

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4331742号
(P4331742)

(45) 発行日 平成21年9月16日(2009.9.16)

(24) 登録日 平成21年6月26日(2009.6.26)

(51) Int. Cl. F I
G06F 13/14 (2006.01) G O 6 F 13/14 3 3 O B
G06F 3/06 (2006.01) G O 6 F 3/06 3 O 1 Z

請求項の数 17 (全 33 頁)

(21) 出願番号	特願2006-289971 (P2006-289971)	(73) 特許権者	000005108
(22) 出願日	平成18年10月25日(2006.10.25)		株式会社日立製作所
(65) 公開番号	特開2008-108050 (P2008-108050A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成20年5月8日(2008.5.8)	(74) 代理人	100075513
審査請求日	平成20年9月3日(2008.9.3)		弁理士 後藤 政喜
		(74) 代理人	100114236
			弁理士 藤井 正弘
		(74) 代理人	100120260
			弁理士 飯田 雅昭
		(72) 発明者	永井 崇之
			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】 1/Oの割り振り比率に基づいて性能を管理する計算機システム、計算機及び方法

(57) 【特許請求の範囲】

【請求項1】

ホスト計算機と、前記ホスト計算機に接続される複数の仮想ストレージ装置と、前記複数の仮想ストレージ装置に接続されるストレージ装置と、前記ホスト計算機、前記仮想ストレージ装置及び前記ストレージ装置に接続される管理計算機と、を備える計算機システムであって、

前記ストレージ装置は、

第1ネットワークを介して前記仮想ストレージ装置に接続される第1インターフェースと、第2ネットワークを介して前記管理計算機に接続される第2インターフェースと、前記第1インターフェース及び前記第2インターフェースに接続される第1プロセッサと、前記第1プロセッサに接続される第1メモリと、前記ホスト計算機に論理ボリュームとして提供される記憶領域と、を備え、

前記各仮想ストレージ装置は、

前記第1ネットワークを介して前記ホスト計算機及び前記ストレージ装置に接続される第3インターフェースと、前記第2ネットワークを介して前記管理計算機に接続される第4インターフェースと、前記第3インターフェース及び前記第4インターフェースに接続される第2プロセッサと、前記第2プロセッサに接続される第2メモリと、を備え、

前記論理ボリュームに対応付けられる仮想記憶領域を前記ホスト計算機に提供し、

前記ホスト計算機は、

前記第1ネットワークを介して前記仮想ストレージ装置に接続される第5インターフェ

ースと、前記第2ネットワークを介して前記管理計算機に接続される第6インターフェースと、前記第5インターフェース及び前記第6インターフェースに接続される第3プロセッサと、前記第3プロセッサに接続される第3メモリと、を備え、

一つの前記論理ボリュームに書き込まれるべきデータを含む複数のデータI/Oを、定められた比率で、前記一つの論理ボリュームと対応付けられた、前記各仮想ストレージ装置が提供する各々の前記仮想記憶領域に分散し、

前記管理計算機は、

前記第2ネットワークを介して前記ホスト計算機、前記仮想ストレージ装置及び前記ストレージ装置に接続される第7インターフェースと、前記第7インターフェースに接続される第4プロセッサと、前記第4プロセッサに接続される第4メモリと、を備え、

前記各仮想記憶領域に対するデータI/Oの量が、前記各仮想記憶領域に設定される閾値よりも大きくなった場合に、その旨を通知する警告を発するものであり、

前記各仮想記憶領域のうち第一仮想記憶領域についての閾値が入力された後に、前記第一仮想記憶領域に対応付けられている前記一つの論理ボリュームと対応付けられている、前記各仮想記憶領域のうち他の仮想記憶領域についての閾値を、前記第一仮想記憶領域についての閾値と前記定められた比率に基づいて算出することを特徴とする計算機システム。

【請求項2】

前記管理計算機は、前記各仮想記憶領域のうち第二仮想記憶領域を経由するI/O量が、前記第一仮想記憶領域を経由するI/O量より多いことが、前記定められた比率に基づいて予測される場合、前記第二仮想記憶領域に設定される前記閾値が前記第一仮想記憶領域に設定される前記閾値より大きくなるように、前記閾値を算出することを特徴とする請求項1に記載の計算機システム。

【請求項3】

前記管理計算機は、前記論理ボリュームと前記仮想記憶領域との対応付けに基づいて、前記比率を定め、前記定められた比率を前記ホスト計算機に送信することを特徴とする請求項1に記載の計算機システム。

【請求項4】

前記管理計算機は、前記各仮想記憶領域のうち第三仮想記憶領域の数が、前記各仮想記憶領域のうち第四仮想記憶領域の数より多い場合、第三仮想記憶領域を経由して実行されるI/O量が第四仮想記憶領域を経由して実行されるI/O量より少なくなるように、前記比率を定めることを特徴とする請求項3に記載の計算機システム。

【請求項5】

前記管理計算機は、前記論理ボリュームと前記仮想記憶領域との対応付けが変更されると、前記変更された対応付けに基づいて、前記比率を定めることを特徴とする請求項3に記載の計算機システム。

【請求項6】

前記管理計算機は、前記変更された対応付けに基づいて前記比率を定めると、前記定められた比率に基づいて、前記閾値を算出することを特徴とする請求項5に記載の計算機システム。

【請求項7】

ホスト計算機、前記ホスト計算機に接続される複数の仮想ストレージ装置、及び、前記複数の仮想ストレージ装置に接続されるストレージ装置に接続される管理計算機であって、

前記ストレージ装置は、前記ホスト計算機に論理ボリュームとして提供される記憶領域を含み、

前記各仮想ストレージ装置は、前記論理ボリュームに対応付けられる仮想記憶領域を前記ホスト計算機に提供し、

前記ホスト計算機は、一つの前記論理ボリュームに書き込まれるべきデータを含む複数のデータI/Oを、定められた比率で、前記一つの論理ボリュームと対応付けられた、前

10

20

30

40

50

記各仮想ストレージ装置が提供する各々の前記仮想記憶領域に分散し、

前記管理計算機は、ネットワークを介して前記ホスト計算機、前記仮想ストレージ装置及び前記ストレージ装置に接続されるインターフェースと、前記インターフェースに接続されるプロセッサと、前記プロセッサに接続されるメモリと、を備え、

前記プロセッサは、前記各仮想記憶領域に対するデータI/Oの量が、前記各仮想記憶領域に設定される閾値よりも大きくなった場合に、その旨を通知する警告を発するものであり、前記各仮想記憶領域のうち第一仮想記憶領域についての閾値が入力された後に、前記第一仮想記憶領域に対応付けられている前記一つの論理ボリュームと対応付けられている、前記各仮想記憶領域のうち他の仮想記憶領域についての閾値を、前記第一仮想記憶領域についての閾値と前記定められた比率に基づいて算出することを特徴とする管理計算機

10

【請求項 8】

前記プロセッサは、前記各仮想記憶領域のうち第二仮想記憶領域を経由するI/O量が、前記第一仮想記憶領域を経由するI/O量より多いことが、前記定められた比率に基づいて予測される場合、前記第二仮想記憶領域に設定される前記閾値が前記第一仮想記憶領域に設定される前記閾値より大きくなるように、前記閾値を算出することを特徴とする請求項 7 に記載の管理計算機。

【請求項 9】

前記プロセッサは、前記論理ボリュームと前記仮想記憶領域との対応付けに基づいて、前記比率を定め、前記定められた比率を前記ホスト計算機に送信することを特徴とする請求項 7 に記載の管理計算機。

20

【請求項 10】

前記プロセッサは、前記各仮想記憶領域のうち第三仮想記憶領域の数が、前記各仮想記憶領域のうち第四仮想記憶領域の数より多い場合、第三仮想記憶領域を経由して実行されるI/O量が第四仮想記憶領域を経由して実行されるI/O量より少なくなるように、前記比率を定めることを特徴とする請求項 9 に記載の管理計算機。

【請求項 11】

前記プロセッサは、前記論理ボリュームと前記仮想記憶領域との対応付けが変更されると、前記変更された対応付けに基づいて、前記比率を定めることを特徴とする請求項 9 に記載の管理計算機。

30

【請求項 12】

前記プロセッサは、前記変更された対応付けに基づいて前記比率を定めると、前記定められた比率に基づいて、前記閾値を算出することを特徴とする請求項 11 に記載の管理計算機。

【請求項 13】

ホスト計算機と、前記ホスト計算機に接続される複数の仮想ストレージ装置と、前記複数の仮想ストレージ装置に接続されるストレージ装置と、前記ホスト計算機、前記仮想ストレージ装置及び前記ストレージ装置に接続される管理計算機と、を備える計算機システムの制御方法であって、

前記ストレージ装置は、

40

第 1 ネットワークを介して前記仮想ストレージ装置に接続される第 1 インターフェースと、第 2 ネットワークを介して前記管理計算機に接続される第 2 インターフェースと、前記第 1 インターフェース及び前記第 2 インターフェースに接続される第 1 プロセッサと、前記第 1 プロセッサに接続される第 1 メモリと、前記ホスト計算機に論理ボリュームとして提供される記憶領域と、を備え、

