

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4147198号
(P4147198)

(45) 発行日 平成20年9月10日(2008.9.10)

(24) 登録日 平成20年6月27日(2008.6.27)

(51) Int.Cl.	F 1		
G 0 6 F 12/08 (2006.01)	G 0 6 F	12/08	5 5 7
G 0 6 F 3/06 (2006.01)	G 0 6 F	12/08	5 5 1 Z
G 0 6 F 13/10 (2006.01)	G 0 6 F	12/08	5 1 1 Z
	G 0 6 F	12/08	5 4 3 B
	G 0 6 F	3/06	3 0 2 A
請求項の数 20 (全 36 頁) 最終頁に続く			

(21) 出願番号	特願2004-84229 (P2004-84229)	(73) 特許権者	000005108
(22) 出願日	平成16年3月23日(2004.3.23)		株式会社日立製作所
(65) 公開番号	特開2005-275525 (P2005-275525A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成17年10月6日(2005.10.6)	(74) 代理人	110000350
審査請求日	平成18年10月27日(2006.10.27)		ポレール特許業務法人
		(74) 代理人	100068504
			弁理士 小川 勝男
		(74) 代理人	100086656
			弁理士 田中 恭助
		(74) 代理人	100094352
			弁理士 佐々木 孝
		(72) 発明者	渡辺 恭男
			神奈川県川崎市麻生区王禅寺1099番地
			株式会社日立製作所 システム開発研究所内
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】

クラスタ構成のストレージシステムにおいて、
 ホストから入出力要求に従い、読み書きされるデータを格納するキャッシュメモリ、及び
 該キャッシュメモリに格納されるデータを保持するデバイスを有する複数のストレージ制
 御部と、

読み書きされるデータを扱う論理デバイス及びキャッシュメモリを持つ外部ストレージを
 該ストレージ制御部に接続する手段と、

該複数のストレージ制御部のキャッシュメモリの使用状況を監視して把握する手段と、
 該把握手段によって得られるキャッシュメモリの使用状況に関する情報を参照して、キャ
 ッシュメモリの使用量を均等化させるように、いずれかの該ストレージ制御部を選択する
 手段と、を有し、

該選択手段により選択された該ストレージ制御部により該接続手段を介して該外部ストレ
 ージの論理デバイスを制御することを特徴とするストレージシステム。

【請求項2】

前記把握手段は、複数のストレージ制御部の該キャッシュメモリ毎にそのダーティデータ
 量(第1のダーティデータ量)を取得し、

前記選択手段は、得られた該ダーティデータ量に従い、ダーティデータ量のより少ないス
 トレージ制御部を選択して、外部ストレージの論理デバイスを制御することを特徴とする
 請求項1記載のストレージシステム。

【請求項 3】

更に、前記把握手段は、該外部ストレージの論理デバイスに格納されるデータであって、いずれかのストレージ制御部が有するキャッシュメモリに格納されているパーティデータの量（第2のパーティデータ量）を把握し、

前記選択手段は、該第1のパーティデータ量及び該第2のパーティデータ量に従って、該外部ストレージの論理デバイスを制御するためのストレージ制御部を選択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 4】

更に、該外部ストレージの論理デバイスに格納されるデータについて、他のストレージへ非同期的リモートコピーを行う手段と、

該ストレージ制御部の該キャッシュメモリに保持されていて、未だ該他のストレージへ送信されていないデータ（サイドファイルデータ）の量を把握する手段と、を有し、

前記選択手段は、該把握手段によって取得されたサイドファイルデータ量を参照して、該外部ストレージの論理デバイスを制御するためのストレージ制御部を選択することを特徴とする請求項2又は3記載のストレージシステム。

【請求項 5】

前記把握手段により取得されたキャッシュメモリの使用状況、又はパーティデータ量、若しくはサイドファイルデータ量を管理者へ提示するために管理端末又は管理サーバへ転送する手段を有し、

前記選択手段は、該管理端末又は管理サーバから指定されたストレージ制御部を選択することを特徴とする請求項1乃至4のいずれか記載のストレージシステム。

【請求項 6】

クラスタ構成のストレージシステムにおいて、

ホストから入出力要求に従い、読み書きされるデータを一時的に格納するキャッシュメモリ、及び該キャッシュメモリのデータを保持するデバイスを有する複数のストレージノードと、

該ホストから入出力要求に従い、読み書きされるデータを格納するキャッシュメモリ及びデバイスを有する外部ストレージと、該ストレージノードと接続するインターフェースと、

、

該複数のストレージノードの該キャッシュメモリの使用状況を監視して把握する手段と、該把握手段によって得られるキャッシュメモリの使用状況に関する情報を参照して、キャッシュメモリの使用量を均等化させるように、あるストレージノードを選択する手段と、を有し、

該選択手段により選択された該ストレージノードにより該インターフェースを介して外部ストレージのデバイスを制御することを特徴とするストレージシステム。

【請求項 7】

前記把握手段は、該キャッシュメモリ毎のパーティデータ量（第1のパーティデータ量）を取得し、

前記選択手段は、得られた該パーティデータ量に従い、パーティデータ量のより少ないストレージノードを選択して、外部ストレージのデバイスを制御することを特徴とする請求項6記載のストレージシステム。

【請求項 8】

更に、前記複数のストレージノードを含む第2のストレージから第3のストレージに対して非同期的リモートコピーを行う手段と、

該ストレージノードの該キャッシュメモリに保持されていて、未だ該第3のストレージへ送信されていないデータ（サイドファイルデータ）の量を把握する手段と、を有し、

該選択手段は、該把握手段によって取得されたサイドファイルデータ量を参照して、ストレージノードを選択することを特徴とする請求項7記載のストレージシステム。

【請求項 9】

前記把握手段により取得されたキャッシュメモリの使用状況、又はパーティデータ量、若

10

20

30

40

50

しくはサイドファイルデータ量を管理者へ提示するために管理端末又は管理サーバへ転送する手段を有し、

前記選択手段は、該管理端末又は管理サーバから指定されたストレージノードを選択することを特徴とする請求項 6 乃至 8 のいずれか記載のストレージシステム。

【請求項 10】

クラスタ構成のストレージシステムにおいて、

ホストとのインターフェースを有し、該ホストとの間で取り交わす入出力プロトコルをシステム内部のプロトコルに変換する第 1 のプロトコル変換部と、ホストとの間で読み書きされるデータを格納するキャッシュメモリ及びディスク装置を有する第 1 のストレージとのインターフェースを有し、該第 1 のストレージとの間で取り交わす入出力プロトコルをシステム内部のプロトコルに変換する第 2 のプロトコル変換部を含む複数のプロトコル変換部と、

該ホストと該第 1 のプロトコル変換部の間で読み書きされるデータを格納するキャッシュメモリ及びディスク装置を有し、かつ該キャッシュメモリと該ディスク装置へのアクセスを制御する制御部を有する複数のストレージ制御部と、

該複数のプロトコル変換部と該複数のストレージ制御部を管理する構成管理部と、

該構成管理部を介して接続されて該ストレージ制御部と通信を行う管理端末と、

該プロトコル変換部と該ストレージ制御部と該構成管理部とを接続する相互結合網と、

該第 1 のストレージのディスク装置及び該ストレージ制御部に保持されるディスク装置を該ストレージ制御部のディスク装置として該ホストへ提示する手段と、

該ホストから受け付けた入出力要求のアクセス対象が該ストレージ制御部のディスク装置もしくは該第 1 のストレージのディスク装置である場合、該入出力要求を該ストレージ制御部で処理する手段と、

該ストレージ制御部の該キャッシュメモリ上にあつて、該ストレージ制御部のディスク装置若しくは該第 1 のストレージのディスク装置に未反映なライトデータの合計量である第 1 のパーティデータ量を含むキャッシュ使用量情報を取得する手段と、

該管理端末に対して該第 1 のパーティデータ量の情報を提示すると共に、該第 1 のストレージのディスク装置のための処理を行うストレージ制御部の選択指示を該管理端末から受け付ける手段と、

を有することを特徴とするストレージシステム。

【請求項 11】

該第 1 のストレージのディスク装置の入出力処理を行う該ストレージ制御部が該管理端末から指定されなかった場合、

該第 1 のストレージのディスク装置の入出力処理を行う該ストレージ制御部を、該第 1 のパーティデータ量情報を用いて決定することを特徴とする請求項 10 記載のストレージシステム。

【請求項 12】

第 1 のストレージのディスク装置の入出力処理を行うストレージ制御部を、あるストレージ制御部（第 1 ストレージ制御部）から他のストレージ制御部（第 2 ストレージ制御部）に変更する変更手段を有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。

【請求項 13】

該第 1 ストレージ制御部のキャッシュメモリに保持した該第 1 のストレージの該ディスク装置のデータを検索し、第 1 のストレージのディスク装置へ未反映なデータは、第 1 のストレージのディスク装置へ書き込んで前記キャッシュメモリ領域を解放し、反映済みのデータについては前記キャッシュメモリ領域を解放することを特徴とする請求項 12 記載のストレージシステム。

【請求項 14】

第 1 のプロトコル変換部が、第 1 のストレージの該ディスク装置に対応する該第 1 ストレージ制御部のディスク装置への前記ホストからの入出力要求に対して、該第 1 ストレージ制御部で変更処理が終わった部分への入出力要求は該第 2 ストレージ制御部へ行い、変更

10

20

30

40

50

処理が終わっていない部分への要求は該第 1 ストレージ制御部へ振り分けることを特徴とする請求項 13 記載のストレージシステム。

【請求項 15】

前記第 1 ストレージ制御部の前記キャッシュメモリ上において、前記第 1 のストレージのディスク装置に未反映なライトデータの量である第 2 のダーティデータ量情報を取得する手段を有し、

該管理端末に、該第 1 のダーティデータ量情報および該第 2 のダーティデータ量情報を提示し、かつ前記変更処理の対象となる前記第 1 のストレージのディスク装置および変更処理による変更先である該第 2 ストレージ制御部の指定を該管理端末から受け付けることを特徴とする請求項 12 記載のストレージシステム。

10

【請求項 16】

前記変更処理の対象となる前記第 1 のストレージのディスク装置および前記変更処理の変更先である該第 2 ストレージ制御部が該管理端末から指定されなかった場合、該第 1 のダーティデータ量および前記第 2 のダーティデータ量を用いて、変更処理の対象となる前記第 1 のストレージのディスク装置および変更先である該第 2 ストレージ制御部を決定する選択処理を行うことを特徴とする請求項 15 記載のストレージシステム。

【請求項 17】

更に、第 3 のストレージシステムとのインターフェースを有し、該第 3 のストレージシステムとの間で取り交わす入出力プロトコルをシステム内部のプロトコルに変換する第 3 のプロトコル変換部と、

20

該第 3 のストレージシステムのディスク装置に前記第 1 のストレージのディスク装置の複製を作成し、前記第 1 のストレージのディスク装置に対するライトデータを該ストレージ制御部のキャッシュメモリに格納し、該ライトデータを該第 3 のストレージシステムのディスク装置に送信する手段と、

該ストレージ制御部のキャッシュメモリ上において、該第 3 のストレージシステムに未送信なライトデータの合計量である第 1 のサイドファイル量情報を取得する手段と、を有し、

前記管理端末に該第 1 のサイドファイル量情報を提示し、第 1 のストレージのディスク装置の処理を行う該ストレージ制御部の指定を該管理端末から受け付けることを特徴とする請求項 10 記載のストレージシステム。

30

【請求項 18】

クラスタ構成のストレージシステムにおいてホストからの入出力要求を処理する方法であって、

データを格納するデバイス、及び該デバイスに格納されるデータを一時的に格納するキャッシュメモリとを有する複数のストレージ部分で、ホストからの入出力要求を処理するステップと、

データを格納するデバイス及びキャッシュメモリを有する外部ストレージの該デバイスを、あるストレージ部分によって制御するステップと、

該複数のストレージ部分におけるキャッシュメモリの使用状況を把握するステップと、得られたキャッシュメモリの使用状況に関する情報を参照して、キャッシュメモリの使用量を均等化させるように、あるストレージ部分を選択するステップと、

40

選択された該ストレージ部分を用いて該外部ストレージに対する該ホストからの入出力要求を処理するステップと、

を有することを特徴とする入出力処理方法。

【請求項 19】

クラスタ構成のストレージシステムにおいて

ホストとのインターフェースを有し、該ホストとの間で取り交わす入出力プロトコルをシステム内部のプロトコルに変換する第 1 のプロトコル変換部と、ホストとの間で読み書きされるデータを格納するキャッシュメモリ及びディスク装置を有する外部ストレージとのインターフェースを有し、該外部ストレージとの間で取り交わす入出力プロトコルをシス

50

テム内部のプロトコルに変換する第2のプロトコル変換部を含む複数のプロトコル変換部と、

該ホストと該第1のプロトコル変換部の間で読み書きされるデータを格納するキャッシュメモリ及びディスク装置を有し、該キャッシュメモリ及び該ディスク装置へのアクセスを制御すると共に、いずれのストレージ制御部も該外部ストレージと交信可能である複数のストレージ制御部と、

該複数のプロトコル変換部と該複数のストレージ制御部を管理する構成管理部と、

該構成管理部を介して接続されて該ストレージ制御部と通信を行う管理端末と、

該プロトコル変換部と該ストレージ制御部と該構成管理部とを接続する相互結合網と、

該ストレージ制御部の該キャッシュメモリ内に在って、該ストレージ制御部のディスク装置及び該外部ストレージのディスク装置に未反映なデータの合計量であるダーティデータ量を含むキャッシュ使用量情報をメモリに保持して管理する手段と、

該管理端末からの指示に基づき、接続対象となる外部ストレージを探索し、正常な応答のあった外部ストレージに関するデバイス管理情報を該構成管理部内のメモリに登録する外部デバイス定義処理手段と、

論理デバイスの定義の対象が外部ストレージである場合、該ダーティデータ量を含むキャッシュ使用量情報を該管理端末に提供する手段と、

該管理端末からの指示により、該外部ストレージを割当てするためのストレージ制御部に関する情報を受付け、又は該キャッシュ使用量情報の中からダーティデータ量がより少ないストレージ制御部を該外部ストレージに割り当てるようにストレージ制御部を選択する選択手段と、

選択された該ストレージ制御部は、該ストレージ制御部内のディスク装置に関する物理デバイス情報、及び該外部ストレージに関するデバイス管理情報に従って論理デバイスの定義を設定する論理デバイス定義処理手段と、

該管理端末からの指示に基づき、該論理デバイス定義処理手段により定義された論理デバイスに対してパスの設定を行うパス定義処理手段と、

該ホストから受け付けた入出力要求のアクセス対象が、該ストレージ制御部のディスク装置もしくは該外部ストレージのディスク装置である場合、該入出力要求を該ストレージ制御部で処理する手段と、

を有することを特徴とするストレージシステム。

【請求項20】

前記ストレージシステムは、更に、

該管理端末から外部ストレージの割当て変更の指示を受付ける手段と、

変更される先の該ストレージ制御部に関する予測ダーティデータ量を算出して、予測キャッシュ使用量の情報を該管理端末へ提供する手段と、を有し

該構成管理部は、該管理端末から変更指示を受け、変更先のストレージ制御部に割当て変更の指示を発行する手段を有し、

割当て変更の指示を受けた変更先のストレージ制御部は、指示された外部ストレージに対して論理デバイスの登録を行う手段を有し、

変更元のストレージ制御部は、外部ストレージのデータについてキャッシュミス化処理を行う手段を有する、

請求項19に記載のストレージシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージシステムに係り、特にクラスタ構成を有するストレージシステムにおける負荷分散のための制御、及びそのためのキャッシュメモリの制御に関する。

