

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4121255号
(P4121255)

(45) 発行日 平成20年7月23日(2008.7.23)

(24) 登録日 平成20年5月9日(2008.5.9)

(51) Int.Cl. F I
G O 6 F 12/00 (2006.01) G O 6 F 12/00 5 O 1 B

請求項の数 3 (全 22 頁)

(21) 出願番号	特願2001-177180 (P2001-177180)	(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13年6月12日(2001.6.12)	(74) 代理人	100093492 弁理士 鈴木 市郎
(65) 公開番号	特開2002-366398 (P2002-366398A)	(74) 代理人	100078134 弁理士 武 顕次郎
(43) 公開日	平成14年12月20日(2002.12.20)	(72) 発明者	平川 裕介 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16年10月6日(2004.10.6)	(72) 発明者	荒川 敬史 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

最終頁に続く

(54) 【発明の名称】 クラスタ構成記憶システム

(57) 【特許請求の範囲】

【請求項1】

複数の記憶システムノードを1つの記憶システムとして運用可能とするクラスタ構成記憶システムにおいて、

前記記憶システムノード内部及び前記記憶システムノード相互間のアクセス情報を採取する手段と、採取したアクセス情報を保守員に提示する手段と、保守員からのポートの設定指示を受け付ける手段と、前記指示に基づいて前記クラスタ構成記憶システム内のポート設定を変更する手段とを有し、

前記クラスタ構成記憶システム内のポート設定を変更する手段は、論理ポートと論理ボリュームとの間の仮想的なパスである論理パス毎の単位時間当りのデータ転送量またはデータ転送時間を示すアクセス頻度を保持するアクセス情報を参照し、全論理パスの中でデータ転送量が規定値以上の論理パスの1つを選択し、選択した論理パスが記憶システムノード間のデータ転送を必要とするか否かをチェックし、記憶システムノード間のデータ転送を必要とした場合、

前記論理パスを経て前記ホストコンピュータがアクセスしている論理ボリュームを格納している記憶装置を有する第1の記憶システムノードに未使用ポートがあり、その未使用ポートが当該ホストコンピュータに接続可能である場合、あるいは、当該ホストコンピュータに接続された未使用ポートがある場合に、前記ホストコンピュータが使用するノードを、それまで使用していた第2の記憶システムノードのポートから前記第1の記憶システムノードの前記未使用ポートとするように、ポート設定を変更することを特徴とするクラ

10

20

スタ構成記憶システム。

【請求項 2】

複数の記憶システムノードを 1 つの記憶システムとして運用可能とするクラスタ構成記憶システムにおいて、

前記記憶システムノード内部及び前記記憶システムノード相互間のアクセス情報を採取する手段と、採取したアクセス情報に基づいて、ポートの設定を決定する手段と、前記クラスタ構成記憶システム内のポート設定を変更する手段とを有し、

前記クラスタ構成記憶システム内のポート設定を変更する手段は、論理ポートと論理ボリュームとの間の仮想的なパスである論理パス毎の単位時間当りのデータ転送量またはデータ転送時間を示すアクセス頻度を保持するアクセス情報を参照し、全論理パスの中でデータ転送量が規定値以上の論理パスの 1 つを選択し、選択した論理パスが記憶システムノード間のデータ転送を必要とするか否かをチェックし、記憶システムノード間のデータ転送を必要とした場合、

前記論理パスを経て前記ホストコンピュータがアクセスしている論理ボリュームを格納している記憶装置を有する第 1 の記憶システムノードに未使用ポートがあり、その未使用ポートが当該ホストコンピュータに接続可能である場合、あるいは、当該ホストコンピュータに接続された未使用ポートがある場合に、前記ホストコンピュータが使用するノードを、それまで使用していた第 2 の記憶システムノードのポートから前記第 1 の記憶システムノードの前記未使用ポートとするように、ポート設定を変更することを特徴とするクラスタ構成記憶システム。

【請求項 3】

複数の記憶システムノードを 1 つの記憶システムとして運用可能とするクラスタ構成記憶システムにおいて、

前記記憶システムノード内部及び前記記憶システムノード相互間のアクセス情報を採取する手段と、前記クラスタ構成記憶システムを利用するホストコンピュータにクラスタ構成記憶システム内のポート情報及び前記アクセス情報を提供する手段と、前記ホストコンピュータからのポートの設定指示を受け付ける手段と、前記指示に基づいてクラスタ構成記憶システム内のポート設定を変更する手段とを有し、

前記クラスタ構成記憶システム内のポート設定を変更する手段は、論理ポートと論理ボリュームとの間の仮想的なパスである論理パス毎の単位時間当りのデータ転送量またはデータ転送時間を示すアクセス頻度を保持するアクセス情報を参照し、全論理パスの中でデータ転送量が規定値以上の論理パスの 1 つを選択し、選択した論理パスが記憶システムノード間のデータ転送を必要とするか否かをチェックし、記憶システムノード間のデータ転送を必要とした場合、

前記論理パスを経て前記ホストコンピュータがアクセスしている論理ボリュームを格納している記憶装置を有する第 1 の記憶システムノードに未使用ポートがあり、その未使用ポートが当該ホストコンピュータに接続可能である場合、あるいは、当該ホストコンピュータに接続された未使用ポートがある場合に、前記ホストコンピュータが使用するノードを、それまで使用していた第 2 の記憶システムノードのポートから前記第 1 の記憶システムノードの前記未使用ポートとするように、ポート設定を変更することを特徴とするクラスタ構成記憶システム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、記憶システムに係り、特に、複数の記憶システムを 1 つの記憶システムとして運用可能とするクラスタ構成の記憶システムに関する。

【0002】

【従来の技術】

記憶システムに関する従来技術として、例えば、特開平 11 - 167521 号公報等に記載された技術が知られている。この従来技術は、記憶システムを使用する上位装置（ホス

10

20

30

40

50

トコンピュータ)に対するインタフェース(ホストアダプタ、CHA)、記憶システム内の磁気ディスク装置等の記憶装置に対するインタフェース(ディスクアダプタ、DKA)、キャッシュメモリ(CACHE)、管理メモリ(SM)の相互間をコモンバス方式で接続して構成したものである。

【0003】

図15は従来技術による記憶システムの構成例を示すブロック図であり、以下、図15を参照して従来技術について説明する。図15において、100は記憶システム、110はホストアダプタ(CHA)、120はディスクアダプタ(DKA)、130はキャッシュメモリ(CACHE)、140は管理メモリ(SM)、150は記憶装置(HDD)、160はコモンバス、170は接続線である。

10

【0004】

従来技術による記憶システム100は、図15に示すように、ホストアダプタ110、ディスクアダプタ120、キャッシュメモリ130、管理メモリ140、記憶装置150、コモンバス160、接続線170より構成される。ホストアダプタ110、ディスクアダプタ120、キャッシュメモリ130、管理メモリ140は、コモンバス160により相互間が接続されている。コモンバス160は、コモンバス160の障害時のために2重化されている。ディスクアダプタ120と記憶装置150との間は、1つのディスクアダプタ120あるいは1つの接続線170の障害時にも記憶装置150を使用できるように、1つの記憶装置150に2つのディスクアダプタ120が異なる接続線170で接続されている。

20

【0005】

ホストアダプタ110は、図示しないホストコンピュータとキャッシュメモリ130との間のデータ転送を制御する。ディスクアダプタ120は、キャッシュメモリ130と記憶装置150との間のデータ転送を制御する。キャッシュメモリ130は、ホストコンピュータから受信したデータ、あるいは、記憶装置150から読み取ったデータを一時的に蓄えるメモリである。管理メモリ140は、全てのホストアダプタ110とディスクアダプタ120とが共有するメモリである。また、図示していないが、記憶システム100の設定、監視、保守等を行なうために保守端末(SVP)が全てのホストアダプタ110とディスクアダプタ120とに専用線を用いて接続されている。

【0006】

前述したような構成を持つ記憶システム100のシステム構成を拡張する場合、ホストアダプタ110、ディスクアダプタ120、キャッシュメモリ130、管理メモリ140、記憶装置150等の構成要素が新たに追加される。例えば、ホストコンピュータとの接続数を増やす場合、ホストアダプタ110が新たにコモンバス160に接続される。また、記憶システム100の記憶容量を増やす場合、記憶装置150を追加する、あるいは、ディスクアダプタ120を新たにコモンバス160に接続して記憶装置150を追加する。

30

【0007】

【発明が解決しようとする課題】

前述した従来技術は、記憶システムのシステム拡張に当たり、ホストアダプタ、ディスクアダプタ、キャッシュメモリ、管理メモリ及び記憶装置等の記憶システムの構成要素を増設することにより対応しているため、記憶システムの拡張性が記憶システムの構成要素の最大搭載数に制限されている。この結果、前述の従来技術は、大規模な記憶システムの要求に応じるために、記憶システムの最大搭載数を大きくすると、小規模な記憶システムの要求時にコスト及び設置面積が大きくなってしまふという問題点を有している。

40

【0008】

前述したような問題点を解決する方法として、クラスタ構成の記憶システムが考えられる。クラスタ構成記憶システムは、複数の前述したような記憶システムを接続した構成で、ホストコンピュータからは1つの記憶システムとして運用可能とした記憶システムである。以下、クラスタ構成の記憶システムを構成する記憶システムを記憶システムノードと呼ぶ。クラスタ構成の記憶システムは、小規模な記憶システムの要求時には、少数の記憶シ

50

ステムノードでクラスタ構成記憶システムを構成し、記憶システムの規模を拡大する場合、クラスタ構成記憶システムに記憶システムノードを増設していくことにより対応することができる。このように、クラスタ構成記憶システムは、小規模なシステムから大規模なシステムまで対応することができ、また、ホストコンピュータからは1つの記憶システムとして運用可能であるため、管理が容易になるという利点を有している。

【0009】

しかし、クラスタ構成の記憶システムは、ホストコンピュータからのアクセス命令を受信する記憶システムノードとアクセス対象のデータを保持する記憶システムノードとが異なる場合、記憶システムノード相互間のデータ転送が必要となり、アクセス性能が低下するという問題点を有している。