前記ホスト計算機に論理ボリュームとして提供される記憶領域を含み、

前記各仮想ストレージ装置は、

前記第 1 ネットワークを介して前記ホスト計算機及び前記ストレージ装置に接続される第 3 インターフェースと、前記第 2 ネットワークを介して前記管理計算機に接続される第 4 インターフェースと、前記第 3 インターフェース及び前記第 4 インターフェースに接続

50

される第2プロセッサと、前記第2プロセッサに接続される第2メモリと、を備え、
 前記論理ボリュームに対応付けられる仮想記憶領域を前記ホスト計算機に提供し、
 前記ホスト計算機は、
 前記第1ネットワークを介して前記仮想ストレージ装置に接続される第5インターフェースと、前記第2ネットワークを介して前記管理計算機に接続される第6インターフェースと、前記第5インターフェース及び前記第6インターフェースに接続される第3プロセッサと、前記第3プロセッサに接続される第3メモリと、を備え、
 一つの前記論理ボリュームに書き込まれるべきデータを含む複数のデータI/Oを、定められた比率で、前記一つの論理ボリュームと対応付けられた、前記各仮想ストレージ装置が提供する各々の前記仮想記憶領域に分散し、
 前記管理計算機は、
 前記第2ネットワークを介して前記ホスト計算機、前記仮想ストレージ装置及び前記ストレージ装置に接続される第7インターフェースと、前記第7インターフェースに接続される第4プロセッサと、前記第4プロセッサに接続される第4メモリと、を備え、
 前記方法は、
 前記各仮想記憶領域に対するデータI/Oの量が、前記各仮想記憶領域に設定される閾値よりも大きくなった場合に、その旨を通知する警告を発するものであり、
 前記各仮想記憶領域のうち第一仮想記憶領域についての閾値が入力された後に、前記第一仮想記憶領域に対応付けられている前記一つの論理ボリュームと対応付けられている、
 前記各仮想記憶領域のうち他の仮想記憶領域についての閾値を、前記第一仮想記憶領域についての閾値と前記定められた比率に基づいて算出することを特徴とする方法。

10

20

【請求項14】

前記方法は、前記論理ボリュームと前記仮想記憶領域との対応付けに基づいて、前記比率を定め、前記定められた比率を前記ホスト計算機に送信することを特徴とする請求項13に記載の方法。

【請求項15】

前記方法は、前記各仮想記憶領域のうち第三仮想記憶領域の数が、前記各仮想記憶領域のうち第四仮想記憶領域の数より多い場合、前記第三仮想記憶領域を経由して実行されるI/O量が前記第四仮想記憶領域を経由して実行されるI/O量より少なくなるように、前記比率を定めることを特徴とする請求項14に記載の方法。

30

【請求項16】

前記方法は、前記論理ボリュームと前記仮想記憶領域との対応付けが変更されると、前記変更された対応付けに基づいて、前記比率を定めることを特徴とする請求項14に記載の方法。

【請求項17】

前記方法は、前記変更された対応付けに基づいて前記比率を定めると、前記定められた比率に基づいて、前記閾値を算出することを特徴とする請求項16に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本願明細書で開示される技術は、計算機システムに用いられるホストコンピュータとストレージシステムを管理するソフトウェアに関し、特に、管理ソフトウェアを用いてホストコンピュータ及びストレージ装置の性能を管理する方法に関する。

40

【背景技術】

【0002】

一般にコンピュータシステムは、業務を実行するホストコンピュータと、ホストコンピュータの指示に従ってデータを読み書きするストレージ装置とによって構成される。ストレージ装置は、データの格納及び読み出しを行う複数の磁気ディスクを備える。ストレージ装置は、ホストコンピュータに対して記憶領域を論理ボリュームという形で提供する。そして、ホストコンピュータ及びストレージ装置の構成及び性能を管理するための管理ソ

50

ソフトウェアが存在するのが一般的である。管理ソフトウェアは、ホストコンピュータ及びストレージ装置から構成情報及び性能情報を定期的を取得してプログラム内部に保持する。そして、管理ソフトウェアは、コンピュータシステムを管理する管理者からの命令に応じて、構成情報及び性能情報を表示する。

【0003】

一方、例えば特許文献1に示すように、ストレージ装置は、ホストコンピュータからのデータの読み書きを受け付ける代わりに、他のストレージ装置からのデータの読み書きを受け付けることができる。例えば、ホストコンピュータは、ストレージ装置Aが提供する仮想的な論理ボリュームXに対してデータの読み書きを実行する。ストレージ装置Aは、論理ボリュームXに対する読み書き指示を受けると、ストレージ装置B上の論理ボリュームYに対してデータの読み書きを実行する。結果として、ホストコンピュータは仮想的な論理ボリュームXに対して読み書きを実行していると認識しているが、実際は論理ボリュームYに対して読み書きを実行していることとなる。このような機能を「ストレージ仮想化機能」と呼び、論理ボリュームXのような仮想的な論理ボリュームを「仮想ボリューム」、論理ボリュームYのような仮想でない論理ボリュームを「実ボリューム」と呼ぶ。

10

【特許文献1】特開2005-11277号公報

【発明の開示】

【発明が解決しようとする課題】

【0004】

ストレージ仮想化機能を備える計算機システムにおいて、1台のストレージ装置が多数の仮想ボリュームを持つ場合、仮想ボリュームに対するデータの読み書き量が急激に増加したとき、データ処理負担が増大してホストコンピュータに対する応答性能が低下する恐れがある。また、仮想ボリュームを持つストレージ装置がダウンしたとき、ホストコンピュータから実ボリュームに対する読み書きが一切行えなくなるという問題がある。

20

【0005】

ストレージ仮想化機能を用いた計算機システムにおいて、負荷集中や装置障害によって特定の仮想ストレージ装置の入出力処理能力が低下したとき、複数の仮想ストレージ装置間でストレージ処理能力と入出力負荷とのバランスを再調整することによって、仮想ストレージ装置の入出力処理能力を回復させることができる。ここで、仮想ストレージ装置とは、仮想ボリュームを持つストレージ装置である。このような機能を具備するストレージシステムを「クラスタ型ストレージ」と呼ぶ。クラスタ型ストレージにおいては、ストレージ管理者の介在なしに、自動的にストレージ装置間の負荷バランスを調整することができる。

30

【0006】

クラスタ型ストレージにおいては、一つの実ボリュームに対応する仮想ボリュームを含む複数の仮想ストレージ装置が存在する。すなわち、ホストコンピュータから仮想ボリュームを経由して一つの実ボリュームに至るデータの経路が複数存在する。この場合、仮想ボリュームへのデータの読み書きを複数経路に振り分ける「交代パスソフト」をホストコンピュータ上に置くことによって、管理者による設定に応じて、各経路に所定の割合でデータの読み書き量を割り振ることが可能となる。このとき、複数の経路のうち、正常時に主にデータが流れる主経路と、正常時は少量のデータのみが流れ、主経路の入出力処理能力低下が発生したときに新たな主経路となる副経路とが存在しうる。

40

【0007】

管理ソフトウェアは、監視対象となるストレージ装置及びホストコンピュータから定期的に構成情報及び性能情報を取得し、管理ソフトウェアが管理するデータベースに蓄積する。そして、管理者からの要求に応じて、蓄積した構成情報及び性能情報を表示する。性能情報は、例えば、監視対象装置の論理ボリューム及びポート等の構成要素ごとに取得された、単位時間当たりの受信又は送信したデータ量等である。以下、上記の論理ボリューム及びポート等、各経路の構成要素を「デバイス」と呼ぶ。

【0008】

50

管理者は、ストレージ管理ソフトウェアを使用して性能を監視する場合、監視対象デバイスに対して、あらかじめ閾値を設定しておく。そして、例えばある論理ボリュームに対するホストコンピュータからの読み書き量が閾値を越えたとき、管理ソフトウェアは管理者に対し閾値を超えたことを通知して注意を喚起することができる。

【0009】

しかし、管理者が閾値を設定する際、交代パスソフトによる仮想ボリュームに対するデータ振り分け比率を考慮せずに閾値を設定してしまう可能性がある。例えば、データ流量の少ない副経路に大きな閾値を設定してしまうと、閾値が本来の意味をなさなくなる恐れがある。このように、不適切な閾値の設定による適切な性能管理の障害を防ぐことが、第1の課題である。

【0010】

また、クラスタ型ストレージは、ストレージ装置間の負荷バランスを取るために、ストレージ装置の機器構成を動的に変化させ、主経路と副経路を入れ替える機能を持つ。しかし、交代パスソフト上での各経路に対するデータ振り分け比率と、管理ソフトウェア上に設けられた閾値は、管理者による操作がないと変化しない。従って、管理者による操作がなければ、各経路の性質と流量バランスが合致しない状態で計算機システムを引き続き使い続けることとなる。このことが、交代パスソフト及び管理ソフトウェアによる適切な性能管理を阻害することとなる。上記のような、適切な性能管理の障害を防ぐことが、第2の課題である。