【背景技術】

【0002】

企業のITシステムにおいては大容量かつ高性能なストレージシステムが要求されてい

10

20

30

40

50

る。

この要求に応えるために、小容量のストレージを複数台導入することによって、大容量のデータを取り扱うことが実現できる。しかしストレージの台数の増加に伴い、障害・保守などによるストレージの管理コストの増加が問題となる。一方、1台のストレージにより大容量のシステムを提供するという方法もあるが、従来のような、メモリ、制御メモリ、内部転送機構といった計算機資源を共有するような形態のストレージにおいては、コストおよび技術的な要因により、現在求められているような大容量かつ高性能なストレージを実現するのは困難になってきている。

【 0 0 0 3 】

このような問題を解決するために、例えば特許文献1（USP6, 256, 740明細書）には、ストレージへのクラスタ技術の適用が開示されている。クラスタ技術はこれまで主にサーバなどホストコンピュータの分野で、大量の処理能力を実現する実装技術として用いられてきた。この技術をストレージに適用することで、大規模なストレージを比較的 low コストで実装することが可能となる。このようなストレージをクラスタ構成ストレージと呼ぶ。

クラスタ構成ストレージは、比較的小構成の複数のストレージノードを相互結合網で結合し、大容量の単一ストレージを実現する。クラスタ構成ストレージに対する入出力要求は、対象データを保持するストレージノードに振り分けられ、各ストレージノードにおいて処理される。ストレージノードは、通常のストレージと同様に、ホストインターフェース、ディスク装置、制御プロセッサ、メモリ、制御メモリ、ディスクキャッシュなどが搭載

され、これらの部位がストレージノード内で内部ネットワークにより結合される。各ストレージノードでは、これらの部位を用いてディスク装置への入出力要求処理を行う。一般に、大規模なキャッシュメモリや制御メモリを共有メモリ空間として管理しようとする

【 0 0 0 4 】

また、特許文献2（特開平10-283272号公報）には、アクセスインターフェースの異なるI/Oサブシステム間において、オープン系のI/Oサブシステムで発生したデータを、そのサブシステムが直接接続されていないメインフレーム系のI/Oサブシステムに転送して、メインフレーム系のI/Oサブシステムの記憶装置内にバックアップする複合計算機システムが開示されている。メインフレーム系のディスク制御装置は自記憶装置のアドレスがオープン系のI/Oサブシステムに割当てられているかを示すテーブルを有し、このテーブルを参照することによりオープン系のI/Oサブシステムに対するアクセスを許容する。

【 0 0 0 5 】

【特許文献1】USP6, 256, 740明細書

【 0 0 0 6 】

【特許文献2】特開平10-283272号公報

【 発明の開示 】

【 発明が解決しようとする課題 】

【 0 0 0 7 】

上記特許文献2には、オープン系のI/Oサブシステムがメインフレーム系I/Oシステムの記憶装置を外部記憶装置として使用するための技術が開示されているが、その記憶装置をクラスタ構成ストレージとして使用するための具体的な開示はない。また、オープン系のI/Oサブシステムとメインフレーム系I/Oシステムを2つのストレージサブシステムと想定した場合、これらサブシステムにおけるアクセス頻度に伴う負荷の偏りを如何に防止するかの示唆もされていない。

【0008】

ところで、クラスタ構成ストレージに対して外部に位置付けられる外部ストレージを接続し、クラスタ構成ストレージが提供する機能を外部ストレージで使用する事が考えられている。外部ストレージの外部デバイスに、ストレージノードの機能を提供するために、いずれかのストレージノードのキャッシュメモリを使用する。そのため、あるストレージノードに、多数の外部デバイスに関する処理が集中する場合や、アクセス頻度の高い外部デバイスの処理を行う場合には、外部デバイスのためにキャッシュメモリが大量に使用され、そのストレージノード内のデバイスに対するアクセス性能が低下するおそれがある。このキャッシュメモリの使用量の不均衡に伴うアクセス性能の低下に対する対策についても、上記特許文献のいずれにも記載されていない。

10

【0009】

本発明の目的は、クラスタ構成ストレージにおいて、あるストレージノードに負荷が集中することを防止し、アクセス性能を向上させ得るストレージシステムを提供することにある。

本発明の他の目的は、外部ストレージ接続を行うクラスタ構成ストレージにおいて、各ストレージノードにおけるキャッシュ使用量を均等化させるストレージシステム及びキャッシュメモリの制御方法を提供することにある。

【課題を解決するための手段】

【0010】

本発明は、クラスタ構成のストレージシステムにおいて、ホストから入出力要求に従い、読み書きされるデータを格納するキャッシュメモリ及びキャッシュメモリに格納されるデータを保持するデバイスを有する複数のストレージ制御部と、読み書きされるデータを扱う論理デバイス及びキャッシュメモリを持つ外部ストレージをストレージ制御部に接続する手段と、複数のストレージ制御部のキャッシュメモリの使用状況を監視して把握する手段と、この把握手段によって得られるキャッシュメモリの使用状況に関する情報を参照して、キャッシュメモリの使用量を均等化させるように、いずれかのストレージ制御部を選択する手段と、を有し、選択手段により選択されたストレージ制御部により接続手段を介して外部ストレージの論理デバイスを制御するように構成したものである。

20

好ましくは、上記把握手段は、複数のストレージ制御部のキャッシュメモリ毎にそのパーティデータ量（第1のパーティデータ量）を取得し、選択手段は、得られたパーティデータ量に従い、パーティデータ量のより少ない、例えば最も少ないストレージ制御部を選択して、外部ストレージの論理デバイスを制御する。

30

更に、前記把握手段は、外部ストレージの論理デバイスに格納されるデータであって、いずれかのストレージ制御部が有するキャッシュメモリに格納されているパーティデータの量（第2のパーティデータ量）を把握し、選択手段は、第1のパーティデータ量及び第2のパーティデータ量に従って、外部ストレージの論理デバイスを制御するためのストレージ制御部を選択する。

更に、外部ストレージの論理デバイスに格納されるデータについて、他のストレージへ非同期のリモートコピーを行う手段と、ストレージ制御部のキャッシュメモリに保持されていて、未だ他のストレージへ送信されていないデータ（サイドファイルデータ）の量を把握する手段とを有し、選択手段は、把握手段によって取得されたサイドファイルデータ量を参照して、外部ストレージの論理デバイスを制御するためのストレージ制御部を選択する。

40

更に、前記把握手段により取得されたキャッシュメモリの使用状況、又はパーティデータ量、若しくはサイドファイルデータ量を管理者へ提示するために管理端末又は管理サーバへ転送する手段を有し、前記選択手段は、管理端末又は管理サーバから指定されたストレージ制御部を選択する。

【0011】

本発明に係るストレージシステムは、また、クラスタ構成のストレージシステムにおいて、ホストから入出力要求に従い、読み書きされるデータを一時的に格納するキャッシュ

50

メモリ、及びキャッシュメモリに格納されるデータを保持するデバイスを有する複数のストレージノードと、ホストから入出力要求に従い、読み書きされるデータを格納するキャッシュメモリ及びデバイスを有する外部ストレージと該ストレージノードと接続するインターフェースと、複数のストレージノードのキャッシュメモリの使用状況を監視して把握する手段と、把握手段によって得られるキャッシュメモリの使用状況に関する情報を参照して、キャッシュメモリの使用量を均等化させるように、あるストレージノードを選択する手段とを有し、選択手段により選択されたストレージノードによりインターフェースを介して外部ストレージのデバイスを制御するように構成される。

好ましくは、上記把握手段は、キャッシュメモリ毎のパーティデータ量（第1のパーティデータ量）を取得し、選択手段は、得られた該パーティデータ量に従い、パーティデータ量のより少ないストレージノードを選択して、外部ストレージを制御する。

10

更に、上記複数のストレージノードを含む第2のストレージから第3のストレージに対して非同期のリモートコピーを行う手段と、ストレージノードのキャッシュメモリに保持されていて、未だ該第3のストレージへ送信されていないデータ（サイドファイルデータ）の量を把握する手段とを有し、選択手段は、把握手段によって取得されたサイドファイルデータ量を参照して、ストレージノードを選択する。

また、上記把握手段により取得されたキャッシュメモリの使用状況、又はパーティデータ量、若しくはサイドファイルデータ量を管理者へ提示するために管理端末又は管理サーバへ転送する手段を有し、選択手段は、管理端末又は管理サーバから指定されたストレージノードを選択する。

20

【0012】

本発明は、また、クラスタ構成のストレージシステムにおいてホストからの入出力要求を処理する方法として把握される。すなわち、本発明は、データを格納するデバイス及びデバイスに格納されるデータを一時的に格納するキャッシュメモリとを有する複数のストレージ部分でホストから入出力要求を処理するステップと、データを格納するデバイス及びキャッシュメモリを有する外部ストレージのデバイスを、あるストレージ部分によって制御するステップと、複数のストレージ部分におけるキャッシュメモリの使用状況を把握するステップと、得られたキャッシュメモリの使用状況に関する情報を参照して、キャッシュメモリの使用量を均等化させるように、あるストレージ部分を選択するステップと、選択されたストレージ部分を用いて外部ストレージに対するホストからの入出力要求を処理するステップと、を有する入出力処理方法である。

30

【0013】

好ましい一例では、クラスタ構成ストレージにおいて、ストレージノードを、上位ホストインターフェースを制御するプロトコル変換部と、ディスク装置およびディスクキャッシュを制御するストレージ制御部に分割し、複数のプロトコル変換部と複数のストレージ制御部を相互結合機構で各部位が他の全ての部位と交信可能なように接続する。プロトコル変換部には、ホストインターフェースとそれを制御する制御プロセッサを搭載し、ホストインターフェースを介して受信した入出力要求を対象デバイス（ディスク装置）が属するストレージ制御部へ振り分ける処理を行う。ストレージ制御部には、ディスク装置とディスクキャッシュと制御プロセッサとメモリ、さらにデバイスへの入出力要求処理に必要な制御情報を格納する制御メモリを搭載し、プロトコル変換部より送信され、内部結合機構を介して転送された入出力要求を受信し、対象デバイスへの入出力処理を実行する。また、ストレージ制御部の制御プロセッサでは、データ複製やデータ再配置などのデータ連携機能を実現する各種処理を実行する。ストレージが提供する論理デバイス（以下、上位論理デバイス）と各ストレージ制御部が提供する論理デバイス（以下、下位論理デバイス）の対応は、同じく内部結合機構に接続された構成管理部で管理する。構成管理部では、こうしたデバイス管理に加えて、プロトコル変換部やストレージ制御部などの障害管理も行う。内部結合機構はプロトコル変換部、ストレージ制御部、構成管理部を接続するストレージ内部のネットワークで、各デバイスのアクセスデータや、各部位間の制御情報のやりとりに用いられる。

40

50

この例では、キャッシュメモリおよび制御メモリが各ストレージ制御部内の制御プロセッサ間でだけ共有されるため、メモリ帯域やメモリと制御プロセッサ間を接続するバックプレーン帯域を抑制し、製造コストを削減することができる。また、内部結合機構により、任意のプロトコル変換部が有する任意のホストインターフェースから、任意のストレージ制御部内の任意の論理デバイスへアクセスできる。

【0014】

上記のようなクラスタ構成ストレージで、外部ストレージ接続を行う場合、外部ストレージをホストインターフェースに接続し、外部ストレージ内デバイス（外部デバイス）と下位論理デバイスの対応を構成管理部で管理する。また、ストレージ制御部でも下位論理デバイスと物理デバイスの対応と共に外部デバイスとの対応を管理できる情報を追加する

10

外部ストレージである第一のストレージを上記のようなクラスタ構成を採る第二のストレージに接続し、外部デバイスを第二のストレージデバイスとして統合する場合、外部デバイスを特定のストレージ制御部の下位論理デバイスと対応づけてから下位論理デバイスを上位論理デバイスへ対応づける。これを、外部デバイスのストレージ制御部への割り当てと呼ぶ。そして、ホストからの上位論理デバイスに対する入出力要求を、ホストインターフェースを介して受信したプロトコル変換部で、上位論理デバイスが対応する下位論理デバイスが属するストレージ制御部へ入出力要求を転送する。

【0015】

上記のような構成のクラスタ構成ストレージにおいて外部ストレージ接続を行う場合、プロトコル変換部に接続された外部ストレージは、任意のストレージ制御部と通信可能である。従って、外部デバイスに関する処理は任意のストレージ制御部により行うことが可能である。この際、クラスタ構成ストレージのシステム全体としての性能を向上させるためには、複数のストレージ制御部の中から、外部デバイスを割り当てるストレージ制御部を適切に選択する。この選択基準としてストレージ制御部におけるキャッシュ使用量を用い、複数のストレージ制御部のうち最もキャッシュ使用量の少ないものを選択する。また、ストレージ制御部のキャッシュ使用状況は時間と共に変化するため、複数ストレージ制御部間にキャッシュ使用量の偏りが発生するが、複数ストレージ制御部間の、時間経過に対応したキャッシュ使用量の均等化を行うために、外部デバイスを割り当てるストレージ制御部を動的に変更する。

20

30

【発明の効果】

【0016】

本発明によれば、クラスタ構成ストレージを構成する複数のストレージノードを有するストレージシステムにおいて、あるストレージノードに負荷が集中することを防止でき、アクセス性能を向上できる。また、各ストレージノードにおけるキャッシュ使用量を均等化させることができ、ストレージノードの負荷を分散させることができる。

【発明を実施するための最良の形態】

【0017】

以下、本発明の複数の実施形態について説明される。最初に各実施形態の概要と参照される図面との関係について述べておく。

40

第1の実施形態は、図1～図16を参照し、第1のストレージである外部ストレージ内のデバイスを、第2のストレージであるクラスタ構成ストレージの論理デバイスとして定義する際に、ダーティ属性を持つキャッシュ量に基づき、外部デバイスの割り当て対象とするストレージ制御部を選択し、割り当てる場合の例である。

【0018】

第2の実施形態は、図17～図19を参照し、当該外部デバイスへの入出力処理を受け付けながら、外部デバイスの割り当て対象ストレージ制御部を変更する場合の例である。第3の実施形態は、図20～図23を参照し、第1のストレージである外部ストレージ内のデバイスに対して後述する非同期リモートコピー機能を適用することを前提とし、外部デバイスをクラスタ構成ストレージである第2のストレージの論理デバイスとして定義す

50

る際に、サイドファイル属性を持つキャッシュ量に基づき、割り当て対象とするストレージ制御部を選択し、割り当ての場合の例である。

【0019】

第4の実施形態は、図24を参照し、第1のストレージである外部ストレージ内のデバイスに対して後述する非同期リモートコピー機能を適用することを前提とし、外部デバイスへのホスト入出力処理を受け付けながら、外部デバイスの割り当て対象ストレージ制御部を変更する場合の例である。

第5の実施形態は、図25を参照し、ストレージノード自体を選択或いは切り替える例である。第1乃至第4の実施形態においては、ストレージノードをプロトコル変換部とストレージ制御部に分割し、ストレージ制御部を選択或いは切り替えるシステムであるのに対し、第5の実施形態ではプロトコル変換部とストレージ制御部の分割を行わないシステムを前提とし、負荷の分散を考慮してストレージノードそのものを選択又は切り替える。