10

【0010】

本発明の目的は、前述したようなクラスタ構成の記憶システムの問題点を解決し、アクセス性能を向上させることができるクラスタ構成の記憶システムを提供することにある。

【0011】

【課題を解決するための手段】

本発明によれば前記目的は、複数の記憶システムノードを1つの記憶システムとして運用可能とするクラスタ構成記憶システムにおいて、前記記憶システムノード内部及び前記記憶システムノード相互間のアクセス情報を採取する手段と、採取したアクセス情報を保守員に提示する手段と、保守員からのポートの設定指示を受け付ける手段と、前記指示に基づいて前記クラスタ構成記憶システム内のポート設定を変更する手段とを有し、前記クラスタ構成記憶システム内のポート設定を変更する手段は、論理ポートと論理ボリュームとの間の仮想的なパスである論理パス毎の単位時間当りのデータ転送量またはデータ転送時間を示すアクセス頻度を保持するアクセス情報を参照し、全論理パスの中でデータ転送量が規定値以上の論理パスの1つを選択し、選択した論理パスが記憶システムノード間のデータ転送を必要とするか否かをチェックし、記憶システムノード間のデータ転送を必要とした場合、前記論理パスを経て前記ホストコンピュータがアクセスしている論理ボリュームを格納している記憶装置を有する第1の記憶システムノードに未使用ポートがあり、その未使用ポートが当該ホストコンピュータに接続可能である場合、あるいは、当該ホストコンピュータに接続された未使用ポートがある場合に、前記ホストコンピュータが使用するノードを、それまで使用していた第2の記憶システムノードのポートから前記第1の記憶システムノードの前記未使用ポートとするように、ポート設定を変更することにより達成される。

20

30

【0012】

前述において、アクセス情報は、例えば、記憶システムノード内部及び記憶システムノード相互間の単位時間内のデータ転送量またはデータ転送に必要な時間等でであり、アクセス命令を受信した論理ポートとデータの管理単位である論理ボリュームとの組み合わせ単位に採取する。このようにアクセス情報を採取することにより、どの論理ポートからどの論理ボリュームへのアクセスが記憶システムノード間のデータ転送を多く必要としているかが判る。そして、保守員は、保守端末を用いて前記アクセス情報を参照することができ、記憶システムノード相互間のデータ転送の頻度が高い論理ポートと論理ボリュームとの組み合わせを認識でき、その論理ポートの設定変更、その論理ボリュームの再配置の検討が可能となる。さらに、前記クラスタ構成記憶システムは、保守員からの指示により、ポート設定変更、記憶システムノード相互間での論理ボリュームの再配置を行い、前記記憶システムノード相互間のデータ転送の頻度をより小さくすることができ、これにより、クラスタ構成記憶システムのアクセス性能の向上を図ることができる。

40

【0013】

また、本発明は、前記クラスタ構成記憶システム、あるいは、前記クラスタ構成記憶システム内の保守端末に、前記アクセス情報に基づいてポート設定変更及び論理ボリュームの再配置を検討する機能を持たせることができ、これにより、保守員の負担を減らすようにすることができる。

50

【 0 0 1 4 】

また、本発明は、クラスタ構成記憶システムに、前記クラスタ構成記憶システムを使用するホストコンピュータに対して前記アクセス情報を提供する機能、前記ホストコンピュータからポート設定指示、論理ボリュームの再配置指示を受け付ける機能を持たせることにより、ホストコンピュータあるいはホストコンピュータの管理者が、前記アクセス情報及び運用状態に基づいて、ポート設定変更、論理ボリュームの再配置を検討し、クラスタ構成記憶システムにポート設定指示及び論理ボリュームの再配置指示を行うようにすることができる。これにより、保守員には判断できない、高度な条件下でのポート設定変更、論理ボリュームの再配置を行うことが可能となり、クラスタ構成記憶システムのアクセス性能を向上させることができる。

10

【 0 0 1 5 】

【発明の実施の形態】

以下、本発明による記憶システムの実施形態を図面により詳細に説明する。以下に説明する本発明の実施形態は、複数の記憶システムを有するクラスタ構成記憶システムがアクセス情報を採取し、クラスタ構成の記憶システム内の保守端末を通じて保守員に提示し、このアクセス情報に基づく保守員の再配置指示により、アクセス命令を受信する論理ポートの変更、あるいは、クラスタ構成の記憶システム内でデータの記憶装置への再配置を行うものである。なお、以下の説明において、クラスタ構成の記憶システムを構成する記憶システムを記憶システムノードと呼ぶこととする。また、保守端末は、クラスタ構成の記憶システム内に備えられていても、外部に備えられていてもよく、いずれの場合も、各記憶システムノードに接続可能であればよい。

20

【 0 0 1 6 】

図 1 は本発明の一実施形態によるクラスタ構成の記憶システムの構成を示すブロック図、図 2 は記憶システムノードの構成を示すブロック図、図 3 はポート情報の例を説明する図、図 4 は物理位置情報の例を説明する図、図 5 は未使用記憶容量情報の例を説明する図、図 6 はアクセス情報の例を説明する図である。図 1 ~ 図 6 において、200 はクラスタ構成記憶システム、210 - 1 ~ 210 - n は記憶システムノード 210、220 は接続線、230 は論理ボリューム、240 はポート、310 はデータ転送コントローラ (D T C)、400 はポート情報、500 は物理位置情報、600 は未使用記憶容量情報、700 はアクセス情報であり、他の符号は図 15 の場合と同一である。

30

【 0 0 1 7 】

本発明の一実施形態によるクラスタ構成記憶システム 200 は、接続線 220 により相互に接続されている複数の記憶システムノード 210 - 1 ~ 210 - n (以下の説明では、これらを区別する必要のない場合、あるいは、全体を示す場合、単に 210 と記す) と図示しない保守端末とから構成される。保守端末は、専用線を用いて全ての記憶システムノード 210 と接続されている。クラスタ構成記憶システム 200 は、ポート情報 400、物理位置情報 500、未使用記憶容量情報 600、アクセス情報 700 を有する。これらの情報は、保守端末から参照することができる。クラスタ構成記憶システム 200 の記憶領域は、分割して管理されおり、分割した記憶領域を論理ボリューム 230 と呼ぶ。各論理ボリューム 230 のフォーマット形式及び記憶容量は、保守端末を用いて指定することができる。論理ボリューム 230 は、記憶システムノード 210 内であれば、複数の記憶装置 150 に分割して保持することができる。論理ボリューム 230 のクラスタ構成記憶システム 200 内の物理的な格納位置 (物理アドレス) は、後述する物理位置情報 500 に保存されている。

40

【 0 0 1 8 】

記憶システムノード 210 は、基本的に図 15 により説明したものと同様な内部構造を有し、記憶システムノード 210 と図 15 により説明した記憶システム 100 との差異は、記憶システムノード 210 が他の記憶システムノード 210 との通信のためのデータ転送制御コントローラ 310 を備えている点である。そして、記憶システムノード 210 は、1 つ以上のホストアダプタ 110、1 つ以上のディスクアダプタ 120、1 つ以上のキャ

50

ッシュメモリ130、1つ以上の管理メモリ140、1つ以上の記憶装置150、2つ以上のコンパス160、1つ以上の接続線170、1つ以上のデータ転送制御コントローラ310を備えて構成される。ホストアダプタ110、ディスクアダプタ120、キャッシュメモリ130、管理メモリ140、データ転送制御コントローラ310はコンパス160により相互間が接続されている。コンパス160は、コンパス160の障害時のために2重化されてある。ディスクアダプタ120と記憶装置150とは接続線170によって接続されている。また、図示していないが、クラスタ構成記憶システム200の設定、監視、保守等を行なうために保守端末(SVP)が全てのホストアダプタ110とディスクアダプタ120とに専用線を用いて接続されている。

【0019】

ホストアダプタ110は、図示しないホストコンピュータとキャッシュメモリ130との間のデータ転送を制御する。ホストアダプタ110は、ホストコンピュータとの接続のための複数のポート240を持ち、さらに、1つのポート240は、1つ以上の論理的なポート(以下、論理ポートとよぶ)を持つ。ディスクアダプタ120は、キャッシュメモリ130と記憶装置150との間のデータ転送を制御する。キャッシュメモリ130は、ホストコンピュータから受信したデータあるいは記憶装置150から読み出したデータを一時的に保持するメモリである。管理メモリ140は、クラスタ構成記憶システム200内の全てのホストアダプタ110とディスクアダプタ120とが共有するメモリである。ホストアダプタ110及びディスクアダプタ120は、データ転送制御コントローラ310及び接続線220を用いて、他の記憶システムノード210内のホストアダプタ110、ディスクアダプタ120との通信及び他の記憶システムノード210内のキャッシュメモリ130、管理メモリ140の使用が可能である。

【0020】

ポート情報400は、論理ポートを有する記憶システムノード番号及びポート番号と、その論理ポートを使用してアクセスする論理ボリューム番号とホストコンピュータ番号とを保存する。ポート情報400は、ホストアダプタ110が使用可能なメモリ、例えば、管理メモリ140あるいはホストアダプタ110の内部メモリに保存される。図3に示すポート情報400の例において、ホストコンピュータと接続していないポートの論理ポートの論理ボリューム番号及びホストコンピュータ番号には、その論理ポートが未使用であることを表すために、論理ボリューム番号及びホストコンピュータ番号で使用しない数値が設定される。使用していない論理ポートの論理ボリューム番号には当該論理ポートが未使用であることを表すために、論理ボリューム番号で使用しない数値が設定される。

【0021】

図3に示す例の場合、その論理ポートが使用されていないことを示すため“0”が使用されている。図3に示すポート情報において、論理ポート3の情報は、論理ポート3が記憶システムノード1のポート2内の論理ポートであることを示し、ホストコンピュータ2が論理ボリューム1をアクセスするために使用していることを示している。論理ポート5の情報は、論理ポート5が記憶システムノード2のポート1内の論理ポートであることを示し、記憶システムノード2のポート1はどのホストコンピュータにも接続されていないことを示している。また、論理ポート8の情報は、論理ポート8が記憶システムノード2のポート2内の論理ポートであることを示し、その論理ポートがどの論理ボリュームにも使用されていないことを示している。