【0011】

また、クラスタ型ストレージは、主経路及び副経路上のデバイスに何らかの障害が発生した場合、ストレージ装置の機器構成を動的に変化させ、それまでの副経路を新たな主経路とすることによって、データ入出力経路を確保する機能を持つ。しかし、各経路の流量バランスが変化しても、管理ソフトウェア上に設けられた閾値は、管理者による操作がないと変化しない。従って、管理者による操作がなければ、各経路の流量バランスに合致しない閾値を引き続き使い続けることとなる。このことが、管理ソフトウェアによる適切な性能管理を阻害することとなる。流量バランスの変化による、管理ソフトウェアによる適切な性能管理の障害を防止することが第3の課題である。

【課題を解決するための手段】

【0012】

本願で開示する代表的な発明は、ホスト計算機と、前記ホスト計算機に接続される複数の仮想ストレージ装置と、前記複数の仮想ストレージ装置に接続されるストレージ装置と、前記ホスト計算機、前記仮想ストレージ装置及び前記ストレージ装置に接続される管理計算機と、を備える計算機システムであって、前記ストレージ装置は、第1ネットワークを介して前記仮想ストレージ装置に接続される第1インターフェースと、第2ネットワークを介して前記管理計算機に接続される第2インターフェースと、前記第1インターフェース及び前記第2インターフェースに接続される第1プロセッサと、前記第1プロセッサに接続される第1メモリと、前記ホスト計算機に論理ボリュームとして提供される記憶領域と、を備え、前記各仮想ストレージ装置は、前記第1ネットワークを介して前記ホスト計算機及び前記ストレージ装置に接続される第3インターフェースと、前記第2ネットワークを介して前記管理計算機に接続される第4インターフェースと、前記第3インターフェース及び前記第4インターフェースに接続される第2プロセッサと、前記第2プロセッサに接続される第2メモリと、を備え、前記論理ボリュームに対応付けられる仮想記憶領域を前記ホスト計算機に提供し、前記ホスト計算機は、前記第1ネットワークを介して前記仮想ストレージ装置に接続される第5インターフェースと、前記第2ネットワークを介して前記管理計算機に接続される第6インターフェースと、前記第5インターフェース及び前記第6インターフェースに接続される第3プロセッサと、前記第3プロセッサに接続される第3メモリと、を備え、一つの前記論理ボリュームに書き込まれるべきデータを含む複数のデータI/Oを、定められた比率で、前記一つの論理ボリュームと対応付けられた、前記各仮想ストレージ装置が提供する各々の前記仮想記憶領域に分散し、前記管理計

10

20

30

40

50

算機は、前記第2ネットワークを介して前記ホスト計算機、前記仮想ストレージ装置及び前記ストレージ装置に接続される第7インターフェースと、前記第7インターフェースに接続される第4プロセッサと、前記第4プロセッサに接続される第4メモリと、を備え、前記各仮想記憶領域に対するデータI/Oの量が、前記各仮想記憶領域に設定される閾値よりも大きくなった場合に、その旨を通知する警告を発するものであり、前記各仮想記憶領域のうち第一仮想記憶領域についての閾値が入力された後に、前記第一仮想記憶領域に対応付けられている前記一つの論理ボリュームと対応付けられている、前記各仮想記憶領域のうち他の仮想記憶領域についての閾値を、前記第一仮想記憶領域についての閾値と前記定められた比率に基づいて算出することを特徴とする。

【発明の効果】

10

【0013】

本発明の一実施形態によれば、計算機システム内のデバイスに適切な閾値を設定することによって、計算機システムの性能を適切に管理することができる。

【発明を実施するための最良の形態】

【0014】

以下に図面を参照しながら本発明の実施形態を説明する。

【0015】

最初に、本発明の第1の実施の形態について説明する。

【0016】

本発明の第1の実施の形態では、ストレージ管理ソフトがホストコンピュータ上の交代パスソフトから、各経路に対するデータ割り振り比率をあらかじめ取得する。ストレージ管理者は、管理対象デバイスに対して性能管理のための閾値を設定する際、データ割り振り比率に基づいて、各経路に属するデバイスの閾値を算出し、設定する。

20

【0017】

図1は、本発明の第1の実施の形態の計算機システムの構成を示すブロック図である。

【0018】

第1の実施の形態の計算機システムは、複数の仮想ストレージ装置20000、ストレージ装置25000、ホストコンピュータ10000及び管理サーバ30000を備える。複数の仮想ストレージ装置20000、ストレージ装置25000及びホストコンピュータ10000は、ストレージエリアネットワーク40000によって接続されている。また、複数の仮想ストレージ装置20000、ストレージ装置25000及びホストコンピュータ10000は、管理用ネットワーク45000によって管理サーバ30000に接続されている。

30

【0019】

ストレージエリアネットワーク40000は、例えばファイバーチャネルプロトコルが適用されるネットワークであってもよいし、他の種類のネットワークであってもよい。管理用ネットワーク45000は、いわゆるLAN(Local Area Network)であってもよいし、他の種類のネットワークであってもよい。ストレージエリアネットワーク40000及び管理用ネットワーク45000は、同じネットワークであってもよい。

40

【0020】

図2は、本発明の第1の実施の形態のホストコンピュータ10000の詳細な構成を示すブロック図である。

【0021】

ホストコンピュータ10000は、ストレージエリアネットワーク40000に接続するためのインターフェースであるI/Oポート11000と、管理用ネットワーク45000に接続するためのインターフェースである管理ポート12000と、プロセッサ13000と、メモリ14000とを備える。これらは内部バス等の回路を介して相互に接続される。メモリ14000には、業務アプリケーション14100、オペレーティングシステム14200、交代パスソフト14300及び交代パス管理表14400が格納され

50

る。

【0022】

プロセッサ13000は、メモリ14000上の業務アプリケーション14100、オペレーティングシステム14200及び交代パスソフト14300を動作させる。一方で、プロセッサ13000は、業務アプリケーション14100において行われる業務の状況に応じて、I/Oポート11000を介してストレージエリアネットワーク40000によって接続されたストレージ装置25000上の記憶領域に対しデータ入出力(I/O)を発行する。

【0023】

なお、ストレージ装置25000上の記憶領域に対しI/Oの発行を実際に行うのはプロセッサ13000であるが、以下では便宜上、プロセッサ13000が動作させる業務アプリケーション14100と、オペレーティングシステム14200と、交代パスソフト14300が主体となってI/Oを発行するものとして記述する。

【0024】

業務アプリケーション14100は、オペレーティングシステムから提供された記憶領域を使用し、記憶領域に対しデータ入出力(以下、I/Oと表記)を行う。

【0025】

オペレーティングシステム14200は、交代パスソフトから提供された仮想デバイス(後述)を記憶領域としてアプリケーションに認識させる。また、オペレーティングシステム14200は、業務アプリケーション14100から受けたI/Oを、交代パスソフト14300を介して仮想デバイス16000に発行する。

【0026】

交代パスソフト14300は、ストレージエリアネットワーク40000を介してホストコンピュータ10000に接続されたストレージ装置25000上の論理ボリューム29100をデバイスファイル15000として認識し、複数のデバイスファイル15000を1つの仮想デバイス16000としてオペレーティングシステム14200に認識させる。論理ボリューム29100、デバイスファイル15000及び仮想デバイス16000については後述する(図4及び図6参照)。また、交代パスソフト14300は、オペレーティングシステム14200から仮想デバイス16000に対して発行された複数のI/Oを、複数のデバイスファイル15000に分散して発行する。

【0027】

交代パス管理表14400には、仮想デバイス16000に対するI/Oを、交代パスソフト14300が各デバイスファイル15000に割り振る比率が格納される(図7参照)。

【0028】

プロセッサ13000は、メモリ14000に格納された業務アプリケーション14100、オペレーティングシステム14200及び交代パスソフト14300等のソフトウェアを実行する。以下の説明においてメモリ14000内のソフトウェアが実行する処理は、実際にはプロセッサ13000によって実行される。

【0029】

図3は、本発明の第1の実施の形態の仮想ストレージ装置20000の詳細な構成例を示すブロック図である。

【0030】

仮想ストレージ装置20000は、ストレージエリアネットワーク40000を介してホストコンピュータ10000又はストレージ装置25000に接続するためのインターフェースであるI/Oポート21000と、管理用ネットワーク45000に接続するためのインターフェースである管理ポート21100と、仮想ストレージ装置20000を制御するプロセッサ22000と、データを一時的に記憶するキャッシュメモリ22100と、プロセッサが用いる管理メモリ23000とを備える。これらは内部バス等の回路を介して相互に接続される。

10

20

30

40

50

【0031】

管理メモリ23000には、仮想ストレージ装置20000の管理プログラム23100及び仮想ボリューム管理表23200が格納される。

【0032】

管理プログラム23100は、ストレージ仮想化機能及びストレージクラスタ機能を持つ。ストレージ仮想化機能は、I/Oポート21000を介して接続されたストレージ装置25000内の論理ボリューム29100(図4参照)を、ホストコンピュータ10000に対して仮想ボリューム24000(図6参照)として提供する。すなわち、仮想ボリューム24000は、論理ボリューム29100と対応付けられた仮想的な記憶領域である。ストレージクラスタ機能については後述する。

10

【0033】

仮想ボリューム管理表23100には、仮想ボリューム24000とストレージ装置25000内の論理ボリューム29100との対応関係を示す情報が格納される。

【0034】

プロセッサ22000は、管理メモリ23000に格納された管理プログラム23100を実行する。以下の説明において管理プログラム23100が実行する処理は、実際にはプロセッサ22000によって実行される。