【0020】

第1の実施形態：

まず図1から図16を参照して、第1の実施形態について説明する。

図1は、第1の実施形態における計算機システムのハードウェア構成を示す図である。計算機システムは、1台以上のホストコンピュータ(単にホストという)100と、管理サーバ110と、ファイバチャネルスイッチ120と、ストレージ130と、1台以上の外部ストレージ180a、180b(180と総称する)と、管理端末190を含んで構成される。

【0021】

ホスト100、ストレージ130は、それぞれポート107、ポート141によりファイバチャネルスイッチ120のポート121に接続して、SAN(Storage Area Network)を構成する。さらに、外部ストレージ180aと180bは各々ポート181によりストレージ130に接続し、ストレージ130によりストレージ130のデバイスとしてホスト100へ提供される。また、ホスト100やスイッチ120など全機器がIPネットワーク175を介して管理サーバ110に接続され、管理サーバ110で動作するSAN管理ソフトウェア(図示せず)によって統合管理される。なお、本実施形態では、ストレージ130は、管理端末190を介して管理サーバ110に接続する形態をとるものとする。

【0022】

ホスト100は、CPU101やメモリ102などを有する計算機であり、ディスク装置や光磁気ディスク装置などの記憶装置103に格納されたオペレーティングシステムやアプリケーションプログラムなどのソフトウェアをメモリ102に読み上げ、CPU101がメモリ102からそれらのプログラムを読み出して実行することで、所定の機能を達成する。キーボードやマウスなどの入力装置104やディスプレイ105を具備し、ホスト管理者などからの操作を受け付け、指示された情報を表示することが可能である。

【0023】

管理サーバ110も、記憶装置103に格納されたSAN管理ソフトウェアなどをメモリ112に読み上げ、CPU111がそれを読んで実行することで、計算機システム全体の運用・保守管理といった、所定の機能を達成する。また、インターフェース116からIPネットワーク175を介して、計算機システム内の各機器から構成情報、リソース利用率、性能監視情報などを収集し、それらの情報をストレージ管理者にディスプレイ115で提示し、入力装置114から受信した運用・保守指示を各機器に送信する。同処理は、図示しないSAN管理ソフトウェアにより行われる。

【0024】

ファイバチャネルスイッチ120は、複数のポート121を有する。各ポート121には、ホスト100のポート107、および、ストレージ130のポート141のいずれかが接続される。ファイバチャネルスイッチ120は、インターフェース123を有しており、IPネットワーク175にも接続されている。ファイバチャネルスイッチ120は1

10

20

30

40

50

台以上のホスト100がストレージ130を自由にアクセスできるようにするために使用される。この構成では、物理的には全てのホスト100がファイバチャネルスイッチ120に接続されたストレージ130にアクセスすることが可能である。また、ファイバチャネルスイッチ120はゾーニングと呼ばれる特定ポートから特定ポートへの通信を制限する機能を有し、例えば、特定ストレージ130の特定ポート141へのアクセスを特定ホスト100に制限する場合などに用いられる。接続元ポートと接続先ポートの組み合わせを制御する方法については、ファイバチャネルスイッチ120のポート121に割り当てられたポートIDを用いる方法、各ホスト100のポート107やストレージ130のポート141が保持するWWN(World Wide Name)を用いる方法などがある。

10

【0025】

ストレージ130は、複数のプロトコル変換部140と複数のストレージ制御部150と構成管理部160を、相互結合網170で接続して構成される。

プロトコル変換部140は、複数のポート141と1つ以上の制御プロセッサ142とメモリ144を搭載し、ポート141から受信した入出力要求についてアクセス対象デバイスを特定し、転送制御部144より相互結合網170を介して適当なストレージ制御部150へ入出力要求やデータを転送する処理を行う。その際、制御プロセッサ142は入出力要求に含まれるポートIDおよびLUN(Logical Unit Number)からストレージ130がホスト100へ提供する上位論理デバイス番号を算出し、さらに上位論理デバイスが対応するストレージ制御部150と下位論理デバイス番号を算出して、対象となるストレージ制御部150へ入出力要求を転送する。また、プロトコル変換部140は、外部ストレージ180など別のストレージを接続し、ストレージ制御部150からの入出力要求を外部ストレージ180へ送信し、外部ストレージ180へのリード/ライトを行うことができる。なお、本実施形態では、ポート141としては、SCSI(Small Computer System Interface)を上位プロトコルとしたファイバチャネルインターフェースを想定しているが、SCSIを上位プロトコルとしたIPネットワークインターフェースなど、他のストレージ接続用ネットワークインターフェースであっても構わない。

20

【0026】

ストレージ制御部150a~150c(150と総称する)は、1台以上のディスク装置157と、1つ以上の制御プロセッサ152とこれに対応するメモリ153、ディスクキャッシュ154、制御メモリ155、および転送制御部151を有する。制御プロセッサ152は、相互結合網170を介して転送制御部151により受信した当該ディスク装置157への入出力要求を処理する。制御プロセッサ152は、特に、ストレージ130がホスト100に対してディスク装置157単体ではなく、ディスクアレイのように複数のディスク装置157を1つ又は複数の論理デバイスに見せかけている場合には、その処理や管理などを行う。ディスクキャッシュ154は、ホスト100からのアクセス処理速度を高めるため、頻繁に読み出されるデータを格納したり、あるいはホスト100からのライトデータを一時的に格納したりする。制御メモリ155は、ディスク装置157や複数のディスク装置157を組み合わせる物理デバイスや同ストレージ130に接続した外部ストレージ180のデバイス(以下、外部デバイス)の管理や、外部/物理デバイスと下位論理デバイスの対応関係の管理を行うための情報を格納する。

30

40

【0027】

なお、制御メモリ155に格納した制御情報の消失は、ディスク装置157に格納したデータへのアクセスが不可になる事態を引き起こすため、制御メモリはバッテリーバックアップなどによる不揮発化や媒体障害への耐性向上のため二重化などの高可用化を行うことが望ましい。同様にディスクキャッシュ154を用いたライトアフトを行う場合、ディスクキャッシュ154に保持したディスク装置157未反映のデータ消失を避けるため、ディスクキャッシュ154も記録媒体の二重化や不揮発化により可用性を向上させることが望ましい。また、本実施形態におけるストレージ130は、複数のディスク装置157を

50

まとめて1つまたは複数の物理デバイスを定義し、1つの物理デバイスに1つの下位/上位論理デバイスを割り当て、ホスト100に提供する。もちろん、個々のディスク装置157を1つの物理デバイスおよび1つの上位/下位論理デバイスとしてホスト100に見せるようにしてもよい。

【0028】

構成管理部160は、制御プロセッサ162、メモリ163、制御メモリ164、記憶装置165、転送制御部161、インターフェース166を有する。固定ディスク装置などの記憶装置165に格納した制御プログラムをメモリ163に読み上げ、これを制御プロセッサ162で実行することでストレージ130の構成管理という所定の動作を実現する。インターフェース166で接続した管理端末190からストレージ管理者への構成情報提示、管理者からの保守・運用指示の受領を行い、受領した指示に従い、ストレージ130の構成変更などを行う。ストレージ130の構成情報は制御メモリ164に保持する。制御メモリ164上の構成情報はプロトコル変換部140の制御プロセッサ142やストレージ制御部150の制御プロセッサ152から参照・更新することで、各部位間での構成情報の共有を実現する。なお、構成管理部160が障害などにより動作不可に陥った場合、ストレージ130全体がアクセス不可に陥るため、構成管理部160内の各部位の二重化もしくは、構成管理部160自体の二重化が望ましい。

10

【0029】

相互結合網170は、例えばクロスバススイッチであり、プロトコル変換部140、ストレージ制御部150、構成管理部160を接続し、各部位間のデータおよび制御情報、構成情報のやりとりを実現する。この相互結合網170により、構成管理部160が全部位の構成管理と構成情報を配布したり、プロトコル変換部140の任意のポート141からストレージ制御部150の任意の下位論理デバイスへアクセスすることが可能となる。なお、可用性向上の観点から相互結合網も多重化されていることが望ましい。

20

【0030】

管理端末190は、例えばPC(パーソナルコンピュータ)であり、ストレージ管理プログラムを動作させる機能と、ストレージ管理者による入出力操作のための機能を有し、構成情報参照、構成変更指示、特定機能の動作指示など、ストレージ130の保守運用に関するストレージ管理者又は管理サーバ110からのインターフェースとなる。このために、プログラムやデータを格納する記憶装置194、記憶装置194から読み込まれたプログラムや種々のデータを格納するメモリ193、プログラムを実行するCPU192、及び入出力機能として入力装置195、ディスプレイ196を有する。なお、変形例では、この管理端末190を省略して、ストレージ130を直接管理サーバ110へ接続し、管理サーバ110で動作する管理ソフトウェアで管理してもよい。

30

【0031】

外部ストレージ180は、ストレージ130と同様に、ポート181から受信したディスク装置186への入出力要求を処理する機能を有する。即ち、内部インターフェースにポート185を介して接続された大容量のディスク186、ディスクキャッシュ184、メモリ183、および制御プロセッサ182を有する。なおこの例では、外部ストレージ180をストレージ130より小規模構成としているが、ストレージ130と同じ構成、規模のストレージであってもかまわない。

40

【0032】

次に、ストレージ130のソフトウェア構成について説明する。図2は、ストレージ130および管理端末190の制御メモリやメモリに格納する制御情報とストレージ制御処理プログラムを示したソフトウェアの構成図である。なお、以降の説明では、簡略化のため、プロトコル変換部140をPA(Protocol Adaptor)、ストレージ制御部150をSA(Storage Adaptor)、構成管理部160をMA(Management Adaptor)、管理端末190をST(Service Terminal)と表記する。

【0033】

50

この例において、ストレージ 130 内のデバイス階層は次の通りである。SA 150 内では複数のディスク装置 157 によりディスクアレイが構成されて物理デバイスが形成される。また、PA 140 に接続した外部ストレージ 180 の外部デバイスは、PA 140 で認識された後、MA 160 で管理される。SA 150 では物理デバイスおよび外部デバイスに対して下位論理デバイスが割り当てられる。下位論理デバイスは各 SA 150 内における論理デバイスであり、その番号は各 SA 150 で独立的に管理される。下位論理デバイスは、MA 160 で管理する上位論理デバイスに対応付き、ストレージ 130 のデバイスとしてホスト 100 に提供される。ストレージ 130 のキャッシュ使用量情報としては、ST 190 のメモリ 193 に CA キャッシュ使用量情報 222 および外部デバイスキャッシュ使用量情報 224 を格納する。

10

これら種々の管理情報および各種の処理については、図 3 ~ 図 16 を参照して後述される。

【0034】

図 3 及び図 4 はストレージ 130 内のキャッシュ使用量を管理するためのテーブルの構成を示す図である。

これらテーブルの説明の前に、まずディスクキャッシュ 154 の制御方法について説明する。

ホスト 100 から受信した書き込みデータはディスクキャッシュ 154 に格納された時点で、ストレージ制御部 150 はホスト 100 に対して書き込み完了報告を返し、この書き込み完了の報告後に適当なタイミングで物理/外部デバイスに書き込みデータを反映する処理を行う。この処理をライトアフト処理と呼ぶ。ディスクキャッシュ 154 に格納されたデータであって物理/外部デバイスに未反映のデータをダーティデータという。ダーティデータは、物理/外部デバイスに反映されるまでの間、ディスクキャッシュ 154 上に保持し続ける必要がある。

20

【0035】

しかし、ダーティデータの物理/外部デバイスへの書き込み処理速度とホスト 100 からの書き込み要求の頻度によっては、ダーティデータがディスクキャッシュ 154 上に大量に滞留し、新たな入出力処理要求に対してディスクキャッシュ 154 を割り当てることができなくなる可能性がある。このような状態をなるべく回避し、ストレージ 130 全体としてのアクセス性能を向上させるためには、各 SA 150 内デバイスへの書き込み要求頻度を均等化することが有効である。SA 150 内のデバイスに対する書き込み要求頻度はデバイスに対応するダーティ量の時間平均と相関があるので、各 SA 150 内デバイスへの書き込み要求頻度の均等化のためにダーティ量の時間平均を指標として用いる。

30

【0036】

ディスクキャッシュ 154 は、セグメントと呼ばれる、ディスクキャッシュ 154 を一定量に分割した単位により管理される。ダーティデータ量はダーティ属性を持つセグメントの個数により算出される。図示しないが、SA 150 は、制御メモリ 155 上にダーティ属性をもつセグメントの数を保持するダーティセグメントカウンタ、およびストレージ制御部 150 内で管理されている外部デバイスに対応する、ダーティ属性をもつセグメントの数を保持するダーティセグメントカウンタを持つ。

40

【0037】

外部ストレージ接続においては、ストレージ 130 に接続されたホスト 100 が外部デバイスにアクセスする際、ストレージ 130 の SA 150 上のディスクキャッシュ 154 を外部デバイスに対して使用させることにより、外部デバイスに対するアクセス性能を向上させることが可能である。

【0038】

図 3 は、SA キャッシュ使用量情報 222 の構成を示す。

SA キャッシュ使用量情報 222 は、ストレージ 130 内の SA 150 の特定時間帯におけるダーティ量を管理するために、SA 番号 301 および対応するダーティ量情報 302 を保持する。ダーティ量情報 302 は、SA 150 において、後述する集計時間情報 50

50

1に含まれる時間帯に、一定の時間間隔で各SA150に対応するパーティセグメントカウンタを参照し、それらの平均値を求めることにより算出される。このようにして得られたパーティ量情報302はSA番号301と共にST190に送信され、ST190内のCAキャッシュ使用量情報222に格納される。

【0039】

図4は、外部デバイスキャッシュ使用量情報224の構成を示す。

外部デバイスキャッシュ使用量情報224は、ストレージ130で管理されている外部デバイスのパーティ量に関する情報であり、外部デバイス番号401および対応するパーティ量情報402を保持する。即ち、外部デバイスキャッシュ使用量情報224は、外部デバイス番号401で示される外部デバイスを管理しているストレージ130内のSA150において、外部デバイスに格納されるデータとしてSA150内のキャッシュに格納しているパーティデータの量を保持するものである。パーティ量情報402の値は、SA150において、後述する集計時間情報501に含まれる時間帯に、一定の時間間隔で外部デバイスに対応するパーティセグメントカウンタを参照し、それらの平均値を求めることにより算出される。このようにして得られたパーティ量情報402は外部デバイス番号401と共にST190に送信され、ST190内の外部デバイスキャッシュ使用量情報224に格納される。

【0040】

図5は、管理端末制御情報270の構成を示す。

管理端末制御情報270は、ストレージ管理者により入力装置195から入力され、ST190のメモリ193上に保持される。またSA150のメモリ153上には、管理端末制御情報の複製271が保持される。集計時間情報501は前述のSAキャッシュ使用量情報222および外部デバイスキャッシュ使用量情報224に格納するキャッシュ使用量情報を算出する際に、どの時間帯の時間平均を算出するかを指定するものである。割当変更時間情報502は、ST190が構成管理部160に対して外部デバイスの割り当て変更の指示をいつ行うかを指定する情報である。

【0041】

ストレージ130の構成管理情報としては、SA150の制御メモリ155に下位論理デバイス管理情報201と物理デバイス管理情報202とキャッシュ管理情報203を、MA160の制御メモリ164に上位論理デバイス管理情報204と外部デバイス管理情報205とLUPAS管理情報206を格納する。

