I/F 部 16 とローカル共有メモリ部 22 の間は複数のディスク制御クラスタ 1-1 乃至 1-n に跨る相互結合網 31 を介して接続され、グローバル情報制御部 21 は相互結合網 31 に接続されている。すなわち、相互結合網 31 を介して、全てのチャネル I/F 部 11 及びディスク I/F 部 12 から、グローバル情報制御部 21 へアクセス可能な構成となっている。

40

ホストコンピュータ 3 は、フロントエンドスイッチ 7 を介してディスク制御クラスタに接続され、任意のホストコンピュータ 3 から任意のディスク制御クラスタへアクセス可能な構成となっている。

チャネル I/F 部 11 の具体的な一例を図 12 に示す。

チャネル I/F 部 11 は、ホストコンピュータ 3 との 2 つの I/F（ホスト I/F）202 と、ホストコンピュータ 3 に対する入出力を制御する 2 つのマイクロプロセッサ 201 と、グローバル情報制御部 21 あるいはローカル共有メモリ部 22 へのアクセスを制御するアクセス制御部（メモリアクセス制御部）206 を有し、ホストコンピュータ 3 とグローバル情報制御部 21 あるいはローカル共有メモリ部 22 間のデータ転送、及びマイクロプロセ

50

ッサ 201 とグローバル情報制御部 21 あるいはローカル共有メモリ部 22 間の制御情報の転送を実行する。マイクロプロセッサ 201 及びホスト I/F 202 は内部バス 205 によって接続され、メモリアクセス制御部 206 は 2 つのホスト I/F 202 に直接接続され、また内部バス 205 に接続されている。

ディスク I/F 部 16 の具体的な一例を図 13 に示す。

ディスク I/F 部 16 は、ディスク装置 2 との 2 つの I/F (ドライブ I/F) 203 と、ディスク装置 2 に対する入出力を制御する 2 つのマイクロプロセッサ 201 と、グローバル情報制御部 21 あるいはローカル共有メモリ部 22 へのアクセスを制御するアクセス制御部 (メモリアクセス制御部) 206 を有し、ディスク装置 2 とグローバル情報制御部 21 あるいはローカル共有メモリ部 22 間のデータ転送、及びマイクロプロセッサ 201 とグローバル情報制御部 21 あるいはローカル共有メモリ部 22 間の制御情報の転送を実行する。マイクロプロセッサ 201 及びドライブ I/F 203 は内部バス 205 によって接続され、メモリアクセス制御部 206 は 2 つのドライブ I/F 203 に直接接続され、また、内部バス 205 に接続されている。ディスク I/F 部 16 は R A I D 機能の実行も行う。

1 つのディスク制御クラスタは 1 つの筐体として構成されるか、またはモジュールとして構成されても良いが、それ自体 1 つのディスク制御装置として機能を備えているものである。

ストレージシステムの具体的な一例を図 3 に示す。

ストレージシステム 1 は、2 つのフロントエンドスイッチ 7 と、複数のディスク制御クラスタ 1 - 1 乃至 1 - n と、グローバル情報制御部 21 と、2 つのグローバルスイッチ (G S W) 115 と、アクセスバス 136 と、アクセスバス 137 を有する。

グローバルスイッチ (G S W) 115 は、グローバル情報制御部 21 からのバスと複数のディスク制御クラスタからのバスを接続する接続部である。

ディスク制御クラスタ 1 - 1 乃至 1 - n は、ホストコンピュータ 3 との 2 つのチャンネル I/F 部 11 と、ディスク装置 2 との 2 つのディスク I/F 部 16 と、2 つのローカルスイッチ (L S W) 110 と、2 つのローカル共有メモリ部 22 と、アクセスバス 131 と、アクセスバス 132 と、アクセスバス 136 を有する。

フロントエンドスイッチ 7 は、スイッチ 71 と、スイッチ制御部 72 と、メモリコントローラ 73 と、メモリモジュール 105 を有する。

1 つのホストコンピュータ 3 からは、2 つのフロントエンドスイッチ 7 内のスイッチ 71 にバスを 1 本ずつ接続し、スイッチ 71 からは、ディスク制御クラスタ 1 - 1 乃至 1 - n の各チャンネル I/F 部 11 に 1 本ずつバスを接続する。

【0021】

グローバル情報制御部 21 は、アクセス制御部 101 と管理機能部 102 とメモリモジュール 105 とを有し、ディスク制御クラスタ 1 - 1 乃至 1 - n の管理情報 (例えば、各ディスク制御クラスタが管理する記憶領域の情報や、ディスク制御クラスタ内の各部位の負荷情報、障害情報、及び構成情報等) を格納する。ローカルスイッチ (L S W) 110 は、チャンネル I/F 部 11 からのバスと、ディスク I/F 部 16 からのバスと、ローカル共有メモリ部 22 からのバスを接続する接続部である。

ローカル共有メモリ部 22 は、メモリコントローラ 100 とメモリモジュール 105 とを有し、ディスク制御クラスタの制御情報 (例えば、チャンネル I/F 部 11 及びディスク I/F 部 16 とローカル共有メモリ部 22 との間のデータ転送制御に関する情報、ディスク装置 2 に記録するデータの管理情報等) とディスク装置 2 に記録するデータを格納する。

チャンネル I/F 部 11 内のメモリアクセス制御部 206 には 2 本のアクセスバス 131 を接続し、それらを 2 つの異なる L S W 110 にそれぞれ接続する。

L S W 110 には 2 本のアクセスバス 132 を接続し、それらを 2 つの異なるローカル共有メモリ部 22 内のメモリコントローラ 100 にそれぞれ接続する。

したがって、メモリコントローラ 100 には、2 つの L S W 110 から 1 本ずつ、計 2 本のアクセスバス 132 が接続される。

こうすることにより、1 つのメモリアクセス制御部 206 から 1 つのメモリコントローラ

10

20

30

40

50

100へのアクセスルートが2つとなる。

これにより、1つのアクセスパスまたはLSW110に障害が発生した場合でも、もう1つのアクセスルートによりローカル共有メモリ部22へアクセスすることが可能となるため、耐障害性を向上させることができる。

LSW110には、2つのチャンネルIF部11と、2つのディスクIF部16からそれぞれ1本ずつ、計4本のアクセスパス131が接続される。

また、LSW110には、2つのローカル共有メモリ部22へのアクセスパス132が2本とGSW115へのアクセスパス136が1本接続される。

LSW110には上記のようなアクセスパスが接続されるため、LSW110内では、チャンネルIF部11及びディスクIF部16からの4本のアクセスパスからの要求を、自己ディスク制御クラスタ内のローカル共有メモリ部22への2本のアクセスパスと、GSW115への1本のアクセスパス136に振分ける機能を有する。

GSW115には、各ディスク制御クラスタから1本ずつ、ディスク制御クラスタ数分の本数のアクセスパス136が接続される。

また、GSW115には、2つのグローバル情報制御部21内のアクセス制御部101へのアクセスパス137が1本ずつ、計2本接続される。

こうすることにより、1つのメモリアクセス制御部206から1つのアクセス制御部101へのアクセスルートが2つとなる。

これにより、1つのアクセスパスまたはLSW110またはGSW115に障害が発生した場合でも、もう1つのアクセスルートによりグローバル情報制御部21へアクセスすることが可能となるため、耐障害性を向上させることができる。GSW115を使わずに、アクセスパス136をメモリコントローラ101に直接接続しても本発明を実施する上で問題ない。そうすることにより、GSW115で発生するデータ転送処理のオーバーヘッドを削減することが可能となり、性能が向上する。

GSW115を使わない場合、1つのメモリアクセス制御部206から1つのメモリコントローラ101へのアクセスルートを2つ確保し、耐障害性を向上するためには、LSW110にアクセスパス136を2本接続し、それぞれを異なるメモリコントローラ101へ接続する。

また、GSW115には、2つのフロントエンドスイッチ7内のメモリコントローラ73へのパスが2本接続される。

フロントエンドスイッチ7内のメモリコントローラ73には、スイッチ制御部72とメモリモジュール105が接続される。

スイッチ制御部72はスイッチ71に接続され、メモリモジュール105に格納されたホストコンピュータ3のチャンネルとディスク制御クラスタのチャンネルの接続を示すルーティングテーブルを参照し、スイッチの切換えを制御する。

また、メモリモジュール105は、グローバル情報制御部21内のメモリモジュール105に格納されたディスク制御クラスタ1 - 1乃至1 - n内の各部位の負荷情報、障害情報、および記憶領域情報のコピーを有し、スイッチ制御部72は定期的あるいは必要に際してこれらの情報をもとにルーティングテーブルを変更する。

具体的な例としては、グローバル情報制御部21は、図18に示すようなシステム構成・稼動状況テーブル500をそのメモリモジュール105に格納しており、フロントエンドスイッチ7は、テーブル500のコピーをメモリモジュール105に格納している。

システム構成・稼動状況テーブル500は、ディスク制御クラスタを識別するクラスタ番号511とディスク制御クラスタのチャンネルを識別するチャンネル番号512に対応する論理ボリューム番号513とチャンネル稼動状況514を示す。

本実施例では、クラスタ番号511、チャンネル番号512、及び論理ボリューム番号513を16進数で示している。チャンネル稼動状況514は、低負荷状態を‘0’、中負荷状態を‘1’、高負荷状態を‘2’、障害を‘3’で示している。

また、フロントエンドスイッチ7は、ホスト - 論理ボリューム対応テーブル600をそのメモリモジュール105に格納している。

10

20

30

40

50

ホスト - 論理ボリューム対応テーブル 600 は、個々のホストコンピュータを識別するホスト番号 615 に対応する論理ボリューム番号 513、すなわち各ホストコンピュータに割当てられた論理ボリュームと、論理ボリューム番号 513 にアクセスするためのチャンネル番号 512 と該当論理ボリュームを管理するディスク制御クラスタのクラスタ番号 511 を示す。

ホスト番号 615 は、例えば、ファイバチャネルのプロトコルで使われるワールドワイドネーム (WWN) や、インターネットプロトコルで使われるマックアドレスや IP アドレスを充てることが考えられる。