【0035】

キャッシュメモリ22100は、仮想ストレージ装置20000とホストコンピュータ10000との間でやりとりされるデータ、及び、仮想ストレージ装置20000とストレージ装置25000との間でやりとりされるデータの少なくとも一方を一時的に記憶する。

20

【0036】

図3には一つのI/Oポート21000を示すが、仮想ストレージ装置20000は、複数のI/Oポート21000を備えてもよい。

【0037】

図4は、本発明の第1の実施の形態のストレージ装置25000の詳細な構成例を示すブロック図である。

【0038】

ストレージ装置25000は、ストレージエリアネットワーク40000を介して仮想ストレージ装置20000に接続するためのインターフェースであるI/Oポート26000と、管理用ネットワーク45000に接続するためのインターフェースである管理ポート26100と、ストレージ装置25000を制御するプロセッサ27000と、仮想ストレージ装置20000との間でやりとりされるデータを一時的に記憶するキャッシュメモリ27100と、プロセッサ27000が用いる管理メモリ28000と、ホストコンピュータ10000に提供するデータを格納する一つ以上のディスクボリューム29000とを備える。これらは内部バス等の回路を介して相互に接続される。

30

【0039】

管理メモリ28000には、ストレージ装置25000の管理プログラム28100が格納されている。

40

【0040】

各ディスクボリューム29000は、1つ又は複数の磁気ディスクによって構成されている。ディスクボリューム29000が複数の磁気ディスクによって構成されている場合、それらの磁気ディスクはRAID(Redundant Arrays of Inexpensive Disks)を構成してもよい。各ディスクボリューム29000は、論理的に一つ以上の論理ボリューム29100に分割されている。

【0041】

プロセッサ27000は、管理メモリ28000に格納された管理プログラム28100を実行する。以下の説明において管理プログラム28100が実行する処理は、実際にはプロセッサ27000によって実行される。

50

【 0 0 4 2 】

図 5 は、本発明の第 1 の実施の形態の管理サーバ 3 0 0 0 0 の詳細な構成を示すブロック図である。

【 0 0 4 3 】

管理サーバ 3 0 0 0 0 は、管理用ネットワーク 4 5 0 0 0 に接続するためのインターフェースである管理ポート 3 1 0 0 0 と、プロセッサ 3 2 0 0 0 と、メモリ 3 3 0 0 0 と、出力部 3 4 0 0 0 と、入力部 3 5 0 0 0 とを備える。これらは内部バス等の回路を介して相互に接続される。

【 0 0 4 4 】

メモリ 3 3 0 0 0 には、構成管理プログラム 3 3 1 1 0、装置構成管理表 3 3 2 0 0、装置性能管理表 3 3 3 0 0、ストレージクラスタ管理表 3 3 4 0 0、デバイスグループ管理表 3 3 5 0 0 及び交代バス管理表 3 3 6 0 0 が格納される。

10

【 0 0 4 5 】

構成管理プログラム 3 3 1 1 0 は、仮想ストレージ装置 2 0 0 0 0、ストレージ装置 2 5 0 0 0 及びホストコンピュータ 1 0 0 0 0 から構成情報及び性能情報を定期的を取得して、装置構成管理表 3 3 2 0 0 及び装置性能管理表 3 3 3 0 0 に格納する。また、構成管理プログラム 3 3 1 1 0 は、管理者からの要求に応じ、取得した構成情報及び性能情報を表示すると共に、管理対象デバイスに対する性能管理用の閾値の設定を受け付ける。

【 0 0 4 6 】

装置構成管理表 3 3 2 0 0 には、各装置から取得した構成情報が格納される。少なくとも、装置構成管理表 3 3 2 0 0 には、各仮想ストレージ装置 2 0 0 0 0 の仮想ボリューム管理表 2 3 2 0 0 が保持する情報と同じ情報が格納される。

20

【 0 0 4 7 】

出力部 3 4 0 0 0 は、入力画面（図 1 0 参照）及び処理結果等を管理者に対して出力する。出力部 3 4 0 0 0 は、例えば、任意の画像表示装置であってもよい。

【 0 0 4 8 】

入力部 3 5 0 0 0 は、管理者が指示を入力するために使用される。入力部 3 5 0 0 0 は、例えば、キーボード及びポインティングデバイス等であってもよい。

【 0 0 4 9 】

装置性能管理表 3 3 3 0 0 には、仮想ストレージ装置 2 0 0 0 0、ストレージ装置 2 5 0 0 0 及びホストコンピュータ 1 0 0 0 0 を構成するデバイスの性能情報が格納される。

30

【 0 0 5 0 】

ストレージクラスタ管理表 3 3 4 0 0、デバイスグループ管理表 3 3 5 0 0 及び交代バス管理表 3 3 6 0 0 は、本発明の第 1 の実施の形態において追加されたものである。すなわち、従来の管理サーバ 3 0 0 0 0 は、これらの管理表を含まない。これらの管理表については、後で詳細に説明する。

【 0 0 5 1 】

図 6 は、本発明の第 1 の実施の形態のストレージクラスタ構成の説明図である。

【 0 0 5 2 】

具体的には、図 6 は、ホストコンピュータ 1 0 0 0 0、3 台の仮想ストレージ装置 2 0 0 0 0 及びストレージ装置 2 5 0 0 0 を組み合わせたストレージクラスタ構成を示す。ストレージ装置（SYS 4）2 5 0 0 0 は、ストレージエリアネットワーク 4 0 0 0 0 を介して仮想ストレージ装置（SYS 1）2 0 0 0 0 と接続されている。

40

【 0 0 5 3 】

ここで、SYS 4 及び SYS 1 は、それぞれ、ストレージ装置 2 5 0 0 0 及び 1 台の仮想ストレージ装置 2 0 0 0 0 に付与された計算機システム内で一意の識別子（すなわち装置 ID）である。残りの 2 台の仮想ストレージ装置 2 0 0 0 0 にも、それぞれ識別子 SYS 2 及び SYS 3 が付与される。以下の説明において、ストレージ装置 2 5 0 0 0 は、単に「SYS 1」とも表記される。各仮想ストレージ装置 2 0 0 0 0 も、同様に、識別子によって表示される。

50

【 0 0 5 4 】

なお、ホストコンピュータ 1 0 0 0 0 には、計算機システム内で一意の識別子「H O S T 1」が付与される。

【 0 0 5 5 】

S Y S 4 は、それぞれ識別子「V O L 1」及び「V O L 2」が付与された論理ボリューム 2 9 1 0 0 を含む。以下、各論理ボリューム 2 9 1 0 0 を、単に「V O L 1」及び「V O L 2」とも表記する。

【 0 0 5 6 】

S Y S 1 は、それぞれ識別子「V O L 1 1 1」及び「V O L 1 1 2」が付与された仮想ボリューム 2 4 0 0 0 を含む。S Y S 2 は、それぞれ識別子「V O L 1 2 1」及び「V O L 1 2 2」が付与された仮想ボリューム 2 4 0 0 0 を含む。S Y S 3 は、それぞれ識別子「V O L 1 3 1」及び「V O L 1 3 2」が付与された仮想ボリューム 2 4 0 0 0 を含む。以下、各仮想ボリューム 2 9 1 0 0 を、「V O L 1 1 1」のように識別子によって表記する。

10

【 0 0 5 7 】

論理ボリューム V O L 1 及び V O L 2 は、それぞれ、S Y S 1 上の仮想ボリューム V O L 1 1 1 及び V O L 1 1 2 に対応付けられている。すなわち、S Y S 1 は、V O L 1 1 1 及び V O L 1 1 2 に対する I / O を受け付けると、それらの I / O を、それぞれ、S Y S 4 の V O L 1 及び V O L 2 に対して発行する。

【 0 0 5 8 】

一方、S Y S 1 はストレージエリアネットワーク 4 0 0 0 0 を介して S Y S 2 と接続されている。S Y S 1 上の仮想ボリューム V O L 1 1 1 及び V O L 1 1 2 は、それぞれ、S Y S 2 上の仮想ボリューム V O L 1 2 1 及び V O L 1 2 2 に対応付けられている。すなわち、S Y S 2 は、V O L 1 2 1 及び V O L 1 2 2 に対する I / O を受け付けると、それらの I / O を、それぞれ、S Y S 1 の V O L 1 1 1 及び V O L 1 1 2 に対して発行する。

20

【 0 0 5 9 】

さらに、S Y S 1 はストレージエリアネットワーク 4 0 0 0 0 を介して仮想ストレージ装置 S Y S 3 と接続されている。S Y S 1 上の仮想ボリューム V O L 1 1 1、V O L 1 1 2 は、それぞれ、S Y S 3 上の仮想ボリューム V O L 1 3 1 及び V O L 1 3 2 に対応付けられている。すなわち、S Y S 3 は、V O L 1 3 1 及び V O L 1 3 2 に対する I / O を受け付けると、それらの I / O を、それぞれ、S Y S 1 の V O L 1 1 1 及び V O L 1 1 2 に対して発行する。