【0022】

物理位置情報500は、論理ボリューム230の物理アドレス情報、フォーマット形式、容量、状態情報、再配置先物理アドレス情報、再配置完了位置、その論理ボリュームをアクセスするために使用する論理ポート番号を保持する。物理位置情報500は、ホストアダプタ110から参照可能なメモリ、例えば、管理メモリ140あるいはホストアダプタ110の内部メモリに保存される。図4に示す物理位置情報500において、物理アドレス情報は、クラスタ構成記憶システム200内の論理ボリューム230の物理的な格納位置を示す情報であり、例えば、記憶システムノード番号と記憶システムノード内部での物

10

20

30

40

50

理位置とからなる。図示例では、これらがカンマによって区切られて示されている。状態情報は、正常、再配置中等の論理ボリューム230の論理的な状態を表す。再配置先物理アドレス情報と再配置完了位置とは、状態情報が再配置中のときのみ有効である。再配置先物理アドレス情報は、後述する論理ボリュームの再配置決定処理により求められた再配置先の論理ボリュームの物理アドレスである。再配置完了位置は、データの再配置処理が終了した論理ボリューム内の位置である。

【0023】

前述において、状態情報が正常の場合、論理ボリュームの物理アドレスは、物理アドレス情報が用いられ、状態情報が再配置中の場合、論理ボリュームの物理アドレスは、アクセス対象のデータの論理ボリューム内での位置（アクセス命令内の論理アドレス）により物理アドレス情報、あるいは、再配置先物理アドレス情報のどちらか一方が使用される。例えば、論理ボリュームの先頭からデータの再配置処理を行う場合、論理アドレスが再配置完了位置より前であれば、アクセス対象のデータが既に再配置されているため、論理ボリュームの物理アドレスは、再配置先物理アドレス情報が用いられる。一方、論理アドレスが再配置完了位置より後の場合、アクセス対象のデータが再配置されていないため、論理ボリュームの物理アドレスは、物理アドレス情報が用いられる。

【0024】

1つの論理ボリュームは、複数のホストコンピュータから同時にアクセスすることが可能であるため、1つの論理ボリュームの項目に対し、論理ポート番号が複数存在する。図4に示す例において、論理ボリューム1の情報は、記憶システムノード1内の先頭から0の位置から格納されており、フォーマット形式がOPEN3、容量が3GB、データの再配置処理を行っておらず、論理ポート2、3を通じてホストコンピュータからアクセスされることを示している。また、論理ボリューム2の情報は、記憶システムノード2内の先頭から0の位置から格納されており、フォーマット形式がOPEN6、容量が6GB、データの再配置処理を行っており、再配置先が記憶システムノード1内の先頭から500の位置から格納されている論理ボリュームで、データの再配置処理が論理ボリュームの先頭から300の位置まで終了しており、論理ポート1を通じてホストコンピュータからアクセスされることを示している。

【0025】

未使用記憶容量情報600は、記憶システムノードの未使用の記憶容量を保持する。未使用記憶容量情報600は、例えば、管理メモリ140に保存される。この情報は、記憶システムノード210に新たに論理ボリュームを作成できる否かを調べるために使用される。図5に示す未使用容量情報600は、例えば、記憶システムノード1に20GBの未使用の記憶容量があることを示している。

【0026】

アクセス情報700は、論理パス毎にアクセス頻度を保存する。論理パスとは、論理ポートと論理ボリュームとの間の仮想的なパスであると定義する。以下、論理パスの論理ポートを有する記憶システムノードを論理パスのフロントノード、論理パスの論理ボリュームを格納する記憶システムノードを論理パスのエンドノードと呼ぶ。アクセス頻度は、例えば、論理パスの論理ポートを使用して、所定の時間を単位時間として、単位時間、例えば、60秒または30秒内に論理パスの論理ボリュームをアクセスしている時間、データ量等であってよい。アクセス頻度は、ホストコンピュータからのアクセス命令の実行時にホストアダプタ110あるいはディスクアダプタ120が更新する。アクセス情報700は、ホストアダプタ110あるいはディスクアダプタ120が使用可能なメモリ、例えば、管理メモリ140、ホストアダプタ110の内部メモリ、ディスクアダプタ120の内部メモリに保存される。さらに、アクセス情報700は、予め指定された時間毎あるいは保守員の指示によって保守端末に転送され、予め指定された期間あるいは保守員の指示で保守端末内に保存される。保守端末に保存されたアクセス情報700は、保守端末から参照することができ、保守員は、アクセス情報700に基づいて、後述する論理ポートの設定変更決定処理及び論理ボリュームの再配置決定処理を実行する。

10

20

30

40

50

【 0 0 2 7 】

図 6 に示すアクセス情報 7 0 0 の例において、例えば、アクセス頻度をデータ転送量とした場合、図 6 の論理パス 1 のアクセス頻度は、論理ポート 1 を使用して単位時間内に論理ボリューム 2 にアクセスするデータ量が 2 0 (単位は、任意であるが、例えば、MB 等であってよい)であることを示す。ポート情報 4 0 0 が図 3、物理位置情報 5 0 0 が図 4 に示すようなものである場合、論理ポート 1 は、記憶システムノード 1 にあり、論理ボリューム 2 は記憶システムノード 2 にあるため、記憶システムノード相互間に 2 0 のデータ転送が生じていることが判る。図 6 に示す論理パス 2 のアクセス頻度は、論理ポート 2 を使用して論理ボリューム 1 にアクセスするデータ量が 1 5 であることを示す。ポート情報 4 0 0 が図 3、物理位置情報 5 0 0 が図 4 に示すようなものである場合、論理ポート 2 は、記憶システムノード 1 にあり、論理ボリューム 1 は記憶システムノード 1 にあるため、記憶システムノード内部に 1 5 のデータ転送が生じていることが判る。

10

【 0 0 2 8 】

図 7 はクラスタ構成記憶システム 2 0 0 の動作を説明する図であり、以下、これについて説明する。まず、リード/ライト処理時のクラスタ構成記憶システム 2 0 0 の動作について説明する。

【 0 0 2 9 】

ホストアダプタ 1 1 0 は、ホストコンピュータからアクセス命令 7 1 0 を受信する。このホストコンピュータからのアクセス命令 7 1 0 は、リード(またはライト)の命令、リード(またはライト)対象のデータの論理ボリューム 2 3 0 内での位置(論理アドレス)、データ量等を含んでいる。ホストアダプタ 1 1 0 は、アクセス命令 7 1 0 を受信すると、まず、アクセス対象のデータのクラスタ構成記憶システム 2 0 0 内での物理的な格納位置(物理アドレス)を求める(物理位置算出処理 7 2 0)。アクセス対象の物理アドレスは、アクセス対象を含む論理ボリューム 2 3 0 の物理アドレスとアクセス命令 7 1 0 内の論理アドレスとにより一意に決定される。

20

【 0 0 3 0 】

物理位置算出処理 7 2 0 の後、ホストアダプタ 1 1 0 は、アクセス処理 7 3 0 を行う。一例として、アクセス命令 7 1 0 がリード命令の場合で説明する。ホストアダプタ 1 1 0 は、このホストアダプタ 1 1 0 を有する記憶システムノード内あるいは物理位置算出処理 7 2 0 で求めた物理アドレスに対応するディスクアダプタ 1 2 0 を有する記憶システムノード内のキャッシュメモリ 1 3 0 にアクセス命令 7 1 0 内のデータ量と等しいメモリ領域を確保する。ホストアダプタ 1 1 0 は、物理位置算出処理 7 2 0 で求めた物理アドレスに対応するディスクアダプタ 1 2 0 に対し、そのキャッシュメモリ 1 3 0 へのアクセス対象のデータの読み出しを命令する。ホストアダプタ 1 1 0 から命令を受けたディスクアダプタ 1 2 0 は、記憶装置 1 5 0 からアクセス対象のデータをキャッシュメモリ 1 3 0 に読み出し、ホストアダプタ 1 1 0 に転送完了を報告する(アクセス処理 7 4 0)。ホストアダプタ 1 1 0 は、キャッシュメモリ 1 3 0 からホストコンピュータにデータを送信し、リード処理が完了する。その後、ホストアダプタ 1 1 0 あるいはディスクアダプタ 1 2 0 は、アクセス対象の論理パスに対応するアクセス情報 7 0 0 のアクセス頻度を変更する。このアクセス頻度の更新は、常時行ってもよいし、保守員が保守端末を用いてアクセス情報 7 0 0 の更新の可否を指示してもよい。

30

40

【 0 0 3 1 】

保守員は、定期的あるいは必要に応じて、論理ポートの設定変更決定処理 7 7 0 を行い、論理ポートの設定変更による記憶システムノード間のデータ転送量の削減の可能性を検討する。保守員は、この検討の結果、削減が可能と判断すれば、保守端末を通じて、クラスタ構成記憶システム 2 0 0 に再配置指示 7 5 0 を行って、再配置処理 7 6 0 を実行させる。

【 0 0 3 2 】

また、保守員は、定期的あるいは必要に応じて、論理ボリュームの再配置決定処理 7 8 0 を行い、論理ボリュームの再配置による記憶システムノード間のデータ転送量の削減の可

50

能性を検討する。保守員は、この検討の結果、削減が可能と判断すれば、保守端末を通じて、クラスタ構成記憶システム 200 に論理ボリュームの再配置指示 750 を行って、再配置処理 760 を実行させる。

【0033】

図 8 は前述した物理位置算出処理 720 での処理動作を説明するフローチャートであり、次に、これについて説明する。

【0034】

(1) ホストアダプタ 110 は、まず、ポート情報 400 を用いて、アクセス命令 710 を受信した論理ポート番号から論理ボリューム番号を求める(ステップ 900)。