スイッチ制御部 72 は、ホスト - 論理ボリューム対応テーブル 600 (ルーティングテーブル) を参照し、スイッチ 71 の切り替えを行う。

10

図 3 で LSW110 はチャンネル IF 部 11 及びディスク IF 部 16 とローカル共有メモリ部 22 との接続部であり、GSW115 はディスク制御クラスタ 1 - 1 乃至 1 - n とグローバル情報制御部 21 との接続部である。

図 3 において、GSW115 とグローバル情報制御部 21 をボックスに実装し、フロントエンドスイッチ 7 を別のボックスに実装し、モジュール化した各ディスク制御クラスタ 1 - 1 乃至 1 - n といっしょに、1 つの筐体の中に実装しても良い。また、各ディスク制御クラスタ 1 - 1 乃至 1 - n を別個の筐体として、距離的に離れた場所に分散しても良い。

【0022】

図 3 において、ホストコンピュータ 3 からストレージシステム 1 に記録されたデータを読み出す場合の一例を述べる。

20

まず、ホストコンピュータ 3 は、ストレージシステム 1 に対してデータの読出し要求を発行する。

要求は、フロントエンドスイッチ 7 に受け取られ、フロントエンドスイッチ 7 内のスイッチ制御部 72 は、要求パケットのヘッダを解析する。

要求パケットのヘッダには、要求を発行したホストコンピュータの番号 (ホスト番号) と要求データが記録されている論理ボリューム番号が格納されており、スイッチ制御部 72 はホスト - 論理ボリューム対応テーブル 600 を参照し、要求を発行したホストコンピュータが接続されたスイッチ 71 のポートを、該当する論理ボリューム番号に割当てられているチャンネル番号のスイッチ 71 のポートに接続し、要求パケットをディスク制御クラスタに送る。

30

要求が送られたチャンネルが接続されたチャンネル IF 部 11 内のマイクロプロセッサ 201 は、自ディスク制御クラスタ 1 - 1 内のローカル共有メモリ部 22 にアクセスし、要求されたデータがどのディスク装置 2 内に格納されているかを調べる。

ローカル共有メモリ部 22 には、要求データのアドレスとそのデータが実際に記録されているディスク装置 2 内のアドレスを対応させる変換テーブルが格納されており、要求されたデータがどのディスク装置 2 内に格納されているかを調べることができる。

さらに、要求を受けたチャンネル IF 部 11 内のマイクロプロセッサ 201 は、自ディスク制御クラスタ 1 - 1 内のローカル共有メモリ部 22 にアクセスし、要求されたデータがローカル共有メモリ部 22 内に格納されているかどうかを確認する。

ローカル共有メモリ部 22 にはディスク装置 2 に格納するデータとともにそのデータのディレクトリ情報が格納されており、ローカル共有メモリ部 22 内に要求データが存在するかどうかを確認できる。

40

【0023】

それにより自ディスク制御クラスタ 1 - 1 のローカル共有メモリ部 22 内にデータがあった場合は、ローカル共有メモリ部 22 にアクセスしてそのデータを自身の LSW110 を介してチャンネル IF 部 11 まで転送し、フロントエンドスイッチ 7 を介して、ホストコンピュータ 3 に送る。

自ディスク制御クラスタ 1 - 1 のローカル共有メモリ部 22 内にデータが存在しなかった場合は、チャンネル IF 部 11 のマイクロプロセッサ 201 は、要求データが格納されているディスク装置 2 が接続されているディスク IF 部 16 内のマイクロプロセッサ 201 に

50

対し、要求データを読み出しローカル共有メモリ部22に格納するというデータ要求の処理内容を示す制御情報を発行し、この制御情報の発行を受けたディスクIF部16内のマイクロプロセッサ201は、要求データが格納されているディスク装置2からデータを読み出し、LSW110を介して、自ディスク制御クラスタ1-1内のローカル共有メモリ部22に要求データを転送し格納する。

すなわち、チャンネルIF部11のマイクロプロセッサ201は、上記データ要求の処理内容を示す制御情報を発行し、ローカル共有メモリ部22の制御情報領域(ジョブ制御ブロック)に格納する。

ディスクIF部16のマイクロプロセッサ201は、ローカル共有メモリ部22の制御情報領域をポーリングで監視し、上記発行された制御情報が上記制御情報領域(ジョブ制御ブロック)に存在した場合は、要求データが格納されているディスク装置2からデータを読み出し、LSW110を介して、自ディスク制御クラスタ1-1内のローカル共有メモリ部22に要求データを転送し格納する。

ディスクIF部16のマイクロプロセッサ201は、要求データをローカル共有メモリ部22へ格納した後、前記制御情報を発行したチャンネルIF部11のマイクロプロセッサ201にローカル共有メモリ部22内のデータを格納したアドレスを、ローカル共有メモリ部22内の制御情報を介して伝える。それを受けたチャンネルIF部11のマイクロプロセッサ201は、ローカル共有メモリ部22からデータを読み出し、フロントエンドスイッチ7を介して、ホストコンピュータ3へ送る。

すなわち、ディスクIF部16のマイクロプロセッサ201は、要求データをローカル共有メモリ部22へ格納した後、処理の実行の終了とデータを格納したアドレスを示す制御情報を発行し、ローカル共有メモリ部22の制御情報領域に格納する。

前記制御情報を発行したチャンネルIF部11のマイクロプロセッサ201は、ローカル共有メモリ部22の制御情報領域をポーリングで監視し、ディスクIF部16のマイクロプロセッサ201から発行された制御情報が上記制御情報領域に存在した場合は、ローカル共有メモリ部内のデータを格納したアドレスによりローカル共有メモリ部22からデータを読み出し、チャンネルIF部11まで転送し、さらにフロントエンドスイッチ7を介して、ホストコンピュータ3へ送る。

本実施例によれば、ホストコンピュータ3はフロントエンドスイッチ7のどの接続ポートに接続しても、ストレージシステム1を構成するディスク制御クラスタを意識することなく、その接続ポートにアクセス要求を発行するだけで、データの書き込み及び読み出しを行うことが可能になり、ホストコンピュータ3に対して、複数台のディスク制御クラスタ1-1乃至1-nを1つのストレージシステムに見せることが可能となる。

そして、ディスク制御クラスタが1個だけの小規模な構成からディスク制御クラスタが数十個接続された超大規模な構成まで、ディスク制御クラスタ単体が持つ高信頼・高性能なアーキテクチャで対応可能な、スケーラビリティがあり使い勝手の良い構成のストレージシステムを提供することが可能となる。

【実施例2】

図4、図5、図12、及び図13に、本発明の一実施例を示す。

【0024】

図4に示すように、ディスク制御ユニット1-1乃至1-nとフロントエンドスイッチ7からなるストレージシステム1の構成は、チャンネルIF部11及びディスクIF部16とローカル共有メモリ部22及び相互結合網31の間の接続構成を除いて、実施例1の図1に示す構成と同様である。

チャンネルIF部11及びディスクIF部16とローカル共有メモリ部22の間は、ディスク制御クラスタ内では直接接続されている。

また、複数のディスク制御クラスタ1-1乃至1-n間では、ローカル共有メモリ部22は相互結合網31を介して接続されており、その相互結合網31にグローバル情報制御部21が接続されている。

上記のように、この実施例ではディスク制御ユニット1-1乃至1-n内においてチャネ

10

20

30

40

50

ル I F 部 1 1 及びディスク I F 部 1 6 とローカル共有メモリ部 2 2 を直接接続することにより、実施例 1 で示した相互結合網 3 1 を介して接続する場合に比べ、ローカル共有メモリ部 2 2 へのアクセス時間を短縮することが可能になる。

チャンネル I F 部 1 1 及びディスク I F 部 1 6 の構成は、それぞれ図 1 2、図 1 3 に示す実施例 1 の構成と同様である。

1 つのディスク制御クラスタは 1 つの筐体として構成されるか、またはモジュールとして構成されても良いが、それ自体 1 つのディスク制御装置として機能を備えているものである。

ストレージシステム 1 の具体的な一例を図 5 に示す。

【 0 0 2 5 】

ディスク制御クラスタ 1 - 1 乃至 1 - n 内の構成も、チャンネル I F 部 1 1 及びディスク I F 部 1 6 とローカル共有メモリ部 2 2 の間の接続構成とディスク制御クラスタ 1 - 1 乃至 1 - n と G S W 1 1 5 の接続構成を除いて、実施例 1 の図 3 に示す構成と同様である。

ストレージシステム 1 は、複数のディスク制御クラスタ 1 - 1 乃至 1 - n と、フロントエンドスイッチ 7 と、グローバル情報制御部 2 1 と、2 つのグローバルスイッチ (G S W) 1 1 5 と、アクセスパス 1 3 6 と、アクセスパス 1 3 7 を有する。

ディスク制御クラスタ 1 - 1 乃至 1 - n は、ホストコンピュータ 3 との 2 つのチャンネル I F 部 1 1 と、ディスク装置 2 との 2 つのディスク I F 部 1 6 と、2 つのローカル共有メモリ部 2 2 と、アクセスパス 1 3 3 と、アクセスパス 1 3 6 を有する。

チャンネル I F 部 1 1 内のメモリアクセス制御部 2 0 6 には 2 本のアクセスパス 1 3 3 を接続し、それらを 2 つの異なるメモリコントローラ 1 0 0 にそれぞれ接続する。

したがって、メモリコントローラ 1 0 0 には、2 つのチャンネル I F 部 1 1 と 2 つのディスク I F 部 1 6 から 1 本ずつ、計 4 本のアクセスパス 1 3 3 が接続される。また、G S W 1 1 5 へのアクセスパス 1 3 6 が 1 本接続される。

メモリコントローラ 1 0 0 には上記のようなアクセスパスが接続されるため、メモリコントローラ 1 0 0 内では、チャンネル I F 部 1 1 及びディスク I F 部 1 6 からの 4 本のアクセスパス 1 3 3 からの要求を、メモリモジュール 1 0 5 への 1 本のアクセスパスと、G S W 1 1 5 への 1 本のアクセスパス 1 3 6 に振分ける機能を有する。