【0042】

図6は、上位論理デバイス管理情報204の構成を示す。

各上位論理デバイスにつき、上位論理デバイス番号601から接続ホスト名606までの情報組を保持する。

サイズ602には、上位論理デバイス番号601により特定される上位論理デバイスの容量が格納される。対応のSA番号、下位論理デバイス番号603には、上位論理デバイスが対応する下位論理デバイスの番号と下位論理デバイスが属するSA番号が格納される。上位論理デバイスが未定義の場合、このエントリには無効値が設定される。この下位論理デバイス番号が特定SA150の下位論理デバイス管理情報201のエントリ番号となる。

【0043】

デバイス状態604には、上位論理デバイスの状態を示す情報が設定される。状態としては、「オンライン」、「オフライン」、「未実装」、「障害オフライン」が存在する。「オンライン」は、上位論理デバイスが正常に稼動し、ホスト100からのアクセスが可能な状態であることを示す。「オフライン」は、上位論理デバイスは定義され、正常に稼動しているが、LUPAS未定義などでホスト100からのアクセスはできない状態にあることを示す。「未実装」は、上位論理デバイスが定義されておらずホスト100からのアクセスはできない状態にあることを示す。「障害オフライン」は、上位論理デバイスに障害が発生してホスト100からのアクセスができないことを示す。

10

20

30

40

50

なお、本実施形態では、簡単のため、上位論理デバイスは、製品の出荷時にあらかじめディスク装置 157 上に作成された物理デバイスへ割り当てられた、下位論理デバイスへ割り当てられているものとする。このため、利用可能な上位論理デバイスについてはデバイス状態 604 の初期値は「オフライン」状態、その他は「未実装」状態となる。

【0044】

エントリ 605 のポート番号には、上位論理デバイスが複数のポート 141 のどのポートに接続されているかを表す情報が設定される。各ポート 141 には、ストレージ 130 内で一意な番号が割り振られており、上位論理デバイスが LUN 定義されているポート 141 の番号が記録される。また、同エントリのターゲット ID と LUN は、上位論理デバイスを識別するための識別子である。ここでは、これらの識別子として、SCSI 上でホ

10

スト 100 からデバイスをアクセスする場合に用いられる SCSI-ID、LUN が用いられる。
接続ホスト名 606 は、このデバイスにアクセスが許可されているホスト 100 を識別するホスト名である。ホスト名としては、ホスト 100 のポート 107 に付与された WWN (World Wide Name) など、ホスト 100 もしくはポート 107 を一意に識別可能な値であればよい。同じストレージ 130 には、このほかに、各ポート 141 の WWN などの属性に関する管理情報を保持する。

【0045】

図 7 は、LUPAS 管理情報 206 の構成を示す。

ストレージ 130 内の各ポート 141 につき、有効な LUN 分の情報を保持する。ターゲット ID / LUN 702 は、ポート番号 701 に対応する LUN のアドレスを格納する。対応上位論理デバイス番号 703 には、LUN を割り当てた上位論理デバイスの番号を格納する。接続ホスト名 704 は、ポート 141 の LUN に対してアクセスを許可されているホスト 100 を示す情報である。一つの上位論理デバイスに対して複数ポート 141 の LUN が定義されている場合、それら全 LUN の接続ホスト名 704 の和集合が上位論理デバイス管理情報 203 の接続ホスト名 606 に保持される。

20

【0046】

図 8 は、下位論理デバイス管理情報 201 の構成を示す。

各 SA150 毎に下位論理デバイスにつき、下位論理デバイス番号 801 から対応上位論理デバイス番号 805 までの情報組を保持する。

30

サイズ 802 には、下位論理デバイス番号 801 により特定される下位論理デバイスの容量が格納されている。対応物理 / 外部デバイス番号 803 には、下位論理デバイスが対応する SA150 内の物理デバイス番号、もしくは外部デバイス番号が格納される。物理 / 外部デバイスへ未割り当ての場合、このエントリには無効値が設定される。このデバイス番号は、物理デバイス管理情報 202、もしくは外部デバイス管理情報 205 のエントリ番号となる。

デバイス状態 804 には、下位論理デバイスの状態を示す情報が設定される。下位論理デバイスの状態を示す値は上位論理デバイス管理情報 203 のデバイス状態 604 と同様であるため、説明を省略する。対応上位論理デバイス番号 805 には、下位論理デバイスが対応する上位論理デバイス番号が設定される。

40

【0047】

図 9 は、SA150 内のディスク装置 157 から構成される物理デバイスを管理する物理デバイス管理情報 202 の構成を示す。

各 SA150 毎に物理デバイスにつき、物理デバイス番号 901 からディスク内サイズ 909 の情報組を保持する。サイズ 902 には、物理デバイス番号 901 により特定される物理デバイスの容量が格納されている。対応下位論理デバイス番号 903 には、物理デバイスが対応する SA150 内の下位論理デバイス番号が格納される。下位論理デバイスへ未割り当ての場合、当該エントリには無効値が設定される。

【0048】

デバイス状態 904 には、物理デバイスの状態を示す情報が設定される。状態としては

50

、「オンライン」、「オフライン」、「未実装」、「障害オフライン」が存在する。「オンライン」は、物理デバイスが正常に稼動し、下位論理デバイスに割り当てられている状態であることを示す。「オフライン」は、物理デバイスは定義され、正常に稼動しているが、下位論理デバイスに未割り当てであることを示す。「未実装」は、物理デバイスがディスク装置157上に定義されていない状態にあることを示す。「障害オフライン」は、物理デバイスに障害が発生して下位論理デバイスに割り当てられないことを示す。なお、本実施形態では、簡単のため、物理デバイスは製品の工場出荷時にあらかじめディスク装置157上に作成されているものとする。このため、利用可能な物理デバイスについてはデバイス状態64の初期値は「オフライン」状態、その他は「未実装」状態となる。

【0049】

RAID構成905には、物理デバイスが割り当てられたディスク装置157のRAIDレベル、データディスクとパリティディスク数などRAID構成に関連する情報が保持される。同じように、ストライプサイズ906には、RAIDにおけるデータ分割単位(ストライプ)長が保持される。ディスク番号リスト907には、物理デバイスが割り当てられたRAIDを構成する複数のディスク装置157の番号が保持される。この番号はSA150内でディスク装置157を識別するために付与した一意な値である。ディスク内開始オフセット908とディスク内サイズ909には、物理デバイスデータが各ディスク装置157内のどの領域に割り当てられているかを示す情報である。本実施形態では簡単のため、全物理デバイスについてRAIDを構成する各ディスク装置157内のオフセットとサイズを統一している。

【0050】

図10は、外部デバイス管理情報205の構成を示す。

この管理情報205は、ストレージ130に接続し、ストレージ130の下位論理デバイスと対応づける、外部ストレージ180内のデバイスを管理するために、ストレージ130全体で各外部デバイス毎に、外部デバイス番号1001乃至ターゲットポートID/ターゲットID/LUNリスト1008の組を保持する。

外部デバイス番号1001にはストレージ130で割り当てた、ストレージ130内で一意な値を保持する。サイズ1002には、外部デバイス番号1001により特定される外部デバイスの容量が格納される。対応SA番号、下位論理デバイス番号1003には、外部デバイスが対応するストレージ130内のSA番号と下位論理デバイス番号が格納される。下位論理デバイスへ未割り当ての場合、このエントリには無効値が設定される。

デバイス状態1004には、当該外部デバイスの状態を示す情報が設定されるが、各状態の意味は物理デバイス管理情報202内デバイス状態904と同じである。ストレージ130の初期状態は外部ストレージを接続していないため、デバイス状態1004の初期値は「未実装」となる。

【0051】

ストレージ識別情報1005には、外部デバイスを搭載する外部ストレージ180の識別情報を保持する。ストレージ識別情報としては、同ストレージのベンダ識別情報と各ベンダが一意に割り振る製造シリアル番号の組み合わせ、などが考えられる。外部ストレージ内デバイス番号1006には、外部デバイスが対応する外部ストレージ180内のデバイス番号を保持する。外部デバイスは外部ストレージ180の論理デバイスであるから、本エントリには外部ストレージ180の論理デバイス番号を保持する。

PA番号/イニシエータポート番号リスト1007には、外部デバイスへアクセス可能なストレージ130のポート141とそれが属するPA140の番号のリストが保持される。ターゲットポートID/ターゲットID/LUNリスト1008には、外部デバイスが外部ストレージ180の1つ以上のポート181にLUN定義されている場合、それらのポート181のポートIDおよび外部デバイスが割り当てられたターゲットID/LUNが1つ又は複数個保持される。

【0052】

次に、再び図2に戻り、各部位のメモリ内に格納される情報およびプログラムについて

10

20

30

40

50

説明する。

制御メモリ 155, 164 に格納した各制御情報は各部位の制御プロセッサから参照・更新可能であるが、その際に相互結合網 170 を介したアクセスが必要となる。よって、処理性能向上のため、各制御プロセッサで動作する処理に必要な制御情報の複製を各部位のメモリに保持する。構成変更により制御情報が更新された場合は、相互結合網 170 を介してその旨を他部位に通知し、最新情報を制御メモリから各部位メモリへ取り込ませる。

別の方法としては、例えば、制御メモリに保持した構成情報毎に制御メモリ上に更新有無を示すフラグを設け、各部位の制御プロセッサは処理開始時もしくは各構成情報を参照する毎に同フラグを参照して更新有無をチェックする方法などが考えられる。また、各部位のメモリにはこうした制御情報の複製に加えて、各制御プロセッサで動作する制御プログラムが格納される。

【0053】

この実施例では、外部デバイスを含むデバイスを特定サーバへ割り当ててストレージ 130 において使用可能とする処理、外部デバイスを含むストレージ 130 のデバイスへの入出力要求処理と、割り当てられた外部デバイスの割当を変更する処理を例にしてその制御方法を説明する。

外部デバイスを含むデバイスを特定サーバへ割り当てて使用可能とする処理については、外部デバイス定義、論理デバイス定義、LUパス定義の大きく3つの処理に分けることができる。

【0054】

図 11 は、外部デバイス定義処理 253 の処理フローを示す。

外部デバイス定義処理 253 は、ストレージ 130 管理下に外部ストレージ 180 のデバイスを外部デバイスとして導入する場合の処理である。

まず、管理端末 190 又は管理サーバ 110 からの外部ストレージ 180 接続指示を受け付けた ST 190 が MA 160 に接続指示を送信する (1101)。接続指示には、対象となる外部ストレージ 180 を特定する情報、例えば、外部ストレージ 180 のポート 181 の WWN や外部ストレージデバイスへの Inquiry コマンド送信などにより得られる装置識別情報もしくはその両方の情報と外部ストレージ 180 へ接続するポート 141 番号が付加される。MA 160 では外部ストレージ 180 の接続指示を受信し、指示された全ポート 141 番号に対応する全 PA 140 へ外部ストレージ接続指示を送信する (1102)。

PA 140 では接続指示に付加された外部ストレージ 180 識別情報を用いて接続対象となる外部デバイスを探索する (1103)。具体的には、外部ストレージ識別情報としてポート 181 の WWN をもらう場合は、指定されたポート 141 から外部ストレージ 180 のポート 181 の全 LUN に対して Inquiry コマンドを送信して、正常な応答のあった LUN を外部デバイス登録候補とする。識別情報として装置識別情報しかもらわない場合には、PA 140 の全ポート 141 から検出した全ノードポートのそれぞれに対して (ノードポートログイン時に検出済み)、全 LUN について Inquiry コマンドを送信し、正常に応答のあったデバイスについて、返却情報内の装置識別情報を接続指示に付加された値と比較する。そして、検出した外部デバイス登録候補についての情報リストを PA 140 は MA 160 へ返却する (1104)。このときの情報には外部デバイスについて外部デバイス管理情報 205 を設定するのに必要な情報を含んでいる。

【0055】

MA 160 では、受け取った外部デバイスリストにある外部デバイスに採番し、外部デバイス管理情報 205 へデバイス情報を登録し、同情報に更新のあった旨を各部へ通知する (1105)。外部デバイス管理情報 205 への情報登録について、具体的には、割り当てた外部デバイス番号のエントリにおいて、Inquiry 情報によりサイズ 1002、ストレージ識別情報 1005、外部ストレージ内デバイス番号 1006 を、PA 140 からの情報によりエントリ 1007、1008 をそれぞれ設定する。エントリ 1003 の

10

20

30

40

50

デバイス番号は未割当のため初期値である無効値を設定する。デバイス状態 1 0 0 4 には「オフライン」状態が設定される。

更新通知を受けた各 S A 1 5 0 , P A 1 4 0 では、M A 1 6 0 の制御メモリ 1 6 4 の外部デバイス管理情報 2 0 5 を各々のメモリに取り込む。S T 1 9 0 では同情報のメモリ取込みに加えて、外部デバイス定義処理の完了を要求元である管理端末 1 9 0 又は管理サーバ 1 1 0 に報告する (1 1 0 6) 。

【 0 0 5 6 】

なお、本実施例では、ストレージ 1 3 0 に対して管理端末 1 9 0 又は管理サーバ 1 1 0 が接続指示と共に、導入対象となる外部ストレージ 1 8 0 を指示する形態をとる。しかし、外部ストレージ 1 8 0 の接続指示のみをストレージ 1 3 0 に指示し、ストレージ 1 3 0 で全ポート 1 4 1 から検出した全ストレージの全デバイスを外部デバイスとして登録しても構わない。また、特に明示的な接続指示を与えず、ストレージ 1 3 0 に外部ストレージ 1 8 0 が接続された契機にストレージ 1 3 0 が検出可能な全デバイスを外部デバイスとして登録してもよい。

10

【 0 0 5 7 】

図 1 2 は、論理デバイス定義処理 2 5 5 の処理フロー図である。

論理デバイス定義処理 2 5 5 は、管理端末 1 9 0 又は管理サーバ 1 1 0 からの指示を受けて、ストレージ 1 3 0 が搭載する物理デバイスもしくは外部デバイス定義処理 2 5 3 で定義された外部デバイスに対して下位論理デバイスを定義する処理である。

まず、S T 1 9 0 は論理デバイス定義指示を受け付ける (1 2 0 1) 。この指示には、論理デバイス定義対象となる物理 / 外部デバイス番号と定義する上位論理デバイス番号が付加される。論理デバイス定義対象となるデバイスが外部デバイスでない場合は、さらに、S A 番号と下位論理デバイス番号が付加される。

20

【 0 0 5 8 】

なお、本実施形態では、説明の簡略化のため 1 つの物理 / 外部デバイスに対して 1 つの論理デバイスを割り当てるものとするが、2 つ以上の物理 / 外部デバイスからなるデバイスグループに対して 1 つの論理デバイスを、1 つの物理 / 外部デバイスに対して 2 つ以上の論理デバイスを、2 つ以上の物理 / 外部デバイスからなるデバイスグループに対して 2 つ以上の論理デバイスを定義しても構わない。ただし、それぞれの場合、下位論理デバイス管理情報 2 0 1 に、物理 / 外部デバイス内での当該論理デバイスの開始位置およびサイズなどの付加情報が必要となる。論理デバイス定義対象となるデバイスが物理デバイスの場合は、S T 1 9 0 は M A 1 6 0 に対して論理デバイス定義指示を送信する (1 2 0 7) 。