【0035】

(2) 次に、論理ボリューム番号と物理位置情報 500 とから論理ボリューム 230 の物理アドレスを算出するため、まず、アクセス対象の論理ボリューム 230 の状態情報を参照して、その論理ボリュームが再配置中か否かを調べる(ステップ 910)。

【0036】

(3) ステップ 910 の調べで、その論理ボリュームが再配置中であった場合、アクセス命令 710 の論理アドレスとアクセス対象の論理ボリューム 230 の再配置完了位置とを比較し、アクセス命令 710 の論理アドレスが再配置完了位置より前、すなわち、アクセス位置が再配置済みか否かを判定する(ステップ 920)。

【0037】

(4) ステップ 920 の判定で、アクセス位置が再配置済みであった場合、論理ボリューム 230 の物理アドレスとして、再配置先物理アドレス情報を求め、これを用いることを決定する(ステップ 930)。

【0038】

(5) ステップ 910 の調べで、その論理ボリュームが再配置中でなかった場合、または、ステップ 920 の判定で、アクセス命令 710 の論理アドレスが再配置完了位置より後、すなわち、アクセス位置が再配置済みでなかった場合、論理ボリューム 230 の物理アドレスとして、現在の物理アドレス情報を求め、これを用いることを決定する(ステップ 940)。

【0039】

(6) 次に、ステップ 930 の処理またはステップ 940 の処理で求めた論理ボリューム 230 の物理アドレスにアクセス命令 710 の論理アドレスを加えてアクセス対象の物理アドレスを求める(ステップ 950)。

【0040】

図 9 は図 7 における論理ポートの設定変更決定処理 770 での処理動作を説明するフローチャート、図 10 は論理ポートの設定変更決定処理 770 の中で作成される論理パス集合 B1 ~ B3 の例を説明する図、図 13 は論理ポートの設定変更決定処理 770、後述する論理ボリュームの再配置決定処理 780 の中で作成される再配置指示について説明する図であり、以下、アクセス情報 700 のアクセス頻度をデータ転送量として、これらについて説明する。

【0041】

(1) アクセス情報 700 を参照し、全論理パスの中でデータ転送量が規定値以上の論理パスがあるか否かをチェックし、なければ変更不要として処理を終了する(ステップ 1010、1020)。

【0042】

(2) ステップ 1010 のチェックで、データ転送量が規定値以上の論理パスが複数存在する場合、例えば、データ転送量が大きいものから順に選択して以後の処理を行う。いま、選択した論理パスを第 1 論理パスと呼び、この第 1 論理パスの論理ポートを論理ポート A、論理ボリュームを論理ボリューム A とし、論理ポート A を有するポートをポート A、論理ポート A を有する記憶システムノードをフロントノード A、論理ボリューム A を格納する記憶システムノードをエンドノード A とする。さらに、論理ポート A を用いて論理

10

20

30

40

50

ボリューム A をアクセスするホストコンピュータをホストコンピュータ A とする。これらの情報は、ポート情報 400 及び物理位置情報 500 から得ることができる。そして、第 1 論理パスが記憶システムノード間のデータ転送を必要とするか否かをチェックし、不要であった場合、変更不要として処理を終了する(ステップ 1015、1020)。

【0043】

(3) ステップ 1015 のチェックで、第 1 論理パスが記憶システムノード間のデータ転送を必要とした場合、ポート情報 400 を用いて、エンドノード A に未使用のポートが有り、そのポートとホストコンピュータ A とを新たに接続することができるかを判定し、その未使用ポートとホストコンピュータ A とを新たに接続することができる場合、図 13 に示す再配置指示 750 - C を作成する(ステップ 1025、1035)。

10

【0044】

再配置指示 750 - C は、論理ポート 1 と論理ポート 2 との 2 つのパラメータを持ち、論理ポート 1 を用いて論理ボリューム A をアクセスすることをやめ、論理ポート 2 を用いて論理ボリューム A をアクセスすることをクラスタ構成記憶システム 200 に指示する命令であり、再配置指示 750 - C の論理ポート 1 に論理ポート A を、論理ポート 2 にそのポートの任意の論理ポートを設定する。

【0045】

(4) ステップ 1025 の判定で、その未使用ポートとホストコンピュータ A とを新たに接続することができなかつた場合、ポート情報 400 を用いて、エンドノード A にホストコンピュータ A と接続されている未使用の論理ポートがあるか否かをチェックし、未使用の論理ポートがあつた場合、図 13 に示す再配置指示 750 - C を作成する。ここでは、再配置指示 750 - C の論理ポート 1 に論理ポート A を、論理ポート 2 に当該論理ポートを設定する(ステップ 1030、1035)。

20

【0046】

(5) ステップ 1030 のチェックで、エンドノード A にホストコンピュータ A と接続されている未使用の論理ポートが存在しなかつた場合、全論理パスから論理パスのフロントノードがエンドノード A と等しく、論理パスの論理ボリュームをアクセスするホストコンピュータがホストコンピュータ A と等しい論理パス集合 B を取得する。さらに、論理パス集合 B を 3 つの集合 B1 ~ B3 に分割する。論理パス集合 B1 は、論理パス集合 B のうち、論理パスのエンドノードがフロントノード A と等しい論理パスとする。論理パス集合 B2 は、論理パス集合 B のうち、論理パスのエンドノードがエンドノード A と等しくない論理パスとする。論理パス集合 B3 は、論理パス集合 B のうち、論理パスのエンドノードがエンドノード A と等しい論理パスとする(ステップ 1040)。

30

【0047】

論理パス集合 B1 ~ B3 の例を図 10 に示しており、この図 10 において、第 1 論理パスのフロントノードは 2、エンドノードは 1、論理ボリュームは 1、ホストコンピュータは A とする。論理パス集合 B1 の例は、フロントノード 1、エンドノード 2 の論理パスである。論理パス集合 B2 の例は、フロントノード 1、エンドノード 3 の論理パスである。論理パス集合 B3 の例は、フロントノード 1、エンドノード 1 の論理パスである。前述した論理パス集合 B1 ~ B3 は、この順に、論理パスの設定変更を行ったときに、記憶システムノード相互間のデータ転送量の削減効果が大きい。このことは、後述するたの論理パス集合においても同様である。

40

【0048】

(6) ステップ 1040 の処理の後、論理パス集合 B1 が存在するか否かをチェックし、論理パス集合 B1 が存在する場合、論理パス集合 B1 の中で、例えば、最もデータ転送量が多い論理パス(以下、この論理パスを第 2 論理パスと呼ぶ)を選択する。これは、第 1 論理パスと第 2 論理パスとの論理ポートを入れ替えることにより、第 1 論理パスと第 2 論理パスとに関する記憶システムノード相互間のデータ転送量を削減することができるからである。このような第 2 論理パスが存在する場合、図 13 に示す再配置指示 750 - D を作成する(ステップ 1045、1065)。

50

【 0 0 4 9 】

再配置指示 7 5 0 - D は、論理ポート 1 と論理ポート 2 との 2 つのパラメータを持ち、それまで論理ポート 1 を使用してアクセスしていた論理ボリュームを論理ポート 2 を用いてアクセスし、それまで論理ポート 2 を使用してアクセスしていた論理ボリュームを論理ポート 1 を用いてアクセスすることをクラスタ構成記憶システム 2 0 0 に指示する命令である。再配置指示 7 5 0 - D の論理ポート 1 と論理ポート 2 に第 2 論理パスの論理ポート番号と第 1 論理パスの論理ポート番号を設定する。

【 0 0 5 0 】

(7) ステップ 1 0 4 5 のチェックで、論理パス集合 B 1 が存在しなかった場合、論理パス集合 B 2 が存在するか否かをチェックし、論理パス集合 B 2 が存在した場合、論理パス集合 B 2 の中で、任意の論理パス（以下、この論理パスを第 3 論理パスと呼ぶ）を選択する。これは、第 1 論理パスと第 3 論理パスとの論理ポートを入れ替えることにより、第 1 論理パスに関する記憶システムノード間のデータ転送量を削減できるからである。このような第 3 論理パスが存在する場合、図 1 3 に示す再配置指示 7 5 0 - D を作成する。ここでは、再配置指示 7 5 0 - D の論理ポート 1 と論理ポート 2 に第 3 論理パスの論理ポート番号と第 1 論理パスの論理ポート番号を設定する（ステップ 1 0 5 0、1 0 6 5）。

10

【 0 0 5 1 】

(8) ステップ 1 0 5 0 のチェックで、論理パス集合 B 2 が存在しなかった場合、論理パス集合 B 3 が存在するか否かをチェックし、論理パス集合 B 3 が存在した場合、論理パス集合 B 3 の中で、例えば、最もデータ転送量が少ない論理パス（以下、この論理パスを第 4 論理パスとよぶ）を選択する。但し、第 4 論理パスが第 1 論理パスよりデータ転送量が大きい場合は選択しない。これは、第 1 論理パスと第 4 論理パスとの論理ポートを入れ替えることにより、第 4 論理パスに関する記憶システムノード間のデータ転送量は増えるが、第 1 論理パスに関する記憶システムノード間のデータ転送量を削減することができるからである。このような第 4 論理パスが存在する場合は、図 1 3 に示す再配置指示 7 5 0 - D を作成する。ここでは、再配置指示 7 5 0 - D の論理ポート 1 と論理ポート 2 に第 4 論理パスの論理ポート番号と第 1 論理パスの論理ポート番号を設定する（ステップ 1 0 5 5、1 0 6 5）。

20

【 0 0 5 2 】

(9) ステップ 1 0 5 5 のチェックで、論理パス集合 B 3 が存在しなかった場合、あるいは、論理パス集合 B 3 が存在したとしても第 4 論理パスが存在しなかった場合、論理ポートの設定変更による記憶システムノード間のデータ転送量の削減が不可能と判断して処理を終了する（ステップ 1 0 6 0）。