実施例 1 と同様に G S W 1 1 5 を使わずに、アクセスパス 1 3 6 をアクセス制御部 1 0 1 に直接接続しても本発明を実施する上で問題ない。そうすることにより、G S W 1 1 5 で発生するデータ転送処理のオーバヘッドを削減することが可能となり、性能が向上する。G S W 1 1 5 を使わない場合、1 つのメモリコントローラ 1 0 0 から 1 つのアクセス制御部 1 0 1 へのアクセスルートを 2 つ確保し、耐障害性を向上するためには、メモリコントローラ 1 0 0 にアクセスパス 1 3 6 を 2 本接続し、それぞれを異なるアクセス制御部 1 0 1 へ接続する。

また実施例 1 と同様に、図 5 において、G S W 1 1 5 とグローバル情報制御部 2 1 をボックスに実装し、フロントエンドスイッチ 7 を別のボックスに実装し、モジュール化した各ディスク制御クラスタ 1 - 1 乃至 1 - n といっしょに、1 つの筐体の中に実装しても良い。また、各ディスク制御クラスタ 1 - 1 乃至 1 - n を別個の筐体として、距離的に離れた場所に分散しても良い。

【 0 0 2 6 】

本実施例において、ホストコンピュータ 3 からストレージシステム 1 へのデータの読み出し / 書き込みを行う場合の、ストレージシステム 1 内の各部の動作は、チャンネル I F 部 1 1 及びディスク I F 部 1 6 からローカル共有メモリ部 2 2 へのアクセスが直接になることと、チャンネル I F 部 1 1 及びディスク I F 部 1 6 からグローバル情報制御部 2 1 へのアクセスがメモリコントローラ 1 0 0 を介して行われることを除いて、実施例 1 と同様である。

本実施例によれば、ホストコンピュータ 3 はフロントエンドスイッチ 7 のどの接続ポートに接続しても、ストレージシステム 1 を構成するディスク制御クラスタを意識することなく、その接続ポートにアクセス要求を発行するだけで、データの書き込み及び読出しを行

10

20

30

40

50

うことが可能になり、ホストコンピュータ 3 に対して、複数台のディスク制御クラスタ 1 - 1 乃至 1 - n を 1 つのストレージシステムに見せることが可能となる。

そして、ディスク制御クラスタが 1 個だけの小規模な構成からディスク制御クラスタが数十個接続された超大規模な構成まで、ディスク制御クラスタ単体が持つ高信頼・高性能なアーキテクチャで対応可能な、スケーラビリティがあり使い勝手の良い構成のストレージシステムを提供することが可能となる。

【実施例 3】

図 6、図 7、図 12、及び図 13 に、本発明の一実施例を示す。

【0027】

図 6 に示すように、ディスク制御ユニット 1 - 1 乃至 1 - n とフロントエンドスイッチ 7 からなるストレージシステム 1 の構成は、チャンネル I/F 部 12 及びディスク I/F 部 17 とローカル共有メモリ部 22 の間の接続構成を除いて、実施例 1 の図 1 に示す構成と同様である。

チャンネル I/F 部 12 及びディスク I/F 部 17 とローカル共有メモリ部 22 の間は、ディスク制御クラスタ内では直接接続されている。

また、複数のディスク制御クラスタ 1 - 1 乃至 1 - n 間では、チャンネル I/F 部 12 及びディスク I/F 部 17 が相互結合網 31 を介して接続されており、その相互結合網 31 にグローバル共有メモリ 21 が接続されている。

上記のように、この実施例ではディスク制御ユニット 1 - 1 乃至 1 - n 内においてチャンネル I/F 部 12 及びディスク I/F 部 17 とローカル共有メモリ部 22 を直接接続することにより、実施例 1 で示した相互結合網 31 を介して接続する場合に比べ、ローカル共有メモリ部 22 へのアクセス時間を短縮することが可能になる。

チャンネル I/F 部 12 及びディスク I/F 部 17 の構成は、それぞれ図 12、図 13 に示すチャンネル I/F 部 11 及びディスク I/F 部 16 の構成において、メモリアクセス制御部 206 のアクセスパスを 4 本に増やした構成となる。

ここで、4 本のアクセスパスの内、2 本はアクセスパス 131、もう 2 本がアクセスパス 133 となる。

1 つのディスク制御クラスタは 1 つの筐体として構成されるか、またはモジュールとして構成されても良いが、それ自体 1 つのディスク制御装置として機能を備えているものである。

ストレージシステム 1 の具体的な一例を図 7 に示す。

【0028】

ディスク制御クラスタ 1 - 1 乃至 1 - n 内の構成も、チャンネル I/F 部 12 及びディスク I/F 部 17 とローカル共有メモリ部 22 の間の接続構成を除いて、実施例 1 の図 3 に示す構成と同様である。

ストレージシステム 1 は、複数のディスク制御クラスタ 1 - 1 乃至 1 - n と、フロントエンドスイッチ 7 と、グローバル情報制御部 21 と、2 つのグローバルスイッチ (G S W) 115 と、アクセスパス 136 と、アクセスパス 137 を有する。

ディスク制御クラスタ 1 - 1 乃至 1 - n は、ホストコンピュータ 3 との 2 つのチャンネル I/F 部 12 と、ディスク装置 2 との 2 つのディスク I/F 部 17 と、2 つのローカルスイッチ (L S W) 110 と、2 つのローカル共有メモリ部 22 と、アクセスパス 131 と、アクセスパス 133 と、アクセスパス 136 を有する。

L S W 110 は、チャンネル I/F 部 12 からのパスと、ディスク I/F 部 17 からのパスを接続する接続部である。

【0029】

チャンネル I/F 部 12 及びディスク I/F 部 17 内のメモリアクセス制御部 206 には 2 本のアクセスパス 133 を接続し、それらを 2 つの異なるメモリコントローラ 100 にそれぞれ接続する。したがって、メモリコントローラ 100 には、2 つのチャンネル I/F 部 11 と 2 つのディスク I/F 部 16 から 1 本ずつ、計 4 本のアクセスパス 133 が接続される。

さらに、チャンネル I/F 部 12 及びディスク I/F 部 17 内のメモリアクセス制御部 206 に

10

20

30

40

50

は2本のアクセスパス131を接続し、それらを2つの異なるLSW110にそれぞれ接続する。したがって、LSW110には、2つのチャンネルIF部12と2つのディスクIF部17から1本ずつ、計4本のアクセスパス131が接続される。また、GSW115へのアクセスパス136が1本接続される。実施例1と同様にGSW115を使わずに、アクセスパス136をアクセス制御部101に直接接続しても本発明を実施する上で問題ない。そうすることにより、GSW115で発生するデータ転送処理のオーバーヘッドを削減することが可能となり、性能が向上する。

GSW115を使わない場合、1つのLSW110から1つのアクセス制御部101へのアクセスルートを2つ確保し、耐障害性を向上するためには、LSW110にアクセスパス136を2本接続し、それぞれを異なるメモリコントローラ101へ接続する。

10

また実施例1と同様に、図7において、GSW115とグローバル情報制御部21をボックスに実装し、フロントエンドスイッチ7を別のボックスに実装し、モジュール化した各ディスク制御クラスタ1-1乃至1-nといっしょに、1つの筐体の中に実装しても良い。また、各ディスク制御クラスタ1-1乃至1-nを別個の筐体として、距離的に離れた場所に分散しても良い。

【0030】

本実施例において、ホストコンピュータ3からストレージシステム1へのデータの読み出し/書き込みを行う場合の、ストレージシステム1内の各部の動作は、チャンネルIF部12及びディスクIF部17からローカル共有メモリ部22へのアクセスが直接になることを除いて、実施例1と同様である。

20

本実施例によれば、ホストコンピュータ3はフロントエンドスイッチ7のどの接続ポートに接続しても、ストレージシステム1を構成するディスク制御クラスタを意識することなく、その接続ポートにアクセス要求を発行するだけで、データの書き込み及び読出しを行うことが可能になり、ホストコンピュータ3に対して、複数台のディスク制御クラスタ1-1乃至1-nを1つのストレージシステムに見せることが可能となる。

そして、ディスク制御クラスタが1個だけの小規模な構成からディスク制御クラスタが数十個接続された超大規模な構成まで、ディスク制御クラスタ単体が持つ高信頼・高性能なアーキテクチャで対応可能な、スケーラビリティがあり使い勝手の良い構成のストレージシステムを提供することが可能となる。

【実施例4】

30

図8、図9、図12、及び図13に、本発明の一実施例を示す。

【0031】

図8に示すように、ディスク制御ユニット1-1乃至1-nとフロントエンドスイッチ7からなるストレージシステム1の構成は、実施例1においてローカル共有メモリ22を除いた構成である。

【0032】

このため、実施例1の各ディスク制御ユニット1-1乃至1-nのローカル共有メモリ22に格納する情報を全てグローバル情報制御部21に格納する。

【0033】

複数のディスク制御クラスタ1-1乃至1-n間では、チャンネルIF部11及びディスクIF部16が相互結合網31を介して接続されており、その相互結合網31にグローバル情報制御部21が接続されている。

40

チャンネルIF部11及びディスクIF部16の構成は、それぞれ図12、図13に示す実施例1の構成と同様である。

1つのディスク制御クラスタは1つの筐体として構成されるか、またはモジュールとして構成されても良い。

ストレージシステム1の具体的な一例を図9に示す。

【0034】

ディスク制御クラスタ1-1乃至1-n内の構成も、ローカル共有メモリ部22が無いことを除いて、実施例1の図3に示す構成と同様である。

50

ストレージシステム 1 は、複数のディスク制御クラスタ 1 - 1 乃至 1 - n と、フロントエンドスイッチ 7 と、グローバル情報制御部 2 1 と、2 つのグローバルスイッチ (G S W) 1 1 5 と、アクセスパス 1 3 6 と、アクセスパス 1 3 7 を有する。

ディスク制御クラスタ 1 - 1 乃至 1 - n は、ホストコンピュータ 3 との 2 つのチャンネル I F 部 1 1 と、ディスク装置 2 との 2 つのディスク I F 部 1 6 と、2 つのローカルスイッチ (L S W) 1 1 0 と、アクセスパス 1 3 1 と、アクセスパス 1 3 6 を有する。