30

【 0 0 6 0 】

結局、ホストコンピュータ 1 0 0 0 0 から V O L 1 1 1、V O L 1 2 1 及び V O L 1 3 1 に対して行われた I / O は、すべて S Y S 1 を介して S Y S 4 上の V O L 1 に対して行われる。また、ホストコンピュータ 1 0 0 0 0 から V O L 1 1 2、V O L 1 2 2、V O L 1 3 2 に対して行われた I / O は、すべて S Y S 1 を介して S Y S 4 上の V O L 2 に対して行われる。

【 0 0 6 1 】

例えば、ホストコンピュータ 1 0 0 0 0 が、データを書き込むための I / O 要求を S Y S 1 の V O L 1 1 1 に対して発行すると、S Y S 1 は、そのデータを書き込むための I / O 要求を S Y S 4 の V O L 1 に対して発行する。S Y S 4 は、S Y S 1 からの要求に従って、データを V O L 1 に書き込む。このように、ホストコンピュータ 1 0 0 0 0 が仮想ボリューム V O L 1 1 1 に書き込むために発行した I / O 要求に含まれるデータは、最終的に、V O L 1 1 1 に対応付けられた論理ボリューム V O L 1 に格納される。

40

【 0 0 6 2 】

ホストコンピュータ 1 0 0 0 0 上の交代パスソフト 1 4 3 0 0 は、ストレージエリアネットワーク 4 0 0 0 0 を介して接続された仮想ストレージ装置 2 0 0 0 0 上の仮想ボリューム 2 4 0 0 0 を、デバイスファイル 1 5 0 0 0 として認識している。図 6 の例において、交代パスソフト 1 4 3 0 0 は、それぞれ識別子（すなわちデバイスファイル ID）「D

50

「DEV111」、「DEV121」、「DEV131」、「DEV112」、「DEV122」及び「DEV132」が付与された六つのデバイスファイル15000を認識する。以下、各デバイスファイル15000を、「DEV111」のように識別子によって表記する。

【0063】

具体的には、交代パスソフト14300は、仮想ストレージ装置SYS1上の仮想ボリュームVOL111及びVOL112を、それぞれ、デバイスファイルDEV111及びDEV112として認識している。交代パスソフト14300は、仮想ストレージ装置SYS2上の仮想ボリュームVOL121及びVOL122を、それぞれ、デバイスファイルDEV121及びDEV122として認識している。さらに、交代パスソフト14300は、仮想ストレージ装置SYS3上の仮想ボリュームVOL131及びVOL132を、それぞれ、デバイスファイルDEV131及びDEV132として認識している。

10

【0064】

ホストコンピュータ10000上の交代パスソフト14300は、認識した複数のデバイスファイル15000を、仮想デバイス16000としてオペレーティングシステム14200に認識させる。図6の例において、交代パスソフト14300は、それぞれ識別子(すなわち仮想デバイスID)「DEV1」及び「DEV2」が付与された二つの仮想デバイス16000をオペレーティングシステム14200に認識させる。以下、各仮想デバイス16000を、「DEV1」のように識別子によって表記する。

【0065】

20

具体的には、交代パスソフト14300は、デバイスファイルDEV111、DEV121及びDEV131を、仮想デバイスDEV1としてオペレーティングシステム14200に認識させる。さらに、交代パスソフト14300は、デバイスファイルDEV112、DEV122及びDEV132を、仮想デバイスDEV2としてオペレーティングシステム14200に認識させる。

【0066】

さらに、交代パスソフト14300は、アプリケーション14100からオペレーティングシステム14200を介して仮想デバイス16000に対して行われるI/Oを、各デバイスファイル15000に振り分ける。具体的には、仮想デバイスDEV1に対して行われるI/Oは、交代パスソフト14300によって、デバイスファイルDEV111、DEV121及びDEV131に振り分けられる。仮想デバイスDEV2に対して行われるI/Oは、交代パスソフト14300によって、デバイスファイルDEV112、DEV122、DEV132に振り分けられる。

30

【0067】

各装置は、I/Oポートを介してストレージエリアネットワーク40000に接続される。各I/Oポートには、計算機システム内で一意の識別子(すなわちポートID)が付与される。

【0068】

具体的には、ホストコンピュータHOST1を、ストレージエリアネットワーク40000に接続するI/Oポート11000のポートIDは、「PORT4」である。以下、各ポートは、「PORT4」のようにポートIDによって表記される。

40

【0069】

仮想ストレージ装置SYS1をストレージエリアネットワーク40000に接続するI/Oポート21000のポートIDは、「PORT1」、「PORT21」及び「PORT11」である。PORT1及びPORT21は、それぞれ、ストレージエリアネットワーク40000を介してHOST1及びSYS4と接続される。PORT11は、ストレージエリアネットワーク40000を介してSYS2及びSYS3と接続される。

【0070】

仮想ストレージ装置SYS2をストレージエリアネットワーク40000に接続するI/Oポート21000のポートIDは、「PORT2」及び「PORT12」である。P

50

PORT 2 及び PORT 1 2 は、それぞれ、ストレージエリアネットワーク 4 0 0 0 0 を介して HOST 1 及び SYS 1 と接続される。

【 0 0 7 1 】

仮想ストレージ装置 SYS 3 をストレージエリアネットワーク 4 0 0 0 0 に接続する I / O ポート 2 1 0 0 0 のポート ID は、「 PORT 3 」、「 PORT 1 3 」及び「 PORT 2 3 」である。 PORT 2 及び PORT 1 3 は、それぞれ、ストレージエリアネットワーク 4 0 0 0 0 を介して HOST 1 及び SYS 1 と接続される。 PORT 2 3 は、図 6 の例では、まだストレージエリアネットワーク 4 0 0 0 0 に接続されていない。

【 0 0 7 2 】

ストレージ装置 SYS 4 をストレージエリアネットワーク 4 0 0 0 0 に接続する I / O ポート 2 6 0 0 0 のポート ID は、「 PORT 2 4 」である。 PORT 2 4 は、ストレージエリアネットワーク 4 0 0 0 0 を介して SYS 1 と接続される。

【 0 0 7 3 】

結局、各仮想デバイス 1 6 0 0 0 は、各論理ボリューム 2 9 1 0 0 と 1 対 1 に対応付けられる。そして、各仮想デバイス 1 6 0 0 0 からそれに対応付けられた論理ボリューム 2 9 1 0 0 に至る複数の経路が存在する。すなわち、オペレーティングシステム 1 4 2 0 0 がある仮想デバイス 1 6 0 0 0 に対して実行した I / O は、複数のうちのいずれかの経路を経由して、その仮想デバイス 1 6 0 0 0 に対応付けられた論理ボリューム 2 9 1 0 0 に至る。そして、その論理ボリューム 2 9 1 0 0 に対する I / O が実行される。

【 0 0 7 4 】

なお、図 6 の構成は一例であり、実際には、ホストコンピュータ 1 0 0 0 0 は任意の数の仮想デバイス 1 6 0 0 0 及び任意の数のデバイスファイル 1 5 0 0 0 を含んでもよい。また、各仮想ストレージ装置 2 0 0 0 0 は、任意の数の仮想ボリューム 2 4 0 0 0 を含んでもよい。ストレージ装置 2 5 0 0 0 は、任意の数の論理ボリューム 2 9 1 0 0 を含んでもよい。

【 0 0 7 5 】

図 2 1 は、本発明の第 1 の実施の形態の交代パスソフト 1 4 3 0 0 が実行する I / O 振り分け処理のフローチャートである。

【 0 0 7 6 】

ホスト計算機 1 0 0 0 0 上の業務アプリケーション 1 4 1 0 0 は、オペレーティングシステム 1 4 2 0 0 から提供された記憶領域に対するデータの読み書きが必要になると、オペレーティングシステム 1 4 2 0 0 に対しデータ I / O 要求を行う (ステップ 6 4 0 0 0)。

【 0 0 7 7 】

データ I / O 要求を受けたオペレーティングシステム 1 4 2 0 0 は、交代パスソフト 1 4 3 0 0 に対し、データ I / O 要求が行われた仮想デバイス 1 6 0 0 0 への I / O を行わせる (ステップ 6 4 0 1 0)。

【 0 0 7 8 】

指示を受けた交代パスソフト 1 4 3 0 0 は、交代パス管理表 1 4 4 0 0 を参照し、仮想デバイス 1 6 0 0 0 がどのデバイスファイル 1 5 0 0 0 と関連付けられているかを調べる (ステップ 6 4 0 2 0)。

【 0 0 7 9 】

仮想デバイス 1 6 0 0 0 が、ただ 1 つのデバイスファイル 1 5 0 0 0 と関連付けられていた場合、交代パスソフト 1 4 3 0 0 は、オペレーティングシステム 1 4 2 0 0 を介して業務アプリケーション 1 4 1 0 0 から受けたデータ I / O 要求を、そのままデバイスファイル 1 5 0 0 0 に行わせる (ステップ 6 4 0 3 0)。

【 0 0 8 0 】

一方、仮想デバイス 1 6 0 0 0 が、複数のデバイスファイル 1 5 0 0 0 と関連付けられていた場合、交代パスソフト 1 4 3 0 0 は交代パス管理表 1 4 4 0 0 を参照し、各デバイスファイルに割り振る I / O の比率を確認する (ステップ 6 4 0 4 0)。