30

【 0 0 5 3 】

図 1 1 は図 7 における論理ボリュームの再配置決定処理 7 8 0 での処理動作を説明するフローチャート、図 1 2 は論理ボリュームの再配置決定処理 7 8 0 の中で作成されるノード間パス集合 A 1 ~ A 3 の例を説明する図であり、以下、アクセス情報 7 0 0 のアクセス頻度をデータ転送量であるとして、これらについて説明する。

【 0 0 5 4 】

(1) まず、アクセス情報 7 0 0 から全てのノード間パスのデータ量を計算する。ノード間パスとは、命令を受信する論理ポートを有する記憶システムノードと論理ボリュームとの間の仮想的なパスであると定義する。また、ノード間パスからの命令を受信する論理ポートを有する記憶システムノードをフロントノードと呼び、ノード間パスの論理ボリュームを格納する記憶システムノードをエンドノードと呼ぶ。ノード間パスのデータ量は、データの再配置により変化する記憶システムノード間のデータ転送量であり、ノード間パスのフロントノード内の任意の論理ポートを使用してノード間パスの論理ボリュームにアクセスするデータ転送量の総和から、ノード間パスのエンドノード内の任意の論理ポートを使用してノード間パスの論理ボリュームにアクセスするデータ転送量の総和を引くことにより算出することができる。例えば、図 3 に示すポート情報 4 0 0、図 4 に示す物理位置情報 5 0 0、図 6 に示すアクセス情報 7 0 0 を用いた場合のフロントノード 4、論理ボリ

40

50

ューム5との間のノード間パスのデータ転送量は、論理パス9のデータ転送量と論理パス10とのデータ転送量の和から論理パス6のデータ転送量を引いたものである。ノード間パスのデータ量が正の場合、ノード間パスの論理ボリュームに関する記憶システムノード間のデータ転送量が記憶システム内部のデータ転送量より大きく、論理ボリュームの再配置を行うことにより、記憶システムノード間のデータ転送量を削減することができることを意味する。また、ノード間パスのデータ量が負の場合、ノード間パスの論理ボリュームに関する記憶システムノード間のデータ転送量が記憶システム内部のデータ転送量より小さく、論理ボリュームの再配置を行うことにより、記憶システムノード間のデータ転送量が増えることを意味する(ステップ1205)。

【0055】

(2)次に、全てのノード間パスの中でデータ転送量が規定値以上のノード間パスがあるか否かを調べ、該当するノード間パスがなければ、再配置不要として処理を終了する(ステップ1210、1220)。

【0056】

(3)ステップ1210のチェックで、該当するノード間パスが存在し、それが複数存在する場合、例えば、データ転送量が多いものから順に選択して以後の処理を行う。いま、選択したノード間パスを第1ノード間パスと呼び、この第1ノード間パスのフロントノードをフロントノードA、論理ボリュームを論理ボリュームAとし、論理ボリュームAを格納する記憶システムノードをエンドノードAとする。これらの情報は、ポート情報400及び物理位置情報500から得ることができる。そして、第1ノード間パスが記憶システムノード間のデータ転送を必要とするか否かを判定し、不要であった場合、再配置不要として処理を終了する(ステップ1215、1220)。

【0057】

(4)ステップ1215の判定で、第1ノード間パスが記憶システムノード間のデータ転送を必要とした場合、未使用記憶容量情報600を用いて、フロントノードAの未使用記憶容量が論理ボリュームAの容量以上あるか否かをチェックして、フロントノードAの未使用記憶容量が論理ボリュームAの容量以上であった場合、図13に示す再配置指示750-Aを作成する(ステップ1225、1235)。

【0058】

再配置指示750-Aは、論理ボリューム1と記憶システムノード2との2つのパラメータをもち、記憶システムノード2に論理ボリューム1と同一フォーマット形式及び同一容量の論理ボリューム2を作成し、論理ボリューム1のデータを論理ボリューム2にコピーし、コピー終了後、論理ボリューム1の物理アドレス情報と論理ボリューム2の物理アドレス情報とを入れ替え、論理ボリューム1の記憶領域を開放することをクラスタ構成記憶システム200に指示する命令である。ここでは、再配置指示750-Aの論理ボリューム1に論理ボリュームAを、記憶システムノード2にフロントノードAを設定する。

【0059】

(5)ステップ1225のチェックで、フロントノードAの未使用領域が論理ボリュームAの容量未満であった場合、全ノード間パスからノード間パスの論理ボリュームを格納する記憶システムノード210がフロントノードAと等しく、ノード間パスの論理ボリュームのフォーマット形式及び容量が論理ボリュームAと等しいノード間パス集合Aを取得する。さらに、ノード間パス集合Aを3つの集合A1~A3に分割する。ノード間パス集合A1は、ノード間パス集合Aのうち、ノード間パスのフロントノードがエンドノードAと等しいノード間パスとする。ノード間パス集合A2は、ノード間パス集合Aのうち、ノード間パスのフロントノードがフロントノードAと等しくないノード間パスとする。ノード間パス集合A3はノード間パス集合Aのうち、ノード間パスのフロントノードがフロントノードAと等しいノード間パスとする(ステップ1240)。

【0060】

ノード間パス集合A1~A3の例を図12に示しており、この図12において、第1ノード間パスは、フロントノード2、論理ボリューム1であるとする。ノード間パス集合A1

10

20

30

40

50

の例は、フロントノード 1、論理ボリューム 10 のノード間パスである。ノード間パス集合 A 2 の例は、フロントノード 3、論理ボリューム 12 のノード間パスである。ノード間パス集合 A 3 の例は、フロントノード 2、論理ボリューム 11 のノード間パスである。

【 0 0 6 1 】

(6) ステップ 1 2 4 0 の処理の後、ノード間パス集合 A 1 が存在するか否かをチェックし、ノード間パス集合 A 1 が存在する場合、ノード間パス集合 A 1 の中で、データ転送量が最も大きいノード間パス(以下、このノード間パスを第 2 ノード間パスとよぶ)を選択する。これは、第 1 ノード間パスの論理ボリューム A を第 2 ノード間パスのエンドノードに格納し、第 2 ノード間パスの論理ボリュームをエンドノード A に格納することにより、第 1 ノード間パスと第 2 ノード間パスに関する記憶システムノード間のデータ転送量を削減することができるからである。そして、このような第 2 ノード間パスが存在する場合、図 1 3 に示す再配置指示 7 5 0 - B を作成する(ステップ 1 2 4 5、1 2 8 0)。

10

【 0 0 6 2 】

再配置指示 7 5 0 - B は、論理ボリューム 1 と論理ボリューム 2 との 2 つのパラメータを持ち、論理ボリューム 1 のデータと論理ボリューム 2 のデータとを入れ替え、入れ替え終了後、論理ボリューム 1 の物理アドレス情報と論理ボリューム 2 の物理アドレス情報とを入れ替えることをクラスタ構成記憶システム 2 0 0 に指示する命令である。ここでは、再配置指示 7 5 0 - B の論理ボリューム 1 と論理ボリューム 2 に論理ボリューム A と第 2 ノード間パスの論理ボリュームを設定する。

【 0 0 6 3 】

(7) ステップ 1 2 4 5 のチェックで、ノード間パス集合 A 1 が存在しなかった場合、ノード間パス集合 A 2 が存在するか否かをチェックし、ノード間パス集合 A 2 が存在した場合、ノード間パス集合 A 2 の中で、任意のノード間パス(以下、このノード間パスを第 3 ノード間パスと呼ぶ)を選択する。これは、第 1 ノード間パスの論理ボリューム A を第 3 ノード間パスのエンドノードに格納し、第 3 ノード間パスの論理ボリュームをエンドノード A に格納することにより、第 1 ノード間パスに関する記憶システムノード間のデータ転送量を削減することができるからである。そして、このような第 3 ノード間パスが存在する場合、再配置指示 7 5 0 - B を作成する。ここでは、再配置指示 7 5 0 - B の論理ボリューム 1 と論理ボリューム 2 に論理ボリューム A と第 3 ノード間パスの論理ボリュームを設定する(ステップ 1 2 5 0、1 2 8 0)。

20

30

【 0 0 6 4 】

(8) ステップ 1 2 5 0 のチェックで、ノード間パス集合 A 2 が存在しなかった場合、ノード間パス集合 A 3 が存在するかチェックし、ノード間パス集合 A 3 が存在した場合、ノード間パス集合 A 3 の中で、第 1 ノード間パスのデータ転送量より小さく、データ転送量が最も小さいノード間パス(以下、このノード間パスを第 4 ノード間パスとよぶ)を選択する。これは、第 1 ノード間パスの論理ボリューム A を第 4 ノード間パスのエンドノードに格納し、第 4 ノード間パスの論理ボリュームをエンドノード A に格納することにより、第 4 ノード間パスに関する記憶システムノード間のデータ転送量は増えるが、第 1 ノード間パスに関する記憶システムノード間のデータ転送量を削減することができるからである。このような第 4 ノード間パスが存在する場合、再配置指示 7 5 0 - B を作成する。ここでは、再配置指示 7 5 0 - B の論理ボリューム 1 と論理ボリューム 2 に論理ボリューム A と第 4 ノード間パスの論理ボリュームを設定する(ステップ 1 2 5 5、1 2 8 0)。

40

【 0 0 6 5 】

(9) ステップ 1 2 5 5 のチェックで、ノード間パス集合 A 3 が存在しなかった場合、あるいは、存在しても、第 4 ノード間パスが存在しなかった場合、論理ボリュームの再配置による記憶システムノード間のデータ転送量の削減が不可能であると判断して処理を終了する(ステップ 1 2 6 0)。

【 0 0 6 6 】

次に、クラスタ構成記憶システム 2 0 0 が再配置指示 7 5 0 を受信して実行する再配置処理 7 6 0 について説明する。クラスタ構成記憶システム 2 0 0 は、再配置指示 7 5 0 - A