30

【 0 0 5 9 】

論理デバイス定義対象となるデバイスが外部デバイスの場合は (1 2 0 2) 、S T 1 9 0 はストレージ管理者に S A キャッシュ使用量情報 2 2 2 を提示し、必要に応じて、外部デバイスを割り当てる S A 番号および下位論理デバイス番号を受け付ける (1 2 0 3) 。ステップ 1 2 0 3 において外部デバイスの割当対象 S A 番号および下位論理デバイス番号が指定されていた場合には (1 2 0 4) 、S T 1 9 0 は M A 1 6 0 に対して論理デバイス定義指示を送信する (1 2 0 7) 。

40

ステップ 1 2 0 4 において、外部デバイスの割当対象 S A 番号および下位論理デバイス番号が指定されていない場合には、S T 1 9 0 は外部デバイスの割当対象 S A 1 5 0 を選択するために、S A キャッシュ使用量情報 2 2 2 を参照する (1 2 0 5) 。S T 1 9 0 は S A キャッシュ使用量情報 2 2 2 の中から最もパーティ量の少ない S A 1 5 0 を割り当て対象として選択する (1 2 0 6) 。また、S A 1 5 0 における未使用下位論理デバイス番号を選択する。その後、M A 1 6 0 に対して論理デバイス定義指示を送信する (1 2 0 7) 。

【 0 0 6 0 】

論理デバイス定義指示を受信した M A 1 6 0 は、定義指示情報より対象 S A 1 5 0 を特定し、S A 1 5 0 に定義指示を送信する (1 2 0 8) 。対象 S A 1 5 0 では、指示された物理 / 外部デバイスに対して下位論理デバイスを登録する (1 2 0 9) 。具体的には、下

50

位論理デバイス管理情報 201 の対象デバイスエントリについて、物理 / 外部デバイスのサイズとデバイス番号をエントリ 802、803 に、対応上位論理デバイス番号 805 に上位論理デバイス番号を、デバイス状態 804 に「オンライン」をそれぞれ設定する。また、物理 / 外部デバイスの対応 SA 番号 / 下位論理デバイス番号を設定し、デバイス状態を「オンライン」に更新する。登録が完了したら SA 150 は MA 160 にその旨通知する。

次に MA 160 では、対象上位論理デバイスを下位論理デバイスへ定義し、情報更新の旨を各部位に通知する (1210)。具体的には、上位論理デバイス管理情報 203 のデバイスエントリについてサイズ 602 と対応 SA 番号、下位論理デバイス番号 603 を設定し、デバイス状態を「オフライン」状態に、エントリ 605、606 は未割当のため無効値を設定する。情報更新を通知された PA 140 では更新のあった管理情報をメモリ 143 へ取り込み、ST 160 では情報取込み後、論理デバイス定義処理の完了を要求元へ報告する (1211)。

【0061】

図 13 は、LUパス定義処理 252 の処理フロー図である。

LUパス定義指示を受け付けた ST 190 が MA 160 に対して同指示を転送する (1301)。同指示には、対象上位論理デバイス番号と LU を定義するポート 141 番号と LUN に加えて、同 LU をアクセスするホスト 100 の識別情報 (ポート 107 の WWN など) が付加される。LUパス定義指示を受け付けた MA 160 は、上位論理デバイスに対して LUパス登録を行う (1302)。具体的には、上位論理デバイス管理情報 203 のデバイスエントリのポート番号、ターゲット ID、LUN 605 と接続ホスト名 606 に対応情報を設定し、LUパス管理情報 206 のポート 141 に対応する空きエントリにターゲット ID / LUN 701 を始めとする構成情報を設定する。登録が完了したら、その旨各部位へ通知し、PA 140 で情報取込み、SA 190 で情報取込みと要求元への完了報告を行う。

【0062】

以上の 3 つの処理により、外部デバイスをストレージ 130 のデバイスとして登録し、ストレージ 130 内の SA 150 のキャッシュ使用量の均等化を考慮していずれかの SA 150 へ割り当て、ホスト 100 からのアクセスが可能となる。

【0063】

次に、このように外部デバイスが割り当てられている状態での、ホスト 100 からの入出力要求処理方法について、PA 140、SA 150 でのリードコマンド処理、ライトコマンド処理、ライトアフタ処理の 3 つに分けて説明する。

【0064】

図 14 は、リードコマンド処理 261 の処理フロー図である。

ホスト 100 がストレージ 130 内の外部デバイスを含むデバイスのデータを読み出す場合の処理を説明するものである。

PA 140 は、ホスト 100 が発行したリードコマンドを特定ポート 141 で受信する (1401)。PA 140 は、受信したコマンドの解析を行い、要求データに対応する上位論理デバイス番号を割り出し、上位論理デバイス管理情報から対応する SA 番号、下位論理デバイス番号を算出し (1402)、SA 番号に対応する SA 150 にリード要求を転送する (1403)。SA 150 は転送制御部 151 よりリード要求を受け取り、要求データがディスクキャッシュ 154 上にあるかどうかを、キャッシュ管理情報 203 を参照することにより判定する。

【0065】

要求データがディスクキャッシュ 154 上にある場合 (キャッシュヒット)、SA 150 は該当するデータを要求の発行元である PA 140 に送信する (1410)。SA 150 よりデータを受信した PA 140 は、ポート 141 を通じて、ホスト 100 へデータを送信する (1411、1412)。

【0066】

10

20

30

40

50

一方、要求データがディスクキャッシュ154上に無い場合（キャッシュミス）には、SA150はキャッシュ管理情報203を更新し、ディスクキャッシュ154上に要求データを格納する領域を確保する。要求が外部デバイスに対するものでない場合は、物理デバイスから要求データを読み出し、対応するディスクキャッシュ154領域にデータを格納する。以降の動作フローはキャッシュヒットした場合と同様である（141～1412）。要求が外部デバイスに対するものである場合、SA150は、PA140を介して外部ストレージ180からデータを読み込み（1413～1416）、対応するディスクキャッシュ154領域にデータを格納する（1409）。以降の動作フローはキャッシュヒットした場合と同様である（1410～1412）。以上のようにして、ホスト100からのリード要求に対して、外部デバイスを含むデバイスからデータを読み出して、ホスト

10

【0067】

図15は、ライトコマンド処理262の処理フロー図である。ホスト100がストレージ130内のディスクキャッシュ154にデータを格納する場合の処理を説明するものである。

まず、PA140は、ホスト100が発行したライトコマンドを特定ポート141で受信する（1501）。PA140は、受信したコマンドの解析を行い、要求データに対応する上位論理デバイス番号を割り出し、上位論理デバイス管理情報から対応するSA番号、下位論理デバイス番号を算出し（1502）、SA番号に対応するSA150にライト要求を転送する（1503）。SA150は転送制御部151よりライト要求を受け取り、要求データがディスクキャッシュ154上にあるかどうかを、キャッシュ管理情報203を参照することにより判定する（1504）。

20

【0068】

要求データがディスクキャッシュ154上にあった場合（キャッシュヒット）、SA150は要求元のPA140に対してライト準備完了を通知する（1507）。SA150よりライト準備完了の通知を受信したPA140は、ポート141を介してホスト100にライト準備完了を通知する（1508、1509）。その後、PA140はポート121を介してホスト100からデータを受信し、そのデータを該当するSA150へ送信する（1510、1511）。PA140からデータを受け取ったSA150は対応するディスクキャッシュ154領域にデータを格納し、PA140に対して完了報告を送信する（1512、1513、1514）。SA150より完了報告を受信したPA140はホスト100に対して、完了報告を送信する（1515、1516）。

30

【0069】

一方、要求データがディスクキャッシュ154上に無かった場合（キャッシュミス）、SA150はキャッシュ管理情報203を更新し、ディスクキャッシュ154上に要求データを格納する領域を確保する（1506）。以降の動作フローは、キャッシュヒットした場合と同様である（1507から1516）。以上のようにして、ホスト100からのライト要求に対して、ホスト100からのライトデータをディスクキャッシュ154へ格納する。

【0070】

図16は、ライトアфта処理257の処理フロー図である。この処理は、SA150におけるライトコマンド処理262の結果、ディスクキャッシュ154に格納されたライトデータをディスク装置157もしくは外部ストレージ180へ書き出す処理である。

40

ディスクキャッシュ154上に保持されたライトデータは、キャッシュ管理情報203で管理される。通常、ライトデータやディスクからリードしたデータはなるべく古いものからディスクキャッシュ154より追い出されるように、キューなどで管理されている。その様な従来の方法によって、管理されたデータの中からライトアфта対象データを決定する（1601）。対象データが、物理デバイスへのライトデータの場合は、ディスク装置157の該当する領域にデータを書き込んだ後に、対象データのディスクキャッシュ15

50

4領域を解放する(1603、1604)。対象データが外部デバイスへのライトデータの場合、SA150はPA140を介して外部ストレージ180に対してデータを書き込み(1605から1610)、対象データのディスクキャッシュ154領域を解放する(1604)。

【0071】

第2の実施形態：

次に図17乃至図19を参照して、第2の実施形態を説明する。

この例は、各SA150のダーティ属性を持つキャッシュ量の変動に応じて、論理デバイス定義時に外部デバイスに割り当てられたSA150を、同デバイスへのホスト100からの入出力要求を受け付けつつ、別のSA150に変更するものである。第2の実施形態では、第1の実施形態と実質的に同様のハードウェアおよびソフトウェア構成を前提としているので、以下、第1の実施形態との差異について説明する。

【0072】

図17に、第2の実施形態における上位論理デバイス管理情報204の構成を示す。図6に示した第1の実施形態のものに比べて、図17の例では、デバイスアクセスモード607、切替進捗ポインタ608および変更先SA番号/下位論理デバイス番号609が追加される。

デバイスアクセスモード607には、上位論理デバイスへの入出力要求の処理形態を示す「正常」もしくは「割当変更中」の値をとる。特定SA150の下位論理デバイスに割り当てられ、通常の入出力処理が行われる上位論理デバイスには「正常」の値が設定され、実体が外部デバイスの上位論理デバイスで、かつ外部デバイスが特定SAの下位論理デバイスから別のSA150の下位論理デバイスへ割り当て変更中である上位論理デバイスに対しては「割当変更中」が設定される。

【0073】

切替進捗ポインタ608は、上位論理デバイスが「割当変更中」に使用され、上位論理デバイスのうち、外部デバイスの割当変更処理が未完了な部分の先頭アドレスを示す情報である。切替進捗ポインタ608は、後述するSA150のキャッシュミス化処理の進捗に応じて更新される。変更先SA番号/下位論理デバイス番号609は上位論理デバイスが「割当変更中」に使用され、外部デバイスの割り当て変更先のSA番号および下位論理デバイス番号を保持する。

【0074】

図18は、第2の実施形態における下位論理デバイス管理情報201の構成を示す。図8に示した第1の実施形態のものに比べて、図18の例では、切替進捗ポインタ806が追加される。その他のエントリの内容は上記上位論理デバイス管理情報201と同じである。SA150によるキャッシュミス化処理の進捗に応じて更新するが、2重に管理することで、MA160で保持する上位論理デバイス管理情報203への更新頻度を抑えることも可能となる。

上位論理デバイス管理情報203の更新時には、全PA140のメモリ143へ更新情報の取り込みが必要となるため、性能を考慮すれば更新頻度の抑制は有効である。ただし、その場合、PA140でのリードコマンド処理261、ライトコマンド処理262で古い切替進捗ポインタ608を参照することになるため、既にキャッシュミス化を終えた領域についてもPA140からSA150へ入出力要求が転送されることになる。

よって、SA150のコマンド処理254で、キャッシュミス化処理を終えた領域へのライト要求についてはライトアフタで処理せず、即物理/外部デバイスへ書き込むなどこの処理で利用したディスクキャッシュ154を処理後速やかに解放する論理が必要となる。

【0075】

図19は外部デバイス割り当て変更処理の処理フローである。

この処理は、管理端末190又は管理サーバ110からの外部デバイス割り当て変更指示を受け、外部デバイスに割り当てられたSA150、即ち外部デバイスの入出力要求を処理するSA150を変更する処理である。

10

20

30

40

50

【 0 0 7 6 】

ST190は、SAキャッシュ使用量情報222および外部デバイスキャッシュ使用量情報224を参照し、全SA150に対するダーティ量情報302および全外部デバイスに対するダーティ量情報402を得て、管理端末190又は管理サーバ110に対してそれらを提示し、外部デバイス割り当て変更の指示を受け付ける(1901、1902)。管理端末190又は管理サーバ110はSA150のキャッシュ使用量が均等化されるように、変更対象とする外部デバイスおよび変更先のSA150を指定することができる。ここで、指定しないことも可能である。

変更対象とする外部デバイスおよび変更先のSA150を指定された場合は(1903)、外部デバイスの割り当て変更によるキャッシュ使用量の均等化の効果を示すために、管理端末190又は管理サーバ110に対して割り当て変更後の予測CAキャッシュ使用量情報を提示する(1906)。予測SAキャッシュ使用量情報は、変更対象とする外部デバイスの現在のダーティ量を、変更元のSA150のダーティ量から減じ、変更先のSA150のダーティ量に加えることにより算出する。その後、ST190は、割当変更時間情報502により指定された時間帯に、MA160に対して割り当て変更の指示を行う(1907)。この指示には、変更対象の外部デバイス番号、変更先SA番号の情報が含まれる。割り当て変更の指示を受け付けたMA160は外部デバイスのデバイスアクセスモードを「正常」から「割当変更中」に変更する(1908)。MA160は変更先SA150に対して割当変更指示を発行する(1909)。割当変更指示を受け取った変更先SA150では、指示された外部デバイスに対して下位論理デバイスを登録し、MA160に対して完了報告を送信する(1910)。

【 0 0 7 7 】

変更先SA150から完了報告を受け取ったMA160は、変更元SA150に対して割当変更指示を送信する(1911)。割当変更指示を受け取った変更元SA150は、対応する下位論理デバイスの全データについてディスクキャッシュ154を検索してミス化、即ち、物理/外部デバイスへ未更新のデータは対応デバイスへの書き出し後、更新済みのデータについては、即そのデータに割り当てられたキャッシュ領域を解放する。検索はキャッシュ管理情報203を用いて、下位論理デバイスの先頭から順次行い、検索およびミス化処理が終わった部分については、上位/下位論理デバイス管理情報の切替進捗ポインタを進めることで進捗を管理する。全領域についてミス化処理が完了したら、その旨をMA160へ報告し、MA160でデバイスアクセスモード607を「正常」へ変更して外部デバイス割当変更処理を完了する(1912、1913)。

【 0 0 7 8 】

変更対象とする外部デバイスおよび変更先のSA150が指定されなかった場合は(1903)、ST190が各SA150のキャッシュ使用量が均等化するように、変更対象とする外部デバイスおよび変更先のSA150を選択する(1904)。変更先のSA150としては最もダーティ量の少ないSA150とする。

キャッシュ使用量の均等化の指標としては、例えば、各SA150のキャッシュ使用量の標準偏差を用いることができる。外部デバイス管理情報を用いて、ストレージ130に割り当てられた各外部デバイスを変更対象候補として1つずつ選ぶ。選択された外部デバイスのダーティ量および各SA150のダーティ量が分かっているので、外部デバイスを割り当てるSA150を変更した場合の各SA150の予測ダーティ量を算出でき、各SA150の予測ダーティ量を用いて予測標準偏差を求めることができる。変更対象として選択する外部デバイスは、予測標準偏差が最小となる外部デバイスとする。ただし、均等化されない場合、即ち、予測標準偏差が現在の標準偏差以上の値になる場合は、外部デバイスの割り当て変更は実施しない(1905)。以降の処理は、変更対象とする外部デバイスおよび変更先のSA150を指定された場合と同様である。