10

20

30

40

50

【 0 0 8 1 】

その上で交代パスソフト 1 4 3 0 0 は、オペレーティングシステム 1 4 2 0 0 を介して業務アプリケーション 1 4 1 0 0 から受けたデータ I / O 要求を、定められた I / O 割り振り比率通りに各デバイスファイル 1 5 0 0 0 に行わせる（ステップ 6 4 0 5 0 ）。

【 0 0 8 2 】

以上が、本発明の第 1 の実施の形態における、交代パスソフト 1 4 3 0 0 が実行する I / O 振り分け処理のフローチャートである。

【 0 0 8 3 】

図 7 は、本発明の第 1 の実施の形態のホストコンピュータ 1 0 0 0 0 が保持する交代パス管理表 1 4 4 0 0 の説明図である。

10

【 0 0 8 4 】

交代パス管理表 1 4 4 0 0 は、フィールド 1 4 4 1 0 から 1 4 4 6 0 までの 6 フィールドからなる。

【 0 0 8 5 】

フィールド 1 4 4 1 0 には、ホストコンピュータ 1 0 0 0 0 内で各仮想デバイス 1 6 0 0 0 を識別する仮想デバイス ID が登録される。

【 0 0 8 6 】

フィールド 1 4 4 2 0 には、ホストコンピュータ 1 0 0 0 0 内で各仮想デバイス 1 6 0 0 0 に対応するデバイスファイル 1 5 0 0 0 を識別するデバイスファイル ID が登録される。

20

【 0 0 8 7 】

フィールド 1 4 4 3 0 には、各デバイスファイル 1 5 0 0 0 に対応する仮想ボリューム 2 4 0 0 0 を含む仮想ストレージ装置 2 0 0 0 0 を識別する装置 ID が登録される。

【 0 0 8 8 】

フィールド 1 4 4 4 0 には、各デバイスファイル 1 5 0 0 0 に対応する仮想ボリューム 2 4 0 0 0 を含む仮想ストレージ装置 2 0 0 0 0 がホストコンピュータ 1 0 0 0 0 と接続するために使用するポート 2 1 0 0 0 を識別するポート ID が登録される。

【 0 0 8 9 】

フィールド 1 4 4 5 0 には、各デバイスファイル 1 5 0 0 0 に対応する仮想ボリューム 2 4 0 0 0 を識別するボリューム ID が登録される。

30

【 0 0 9 0 】

フィールド 1 4 4 6 0 には、交代パスソフト 1 4 3 0 0 が各デバイスファイル 1 5 0 0 0 に対し I / O を割り振る際のデータ量の比率を示す値が登録される。

【 0 0 9 1 】

図 7 には、ホストコンピュータ 1 0 0 0 0 が保持する交代パス管理表 1 4 4 0 0 に登録された具体的な値の一例を示している。

【 0 0 9 2 】

この例では、ホストコンピュータ 1 0 0 0 0 上の交代パスソフト 1 4 3 0 0 は、仮想デバイス DEV 1 への I / O を、デバイスファイル DEV 1 1 1、DEV 1 2 1 及び DEV 1 3 1 へ割り振っている。

40

【 0 0 9 3 】

DEV 1 1 1 は、仮想ストレージ装置 SYS 1 上の仮想ボリューム VOL 1 1 1 と、PORT 1 を介して接続されている。同様に、DEV 1 2 1 は、SYS 2 上の VOL 1 2 1 と、PORT 2 を介して接続されている。DEV 1 3 1 は、SYS 3 上の VOL 1 3 1 と、PORT 3 を介して接続されている。

【 0 0 9 4 】

交代パスソフト 1 4 3 0 0 は、仮想デバイス DEV 1 への I / O を、デバイスファイル DEV 1 1 1、DEV 1 2 1 及び DEV 1 3 1 に対し、それぞれ 7 0 %、2 0 % 及び 1 0 % の割合で割り振ることを示している。

【 0 0 9 5 】

50

図8は、本発明の第1の実施の形態の仮想ストレージ装置20000が保持する仮想ボリューム管理表23200の説明図である。

【0096】

仮想ボリューム管理表23200は、フィールド23210から23250までの5フィールドからなる。

【0097】

フィールド23210には、仮想ストレージ装置20000内で各仮想ボリューム24000を識別する仮想ボリュームIDが登録される。

【0098】

フィールド23220には、各仮想ボリューム24000に対応するボリュームを含む装置と接続するために使用されるポートを識別するポートIDが登録される。

10

【0099】

フィールド23230には、各仮想ボリューム24000に対応するボリュームを含む装置を識別する接続先装置IDが登録される。

【0100】

フィールド23240には、各仮想ボリューム24000に対応するボリュームを含む装置が論理ボリューム29100と接続するために使用されるポートを識別する接続先ポートIDが登録される。

【0101】

フィールド23250には、各仮想ボリューム24000に対応するボリュームを識別する接続先ボリュームIDが登録される。

20

【0102】

図8には、仮想ストレージ装置20000が保持する仮想ボリューム管理表23200に登録された具体的な値の一例を示している。この例では、仮想ストレージ装置20000上の仮想ボリュームVOL111及びVOL112は、それぞれ、SYS4上のVOL1及びVOL2と、PORT21及びPORT24を介して接続されている。

【0103】

図9は、本発明の第1の実施の形態の管理サーバ30000が保持する装置性能管理表33300の説明図である。

【0104】

装置性能管理表33300は、フィールド33310から33340までの4フィールドからなる。

30

【0105】

フィールド33310には、管理対象となるデバイスが属するストレージ装置25000、仮想ストレージ装置20000又はホストコンピュータ10000の識別子である装置IDが登録される。

【0106】

フィールド33320には、管理対象デバイスの識別子であるデバイスIDが登録される。

【0107】

フィールド33330には、管理対象デバイスについて実測された性能値が格納される。図9の例では、フィールド33330には、管理対象デバイスの単位時間当たりのI/O量(すなわち、そのデバイスに入出力される単位時間当たりのデータ量)が登録される。このI/O量は、管理対象デバイスが属するストレージ装置25000等から管理サーバ30000が取得したものである。

40

【0108】

フィールド33340には、アラート実行のための性能値の閾値が登録される。この閾値は、管理者から入力されたものである。

【0109】

管理サーバ30000は、管理対象デバイスの単位時間当たりのI/O量33330が

50

アラート実行閾値 33340 を超えた場合、管理者に対しメール等の手段によってアラートを送信する。このように、アラート実行閾値 33340 に登録された値は、管理サーバ 30000 がアラートを送信するか否かを判定するために、管理対象デバイスに設定される。

【0110】

図9には、管理サーバ 30000 が保持する装置性能管理表 33300 に登録された具体的な値の一例を示している。この例では、ストレージ装置 SYS1 内のボリューム VOL111 では、単位時間当たり 200 の I/O 量が発生している。また、VOL111 の単位時間当たりの I/O が 1000 を超えた場合、管理サーバ 30000 は管理者に対しアラートを送信する。

10

【0111】

なお、ここでは、管理サーバ 30000 が管理するデバイスの性能値の例として単位時間当たりの I/O 量を挙げたが、管理サーバ 30000 が管理する性能値はこれ以外の値（例えば、単位時間当たりの I/O 回数）でもよい。また、管理サーバ 30000 は、複数種類の性能値を同時に管理してもよい。

【0112】

図10は、本発明の第1の実施の形態において表示されるアラート実行閾値設定画面の説明図である。

【0113】

図10に示すアラート実行閾値設定画面 71000 は、ストレージ管理者が管理サーバ 30000 を用いてデバイスに対して性能管理用の閾値を設定する際に、管理サーバ 30000 の出力部 34000 に表示される画面の例である。アラート実行閾値設定画面 71000 では、管理者が、閾値を設定するデバイスの ID を指定し（テーブル 71010）、閾値を指定し（テーブル 71020）、閾値を超えたときにその旨を通知するメールアドレスを指定する（テーブル 71030）。管理者は、指定したパラメータを確認の上、論理ボリューム作成を続行する場合は「確認」ボタン 71040 を、論理ボリューム作成を中止する場合は「中止」ボタン 71050 を押下する。「中止」ボタンが押下されると、管理サーバ 30000 の構成管理プログラム 33110 は閾値を設定せずに処理を終了する。「確認」ボタンが押下されると、構成管理プログラム 33110 は、管理者からの閾値設定指示を受け、指示された閾値を装置性能管理表 33300 のフィールド 33340 に登録する。

20

30

【0114】

次に、本発明の第1の実施の形態において追加された管理表 33400 から 33600、及び、本実施の形態の構成管理プログラム 33110 が実行する処理について、図11から図13を参照して説明する。その説明に先立って、従来のシステム構成及びプログラム構成の問題点について説明する。従来のシステム構成及びプログラム構成では、管理者がデバイスに対して性能値の閾値を設定する際、交代パスソフト 14300 によって設定された各経路の I/O 量の比を意識せずに閾値を設定してしまう恐れがある。その結果、設定された閾値が適切に効果を発揮しなくなるという問題点が発生する。