50

を受信すると、記憶システムノード2に論理ボリューム1と同一フォーマット形式及び容量の論理ボリューム2を作成し、未使用記憶容量情報600の記憶システムノード2の未使用容量から論理ボリューム1の容量を減じる。そして、後述する再配置処理760-Bと同様の方法で論理ボリューム1のデータを論理ボリューム2にコピーする。データのコピー終了後、論理ボリューム1の物理アドレス情報を論理ボリューム2の物理アドレス情報に変更し、論理ボリューム1の領域を開放し、未使用記憶容量情報600の論理ボリューム1を格納する記憶システムノード1の未使用容量に論理ボリューム1の容量を加える。以上の処理は、論理ボリューム1へのアクセスを停止せずに行うことができる。

【0067】

図14はクラスタ構成記憶システム200が再配置指示750-Bの受信時に実行する論理ボリュームの再配置処理760-Bでの処理動作を説明するフローチャートであり、以下、これについて説明する。すでに説明したように、再配置指示750-Bは、論理ボリューム1と論理ボリューム2とからなる。ここでは、論理ボリューム1を格納する記憶システムノード210をノードA、論理ボリューム2を格納する記憶システムノード210をノードBとして説明する。

10

【0068】

(1)まず、ノードA及びノードBは、再配置処理の1回の処理単位分のメモリをキャッシュメモリ130(キャッシュメモリA、キャッシュメモリBとする)に確保する。キャッシュメモリA及びキャッシュメモリBは、どの記憶システムノード210のキャッシュメモリ130であってもよい(ステップ1500)。

20

【0069】

(2)次に、論理ボリューム1と論理ボリューム2の状態情報を「再配置中」に設定する。そして、論理ボリューム1と論理ボリューム2との再配置完了位置を先頭位置に初期化し、論理ボリューム1の再配置先物理アドレス情報を論理ボリューム2の物理アドレス情報に設定し、論理ボリューム2の再配置先物理アドレス情報を論理ボリューム1の物理アドレス情報に設定する(ステップ1510、1520)。

【0070】

(3)論理ボリューム1と論理ボリューム2との再配置完了位置を調べ、全領域の再配置が完了しているか否かをチェックし、全領域の再配置が完了してした場合、論理ボリューム1と論理ボリューム2との物理アドレス情報を交換し、状態情報を正常に戻し、ステップ1500で確保したキャッシュメモリ130を開放して処理を終了する(ステップ1530、1570)。

30

【0071】

(4)ステップ1530のチェックで、全領域の再配置が完了していなかった場合、ノードAのディスクアダプタ120が、論理ボリューム1の再配置完了位置が示しているデータ位置から再配置処理の1回の処理単位分のデータを記憶装置150からキャッシュメモリAへ読み込む。同様に、ノードBのディスクアダプタ120が、論理ボリューム2の再配置完了位置が示しているデータ位置から再配置処理の1回の処理単位分のデータを記憶装置150からキャッシュメモリBへ読み込む(ステップ1540)。

【0072】

40

(5)次に、ノードAのディスクアダプタ120が、ステップ1540の処理でキャッシュメモリBに格納されたデータを読み取り、そのデータを論理ボリューム1の再配置完了位置が示しているデータ位置に書き込む。同様に、ノードBのディスクアダプタ120が、ステップ1540の処理でキャッシュメモリAに格納されたデータを読み取り、そのデータを論理ボリューム2の再配置完了位置が示しているデータ位置に書き込む(ステップ1550)。

【0073】

(6)ステップ1550の処理の後、1回の処理単位分だけ論理ボリューム1と論理ボリューム2との再配置完了位置を進めて、ステップ1530からの処理に戻って処理を繰り返す(ステップ1560)。

50

【 0 0 7 4 】

前述で説明した再配置処理は、論理ボリューム 1 及び論理ボリューム 2 へのアクセスを停止することなく行うことができる。

【 0 0 7 5 】

前述では、クラスタ構成記憶システム 2 0 0 が再配置指示 7 5 0 - B の受信時に実行する再配置処理について説明したが、前述の例において、別の再配置指示を受信した場合、クラスタ構成記憶システム 2 0 0 は、次のように動作する。

【 0 0 7 6 】

クラスタ構成記憶システム 2 0 0 は、再配置指示 7 5 0 - C を受信すると、ポート情報 4 0 0 の論理ポート 2 の論理ボリューム番号に論理ポート 1 の論理ボリューム A を設定し、論理ポート 1 の論理ボリューム番号を未使用に変更する。クラスタ構成記憶システム 2 0 0 は、変更終了後、保守端末を通じて、保守員に変更終了を知らせる。以上の処理中、論理ボリューム A へのアクセスを停止する必要がある。

10

【 0 0 7 7 】

クラスタ構成記憶システム 2 0 0 は、再配置指示 7 5 0 - D を受信すると、ポート情報 4 0 0 の論理ポート 1 の論理ボリューム番号を論理ポート 2 の論理ボリューム A に変更し、論理ポート 2 の論理ボリューム番号を論理ポート 1 の論理ボリューム B に変更する。クラスタ構成記憶システム 2 0 0 は、変更終了後、保守端末を通じて、保守員に変更終了を知らせる。以上の処理中、論理ボリューム A 及び論理ボリューム B へのアクセスを停止する必要がある。

20

【 0 0 7 8 】

保守員は、新たに論理ポートに論理ボリュームを設定する場合、保守端末を用いて全論理パスのアクセス情報 7 0 0 を参照し、極力記憶システムノード相互間のデータ転送量が増加しないように、論理ポートの設定を行って論理ボリュームを作成する。

【 0 0 7 9 】

次に、ホストコンピュータ A が未使用の論理ポートを用いて既存の論理ボリューム A を使用する場合について説明する。

【 0 0 8 0 】

まず、論理ボリューム A を格納している記憶システムノード A を調べる。次に、記憶システムノード A に未使用のポートが存在するか否かを調べる。記憶システムノード A に未使用のポートが存在し、ホストコンピュータ A とそのポートとを新たに接続することができる場合、記憶システムノード相互間のデータ転送が生じないため、ホストコンピュータ A とその未使用ポートとを接続し、そのポートの任意の論理ポートに論理ボリューム A を設定する。条件を満たすポートが存在しない場合、記憶システムノード A 内のホストコンピュータ A と接続されたポートに未使用の論理ポートが存在するか否かを調べる。条件を満たす論理ポートが存在する場合、記憶システムノード相互間のデータ転送が生じないため、その論理ポートに論理ボリューム A を設定する。

30

【 0 0 8 1 】

前述で条件を満たす論理ポートが存在しない場合、必ず記憶システムノード間のデータ転送が生じるため、任意の未使用の論理ポートに論理ボリューム A を設定する。この場合、保守員が前述した検討を行わずに、保守端末から論理ボリューム A 及び使用してよいポートを指定してもよく、また、保守端末が前述の検討を行って、処理結果を保守端末上に表示させ、処理結果の論理ポートに論理ボリューム A を設定してもよい。さらに、保守員の代わりに、保守端末が自動的に検討結果の論理ポートに論理ボリューム A を設定してもよい。

40

【 0 0 8 2 】

次に、ホストコンピュータ A が新たに作成する論理ボリューム A を使用する場合について説明する。

【 0 0 8 3 】

まず、論理ボリューム A の容量以上の未使用記憶容量を有する記憶システムノード集合 A

50

を調べる。次に、記憶システムノード集合 A 内に未使用のポートを有する記憶システムノード B を調べる。そして、記憶システムノード B が存在し、ホストコンピュータ A と記憶システムノード B の有する未使用のポート B とを新たに接続することができる場合、記憶システムノード相互間のデータ転送が生じないため、記憶システムノード B に論理ボリューム A を作成し、ホストコンピュータ A とポート B とを接続し、ポート B の任意の論理ポートに論理ボリューム A を設定する。条件を満たす記憶システムノードが存在しない場合、記憶システムノード集合 A 内に、ホストコンピュータ A と接続され、ポート内に未使用の論理ポートを有する記憶システムノード C を調べる。記憶システムノード C が存在する場合、記憶システムノード相互間のデータ転送が生じないため、記憶システムノード C に論理ボリューム A を作成し、記憶システムノード C の未使用の論理ポートに論理ボリューム A を設定する。記憶システムノード C が存在しない場合、必ず記憶システムノード相互間のデータ転送が生じるため、記憶システムノード集合 A の任意の記憶システムノードに論理ボリューム A を作成し、任意の未使用の論理ポートに論理ボリューム A を設定する。この場合、保守員が前記検討を行わずに、保守端末に作成する論理ボリュームのフォーマット形式及び容量と使用してよいポートを指定してもよく、また、保守員が保守端末に前記検討を行わせ、検討結果を保守端末上に表示させて、検討結果の記憶システムノードに論理ボリューム A を作成し、検討結果の論理ポートに論理ボリューム A を設定してもよい。さらに、保守員の代わりに、保守端末が自動的に、検討結果の記憶システムノードに論理ボリューム A を作成し、検討結果の論理ポートに論理ボリューム A を設定してもよい。

10

【 0 0 8 4 】

20

前述で説明した本発明の実施形態は、複数の記憶システムを 1 つの記憶システムとして運用可能とするクラスタ構成記憶システム 2 0 0 において、論理パス毎にアクセス情報 7 0 0 を採取し、保守端末を用いて保守員にアクセス情報 7 0 0 を提示し、論理ポートの設定変更決定処理 7 7 0 及び論理ボリュームの再配置決定処理 7 8 0 に基づく保守員の判断により、論理ポートの設定変更処理及び論理ボリュームの再配置処理を行わせることができるので、クラスタ構成記憶システム 2 0 0 を構成する記憶システムノード相互間の通信負荷を抑えることができ、クラスタ構成記憶システム 2 0 0 のアクセス性能の向上を図ることができる。

【 0 0 8 5 】

前述した本発明の実施形態は、論理ポートの設定変更及び論理ボリュームの再配置処理の両方を行うものとして説明したが、本発明は、前述した本発明の実施形態に比較してその効果が小さくなるが、いずれか一方の処理を行うだけでもよい。そして、この場合、システムの構成を簡易化することができるので、システム全体のコストの低減を図ることができる。