【 0 0 7 9 】

第2の実施形態においては、第1の実施形態に比べてリードコマンド処理261、ライトコマンド処理262が一部変わる。

即ち、図14のリードコマンド処理261の1402において、PA140は受信したコマンドの解析を行い、要求データに対応する上位論理デバイス番号を割り出した後、デバイスアクセスモードをチェックする。デバイスアクセスモードが「正常」の場合には、上位論理デバイス管理情報から対応するSA番号、下位論理デバイス番号を算出する。その後、SA番号に対応するSA150にリード要求を転送する(1403)。

デバイスアクセスモードが「割当変更中」の場合には、切替進捗ポイントを参照し、その要求のアクセス領域が切替進捗ポイントの前後どちらにあるかを判定する。切替進捗ポイントよりも前、即ち変更先SA150へ切替済みの領域については、変更先SA150を要求先SA150として決定する。もし、そうでない場合、即ち変更先SA150への切替が完了していない領域に関しては、変更元SA150を要求先SA150として決定する。図15のステップ1502においても、同様の判定を行い要求先SA150を決定する。

10

【0080】

第3の実施形態：

次に図20乃至図23を参照して、第3の実施形態を説明する。

この例は、第1のストレージである外部ストレージ180aのデバイス(外部デバイス)に対して、ストレージ130が非同期リモートコピー機能を適用し、ストレージ130が有する非同期リモートコピー機能によって、外部ストレージ180aのデバイス内に格納されるデータを他の外部ストレージ180bのデバイスにコピーすることを前提とする。この際、本実施形態においては、サイドファイル属性を持つキャッシュのキャッシュ使用量に基づき、外部ストレージ180aのデバイスを自身の論理デバイスとして管理し、外部ストレージ180aのデバイスについてリモートコピー処理を実行するストレージ130内のストレージ制御部を選択する。

20

【0081】

第3の実施形態を説明する前に、まず、リモートコピー機能について説明する。

リモートコピーは、正サイトにあるストレージシステムのデバイス(正デバイス)のバックアップを副サイトにあるストレージシステムのデバイス(副デバイス)にリアルタイムで作成する機能である。正デバイスのコピーを遠隔地にある副デバイスに保持することにより、テロや災害等によるデータの損失を最小限に抑えることが可能である。リモートコピーの操作には、形成コピーと更新コピーの2種類がある。形成コピーはホスト100からのリード/ライトとは別に正デバイスと副デバイスのデバイスペアを同期させる操作である。形成コピーにより、デバイスペアが同期された後の正デバイスに対するライトは、更新コピーにより副デバイスに対しても適用され、同期状態が維持される。リモートコピーは、更新コピーの方式により、同期リモートコピーと非同期リモートコピーに大別される。

30

【0082】

非同期リモートコピーは、特に、ホスト100からライトデータを受信しキャッシュメモリに書き込んだ時点で、ホスト100に対して完了報告を送信する。非同期リモートコピー対象のライトデータのうち、正ストレージシステムのキャッシュ上に保持され、副ストレージシステムに未送信のデータはサイドファイルと呼ばれる。サイドファイルは副ストレージシステムに対して送信されるまで、キャッシュメモリ上に保持し続ける必要がある。しかし、サイドファイルの副ストレージシステムへの送信速度とホスト100からのライトの頻度によっては、サイドファイルがキャッシュメモリ上に大量に滞留し、新たなリード/ライトに対してキャッシュメモリを割り当てることができなくなる可能性がある。

40

【0083】

サイドファイルは、正ストレージシステムのSA150の制御メモリ155上のキャッシュ管理情報203により、サイドファイル属性を持つセグメントとして管理される。サイドファイル量はサイドファイル属性を持つセグメントの個数により算出される。キャッシュ管理情報203としては、ライトデータがディスクキャッシュ154に書き込まれた

50

時間を示すタイムスタンプ情報、ライトデータに対応する上位論理デバイス番号、ライトデータの上位論理デバイス内での位置情報、ライトデータのサイズのような制御情報が格納される。タイムスタンプ情報は、サイドファイルが副ストレージシステムに送信される際の付加情報としてライトデータと共に送信される。互いに依存関係のある複数正デバイスにより構成されるデバイスグループに対するライトを副デバイスに適用する際には、デバイスグループ内のデータの整合性を維持するために、ライトの順序を保障する必要がある。副ストレージシステムにおいては、タイムスタンプ情報に基づきライトの順序を維持しながら、複数の副デバイスに対する更新が行われる。

【 0 0 8 4 】

第三の実施形態においては、例えば第一の外部ストレージ 1 8 0 a 内の外部デバイスと対応付けられているストレージ 1 3 0 内の上位論理デバイスと、遠隔地にある第2の外部ストレージ 1 8 0 b 内の論理デバイスとの間で非同期リモートコピーを行う。即ち、ホスト 1 0 0 からストレージ 1 3 0 内の上位論理デバイスへの書き込み要求があった際には、上位論理デバイスを管理しているストレージ 1 3 0 内のストレージ制御部は、キャッシュにライトデータを書き込んだ後ホスト 1 0 0 に完了報告を返送し、その後ライトデータを第2の外部ストレージ 1 8 0 b に転送すると共に、前記上位論理デバイスと対応付けられている第1の外部ストレージ 1 8 0 a 内の外部デバイスにもライトデータ書き込む。

【 0 0 8 5 】

尚、本実施形態においては、リモートコピー先のストレージは第2の外部ストレージ 1 8 0 b であるものと仮定して説明するが、リモートコピー先のストレージは外部ストレージ 1 8 0 に限られるものではない。即ち、リモートコピー先のストレージは、ストレージ 1 3 0 のストレージ制御部によってストレージ 1 3 0 の上位論理デバイスとして管理される外部デバイスを有する外部ストレージであっても良いが、ストレージ 1 3 0 とネットワークを介して接続されるストレージであれば外部ストレージに限られることはなく、ホスト 1 0 0 からストレージ 1 3 0 を介することなくアクセスされ得るデバイスを有するストレージであっても良い。

第3の実施形態も、第1の実施形態と実質的に同様のハードウェアおよびソフトウェア構成を前提としているので、以下、第1の実施形態との差異について説明する。

【 0 0 8 6 】

図 2 0 は SA キャッシュ使用量情報 2 2 2 を示す図である。
この図 2 0 は、図 3 に比べて、ストレージ 1 3 0 内の SA 1 5 0 のサイドファイル量に関する情報を管理するサイドファイル量情報 3 0 3 が追加されている。サイドファイル量情報の値は、SA 1 5 0 において、集計時間情報 5 0 1 に含まれる時間帯に、一定の時間間隔で SA 1 5 0 に対応するサイドファイルセグメントカウンタを参照し、それらの平均値を求めることにより算出される。

【 0 0 8 7 】

図 2 1 は外部デバイスキャッシュ使用量情報 2 2 4 を示す図である。
この図 2 1 では、図 4 に比べて、ストレージ 1 3 0 内の外部デバイスの特定の時間帯におけるサイドファイル量を管理するサイドファイル量情報 4 0 3 が追加されている。サイドファイル量情報 4 0 3 の値は、SA 1 5 0 において、集計時間情報 5 0 1 に含まれる時間帯に、一定の時間間隔で当該外部デバイスに対応するサイドファイルセグメントカウンタを参照し、それらの平均値を求めることにより算出される。

【 0 0 8 8 】

第3の実施形態においては、論理デバイス定義処理 2 5 5 が部分的に変わる。
第1の実施形態においては、図 1 2 のステップ 1 2 0 6 において、最もダーティ量の少ない SA 1 5 0 を割り当て対象として選択したが、第3の実施形態においては、最もサイドファイル量の少ない SA 1 5 0 を割り当て対象とする。この処理により、外部デバイスを、ストレージ 1 3 0 内の SA 1 5 0 のキャッシュ使用量の均等化を考慮して、いずれかの SA 1 5 0 へ割り当てることが可能となる。

【 0 0 8 9 】

外部デバイス定義処理 2 5 3 および L U パス定義処理 2 5 2 は、第 1 の実施形態と同様であり、上記の論理デバイス定義処理 2 5 5 と合わせた 3 つの処理によりホスト 1 0 0 からのアクセスが可能となる。

【 0 0 9 0 】

次に、第 1 の外部ストレージ 1 8 0 a 内の外部デバイスを実体とする上位論理デバイスである正デバイスと、第 2 の外部ストレージ 1 8 0 b 内の論理デバイスである副デバイス間で非同期リモートコピーを適用する場合における、リードコマンド処理 / ライトコマンド処理について述べる。尚、リードコマンド処理 2 6 1 は第 1 の実施形態と同様なので、その説明を省略する。

【 0 0 9 1 】

次に、図 2 2 を参照して、ライトコマンド処理 2 6 2 について説明する。第 3 の実施形態におけるライトコマンド処理は、図 1 5 におけるライトコマンド処理 2 6 2 とほぼ同じであるが、ステップ 2 2 1 7 から 2 2 1 9 が追加されている。ステップ 2 2 1 7 においては、ライトデータに対応するデバイスが非同期リモートコピー機能を使用中かを判定し、使用中なら、ディスクキャッシュ 1 5 4 上にライトデータ用のサイドファイル領域を確保し、管理メモリ上のキャッシュ管理情報を更新する (2 2 1 8)。ステップ 2 2 1 9 においては、ライトデータをディスクキャッシュ 1 5 4 上のサイドファイル領域に格納し、キャッシュ管理情報にタイムスタンプ情報等の制御情報を格納する。

【 0 0 9 2 】

図 2 3 はサイドファイル送信処理を説明するための図である。この処理は、S A 1 5 0 におけるライトコマンド処理 2 6 2 の結果、ディスクキャッシュ 1 5 4 上のサイドファイル領域に格納されたライトデータを第 2 の外部ストレージ 1 8 0 b に書き出す処理である。

S A 1 5 0 は、管理メモリ 1 5 5 上のキャッシュ管理情報 2 0 3 から、サイドファイルが書き込まれた順番にサイドファイルを送信するように、サイドファイルを決定する (2 3 0 1)。S A 1 5 0 は該当する P A 1 4 0 にサイドファイルデータをタイムスタンプ情報等の制御情報と共に送信する (2 3 0 2)。サイドファイルデータと制御情報を受信した P A 1 4 0 は、該当する外部ストレージ 1 8 0 に対して、サイドファイルデータと制御情報を送信する (2 3 0 3)。その後、S A 1 5 0 は対象データのサイドファイル領域を解放する (2 3 0 4)。

【 0 0 9 3 】

第 4 の実施形態：

この例は、非同期リモートコピー機能が適用されている外部デバイスに対して割り当てられた S A 1 5 0 をその後の各 S A 1 5 0 のサイドファイル属性を持つキャッシュ量に応じて、そのデバイスへのホスト 1 0 0 からの入出力要求を受け付けつつ、変更するものである。

第 4 の実施形態も、第 1 の実施形態と実質的に同様のハードウェアおよびソフトウェア構成を前提としているので、以下、第 1 の実施形態との差異について説明する。

第 4 の実施形態においては、第 2 の実施形態で説明した、図 1 7 の上位論理デバイス管理情報 2 0 3 および図 1 8 の下位論理デバイス管理情報 2 0 1 を用いる。

【 0 0 9 4 】

図 2 4 は外部デバイス割当変更処理を説明する図である。

管理端末 1 9 0 又は管理サーバ 1 1 0 からの外部デバイス割当て変更指示を受け、外部デバイスに割り当てられた S A 1 5 0、即ち、当該外部デバイスの入出力要求を処理する S A 1 5 0 を変更する処理である。図 2 4 は第 2 の実施形態で説明した図 1 9 とほぼ同様であるが、ステップ 2 4 0 1、2 4 0 2、2 4 0 4、2 4 0 5、2 4 0 6 が変わっている。

ステップ 1 9 0 1 においてはダーティ量情報 3 0 2 およびダーティ量情報 4 0 2 を参照していたのに対して、ステップ 2 4 0 1 においては、S T 1 9 0 は S A キャッシュ使用量情

10

20

30

40

50

報 2 2 2 および外部デバイスキャッシュ使用量情報 2 2 4 を参照して、全 S A 1 5 0 に対するサイドファイル量情報 3 0 3 および全外部デバイスに対するサイドファイル量情報 4 0 3 を得る。この変更に伴い、ステップ 2 4 0 2、2 4 0 4、2 4 0 5、2 4 0 6 においては、ダーティ量情報をサイドファイル量情報に置き換えて、ステップ 1 9 0 2、1 9 0 4、1 9 0 5、1 9 0 6 と同様の処理を行う。

【 0 0 9 5 】

第 4 の実施形態におけるリードコマンド処理 2 6 1 は、第 2 の実施形態で説明したものと同様である。第 4 の実施形態におけるライトコマンド処理 2 6 2 としては、第 3 の実施形態で説明したライトコマンド処理 2 6 2 を一部変更したものをを用いる。変更箇所は図 2 2 のステップ 2 2 0 2 であり、変更内容は第 2 の実施形態で説明したリードコマンド処理 2 6 1 に対する変更と同様である。

10

【 0 0 9 6 】

第 5 の実施形態：

上記した実施形態では、ストレージ 1 3 0 は P A 1 4 0 と S A 1 5 0 と M A 1 6 0 が相互結合網により結合された構成を持つクラスタ構成ストレージを例として挙げた。しかし本発明は、上記したクラスタ構成ストレージに限定されず、他の構成からなるクラスタ構成ストレージに関しても適用される。

【 0 0 9 7 】

図 2 5 は、本発明の他の構成による計算機システム例を示す。

この例において、ストレージ 2 5 3 0 は、複数のストレージノード 2 5 5 0 とストレージノード 2 5 5 0 間でのデータ連携を行うための相互結合網 2 5 7 0 から構成されるクラスタ構成ストレージである。各ストレージノード 2 5 5 0 は、通常のストレージと同様に、1 または複数のディスク装置 2 5 5 7 と、ディスクキャッシュ 2 5 5 4 と、制御プロセッサ 2 5 5 2 と、メモリ 2 5 5 3 と、ポート 2 5 5 1 から構成される。

20

本構成例においては、前述したストレージ 2 5 3 0 が提供する高機能処理を担当する部位である S A 1 5 0 と入出力要求を振り分ける P A 1 4 0 が、共にストレージノード 2 5 5 0 として構成されている点が第 1 から第 4 の実施形態とは異なる。この例においても、複数のストレージノード 2 5 5 0 はそれぞれ相互結合網 2 5 7 0 により接続されているため、任意のストレージノード 2 5 5 0 に対して外部デバイスを割り当てることができる。なお、この例において、第 1 乃至第 4 の実施形態における構成管理部 1 6 0 に相当する処理はストレージ 2 5 3 0 内のいずれかのストレージノード 2 5 5 0 が行う。

30

【 0 0 9 8 】

管理サーバ 2 5 1 0 は、第 1 乃至第 4 の実施形態における管理端末 1 9 0 の役割も兼ねており、IP ネットワーク 2 5 7 5 により計算機システム内の各機器とインターフェース 2 5 1 6 を介して交信し、計算機システム内の各機器から構成情報、リソース利用率、性能監視情報などを収集する。またそれらの情報をディスプレイ 2 5 1 5 に表示してストレージ管理者に提示する。さらに入力装置から入力され受信した運用・保守に関する指示を各機器に送信する。また、第 1 乃至第 4 の実施形態における S A 1 5 0 と同様に、ストレージノード 2 5 5 0 はディスクキャッシュ 2 5 5 4 上のダーティ量やサイドファイル量に関する情報を採取する。管理サーバ 2 5 1 0 は各ストレージノード 2 5 5 0 から、ダーティ量やサイドファイル量に関する情報を収集し、メモリ 2 5 1 2 上に保持する。さらに、管理サーバ 2 5 1 0 は外部ストレージ 2 5 8 0 を含めた計算機システム全体の管理も行う。