L S W 1 1 0 は、チャンネル I F 部 1 1 からのパスとディスク I F 部 1 6 からのパスを接続する接続部である。

チャンネル I F 部 1 1 及びディスク I F 部 1 6 内のメモリアクセス制御部 2 0 6 には 2 本のアクセスパス 1 3 1 を接続し、それらを 2 つの異なる L S W 1 1 0 にそれぞれ接続する。したがって、L S W 1 1 0 には、2 つのチャンネル I F 部 1 1 と 2 つのディスク I F 部 1 6 から 1 本ずつ、計 4 本のアクセスパス 1 3 1 が接続される。また、G S W 1 1 5 へのアクセスパス 1 3 6 が 1 本接続される。

実施例 1 と同様に G S W 1 1 5 を使わずに、アクセスパス 1 3 6 をアクセス制御部 1 0 1 に直接接続しても本発明を実施する上で問題ない。そうすることにより、G S W 1 1 5 で発生するデータ転送処理のオーバーヘッドを削減することが可能となり、性能が向上する。G S W 1 1 5 を使わない場合、1 つの L S W 1 1 0 から 1 つのアクセス制御部 1 0 1 へのアクセスルートを 2 つ確保し、耐障害性を向上するためには、L S W 1 1 0 にアクセスパス 1 3 6 を 2 本接続し、それぞれを異なるアクセス制御部 1 0 1 へ接続する。

また実施例 1 と同様に、図 9 において、G S W 1 1 5 とグローバル情報制御部 2 1 をボックスに実装し、フロントエンドスイッチ 7 を別のボックスに実装し、モジュール化した各ディスク制御クラスタ 1 - 1 乃至 1 - n といっしょに、1 つの筐体の中に実装しても良い。また、各ディスク制御クラスタ 1 - 1 乃至 1 - n を別個の筐体として、距離的に離れた場所に分散しても良い。

【 0 0 3 5 】

本実施例において、ホストコンピュータ 3 からストレージシステム 1 へのデータの読み出し / 書き込みを行う場合の、ストレージシステム 1 内の各部の動作は、実施例 1 の処理においてローカル共有メモリ 2 2 における処理を全てグローバル情報制御部 2 1 で行うことを除いて、実施例 1 と同様である。

本実施例によれば、ホストコンピュータ 3 はフロントエンドスイッチ 7 のどの接続ポートに接続しても、ストレージシステム 1 を構成するディスク制御クラスタを意識することなく、その接続ポートにアクセス要求を発行するだけで、データの書き込み及び読出しを行うことが可能になり、ホストコンピュータ 3 に対して、複数台のディスク制御クラスタ 1 - 1 乃至 1 - n を 1 つのストレージシステムに見せることが可能となる。

そして、ディスク制御クラスタが 1 個だけの小規模な構成からディスク制御クラスタが数十個接続された超大規模な構成まで、ディスク制御クラスタ単体が持つ高信頼・高性能なアーキテクチャで対応可能な、スケーラビリティがあり使い勝手の良い構成のストレージシステムを提供することが可能となる。

[実施例 5]

図 1 0 に、本発明の一実施例を示す。

以下の実施例において、相互結合網はスイッチを利用したものを例にして説明してあるが、相互に接続され制御情報やデータが転送されれば良いのであり、例えばバスで構成されても良い。

図 1 0 に示すように、ストレージシステム 1 は複数のディスク制御クラスタ 1 - 1 乃至 1 - n とフロントエンドスイッチ 7 から構成される。

ディスク制御クラスタ 1 - 1 乃至 1 - n は、ホストコンピュータ 3 とのインターフェース部 (チャンネル I F 部) 1 3 と、ディスク装置 2 とのインターフェース部 (ディスク I F 部) 1 8 と、メモリ 1 : 2 5 とメモリ 2 : 2 6 を有するローカル共有メモリ部 2 2 を有し、チャンネル I F 部 1 3 及びディスク I F 部 1 8 とメモリ 2 の間は、ディスク制御クラスタ内部では直接接続される。

また、チャンネルＩＦ部１３及びディスクＩＦ部１８は複数のディスク制御クラスタ１－１乃至１－ｎに跨る相互結合網１：３２を介して接続され、グローバル情報制御部２１は相互結合網１：３２に接続される。すなわち、相互結合網１：３２を介して、全てのチャンネルＩＦ部１３及びディスクＩＦ部１８から、グローバル情報制御部２１へアクセス可能な構成となっている。

また、チャンネルＩＦ部１３及びディスクＩＦ部１８とメモリ１の間は、複数のディスク制御クラスタ１－１乃至１－ｎに跨る相互結合網２：３３を介して接続される。

ホストコンピュータ３は、フロントエンドスイッチ７を介してディスク制御クラスタに接続され、任意のホストコンピュータ３から任意のディスク制御クラスタへアクセス可能な構成となっている。

10

チャンネルＩＦ部１３の具体的な一例を図１４に示す。

チャンネルＩＦ部１３は、ホストコンピュータ３との２つのＩＦ（ホストＩＦ）２０２と、ホストコンピュータ３に対する入出力を制御する２つのマイクロプロセッサ２０１と、グローバル情報制御部２１あるいはメモリ２：２６へのアクセスを制御するアクセス制御部１（メモリアクセス制御部１）２０７と、メモリ１：２５へのアクセスを制御するアクセス制御部２（メモリアクセス制御部２）２０８とを有し、ホストコンピュータ３とメモリ１間のデータ転送、及びマイクロプロセッサ２０１とグローバル情報制御部２１あるいはメモリ２間の制御情報の転送を実行する。

マイクロプロセッサ２０１及びホストＩＦ２０２は内部バス２０５によって接続され、メモリアクセス制御部１：２０７は内部バス２０５に接続され、メモリアクセス制御部２：２０８は２つのホストＩＦ２０２に直接接続され、また内部バス２０５に接続されている。

20

ディスクＩＦ部１８の具体的な一例を図１５に示す。

ディスクＩＦ部１８は、ディスク装置２との２つのＩＦ（ドライブＩＦ）２０３と、ディスク装置２に対する入出力を制御する２つのマイクロプロセッサ２０１と、グローバル情報制御部２１あるいはメモリ２へのアクセスを制御するアクセス制御部１（メモリアクセス制御部１）２０７と、メモリ１へのアクセスを制御するアクセス制御部２（メモリアクセス制御部２）２０８とを有し、ディスク装置２とメモリ１間のデータ転送、及びマイクロプロセッサ２０１とグローバル情報制御部２１あるいはメモリ２間の制御情報の転送を実行する。

30

マイクロプロセッサ２０１及びドライブＩＦ２０３は内部バス２０５によって接続され、メモリアクセス制御部１：２０７は内部バス２０５に接続され、メモリアクセス制御部２：２０８は２つのドライブＩＦ２０３に直接接続され、また内部バス２０５に接続されている。ディスクＩＦ部１８はＲＡＩＤ機能の実行も行う。

１つのディスク制御クラスタは１つの筐体として構成されるか、またはモジュールとして構成されても良いが、それ自体１つのディスク制御装置として機能を備えているものである。

ストレージシステムの具体的な一例は、チャンネルＩＦ部１３及びディスクＩＦ部１８とメモリ２：２６と相互結合網１：３２とグローバル共有メモリ２１との接続構成は、実施例３の図７に示す構成と同様になる。また、チャンネルＩＦ部１３及びディスクＩＦ部１８とメモリ１と相互結合網２：３３との接続構成は、実施例１の図３に示す構成においてグローバル情報制御部２１を除いた構成と同様になる。

40

フロントエンドスイッチ７は、スイッチ７１と、スイッチ制御部７２と、メモリコントローラ７３と、メモリモジュール１０５を有する。

１つのホストコンピュータ３からは、２つのフロントエンドスイッチ７内のスイッチ７１にパスを１本ずつ接続し、スイッチ７１からは、ディスク制御クラスタ１－１乃至１－ｎの各チャンネルＩＦ部１１に１本ずつパスを接続する。

【００３６】

グローバル情報制御部２１は、アクセス制御部１０１と管理機能部１０２とメモリモジュール１０５とを有し、ディスク制御クラスタ１－１乃至１－ｎの管理情報（例えば、各デ

50

ィスク制御クラスタが管理する記憶領域の情報や、ディスク制御クラスタ内の各部位の負荷情報、障害情報、及び構成情報等)を格納する。メモリ1は、ディスク装置2に記録するデータを一時的に格納する。また、メモリ2は、ディスク制御クラスタの制御情報(例えば、チャンネルIF部13及びディスクIF部18とメモリ1:25との間のデータ転送制御に関する情報、ディスク装置2に記録するデータの管理情報等)を格納する。

フロントエンドスイッチ7は、そのメモリモジュール105内に実施例1と同様の情報、テーブルを格納しており、同様の制御を行う。

また、グローバル情報制御部21も、そのメモリモジュール105内に実施例1と同様の情報、テーブルを格納している。

図10において、相互結合網1:32を形成するディスク制御クラスタ外のスイッチ及び相互結合網2:33を形成するディスク制御クラスタ外のスイッチとグローバル情報制御部21をボックスに実装し、フロントエンドスイッチ7を別のボックスに実装し、モジュール化した各ディスク制御クラスタ1-1乃至1-nといっしょに、1つの筐体の中に実装しても良い。また、各ディスク制御クラスタ1-1乃至1-nを別個の筐体として、距離的に離れた場所に分散しても良い。

【0037】

図10において、ホストコンピュータ3からストレージシステム1に記録されたデータを読み出す場合の一例を述べる。

まず、ホストコンピュータ3は、ストレージシステム1に対してデータの読出し要求を発行する。

要求は、フロントエンドスイッチ7に受け取られ、フロントエンドスイッチ7内のスイッチ制御部72は、要求パケットのヘッダを解析する。

要求パケットのヘッダには、要求を発行したホストコンピュータの番号(ホスト番号)と要求データが記録されている論理ボリューム番号が格納されており、スイッチ制御部72はホスト-論理ボリューム対応テーブル600を参照し、要求を発行したホストコンピュータが接続されたスイッチ71のポートを、該当する論理ボリューム番号に割当てられているチャンネル番号のスイッチ71のポートに接続し、要求パケットをディスク制御クラスタに送る。