30

【 0 0 8 6 】

次に、前述で説明した本発明の実施形態の変形例について説明する。この変形例は、保守員が論理ポートの設定変更決定処理 7 7 0 を行うのではなく、クラスタ構成記憶システム 2 0 0 の保守端末、クラスタ構成記憶システム 2 0 0 内の 1 つのホストアダプタ 1 1 0 あるいはディスクアダプタ 1 2 0 がアクセス情報 7 0 0 を参照して論理ポートの設定変更決定処理 7 7 0 を行い、論理ポートの設定変更による記憶システムノード相互間のデータ転送量の削減が可能な場合に、処理結果のポートの設定変更をホストコンピュータの管理者あるいは保守員に提案するようにしたものである。ホストコンピュータの管理者あるいは保守員は、運用上問題であると判断したら、提案されたポート設定を実行する。

40

【 0 0 8 7 】

次に、前述で説明した本発明の実施形態の他の変形例について説明する。この他の変形例は、保守員が論理ボリュームの再配置決定処理 7 8 0 を行うのではなく、クラスタ構成記憶システム 2 0 0 の保守端末、クラスタ構成記憶システム 2 0 0 内の 1 つのホストアダプタ 1 1 0 あるいはディスクアダプタ 1 2 0 がアクセス情報 7 0 0 を参照して論理ボリュームの再配置決定処理 7 8 0 を行い、論理ボリュームの再配置処理を自動的に実行するようにしたものである。

50

【 0 0 8 8 】

前述した本発明の実施形態の2つの変形例によれば、保守員の負担を低減することが可能となる。

【 0 0 8 9 】

次に、前述で説明した本発明の実施形態のさらに他の変形例について説明する。このさらに他の変形例は、クラスタ構成記憶システム200がアクセス情報700を採取し、クラスタ構成記憶システム200を使用するホストコンピュータにポート情報400、物理位置情報500、未使用記憶容量情報600、アクセス情報700を提供し、ホストコンピュータからの再配置指示750により、再配置処理760を行うようにしたものである。そして、ホストコンピュータあるいはホストコンピュータの管理者は、前述で説明した実施形態の場合と同様に、論理ポートの設定変更決定処理770と論理ボリュームの再配置決定処理780と運用状況とによりポートの設定変更及び論理ボリュームの再配置を行うか否かを決定する。これにより、保守員には判断することが困難な高度な条件下でのポートの設定変更及びデータの再配置を行うことが可能となる。例えば、負荷が低いときに論理ボリュームの再配置を行う等の状況に応じた運用が可能となる。

【 0 0 9 0 】

【発明の効果】

以上説明したように本発明によれば、クラスタ構成記憶システムを構成する記憶システムノード相互間のデータ転送の頻度、データ転送量を削減することができ、クラスタ構成記憶システムのアクセス性能を向上させることができる。

【図面の簡単な説明】

【図1】本発明の一実施形態によるクラスタ構成の記憶システムの構成を示すブロック図である。

【図2】記憶システムノードの構成を示すブロック図である。

【図3】ポート情報の例を説明する図である。

【図4】物理位置情報の例を説明する図である。

【図5】未使用記憶容量情報の例を説明する図である。

【図6】アクセス情報の例を説明する図である。

【図7】クラスタ構成記憶システムの動作を説明する図である。

【図8】物理位置算出処理での処理動作を説明するフローチャートである。

【図9】図7における論理ポートの設定変更決定処理での処理動作を説明するフローチャートである。

【図10】論理ポートの設定変更決定処理の中で作成される論理パス集合B1～B3の例を説明する図である。

【図11】図7における論理ボリュームの再配置決定処理での処理動作を説明するフローチャートである。

【図12】論理ボリュームの再配置決定処理の中で作成されるノード間パス集合A1～A3の例を説明する図である。

【図13】論理ポートの設定変更決定処理、論理ボリュームの再配置決定処理の中で作成される再配置指示について説明する図である。

【図14】クラスタ構成記憶システムが再配置指示の受信時に実行する論理ボリュームの再配置処理での処理動作を説明するフローチャートである。

【図15】従来技術による記憶システムの構成例を示すブロック図である。

【符号の説明】

- 100 記憶システム
- 110 ホストアダプタ (CHA)
- 120 ディスクアダプタ (DKA)
- 130 キャッシュメモリ (CACHE)
- 140 管理メモリ (SM)
- 150 記憶装置 (HDD)

10

20

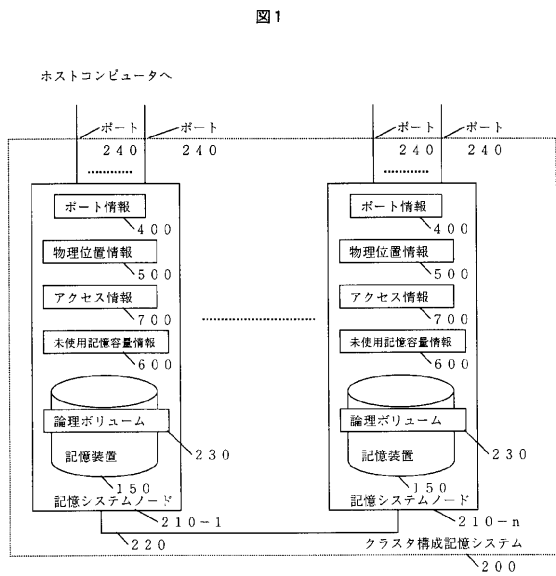
30

40

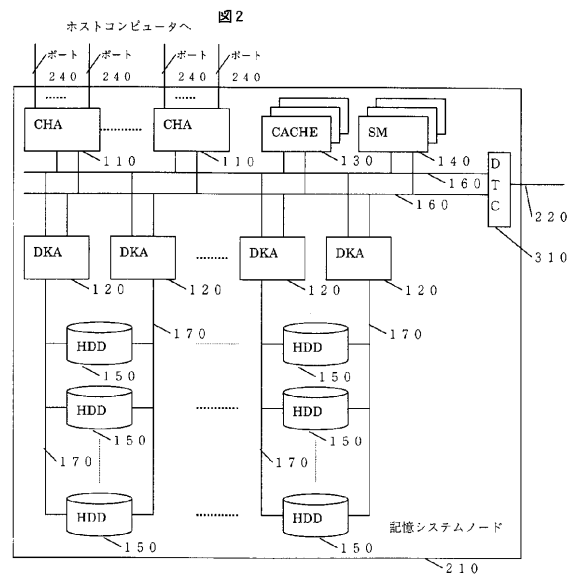
50

- 1 6 0 コモンパス
- 1 7 0 接続線
- 2 0 0 クラスタ構成記憶システム
- 2 1 0 - 1 ~ 2 1 0 - n 記憶システムノード
- 2 2 0 接続線
- 2 3 0 論理ボリューム
- 2 4 0 ポート
- 3 1 0 データ転送コントローラ (D T C)
- 4 0 0 ポート情報
- 5 0 0 物理位置情報
- 6 0 0 未使用記憶容量情報
- 7 0 0 アクセス情報

【 図 1 】



【 図 2 】



【 図 3 】

図 3

論理ポート番号	記憶システムノード番号	ポート番号	論理ボリューム番号	ホストコンピュータ番号
1	1	1	2	1
2	1	1	1	1
3	1	2	1	2
4	1	2	3	2
5	2	1	0	0
6	2	1	0	0
7	2	2	4	2
8	2	2	0	2
9	3	1	5	1
10	3	1	6	1
11	3	2	6	3
12	3	2	0	3
13	4	1	5	1
14	4	1	0	1
15	4	2	5	4
16	4	2	0	4

400 ポート情報

【 図 5 】

図 5

記憶システムノード番号	未使用容量 (GB)
1	20
2	40
3	10
4	0

600 未使用記憶容量情報

【 図 4 】

図 4

論理ボリューム番号	物理アドレス情報	フォーマット形式	容量 (GB)	状態情報	再配置先物理アドレス情報	再配置完了位置	論理ポート番号
1	1,0	OPEN3	3	正常	0,0	0	2,3
2	2,0	OPEN6	6	再配置中	1,500	300	1
3	1,500	OPEN6	6	再配置中	2,0	300	4
4	2,1000	OPEN3	3	正常	0,0	0	7
5	3,0	OPEN3	3	正常	0,0	0	8,13,15
6	3,500	OPEN3	3	正常	0,0	0	10,11

500 物理位置情報

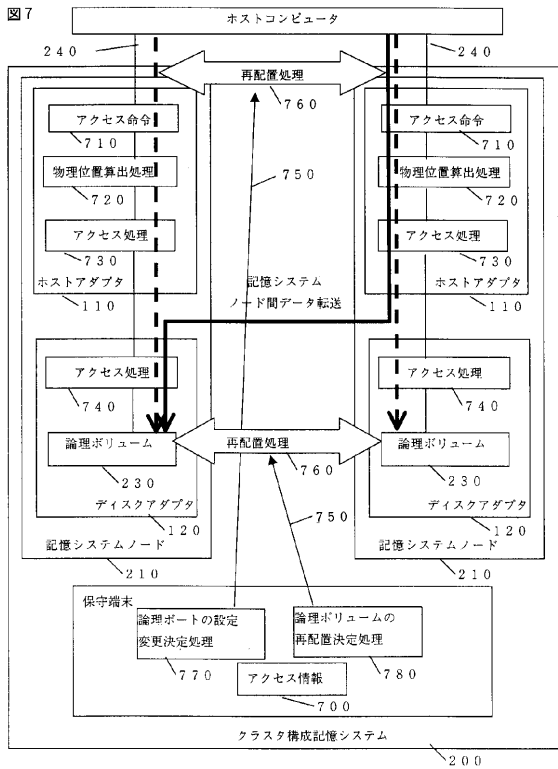
【 図 6 】

図 6

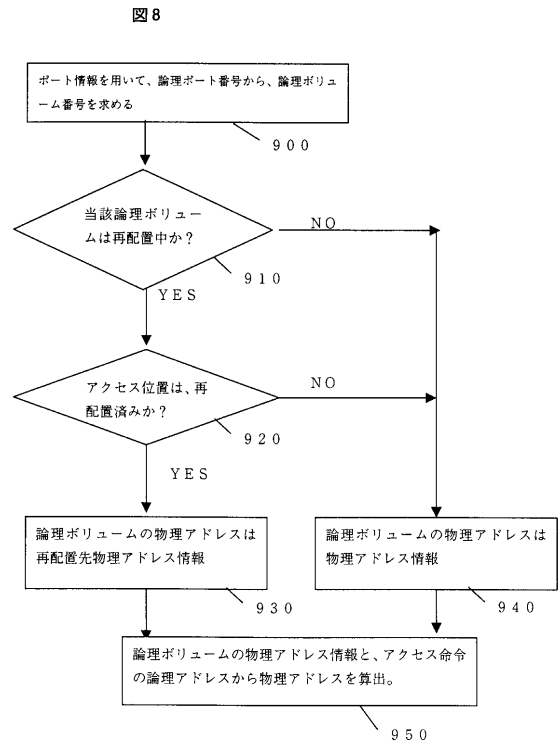
論理バス番号	論理ポート番号	論理ボリューム番号	アクセス頻度
1	1	2	20
2	2	1	15
3	3	1	15
4	4	3	10
5	7	4	13
6	9	5	10
7	10	6	5
8	11	6	2
9	13	5	11
10	15	6	6