40

【 0 0 9 9 】

ファイバチャネルスイッチ 2 5 2 0 は、ホスト 2 5 0 0 のポート 2 5 0 7 およびストレージ 2 5 3 0 のポート 2 5 5 1 に加え、外部ストレージのポート 2 5 8 1 にも接続されている。その他の機器は第 1 乃至第 4 の実施形態で説明したものと同等の役割を果たす。

【 0 1 0 0 】

次に、第 5 の実施形態と上記した第 1 乃至第 4 の実施形態との処理の類似点および相違点について説明する。この第 5 の実施形態では、ハードウェア構成が前述の第 1 乃至第 4

50

の実施形態とは異なるので、それによりリードライト処理も相違してくる。

具体的には、あるストレージノード2550がホスト2500からの入出力要求を受領した際、その入出力要求がそのストレージノード2550にて処理可能な入出力要求であった場合、そのストレージノード2550において処理し、他のストレージノード2550において処理すべき入出力要求であった場合には、相互結合網2570を介して他のストレージノード2550に入出力要求を転送する。さらに、第1乃至第4の実施形態におけるSA150と同様に、ストレージノード2550においては、外部ストレージ2580のデバイスはストレージ2530の論理デバイスとして定義され、ホスト2500からの入出力要求に対して、ストレージノード内部のディスク装置に対するアクセスか、外部ストレージに対するアクセスかを識別して、要求を振り分ける。

10

【0101】

この実施形態においては、第1の実施形態のように、外部デバイスの処理を行うストレージノード2550を選択する際に、ストレージノード2550のダーティ量を用いて、ストレージシステム全体としてのキャッシュ使用量の均等化を行うことができる。

具体的な処理内容は、ストレージノード2550がPA140、SA150の両方の処理を行う点を除いては、第1の実施形態で説明したものと同様である。ただし、以下の点で処理が異なる。

【0102】

第一に、予め、全てのストレージノード2550から外部ストレージ2580にアクセスできるように、ファイバチャネルスイッチ2520のゾーニング設定を変更しておくものとする。これは、どのストレージノードも外部ストレージへの経路としての役割を果たすことができることを意味し、第2から第4の実施形態の構成例においても同様である。

20

第二に、論理デバイス定義処理255において、SA150としての役割を果たすストレージノード2550が同時にPA140としての役割も果たすように設定を行う。これは第3の実施形態の構成例においても同様である。

【0103】

さらに、この例においては、第2の実施形態のように、各ストレージノード2550のダーティ属性を持つキャッシュ量の変動に応じて、論理デバイス定義時に外部デバイスに割り当てられたストレージノード2550を、同デバイスへのホスト2500からの入出力要求を受け付けつつ、別のストレージノード2550に変更することができる。

30

具体的な処理内容は、ストレージノード2550がPA140、SA150の両方の処理を行う点を除いては、第2の実施形態で説明したものと同様である。ただし、次の点で処理が異なる。即ち、外部デバイス割当変更処理256において、SA150としての役割を果たすストレージノード2550が同時にPA140としての役割も果たすように設定を行う。これは第4の実施形態の構成例においても同様である。

【0104】

また、第3の実施形態のように、外部デバイスを実体とする、ストレージノード2550内の論理デバイスに対して非同期リモートコピーを適用する場合においても、各ストレージノード2550におけるサイドファイル量に関する情報を用いて、非同期リモートコピーの処理を行うストレージノード2550を適切に選択することにより、各ストレージノード2550のサイドファイル量の均等化を行うことが可能である。具体的な処理内容は、ストレージノード2550がPA140、SA150の両方の処理を行う点を除いては、第3の実施形態で説明したものと同様である。

40

【0105】

さらに、本構成例においては、第4の実施形態のように、非同期リモートコピー機能が適用されている外部デバイスに対して割り当てられたストレージノード2550をその後の各ストレージノード2550のサイドファイル属性を持つキャッシュ量に応じて、同デバイスへのホスト2500からの入出力要求を受け付けつつ、変更することができる。具体的な処理内容は、ストレージノード2550がPA140、SA150の両方の処理を行う点を除いては、第4の実施形態で説明したものと同様である。

50

この他にも本発明の趣旨を逸脱しない範囲で種々変形して実施し得る。

【0106】

尚、上記第5の実施形態に関する態様を整理すると、例えば以下のように列挙できる。

(1) ホストとのインターフェースを有し、ホストから入出力要求に従って読み書きされるデータを格納するキャッシュメモリ及びディスク装置を有すると共に、第1のストレージシステムとのインターフェースを有して第1のストレージシステムに対して入出力要求を行う複数のストレージノードと、ストレージノードの間を互いに接続する相互結合網と、ストレージノードと通信を行う管理サーバとを有するストレージシステムにおいて、第1のストレージシステムのディスク装置及びストレージノードに保持するディスク装置をストレージシステムが持つディスク装置としてホストへ提示する手段と、ホストから受け付けた入出力要求のアクセス対象であるストレージシステムのディスク装置がストレージノードのディスク装置もしくは第1のストレージシステムのディスク装置である場合、入出力要求をストレージノードで処理する手段と、ストレージノードのキャッシュメモリ上において、ストレージノードのディスク装置と第1のストレージシステムのディスク装置に未反映なライトデータの合計量である第1のダーティ量情報を取得する手段と、管理サーバに第1のダーティ量情報を提示し、第1のストレージシステムのディスク装置の処理を行う該ストレージノードの指定を受け付ける手段と、を有するストレージシステム。

10

(2) 上記(1)のストレージシステムにおいて、第1のストレージシステムのディスク装置の処理を行う第2のストレージシステム内のストレージノードが管理サーバで指定されなかった場合、第1のストレージシステムのディスク装置の処理を行うストレージノードを、第1のダーティ量情報を用いて決定するストレージシステム。

20

(3) 上記(2)のストレージシステムにおいて、第1のストレージシステムのディスク装置の処理を行うストレージノードを、第1のストレージノードから第2のストレージノードに変更する切替処理手段を有するストレージシステム。

(4) 上記(3)のストレージシステムにおいて、ストレージノードのキャッシュメモリ上において第1のストレージノードのディスク装置のデータを検索し、第1のストレージシステムのディスク装置へ未反映なデータは第1のストレージシステムのディスク装置へ書き込んでキャッシュメモリ領域を解放し、反映済みのデータはキャッシュメモリ領域を解放するストレージシステム。

(5) 上記(4)のストレージシステムにおいて、第3のストレージノードが第1のストレージノードのディスク装置へのホストからの入出力要求に対して、第1のストレージノードのディスク装置で切替処理が終わった部分への入出力要求は第2のストレージノードへ、切替処理が終わっていない部分への要求は第1のストレージノードへ振り分けるストレージシステム。

30

(6) 上記(3)のストレージシステムにおいて、ストレージノードのキャッシュメモリ上において、第1のストレージシステムのディスク装置に未反映なライトデータの量である第2のダーティ量情報を取得する手段を有し、上記管理サーバへ第1のダーティ量情報および第2のダーティ量情報を提示し、上記切替処理の対象となる第1のストレージシステムのディスク装置および切替処理の切替先である第2のストレージノードの指定を受け付ける手段を有するストレージシステム。

40

(7) 上記(6)のストレージシステムにおいて、上記切替処理の対象となる第1のストレージシステムのディスク装置および切替処理の切替先である第2のストレージノードの指定が無かった場合、第1のダーティ量および第2のダーティ量を用いて、切替処理の対象となる第1のストレージシステムのディスク装置および切替処理の切替先である第2のストレージノードを決定するストレージシステム。

(8) 上記(1)のストレージシステムにおいて、第3のストレージシステムとのインターフェースを有し、第3のストレージシステムに対してコピー処理を行う第4のストレージノードを含み、第3のストレージシステムのディスク装置にストレージノードにおけるディスク装置の複製を作成し、このストレージノードにおけるディスク装置に対するライトデータを逐次そのストレージノードのキャッシュメモリに格納し、書き込みデータを第

50

3のストレージシステムのディスク装置に送信する機能を有するストレージシステムであって、前記ストレージノードのキャッシュメモリ上にあって、第3のストレージシステムに未送信なライトデータの合計量である第1のサイドファイル量情報を取得する手段を有し、管理サーバに第1のサイドファイル量情報を提示し、第1のストレージシステムのディスク装置の処理を行うストレージノードの指定をから受け付けるストレージシステム。

(9)上記(8)のストレージシステムにおいて、第1のストレージシステムのディスク装置の処理を行うストレージノードの指定が無かった場合、第1のストレージシステムのディスク装置の処理を行うストレージノードを、第1のサイドファイル量情報を用いて決定するストレージシステム。

(10)上記(8)のストレージシステムにおいて、第1のストレージシステムのディスク装置をホストに提示する処理を行うストレージノードを、第1のストレージノードから第2のストレージノードに変更するストレージシステム。

(11)上記(10)のストレージシステムにおいて、ストレージノードのキャッシュメモリ上にあって、第3のストレージシステムのディスク装置に未送信なデータ量である第2のサイドファイル量情報を取得する手段を有し、管理サーバに第1のサイドファイル量情報および前記第2のサイドファイル量情報を提示し、切替処理の対象となる第1のストレージシステムのディスク装置および切替処理の切替先である第2のストレージノードの指定を受け付けるストレージシステム。

(12)上記(11)のストレージシステムにおいて、切替処理の対象となる第1のストレージシステムのディスク装置および切替処理の切替先である第2のストレージノードの指定が無かった場合、第1のサイドファイル量情報および第2のサイドファイル量情報を用いて、切替処理の対象となる第1のストレージシステムのディスク装置および切替処理の切替先である第2のストレージノードを決定するストレージシステム。

【図面の簡単な説明】

【0107】

【図1】本発明の第1の実施形態による計算機システムのハードウェア構成を示す図である。

【図2】第1の実施形態におけるソフトウェアの構成を示す図である。

【図3】第1の実施形態におけるストレージ制御部のキャッシュ使用量情報の構成を示す図である。

【図4】第1の実施形態における外部デバイスキャッシュ使用量情報の構成を示す図である。

【図5】第1の実施形態における管理端末制御情報の構成を示す図である。

【図6】第1の実施形態における上位論理デバイス管理情報の構成を示す図である。

【図7】第1の実施形態におけるLUパス管理情報の構成を示す図である。

【図8】第1の実施形態における下位論理デバイス管理情報の構成を示す図である。

【図9】第1の実施形態における物理デバイス管理情報の構成を示す図である。

【図10】第1の実施形態における外部デバイス管理情報の構成を示す図である。

【図11】第1の実施形態における外部デバイス定義処理の処理フロー図である。

【図12】第1の実施形態における論理デバイス定義処理の処理フロー図である。

【図13】第1の実施形態におけるLUパス定義処理の処理フロー図である。

【図14】第1の実施形態におけるリードコマンド処理の処理フロー図である。

【図15】第1の実施形態におけるライトコマンド処理の処理フロー図である。

【図16】第1の実施形態におけるライトアフト処理の処理フロー図である。

【図17】第2の実施形態における上位論理デバイス管理情報の構成を示す図である。

【図18】第2の実施形態における上位論理デバイス管理情報の構成を示す図である。

【図19】第2の実施形態における外部デバイス割当変更処理の処理フロー図である。

【図20】第3実施形態におけるストレージ制御部キャッシュ使用量情報の構成を示す図である。

【図21】第3の実施形態における外部デバイスキャッシュ使用量情報の構成を示す図で

10

20

30

40

50

ある。

【図 2 2】第 3 の実施形態におけるライトコマンド処理の処理フロー図である。

【図 2 3】第 3 の実施形態におけるサイドファイル送信処理の処理フロー図である。

【図 2 4】第 4 の実施形態における外部デバイス割当変更処理の処理フロー図である。

【図 2 5】第 5 の実施形態における計算機システムの構成を示す図である。

【符号の説明】

【 0 1 0 8 】

- 1 0 0 ... ホスト、
- 1 2 0 ... ファイバチャネルスイッチ、
- 1 4 0 ... プロトコル変換部、
- 1 5 4 ... ディスクキャッシュ、
- 1 6 0 ... 構成管理部、
- 1 8 0 ... 外部ストレージ、
- 2 0 1 ... 下位論理デバイス管理情報、
- 2 0 3 ... キャッシュ管理情報、
- 2 0 5 ... 外部デバイス管理情報、
- 2 5 2 ... LUパス定義処理、
- 2 5 5 ... 論理デバイス定義処理、
- 2 5 7 ... ライトアフタ処理、
- 2 6 2 ... ライトコマンド処理、
- 2 5 0 0 ... ホスト、
- 2 5 3 0 ... ストレージ、
- 2 5 5 4 ... ディスクキャッシュ、
- 2 5 7 0 ... 相互結合網。

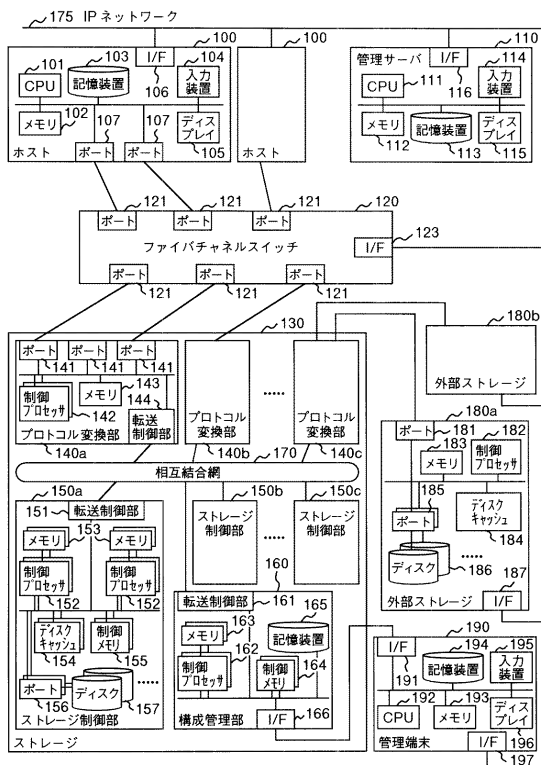
- 1 1 0 ... 管理サーバ、
- 1 3 0 ... ストレージ、
- 1 5 0 ... ストレージ制御部、
- 1 5 7 ... ディスク
- 1 7 0 ... 相互結合網、
- 1 9 0 ... 管理端末、
- 2 0 2 ... 物理デバイス管理情報、
- 2 0 4 ... 上位論理デバイス管理情報、
- 2 0 6 ... LUパス管理情報、
- 2 5 3 ... 外部デバイス定義処理、
- 2 5 6 ... 外部デバイス割当変更処理、
- 2 6 1 ... リードコマンド処理、
- 2 6 3 ... サイドファイル送信処理
- 2 5 1 0 ... 管理サーバ
- 2 5 5 0 ... ストレージノード
- 2 5 8 0 ... 外部ストレージ