【0115】

例えば、図6に示す構成において、VOL111を通る経路に I/O 割り振り比率が 70%、VOL121を通る経路に I/O 割り振り比率が 20%、VOL131を経由する経路に I/O 割り振り比率が 10%と設定されている時、VOL111、VOL121、VOL131にそれぞれ閾値 100 を設定したとする。このとき、VOL121にはVOL111の7分の2、VOL131にはVOL111の7分の1しか I/O が割り振られない。このため、VOL111に閾値 100 が設定されている以上、VOL121及びVOL131を経由する I/O 量は閾値を超えない。すなわち、VOL121及びVOL131に設定された閾値 100 は効果を発揮しないこととなる。

40

【0116】

図11は、本発明の第1の実施の形態の管理サーバ 30000 が保持するストレージク

50

ラスタ管理表 3 3 4 0 0 の説明図である。

【 0 1 1 7 】

ストレージクラスタ管理表 3 3 4 0 0 は、フィールド 3 3 4 1 0 から 3 3 4 7 0 までの 7 フィールドからなる。

【 0 1 1 8 】

フィールド 3 3 4 1 0 には、ホストコンピュータ 1 0 0 0 0 内で各仮想デバイス 1 6 0 0 0 を識別する仮想デバイス ID が登録される。

【 0 1 1 9 】

フィールド 3 3 4 2 0 には、ホストコンピュータ 1 0 0 0 0 内で各仮想デバイス 1 6 0 0 0 に対応するデバイスファイル 1 5 0 0 0 を識別するデバイスファイル ID が登録される。

10

【 0 1 2 0 】

フィールド 3 3 4 3 0 には、各デバイスファイル 1 5 0 0 0 に対応する仮想ボリューム 2 4 0 0 0 を含む仮想ストレージ装置 2 0 0 0 0 の識別子である装置 ID が登録される。

【 0 1 2 1 】

フィールド 3 3 4 4 0 には、各デバイスファイル 1 5 0 0 0 に対応する仮想ボリューム 2 4 0 0 0 の識別子である仮想ボリューム ID が登録される。

【 0 1 2 2 】

フィールド 3 3 4 5 0 には、各仮想ボリューム 2 4 0 0 0 に対応する論理ボリューム 2 9 1 0 0 を含むストレージ装置 2 5 0 0 0 の識別子である装置 ID が登録される。

20

【 0 1 2 3 】

フィールド 3 3 4 6 0 には、各仮想ボリューム 2 4 0 0 0 に対応する論理ボリューム 2 9 1 0 0 の識別子である論理ボリューム ID が登録される。

【 0 1 2 4 】

フィールド 3 3 4 7 0 には、ホストコンピュータ 1 0 0 0 0 から論理ボリューム 2 9 1 0 0 に至る経路が主経路 (M a s t e r ルート) であるか副経路 (S l a v e ルート) であるかを示す値が登録される。

【 0 1 2 5 】

なお、ここでいう主経路 (M a s t e r ルート) とは、ホストコンピュータ 1 0 0 0 0 上の仮想デバイス 1 6 0 0 0 から、1 台の仮想ストレージ装置 2 0 0 0 0 を経由してストレージ装置 2 5 0 0 0 上の論理ボリューム 2 9 1 0 0 に至る情報伝達経路を示す。一方、副経路 (S l a v e ルート) とは、ホストコンピュータ 1 0 0 0 0 上の仮想デバイス 1 6 0 0 0 から、2 台の仮想ストレージ装置 2 0 0 0 0 を経由してストレージ装置 2 5 0 0 0 上の論理ボリューム 2 9 1 0 0 に至る情報伝達経路を示す。

30

【 0 1 2 6 】

図 1 1 は、管理サーバ 3 0 0 0 0 が保持するストレージクラスタ管理表 3 3 4 0 0 に登録された具体的な値の一例を示している。この例では、ホストコンピュータ 1 0 0 0 0 上の交代パスソフト 1 4 3 0 0 は、仮想デバイス D E V 1 への I / O を、デバイスファイル D E V 1 1 1、D E V 1 2 1 及び D E V 1 3 1 へ割り振っている。

【 0 1 2 7 】

デバイスファイル D E V 1 1 1 は、仮想ストレージ装置 S Y S 1 上の仮想ボリューム V O L 1 1 1 を介して、ストレージ装置 S Y S 4 上の論理ボリューム V O L 1 に接続されている。また、この経路は、M a s t e r ルートである。

40

【 0 1 2 8 】

同様に、デバイスファイル D E V 1 2 1 は、仮想ストレージ装置 S Y S 2 上の仮想ボリューム V O L 1 2 1 及び仮想ストレージ装置 S Y S 1 上の仮想ボリューム V O L 1 1 1 を介して、ストレージ装置 S Y S 4 上の論理ボリューム V O L 1 に接続されている。この経路は S l a v e ルートである。デバイスファイル D E V 1 3 1 は、仮想ストレージ装置装置 S Y S 3 上の仮想ボリューム V O L 1 3 1 及び仮想ストレージ装置装置 S Y S 1 上の仮想ボリューム V O L 1 1 1 を介して、ストレージ装置装置 S Y S 4 上の論理ボリューム V

50

OL1に接続されている。この経路はSlaveルートである。

【0129】

図12は、本発明の第1の実施の形態の管理サーバ30000が保持するデバイスグループ管理表33500の説明図である。

【0130】

デバイスグループ管理表33500には、計算機システム内に存在するデバイスグループを構成するデバイスを示す情報が登録される。

【0131】

ここで、デバイスグループについて説明する。既に説明したように、本実施の形態においては、一つの仮想デバイス16000から一つの論理ボリューム29100に至る複数の経路が存在する。それらの各経路に属するデバイス（すなわち、各経路が経由するデバイス）からなるグループが、本実施の形態におけるデバイスグループである。例えば、図6及び図12の例では、DEV1からVOL1に至る三つの経路が存在する。それらの三つの経路は、それぞれ、VOL111、VOL121及びVOL131を經由する。この場合、VOL111、VOL121及びVOL131がデバイスグループを構成する。

10

【0132】

デバイスグループ管理表33500は、フィールド33510及び33520の2フィールドからなる。フィールド33510には、管理サーバ30000によって管理されるデバイスグループの識別子であるグループIDが登録される。フィールド33520には、ストレージクラスタ管理表33400において同じ仮想デバイス24000から同じ論理ボリューム29100へ至る複数の経路を構成するデバイスのうち、グループを構成するデバイスの識別子が登録される。なお、フィールド33520には、登録されたデバイスが主経路を構成するデバイスである場合、その旨も登録される。

20

【0133】

図12には、管理サーバ30000が保持するデバイスグループ管理表33500に登録された具体的な値の一例を示している。この例では、ストレージクラスタ管理表33400に登録された仮想ボリュームVOL111、VOL121及びVOL131は、同じ仮想デバイス16000（この例では、DEV1）から同じ論理ボリューム29100（この例では、VOL1）へ至る経路群においてグループを構成している。そして、これらのデバイスは、G3というIDを持つデバイスグループとして登録されている。

30

【0134】

なお、管理サーバ30000が保持する交代パス管理表33600の内容は、ホストコンピュータ10000が保持する交代パス管理表14400と同じであるため、説明を省略する。

【0135】

図13は、本発明の第1の実施の形態の構成管理プログラム33110が実行する閾値設定処理のフローチャートである。

【0136】

ストレージ管理者は、管理サーバ30000が提供するアラート実行閾値設定画面71000を用いて、各デバイスに対して性能管理用の閾値を設定する（ステップ61000）。

40

【0137】

構成管理プログラム33110は、閾値の入力を受けると、入力された閾値を装置性能管理表33300に書き込む（ステップ61010）。具体的には、設定対象ID71010に入力された値と同一のデバイスID33320に対応するアラート実行閾値33340に、閾値71020に入力された値を登録する。

【0138】

次に、構成管理プログラム33110は、閾値の設定対象デバイスがボリューム（仮想ボリューム24000又は論理ボリューム29100）又はデバイスファイル15000であるか否かを判定する（ステップ61020）。

50

【 0 1 3 9 】

ステップ 6 1 0 2 0 において、設定対象デバイスがボリューム又はデバイスファイル 1 5 0 0 0 のいずれでもない判定された場合、設定対象デバイスは、例えばポート 2 1 0 0 0 等である。この場合、交代パス管理表 3 3 6 0 0 に定義された I / O 割り振り比率のみに基づいて閾値を算出することができない。このため、構成管理プログラム 3 3 1 1 0 は処理を終了する。

【 0 1 4 0 】

一方、ステップ 6 1 0 2 0 において、設定対象デバイスがボリューム又はデバイスファイル 1 5 0 0 0 のいずれかであると判定された場合、構成管理プログラム 3 3 1 1 0 は、次にデバイスグループ管理表 3 3 5 0 0 を参照し、設定対象デバイスと同一のデバイスグループに属するデバイスが存在するか否かを調べる（ステップ 6 1 0 3 0 ）。以下、図 1 4 の説明において、設定対象デバイスと同一のデバイスグループに属するデバイスを、「同一グループのデバイス」と記載する。