要求が送られたチャンネルが接続されたチャンネルIF部11内のマイクロプロセッサ201は、自ディスク制御クラスタ1-1内のメモリ2:26にアクセスし、要求されたデータがどのディスク装置2内に格納されているかを調べる。メモリ2:26には、要求データのアドレスとそのデータが実際に記録されているディスク装置2内のアドレスを対応させる変換テーブルが格納されており、要求されたデータがどのディスク装置2内に格納されているかを調べることができる。

さらに、要求を受けたチャンネルIF部13内のマイクロプロセッサ201は、自ディスク制御クラスタ1-1内のメモリ2:26にアクセスし、要求されたデータがメモリ1:25に格納されているかどうかを確認する。メモリ2:26にはメモリ1:25に格納されているデータのディレクトリ情報が格納されており、メモリ1:25に要求データが存在するかどうかを確認できる。

【0038】

それにより自ディスク制御クラスタ1-1のメモリ1:25にデータがあった場合は、そのデータをチャンネルIF部13まで転送し、フロントエンドスイッチ7を介して、ホストコンピュータ3に送る。

【0039】

自ディスク制御クラスタ1-1のメモリ1:25にデータが存在しなかった場合は、チャンネルIF部13内のマイクロプロセッサ201は要求データが格納されているディスク装置2が接続されているディスクIF部18内のマイクロプロセッサ201に対し、要求データを読み出し、メモリ1:25に格納するというデータ要求の処理内容を示す制御情報を発行し、この制御情報の発行を受けたディスクIF部18内のマイクロプロセッサ201は、要求データが格納されているディスク装置2からデータを読み出し、自ディスク制御ク

10

20

30

40

50

ラスト 1 - 1 内のメモリ 1 : 2 5 に要求データを転送し格納する。

すなわち、チャンネル I F 部 1 3 のマイクロプロセッサ 2 0 1 は、上記データ要求の処理内容を示す制御情報を発行し、メモリ 2 : 2 6 の制御情報領域（ジョブ制御ブロック）に格納する。

ディスク I F 部 1 8 のマイクロプロセッサ 2 0 1 は、メモリ 2 : 2 6 の制御情報領域をポーリングで監視し、上記発行された制御情報が上記制御情報領域（ジョブ制御ブロック）に存在した場合は、要求データが格納されているディスク装置 2 からデータを読み出し、自ディスク制御クラスタ 1 - 1 内のメモリ 1 : 2 5 に要求データを転送し格納する。

ディスク I F 部 1 8 内のマイクロプロセッサ 2 0 1 は、要求データをメモリ 1 : 2 5 へ格納した後、制御情報を発行したチャンネル I F 部 1 3 内のマイクロプロセッサ 2 0 1 に、メモリ 1 : 2 5 内のデータを格納したアドレスを、メモリ 2 : 2 6 内の制御情報を介して伝える。それを受けたチャンネル I F 部 1 3 内のマイクロプロセッサ 2 0 1 は、メモリ 1 : 2 5 からデータを読み出し、フロントエンドスイッチ 7 を介して、ホストコンピュータ 3 へ送る。

10

すなわち、ディスク I F 部 1 8 のマイクロプロセッサ 2 0 1 は、要求データをメモリ 1 : 2 5 へ格納した後、処理の実行の終了とデータを格納したアドレスを示す制御情報を発行し、ローカル共有メモリ部 2 2 の制御情報領域に格納する。

前記制御情報を発行したチャンネル I F 部 1 3 のマイクロプロセッサ 2 0 1 は、メモリ 2 : 2 6 の制御情報領域をポーリングで監視し、ディスク I F 部 1 8 のマイクロプロセッサ 2 0 1 から発行された制御情報が上記制御情報領域に存在した場合、メモリ 1 : 2 5 内のデータを格納したアドレスによりメモリ 1 : 2 5 からデータを読み出し、チャンネル I F 部 1 3 まで転送し、さらにフロントエンドスイッチ 7 を介して、ホストコンピュータ 3 へ送る。

20

【 0 0 4 0 】

制御情報とデータはデータ長が数千倍異なるため、1 回のデータ転送時間がかかなり異なる。このため、同じ相互結合網及びメモリを用いた場合、両者が互いの転送を妨げる。本実施例によれば、制御情報を転送する相互結合網 1 : 3 2 とデータを転送する相互結合網 2 : 3 3 を分けることができるため、両者が互いの転送を妨げることがなくなるため、性能が向上する。

[実施例 6]

図 1 1、図 1 6、図 1 7 に、実施例 1 のストレージシステム 1 におけるディスク制御クラスタの増設手順の一例を示す。

30

図 1 1 に示すように、グローバルスイッチボックス 3 1 0 と、フロントエンドスイッチボックス 3 1 5 と、クラスタ筐体 3 0 2 は、それぞれ別筐体として、筐体 3 0 1 内に実装されている。

グローバルスイッチボックス 3 1 0 内には、G S W 1 1 5 とグローバル情報制御部 2 1 が実装されている。

グローバルスイッチボックス 3 1 0 はコネクタ 3 2 1、コネクタ 3 2 2 をそれぞれ 8 個有し、ディスク制御クラスタを 8 クラスタ接続することができる。図ではディスク制御クラスタを 3 クラスタ接続した場合について示している。

G S W 1 1 5 のアクセスパス 1 3 6 は 1 本ずつコネクタ 3 2 1、コネクタ 3 2 2 に接続される。

40

また、グローバルスイッチボックス 3 1 0 はフロントエンドスイッチボックス 3 1 5 を接続するためのコネクタ 3 2 6、コネクタ 3 2 7 をそれぞれ 2 個有する。グローバル情報制御部 2 1 の 1 つのアクセス制御部 1 0 1 をフロントエンドスイッチ 7 の 2 つのメモリコントローラ 7 3 に接続するための 2 本の接続パスは 2 つのコネクタ 3 2 6 に接続される。また、もう 1 つのアクセス制御部 1 0 1 をフロントエンドスイッチ 7 の 2 つのメモリコントローラ 7 3 に接続するための 2 本の接続パスは 2 つのコネクタ 3 2 7 に接続される。

フロントエンドスイッチボックス 3 1 5 内には、フロントエンドスイッチ 7 が実装されている。

フロントエンドスイッチボックス 3 1 5 はグローバルスイッチボックス 3 1 0 を接続する

50

ためのコネクタ 3 2 8、コネクタ 3 2 9 をそれぞれ 2 個有する。

フロントエンドスイッチ 7 の 1 つのメモリコントローラ 7 3 をグローバル情報制御部 2 1 の 2 つのアクセス制御部 1 0 1 に接続するための 2 本の接続パスは 2 つのコネクタ 3 2 8 に接続される。また、もう 1 つのメモリコントローラ 7 3 をグローバル情報制御部 2 1 の 2 つのアクセス制御部 1 0 1 に接続するための 2 本の接続パスは 2 つのコネクタ 3 2 9 に接続される。

上記個数は一実施例に過ぎず、個数を上記に限定するものではない。

各ディスク制御クラスタ 1 - 1 乃至 1 - 3 は、それぞれクラスタ筐体 3 0 2 に実装されている。クラスタ筐体 3 0 2 はコネクタ 3 2 1、コネクタ 3 2 2 を有し、2 本のアクセスパス 1 3 6 がそれぞれに 1 本ずつ接続されている。

10

グローバルスイッチボックス 3 1 0 に、ケーブル 3 3 1、ケーブル 3 3 2 を介してそれぞれのコネクタ 3 2 1、コネクタ 3 2 2 により 3 つのクラスタ筐体 3 0 2 が接続される。

また、グローバルスイッチボックス 3 1 0 のコネクタ 3 2 6、コネクタ 3 2 7 に、ケーブル 3 3 6、ケーブル 3 3 7 を介して、フロントエンドスイッチボックス 3 1 5 のコネクタ 3 2 8、コネクタ 3 2 9 が接続される。

ここで、2 つのコネクタ 3 2 6 に接続された 2 本のケーブル 3 3 6 は、1 本ずつコネクタ 3 2 8 とコネクタ 3 2 9 に接続される。同様に、2 つのコネクタ 3 2 7 に接続された 2 本のケーブル 3 3 7 は、1 本ずつコネクタ 3 2 8 とコネクタ 3 2 9 に接続される。こうすることにより、アクセス制御部 1 0 1 の 2 本の接続パスは、1 本ずつ 2 つのメモリコントローラ 7 3 に接続される。

20

ストレージシステム 1 において、ディスク制御クラスタを増設する場合は、次の手順による。グローバルスイッチボックス 3 1 0 にディスク制御クラスタを増設するコネクタに余分があれば、そのコネクタにケーブル 3 3 1、ケーブル 3 3 2 を接続する。

余分がなければ、G S W のみを実装したグローバルスイッチボックスを用意し、グローバルスイッチボックスを多段に接続した上でそのコネクタにケーブル 3 3 1、ケーブル 3 3 2 を接続する。

それと共に、図 1 6 に示す G S W 1 1 5 のポートに接続されるディスク制御クラスタを示す、言い換えるとストレージシステム 1 を構成しているディスク制御クラスタを示す G S W ポート - クラスタ対応テーブル 4 0 0 と、図 1 7 に示すディスク制御クラスタが管理する論理ボリュームを示す G S W ポート - クラスタ対応テーブル 4 0 5 とを書き換える。

30

G S W ポート - クラスタ対応テーブル 4 0 0 と G S W ポート - クラスタ対応テーブル 4 0 5 はグローバル情報制御部 2 1 に格納されており、サービスプロセッサ (S V P) により書き換えることが可能である。

S V P は通常ノートパソコンであることが多く、ノートパソコンのディスプレイ上に図 1 6 及び図 1 7 に示すテーブルが表示され、そこで内容を書き換える。

図 1 6 及び図 1 7 は、それぞれディスク制御クラスタの増設前、増設後の G S W ポート - クラスタ対応テーブル 4 0 0 及び G S W ポート - クラスタ対応テーブル 4 0 5 を示している。