10

20

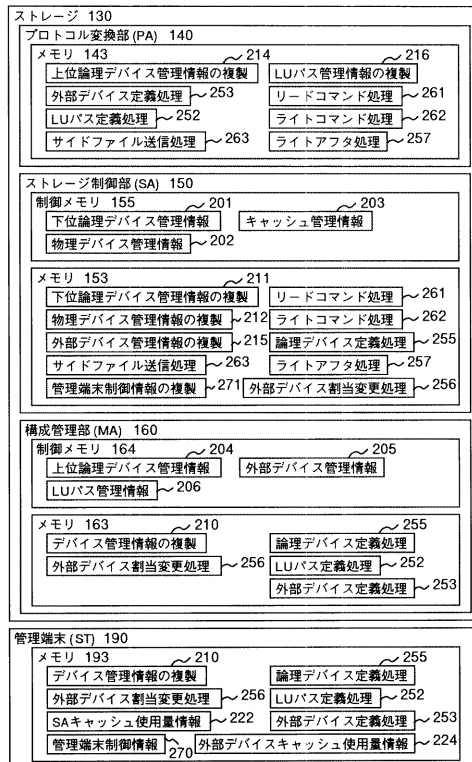
【図 1】

図 1



【図 2】

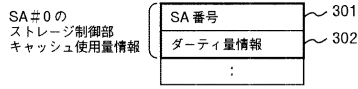
図 2



【図 3】

図 3

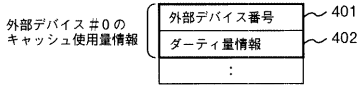
ストレージ制御部キャッシュ使用量情報 222



【図 4】

図 4

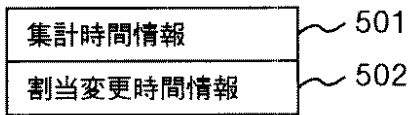
外部デバイスキャッシュ使用量情報 224



【図 5】

図 5

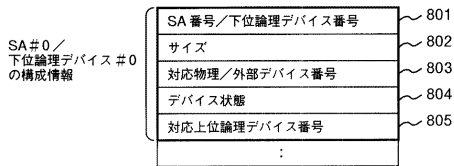
管理端末制御情報 270



【図 8】

図 8

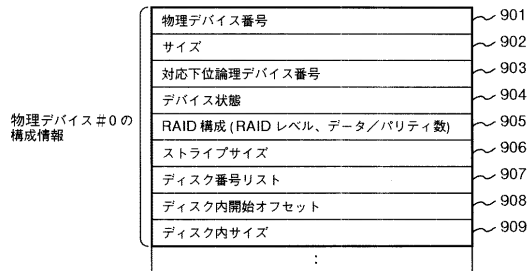
下位論理デバイス管理情報 201



【図 9】

図 9

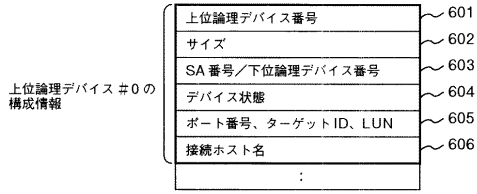
物理デバイス管理情報 202



【図 6】

図 6

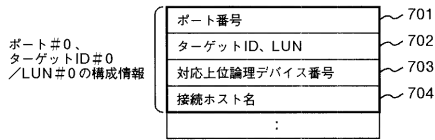
上位論理デバイス管理情報 204



【図 7】

図 7

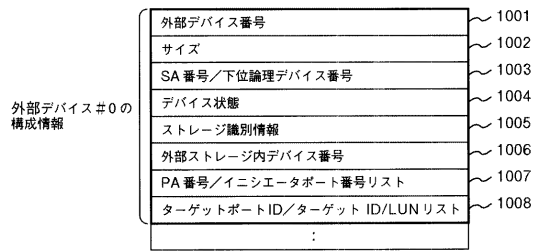
LU バス管理情報 206



【図 10】

図 10

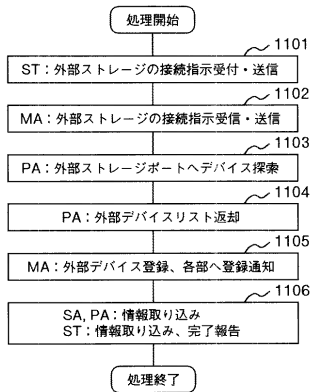
外部デバイス管理情報 205



【図 1 1】

図 1 1

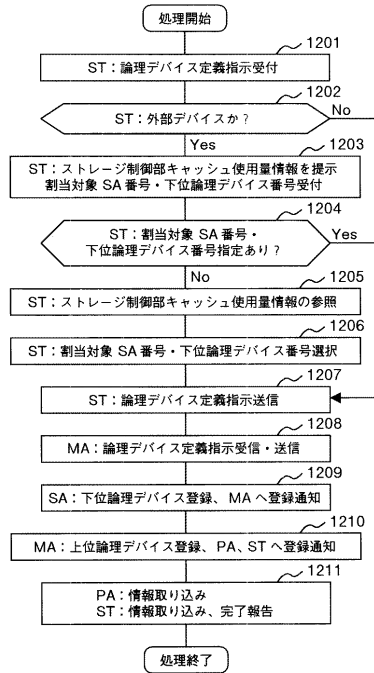
外部デバイス定義処理 253



【図 1 2】

図 1 2

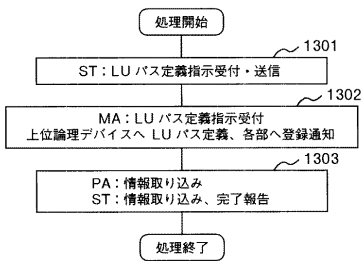
論理デバイス定義処理 255



【図 1 3】

図 1 3

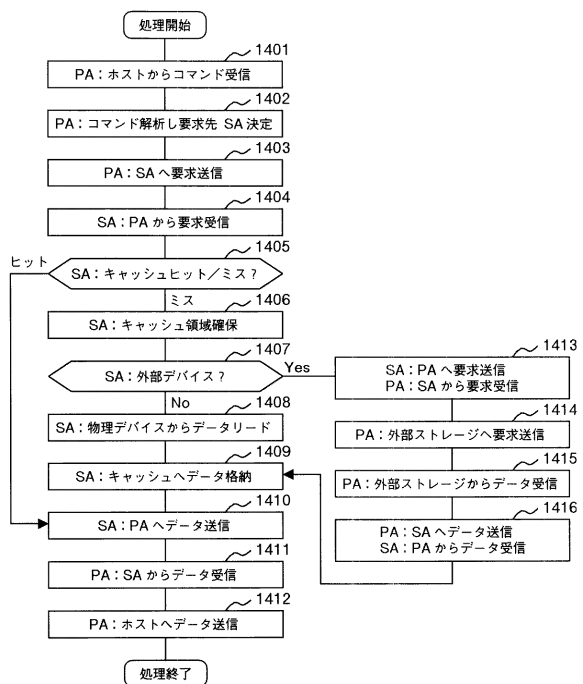
LU パス定義処理 252



【図 1 4】

図 1 4

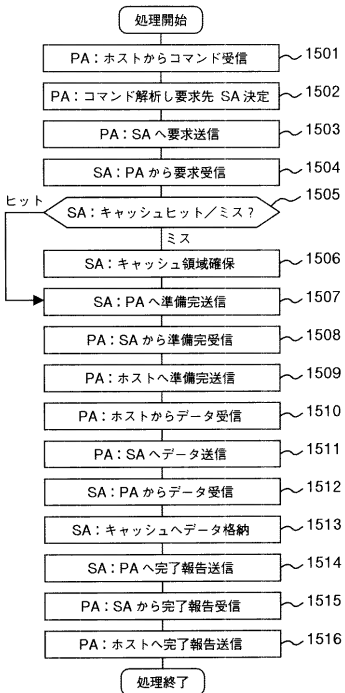
リードコマンド処理 261



【図 15】

図 15

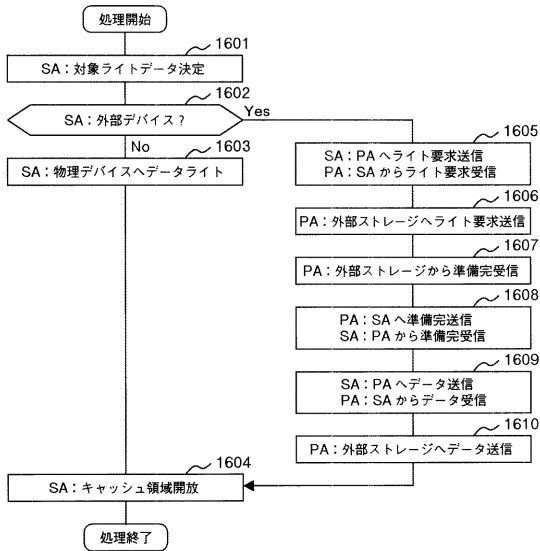
ライトコマンド処理 262



【図 16】

図 16

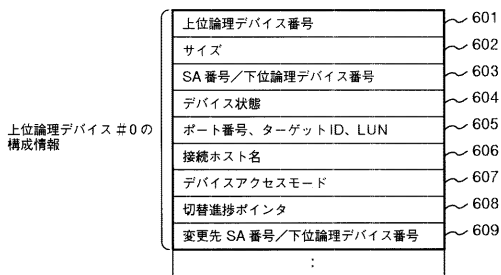
ライトアフタ処理 257



【図 17】

図 17

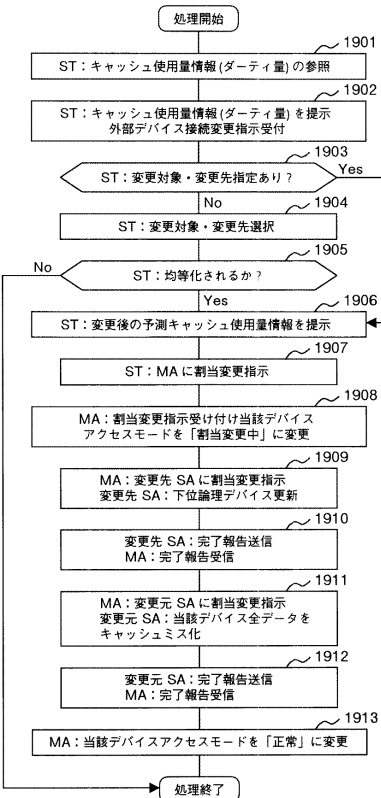
上位論理デバイス管理情報 204



【図 19】

図 19

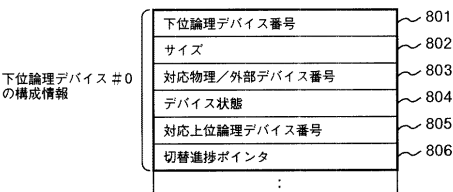
外部デバイス割当変更処理 256



【図 18】

図 18

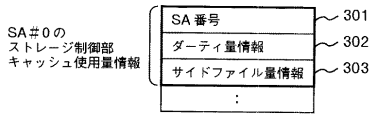
下位論理デバイス管理情報 201



【図 20】

図 20

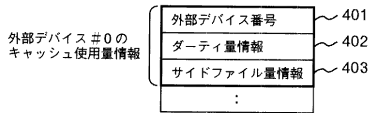
ストレージ制御部キャッシュ使用量情報 222



【図 21】

図 21

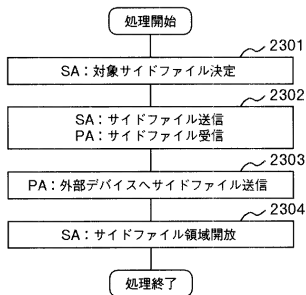
外部デバイスキャッシュ使用量情報 224



【図 23】

図 23

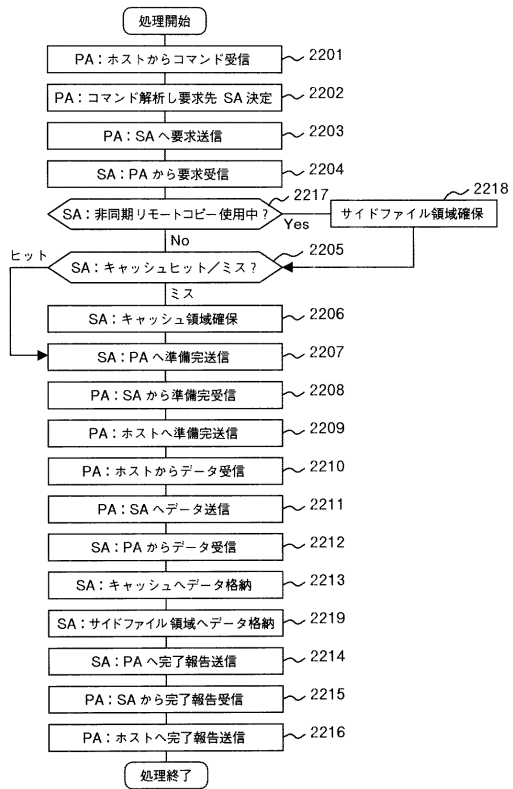
サイドファイル送信処理 263



【図 22】

図 22

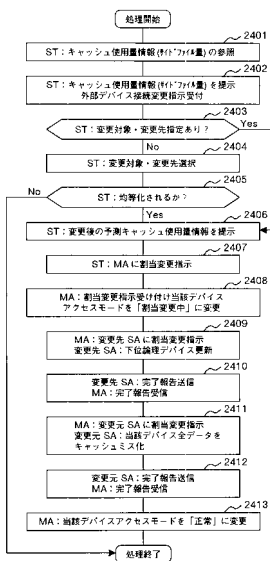
ライトコマンド処理 262



【図 24】

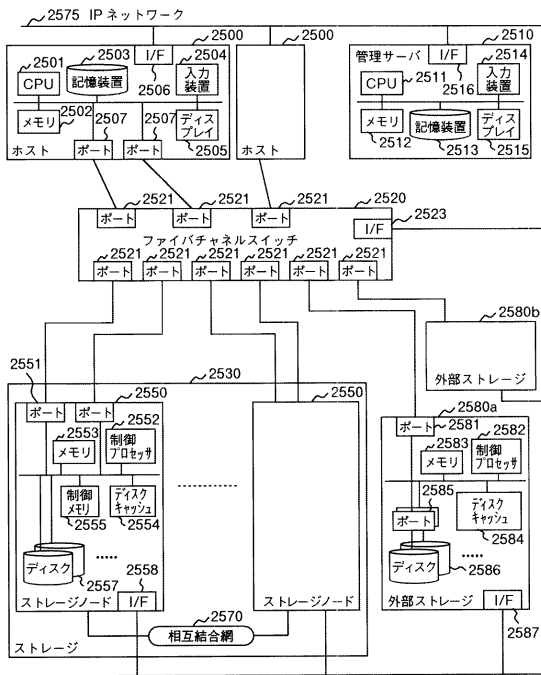
図 24

外部デバイス割当変更処理 256



【図 25】

図 25



フロントページの続き

(51)Int.Cl. F I
G 0 6 F 13/10 3 4 0 A

(72)発明者 山本 康友
神奈川県川崎市麻生区王禅寺 1 0 9 9 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 藤本 和久
神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所 システム開発研究所内

審査官 清木 泰

(56)参考文献 特開 2 0 0 5 - 0 5 5 9 7 0 (J P , A)
特開 2 0 0 3 - 1 3 1 9 4 4 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)
G 0 6 F 1 2 / 0 8 - 1 2 / 1 2
G 0 6 F 3 / 0 6 - 3 / 0 8
G 0 6 F 1 3 / 1 0 - 1 3 / 1 4
G 0 6 F 1 2 / 0 0 - 1 2 / 0 0 , 5 4 9