700 アクセス情報

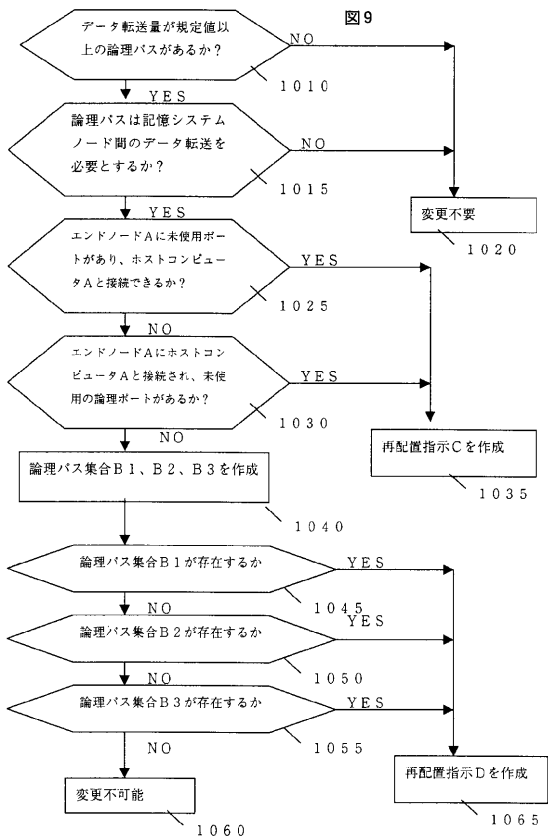
【 図 7 】



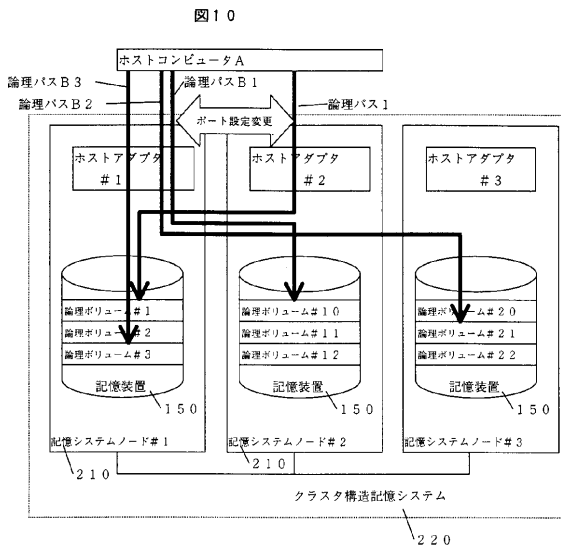
【 図 8 】



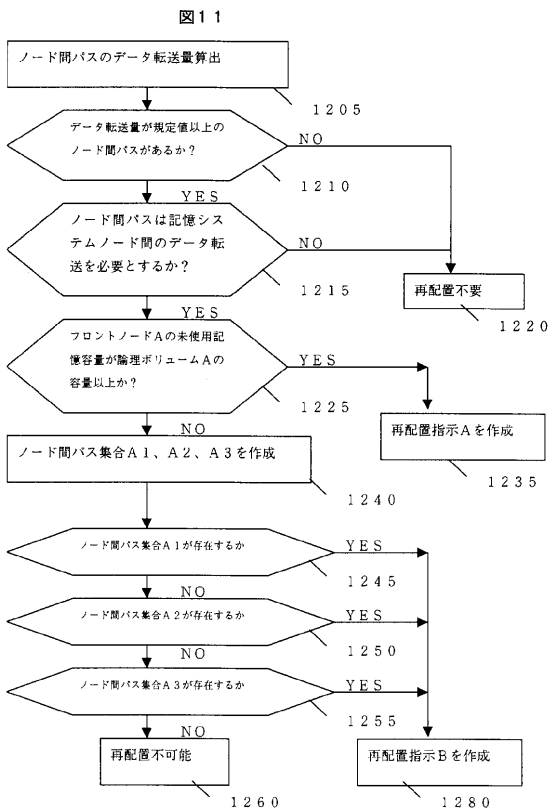
【図9】



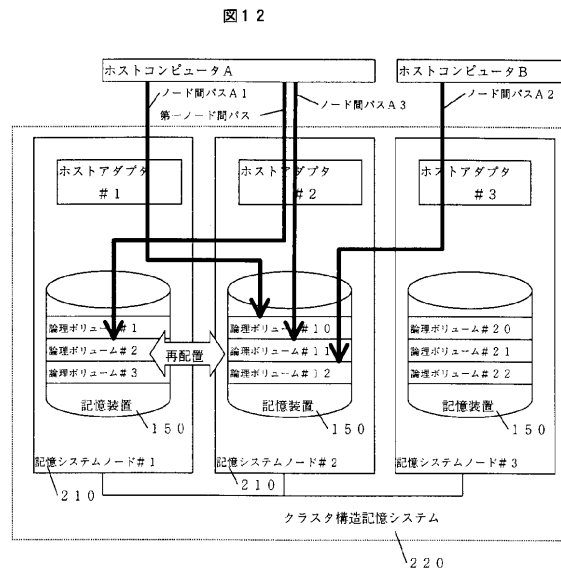
【図10】



【図11】

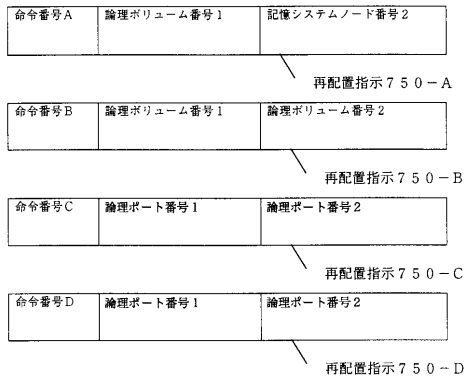


【図12】



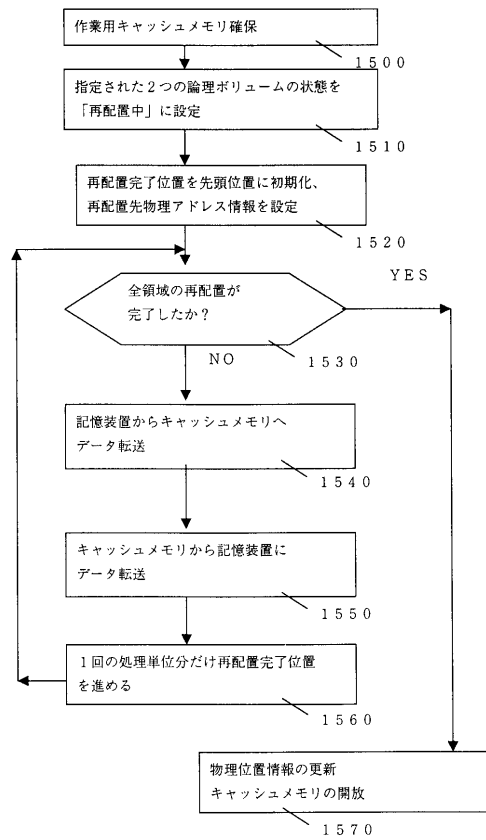
【図 13】

図 13



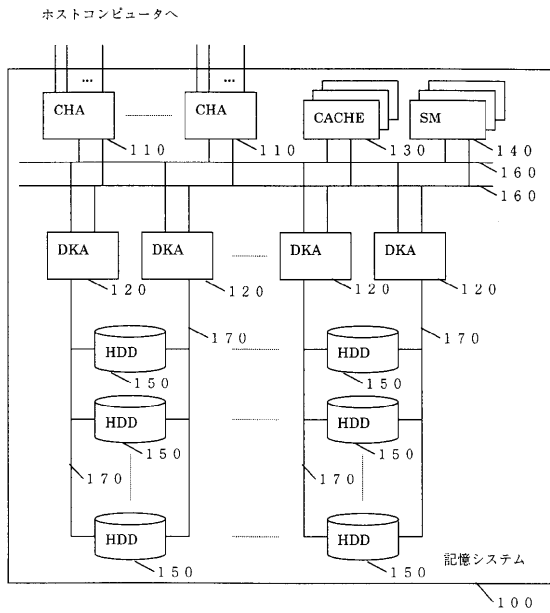
【図 14】

図 14



【図 15】

図 15



フロントページの続き

(72)発明者 大枝 高

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 荒井 弘治

神奈川県小田原市中里322番地2号 株式会社 日立製作所 R A I Dシステム事業部内

審査官 浜岸 広明

(56)参考文献 特開平11-095934(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00