ここでは、ストレージシステム 1 が増設前に 5 台のディスク制御クラスタで構成されており、そこに 1 台のディスク制御クラスタを増設する例を示している。

40

図 1 6 に示すように、G S W ポート番号 4 0 1 の 4 番ポートが未接続となっており、そのポートにクラスタ 5 のケーブルを接続した後、S V P のディスプレイ上でポート番号 4 の行のクラスタ番号 4 0 2 の列の未接続表示を 5 に書き換える。その後、図 1 7 に示すように、クラスタ番号 4 0 2 のクラスタ 5 の行の論理ボリューム番号 4 0 6 の列の未接続表示を 1 6 6 4 0 ~ 2 0 7 3 5 に書き換える。

ここで、論理ボリューム番号 4 0 6 は各クラスタが管理する論理ボリュームの範囲を示している。

増設前の論理ボリューム番号の最大値は 1 6 6 3 9 で、ディスク制御クラスタは 4 0 9 6 個の論理ボリュームを持っているため、ディスク制御クラスタ 5 の管理する論理ボリュームの範囲は 1 6 6 4 0 ~ 2 0 7 3 5 となる。論理ボリューム番号は連続せず飛び飛びにな

50

っていても問題ない。

上記のようにすることで、ストレージシステムに新たにディスク制御クラスタを増設することができる。

〔実施例 7〕

図 18～図 21 に、ストレージシステム 1 において、あるチャンネル I/F 部のホストコンピュータとの 1 つのインタフェース(チャンネル)の負荷が高くなった場合の、グローバル制御情報部及びフロントエンドスイッチ 7 の動作の一例を示す。実施例 1 に示すように、一例として、グローバル情報制御部 21 は、図 18 に示すようなシステム構成・稼動状況テーブル 500 をそのメモリモジュール 105 に格納しており、フロントエンドスイッチ 7 は、テーブル 500 のコピーをメモリモジュール 105 に格納している。

10

グローバル情報制御部 21 内の管理機能部 102 は、各ディスク制御クラスタのローカル共有メモリ部 22 内に格納された該ディスク制御クラスタが管理する記憶領域の情報や、ディスク制御クラスタ内の各部位の負荷情報、障害情報を定期的に参照し、システム構成・稼動状況テーブル 500 を更新する。

また、フロントエンドスイッチ 7 内のスイッチ制御部 72 は、テーブル 500 を定期的に自メモリモジュール 105 にコピーする。

一方、各ディスク制御クラスタのローカル共有メモリ部 22 内に格納された該ディスク制御クラスタが管理する記憶領域の情報、ディスク制御クラスタ内の各部位の負荷情報、あるいは障害情報が更新された時点で、チャンネル I/F 部またはディスク I/F 部内のマイクロプロセッサが、グローバル情報制御部のメモリモジュール 105 内の該当情報を更新し、さらに、その更新の時点で、管理機能部 102 が該当情報をフロントエンドスイッチ 7 のメモリモジュール 105 にコピーする場合もある。これは、特にディスク制御部内のある部位に障害が起こった場合において有効である。

20

システム構成・稼動状況テーブル 500 は、ディスク制御クラスタを識別するクラスタ番号 511 とディスク制御クラスタのチャンネルを識別するチャンネル番号 512 に対応する論理ボリューム番号 513 とチャンネル稼動状況 514 を示す。

本実施例では、クラスタ番号 511、チャンネル番号 512、及び論理ボリューム番号 513 を 16 進数で示している。チャンネル稼動状況 514 は、低負荷状態を '0'、中負荷状態を '1'、高負荷状態を '2'、障害を '3' で示している。

また、フロントエンドスイッチ 7 は、ホスト - 論理ボリューム対応テーブル 600 をそのメモリモジュール 105 に格納している。

30

ホスト - 論理ボリューム対応テーブル 600 は、個々のホストコンピュータを識別するホスト番号 615 に対応する論理ボリューム番号 513、すなわち各ホストコンピュータに割当てられた論理ボリュームと、論理ボリューム番号 513 にアクセスするためのチャンネル番号 512 と該当論理ボリュームを管理するディスク制御クラスタのクラスタ番号 511 を示す。

ホスト番号 615 は、例えば、ファイバチャンネルのプロトコルで使われるワールドワイドネーム(WWN)や、インターネットプロトコルで使われるマックアドレスや IP アドレスを充てることが考えられる。

図 19 に示すように、1 つのホストコンピュータに対し、連続しない複数の論理ボリュームの範囲を割当てても問題ない。

40

スイッチ制御部 72 は、ホスト - 論理ボリューム対応テーブル 600 (ルーティングテーブル)を参照し、スイッチ 71 の切り替えを行う。

図 20 に示すように、チャンネル番号 512 の 2 番チャンネルの稼動状況 551 が高負荷を示す '2' となった場合、2 番チャンネルに割当てられた論理ボリューム 0100～01FF の一部を、稼動状況 552 が低負荷を示す '0' となっている 1F 番チャンネルに割当て替える。

図 21 は、上記割当て替えを行った後のテーブル 500 の内容を示している。

本実施例では、2 番チャンネルに割当てられていた論理ボリューム 0100～01FF の半分の 0180～01FF を 1F 番チャンネルに割当て替えた。これにより、2 番チャンネルの

50

稼動状況 5 5 1 は中負荷を示す ' 1 ' となり、 1 F 番チャネルの稼動状況 5 5 2 も中負荷を示す ' 1 ' となる。

図 2 1 に示すテーブル 5 0 0 がフロントエンドスイッチ 7 のメモリモジュール 1 0 5 にコピーされると、スイッチ制御部 7 2 は、テーブル 5 0 0 を参照し図 1 9 に示すホスト - 論理ボリューム対応テーブル 6 0 0 を変更する。

図 2 2 は変更後のテーブル 6 0 0 を示す。

テーブル 5 0 0 において、論理ボリュームの 0 1 8 0 ~ 0 1 F F は 1 F 番チャネルに割当て返られているので、 2 番のホストコンピュータ 6 5 2 に新たに 1 F 番チャネルが割当てられる。

これにより、 2 番のホストコンピュータ 6 5 2 から論理ボリューム 0 1 8 0 ~ 0 1 F F へのアクセス要求が発行された場合、スイッチは 1 F 番チャネルに切り替えられる。

この場合、アクセス要求を受けたチャネル I F 部のマイクロプロセッサは、ローカル共有メモリ部にアクセスし、自ディスク制御クラスタ内に要求データがないことを知るが、その場合、グローバル情報制御部にアクセスすることにより、要求データのあるディスク制御クラスタを知り、該当ディスク制御クラスタから要求データを読み出すことができる。これにより、ホストコンピュータに意識させること無しに、ストレージシステム 1 の各チャネル I F 部のチャネルの負荷を均等にすることができ、それによりシステム全体の性能向上が可能となる。

【 0 0 4 1 】

また、あるチャネル番号のチャネルが障害で使用できなくなった場合にも、上記の方法で論理ボリュームの割当て替えを行うことにより、別のチャネルから目的のディスク制御クラスタのディスク制御装置に格納されたデータへアクセスすることが、チャネルの障害をホストコンピュータに意識させること無しに、可能となる。

【 0 0 4 2 】

また、あるディスク制御クラスタに繋がるディスク装置上の論理ボリュームの負荷が高くなった場合、負荷を均等にするために該当論理ボリュームを他のディスク制御クラスタのディスク装置にコピーまたは移動する場合があるが、その場合も上記の方法でテーブル 5 0 0 の論理ボリューム番号とチャネル番号を割当て替えることにより、ホストコンピュータに意識させること無しに、論理ボリュームの負荷を均等にすることが可能となる。

[実施例 8]

図 2 4 に、本発明の一実施例を示す。

以下の実施例において、相互結合網はスイッチを利用したものを例にして説明してあるが、相互に接続され制御情報やデータが転送されれば良いのであり、例えばバスで構成されても良い。

図 2 4 に示すように、ストレージシステム 1 は複数のディスク制御クラスタ 1 - 1 乃至 1 - n とフロントエンドスイッチ 7 から構成される。

本実施例のストレージシステム 1 は、グローバル制御部 2 1 とフロントエンドスイッチ 7 が接続されていない点を除いて、実施例 1 の図 1 に示すストレージシステムと同様である。

また、その具体的な一例は、図 3 に示すストレージシステム 1 においてアクセス制御部 1 0 1 とメモリコントローラ 7 3 が接続されていない構成となる。

上記のように、アクセス制御部 1 0 1 とメモリコントローラ 7 3 が接続されないため、グローバル制御情報部 2 1 内の各ディスク制御クラスタが管理する記憶領域の情報、ディスク制御クラスタ内の各部位の負荷情報、および障害情報をチャネル I F 部 1 1 とスイッチ 7 1 の接続バスを介してフロントエンドスイッチ 7 のメモリモジュール 1 0 5 にコピーする処理を行う。

この処理は、フロントエンドスイッチ 7 内のスイッチ制御部 2 1 がチャネル I F 部 1 1 とスイッチ 7 1 の接続バスを介して、チャネル I F 部 1 1 内のマイクロプロセッサに処理要求を出すことにより行う。

例えば、上記接続バスが、インターネットプロトコルを流すことが可能である場合、シン

10

20

30

40

50

ブル・ネットワーク・マネージメント・プロトコル (S N M P) により、上記情報を取得することが可能である。

本実施例によれば、フロントエンドスイッチ 7 とグローバル情報制御部 2 1 を接続する必要がなくなり、ストレージシステムを構成する筐体の実装を簡素化することが可能となる。

【 0 0 4 3 】

【 発明の効果 】

本発明によれば、複数台のディスク制御クラスタを 1 つのシステムとして運用するストレージシステムにおいて、ディスク制御クラスタが 1 個だけの小規模な構成からディスク制御クラスタが数十個接続された超大規模な構成まで、ディスク制御クラスタ単体が持つ高信頼・高性能なアーキテクチャで対応可能な、スケーラビリティがあり使い勝手の良い構成のストレージシステムを提供することが可能となる。

【 図面の簡単な説明 】

【 図 1 】 本発明によるストレージシステムの実施例 1 の構成を示す図。

【 図 2 】 従来の複数のディスク制御装置の構成を示す図。

【 図 3 】 図 1 に示す実施例 1 のストレージシステムの詳細構成を示す図。

【 図 4 】 本発明によるストレージシステムの実施例 2 の構成を示す図。

【 図 5 】 図 4 に示す実施例 2 のストレージシステムの詳細構成を示す図。

【 図 6 】 本発明によるストレージシステムの実施例 3 の構成を示す図。

【 図 7 】 図 6 に示す実施例 3 のストレージシステムの詳細構成を示す図。

【 図 8 】 本発明によるストレージシステムの実施例 4 の構成を示す図。

【 図 9 】 図 8 に示す実施例 4 のストレージシステムの詳細構成を示す図。

【 図 1 0 】 本発明によるストレージシステムの実施例 5 の構成を示す図。

【 図 1 1 】 本発明によるディスク制御クラスタの増設方法を説明するための図。

【 図 1 2 】 本発明によるストレージシステムを構成するチャンネルインターフェース部の構成を示す図。

【 図 1 3 】 本発明によるストレージシステムを構成するディスクインターフェース部の構成を示す図。

【 図 1 4 】 本発明によるストレージシステムを構成するチャンネルインターフェース部の他の構成を示す図。

【 図 1 5 】 本発明によるストレージシステムを構成するディスクインターフェース部の他の構成を示す図。

【 図 1 6 】 グローバル情報制御部内に格納されたストレージシステムの構成情報の一例を示す図。

【 図 1 7 】 グローバル情報制御部内に格納されたストレージシステムの構成情報の他の一例を示す図。

【 図 1 8 】 グローバル情報制御部内に格納されたストレージシステムの構成・稼動状況の一例を示す図。

【 図 1 9 】 フロントエンドスイッチ内に格納されたスイッチ切換え制御のためのテーブルの一例を示す図。

【 図 2 0 】 グローバル情報制御部内に格納されたストレージシステムの構成・稼動状況の他の一例を示す図。

【 図 2 1 】 グローバル情報制御部内に格納されたストレージシステムの構成・稼動状況の他の一例を示す図。

【 図 2 2 】 フロントエンドスイッチ内に格納されたスイッチ切換え制御のためのテーブルの他の一例を示す図。

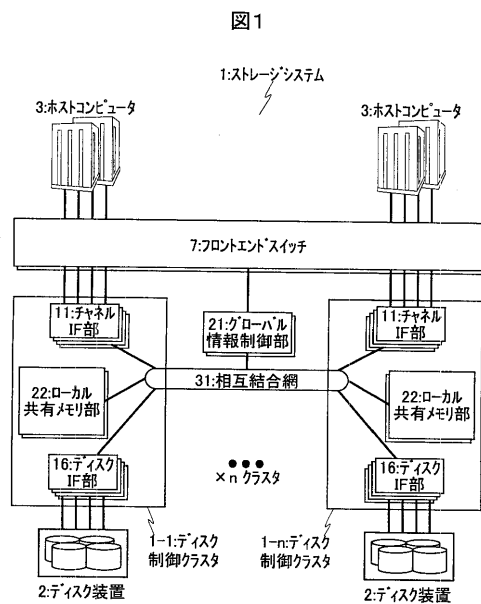
【 図 2 3 】 従来の複数のディスク制御クラスタから成るストレージシステムの構成を示す図。

【 図 2 4 】 本発明によるストレージシステムの実施例 8 の構成を示す図。

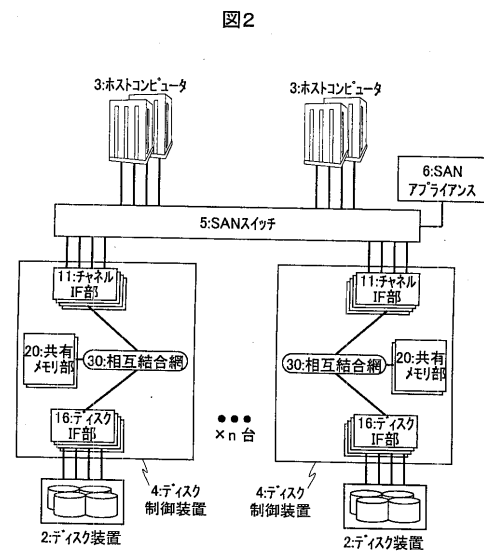
【 符号の説明 】

1 : ストレージシステム、1 - 1、1 - n ... ディスク制御クラスタ、2 ... ディスク装置、3 ... ホストコンピュータ、7 ... フロントエンドスイッチ、11 ... チャンネル I/F 部、16 ... ディスク I/F 部、21 ... グローバル情報制御部、22 ... ローカル共有メモリ部、31 ... 相互結合網。

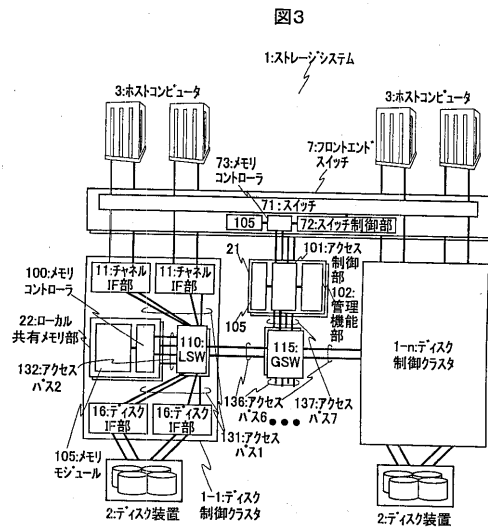
【図 1】



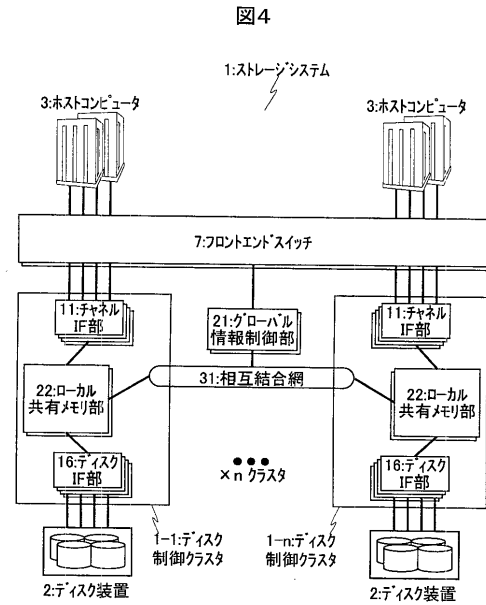
【図 2】



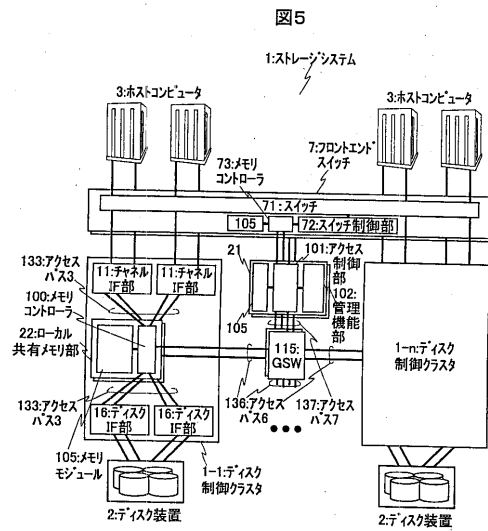
【図 3】



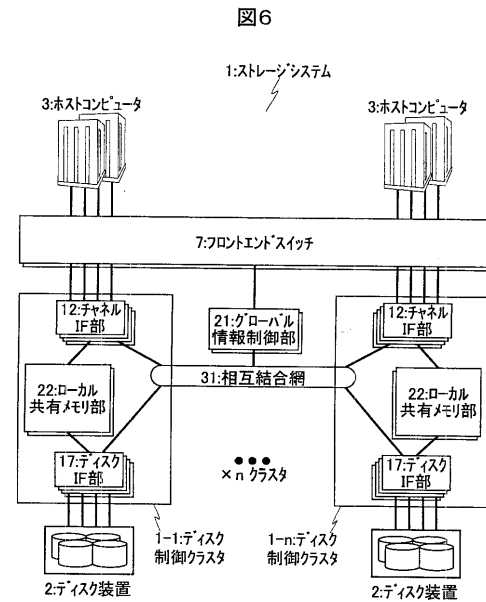
【図 4】



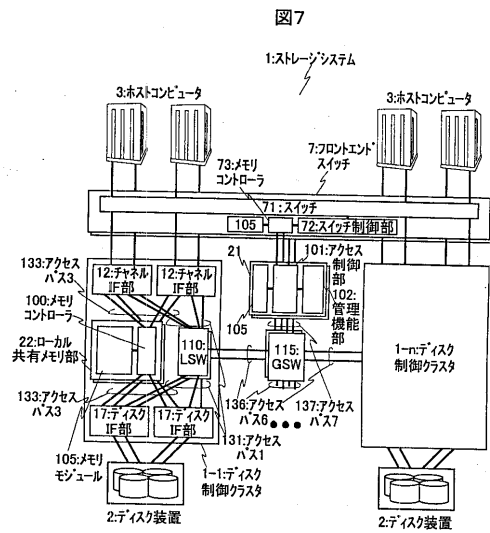
【図 5】



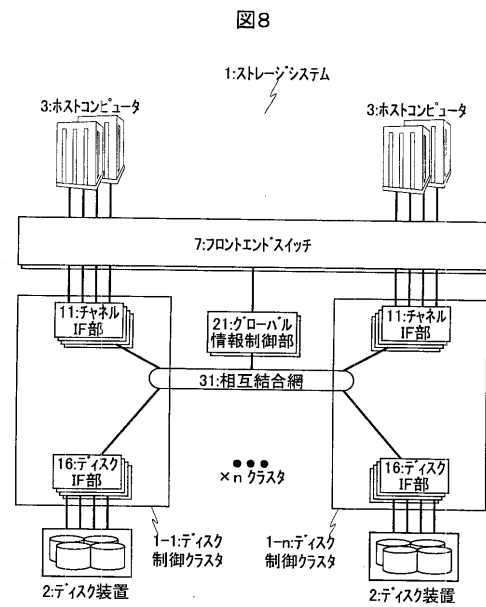
【図 6】



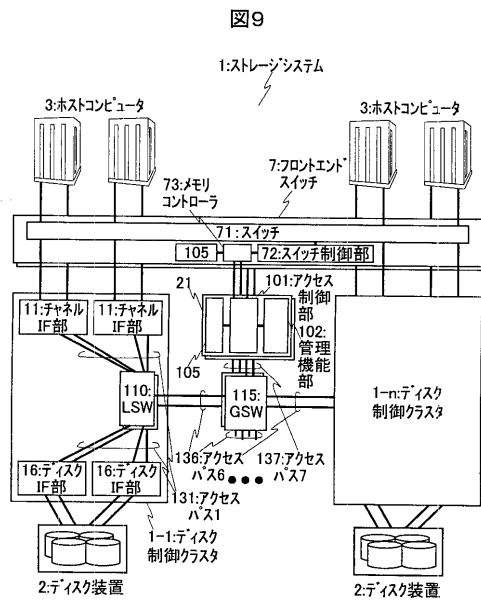
【図 7】



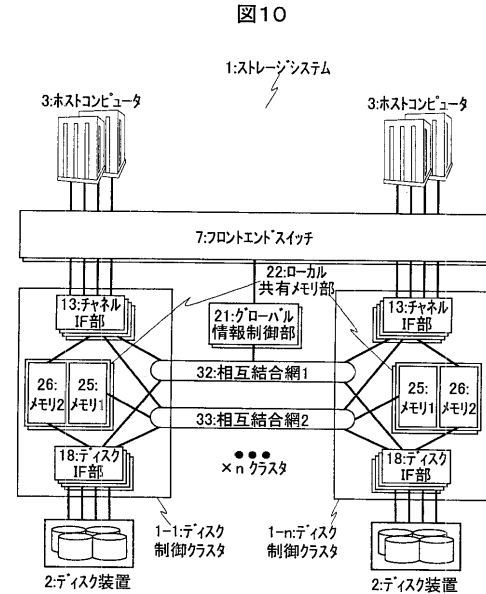
【図 8】



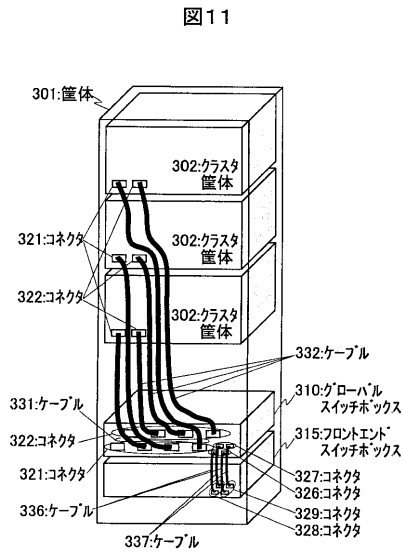
【図 9】



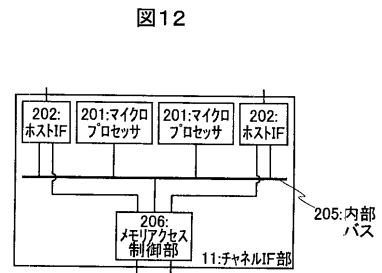
【図 10】



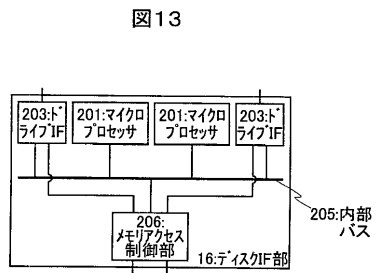
【図 1 1】



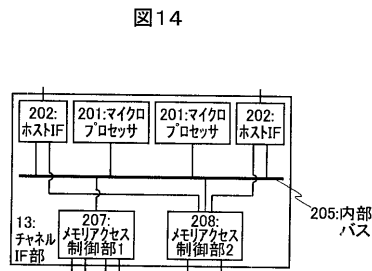
【図 1 2】



【図 1 3】



【図 1 4】



【図 1 6】

図16

(増設前)

GSWホート番号	クラス番号
0	0
1	1
2	3
3	2
4	未接続
5	4
6	未接続
7	未接続

400:GSWホート-クラス対応テーブル

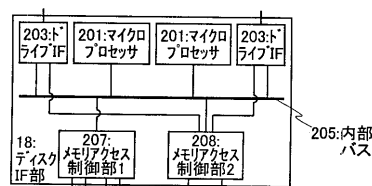
(増設後)

GSWホート番号	クラス番号
0	0
1	1
2	3
3	2
4	5
5	4
6	未接続
7	未接続

400:GSWホート-クラス対応テーブル

【図 1 5】

図15



【図 17】

図17

(増設前)

クラス番号	論理ボリューム番号
0	0~4095
1	4096~6143
2	12288~16383
3	6144~12287
4	16384~16639
5	未実装
6	未実装
7	未実装

405:クラス-論理ボリューム対応テーブル

(増設後)

クラス番号	論理ボリューム番号
0	0~4095
1	4096~6143
2	12288~16383
3	6144~12287
4	16384~16639
5	16640~20735
6	未実装
7	未実装

405:クラス-論理ボリューム対応テーブル

【図 18】

図18

クラス番号	チャネル番号	論理ボリューム番号	チャネル稼動状況
0	0	0000~007F	1
0	1	0080~00FF	1
0	2	0100~01FF	1
0	3	0200~02FF	1
1	4	0300~033F	1
1	5	0340~03FF	1
1	6	0400~04FF	1
1	7	0300~033F	1
⋮	⋮	⋮	⋮
7	1C	1000~10FF	1
7	1D	1100~12FF	1
7	1E	1300~14FF	1
7	1F	1500~17FF	0

500:システム構成・稼動状況テーブル

【図 19】

図19

ホスト番号	論理ボリューム番号	チャネル番号	クラス番号
0	0000~007F	0	0
1	0080~00FF	1	0
1	0340~03FF	5	1
2	0100~01FF	2	0
3	0200~02FF	3	0
4	0300~033F	4	1
5	0400~04FF	6	1
6	0300~033F	7	1
⋮	⋮	⋮	⋮
1C	1000~10FF	1C	7
1D	1100~12FF	1D	7
1E	1300~14FF	1E	7
1F	1500~17FF	1F	7

600:ホスト-論理ボリューム対応テーブル

【図 20】

図20

クラス番号	チャネル番号	論理ボリューム番号	チャネル稼動状況
0	0	0000~007F	1
0	1	0080~00FF	1
0	2	0100~01FF	2
0	3	0200~02FF	1
1	4	0300~033F	1
1	5	0340~03FF	1
1	6	0400~04FF	1
1	7	0300~033F	1
⋮	⋮	⋮	⋮
7	1C	1000~10FF	1
7	1D	1100~12FF	1
7	1E	1300~14FF	1
7	1F	1500~17FF	0

500:システム構成・稼動状況テーブル

【図 2 1】

図21

クラスタ番号	チャネル番号	論理ボリューム番号	チャネル稼動状況
0	0	0000~007F	1
0	1	0080~00FF	1
0	2	0100~017F	1
0	3	0200~02FF	1
1	4	0300~033F	1
1	5	0340~03FF	1
1	6	0400~04FF	1
1	7	0300~033F	1
⋮	⋮	⋮	⋮
7	1C	1000~10FF	1
7	1D	1100~12FF	1
7	1E	1300~14FF	1
7	1F	0180~01FF	1
7	1F	1500~17FF	1

500:システム構成・稼動状況テーブル

【図 2 2】

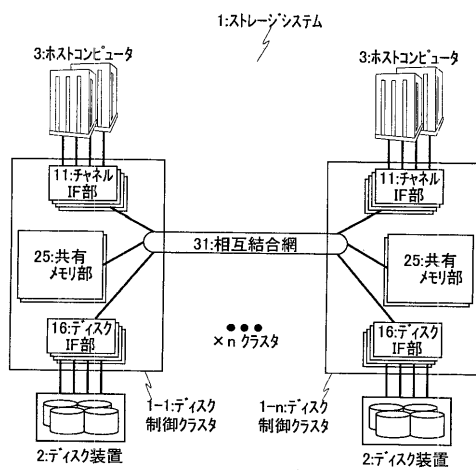
図22

ホスト番号	論理ボリューム番号	チャネル番号	クラスタ番号
0	0000~007F	0	0
1	0080~00FF	1	0
1	0340~03FF	5	1
2	0100~017F	2	0
2	0180~01FF	1F	7
3	0200~02FF	3	0
4	0300~033F	4	1
5	0400~04FF	6	1
6	0300~033F	7	1
⋮	⋮	⋮	⋮
1C	1000~10FF	1C	7
1D	1100~12FF	1D	7
1E	1300~14FF	1E	7
1F	1500~17FF	1F	7

600:ホスト-論理ボリューム対応テーブル

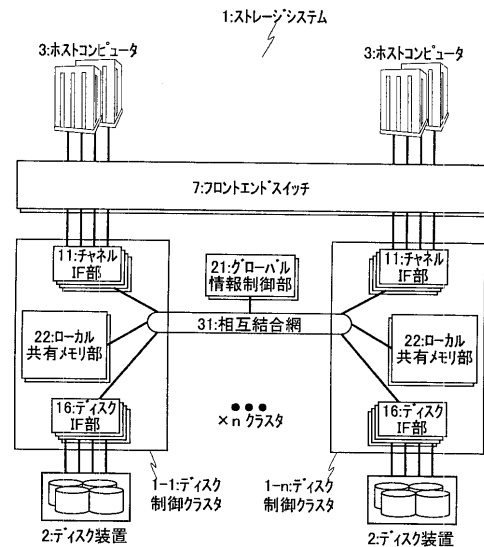
【図 2 3】

図23



【図 2 4】

図24



フロントページの続き

(56)参考文献 特開2000-242434(JP,A)
特開平11-007359(JP,A)
特開2001-256003(JP,A)
特開2001-229042(JP,A)
特開平08-328760(JP,A)
特開2001-290608(JP,A)
特開2001-142648(JP,A)
特開平11-296313(JP,A)
特開平07-056691(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06

G06F 13/10-13/14

G06F 12/00