10

【 0 1 4 1 】

次に、ステップ 6 1 0 3 0 の調査の結果に基づいて、処理が分岐する（ステップ 6 1 0 4 0 ）。

【 0 1 4 2 】

同一グループのデバイスが存在しない場合（ステップ 6 1 0 4 0 ）、I / O 割り振り比率に基づいて設定対象デバイス以外のデバイスの閾値を設定する必要はない。このため、構成管理プログラム 3 3 1 1 0 は処理を終了する。

20

【 0 1 4 3 】

一方、同一グループのデバイスが存在する場合（ステップ 6 1 0 4 0 ）、構成管理プログラム 3 3 1 1 0 は、装置性能管理表 3 3 3 0 0 を参照し、同一グループのデバイスに対して閾値が設定されているか否かを調べる（ステップ 6 1 0 5 0 ）。具体的には、同一グループのデバイスのデバイス ID 3 3 3 2 0 に対応するアラート実行閾値 3 3 3 4 0 に値が登録されている場合、そのデバイスに対して閾値が設定されている。

【 0 1 4 4 】

次に、ステップ 6 1 0 5 0 の結果に基づいて、処理が分岐する（ステップ 6 1 0 6 0 ）。

【 0 1 4 5 】

同一グループのデバイスに閾値が設定されている場合（ステップ 6 1 0 6 0 ）、そのデバイスに対してもう一度閾値を設定する必要はない。このため、構成管理プログラム 3 3 1 1 0 は、同一グループのデバイスに対して閾値を設定せずに処理を終了する。

30

【 0 1 4 6 】

一方、同一グループのデバイスに閾値が設定されていない場合（ステップ 6 1 0 6 0 ）、構成管理プログラム 3 3 1 1 0 は、交代パス管理表 3 3 5 0 0 を参照し、同一グループのデバイスに対して設定されるべき閾値を算出する（ステップ 6 1 0 7 0 ）。

【 0 1 4 7 】

次に、構成管理プログラム 3 3 1 1 0 は、算出した閾値を装置性能管理表 3 3 3 0 0 のアラート実行閾値 3 3 3 4 0 に書き込む（ステップ 6 1 0 8 0 ）。

40

【 0 1 4 8 】

ステップ 6 1 0 7 0 における閾値の算出方法の例を具体的に示す。図 6 に示す構成において、DEV 1 1 1、DEV 1 2 1 及び DEV 1 3 1 に対する I / O 割り振り比率が、それぞれ、70%、20% 及び 10% である場合（図 7 参照）、VOL 1 1 1、VOL 1 2 1 及び VOL 1 3 1 を経由する I / O 量の比率は、100 対 20 対 10 となると予想される。DEV 1 2 1 及び DEV 1 3 1 に対する I / O も、VOL 1 1 1 を経由し、結局、VOL 1 に至る全ての I / O が VOL 1 1 1 を経由するためである。このような場合において、VOL 1 1 1 に閾値 100 が設定された場合、VOL 1 2 1 及び VOL 1 3 1 には、予想される I / O の比率に応じて、それぞれ閾値 20 及び閾値 10 が設定される。

【 0 1 4 9 】

50

このように、各デバイスを経由すると予想されるI/O量の比率に比例した閾値がそれらのデバイスに設定される。

【0150】

その後、各デバイスを経由する実際のI/O量が、そのデバイスに設定された閾値を超えたとき(具体的には、単位時間当たりI/O量33330に登録された値がアラート実行閾値33340に登録された値を超えたとき)、構成管理プログラム33110は、警告を出力部34000に表示する。

【0151】

なお、上述した閾値算出方式は本発明の実現手段を限定するものではない。閾値は、上記以外の方法によって算出されてもよい。ただし、各デバイスを経由するI/O量が異なると予測される場合、高いI/O量が予測されるデバイスの閾値は、低いI/O量が予測されるデバイスの閾値より高くなるように設定される。また、管理者が、高いI/O量が予測されるデバイスに、低いI/O量が予測されるデバイスの閾値より低い閾値を設定しようとした場合、管理サーバ30000は、アラート実行閾値設定画面71000に警告を表示してもよい。

【0152】

以上に、本発明の第1の実施の形態における閾値設定機能について述べたが、ストレージ装置及びホストコンピュータ10000の構成は図14に示す通りであってもよい。その場合の計算機システム構成について述べる。なお、図6の構成と相違する点についてのみ述べる。

【0153】

図14は、本発明の第1の実施の形態の変形例のストレージクラスタ構成の説明図である。

【0154】

図14には、ホストコンピュータ10000と、3台の仮想ストレージ装置20000と、ストレージ装置25000とを組み合わせたストレージクラスタ構成を示す。SYS4は、ストレージエリアネットワーク40000を介してSYS1、SYS2及びSYS3と接続されている。

【0155】

SYS4上の論理ボリュームVOL1は、SYS1上の仮想ボリュームVOL111、SYS2上の仮想ボリュームVOL121、及び、SYS3上の仮想ボリュームVOL131に対応付けられている。一方、SYS4上の論理ボリュームVOL2は、SYS1上の仮想ボリュームVOL112、SYS2上の仮想ボリュームVOL122、及び、SYS3上の仮想ボリュームVOL132に対応付けられている。すなわち、ホストコンピュータ10000からVOL111、VOL121、VOL131に対して行われたI/Oが、それぞれSYS1、SYS2、SYS3を介してSYS4上のVOL1に対して行われる。

【0156】

この場合、VOL111、VOL121及びVOL131が一つのデバイスグループを構成する。一方、VOL112、VOL122及びVOL132が、別のデバイスグループを構成する。

【0157】

図14に示す構成では、仮想デバイス16000から論理ボリューム29100に至る各経路が、一つの仮想ストレージ装置21000のみを経由する。言い換えると、二つの仮想ストレージ装置21000を経由する経路は存在しない。このため、図14に示す構成では、図6に示す構成と異なり、各経路は、主経路又は副経路に分類されない。

【0158】

しかし、例えば、各仮想ストレージ装置21000の性能に相違がある場合がある。例えば、SYS1の性能が、SYS2及びSYS3の性能と比較して高い場合がある。その場合、SYS1上の仮想ボリューム24000を経由する経路のI/O割り振り比率を、

10

20

30

40

50

それ以外の経路の比率より高く設定してもよい。この場合、SYS 1上の仮想ボリューム24000を經由する経路を、主経路として扱うことができる。

【0159】

あるいは、各仮想ストレージ装置21000が仮想ボリューム24000の管理以外の処理も実行している場合、それらの処理の負荷に相違がある場合がある。例えば、SYS 1の負荷が、SYS 2及びSYS 3の性能と比較して低い場合がある。その場合、SYS 1上の仮想ボリューム24000を經由する経路のI/O割り振り比率を、それ以外の経路の比率より高く設定してもよい。この場合、SYS 1上の仮想ボリューム24000を經由する経路を、主経路として扱うことができる。

【0160】

このように、図14に示す構成においても、図13に示す処理を実行することによって、各デバイスの閾値を設定することができる。

【0161】

以上、本発明の第1の実施の形態によれば、管理ソフトウェアを用いる管理者は、管理対象のデバイスに対して、交代パス管理表33600に設定されたI/O割り振り比率に応じて算出された閾値を設定することが可能となる。そのため、ホストコンピュータ10000上で、仮想ストレージ装置20000内の仮想ボリューム24000に対するI/O割り振り比率が定まっている場合も、I/O割り振り比率に応じた適切な閾値を仮想ボリューム24000に対して設定することができる。

【0162】

次に、本発明の第2の実施の形態について説明する。

【0163】

第2の実施の形態では、ストレージ管理ソフトがホストコンピュータ10000上の交代パスソフトからデータ割り振り比率をあらかじめ取得する。そして、複数の仮想ストレージ装置間でストレージ処理能力と入出力負荷とのバランス再調整が行われた結果、機器構成が変化したときに、データ割り振り比率と、関連するデバイスの閾値とが算出され、設定される。

【0164】

第2の実施の形態の計算機システムの構成、及び、各装置が保持する管理情報は、第1の実施の形態と同様である。以下、第2の実施の形態が第1の実施の形態と異なる部分について説明する。

【0165】

仮想ストレージ装置20000は、負荷集中によって特定の仮想ストレージ装置20000の入出力処理能力低下が発生したとき、複数の仮想ストレージ装置20000間でストレージ処理能力と入出力負荷とのバランスを再調整する機能を備える。ここでは例として、図6において仮想ストレージ装置SYS 1の負荷が高まった場合を想定する。

【0166】

図15は、本発明の第2の実施の形態において、負荷バランスが再調整された後の機器構成を示す説明図である。

【0167】

SYS 4は、ストレージエリアネットワーク40000を介してSYS 3と接続されている。SYS 4上の論理ボリュームVOL 1及びVOL 2は、それぞれ、SYS 3上の仮想ボリュームVOL 131及びVOL 132に対応付けられている。一方、SYS 3は、ストレージエリアネットワーク40000を介してSYS 2と接続されている。SYS 3上の仮想ボリュームVOL 131及びVOL 132は、それぞれ、SYS 2上の仮想ボリュームVOL 121及びVOL 122に対応付けられている。SYS 3は、さらに、ストレージエリアネットワーク40000を介してSYS 1と接続されている。SYS 3上の仮想ボリュームVOL 131及びVOL 132は、それぞれ、SYS 1上の仮想ボリュームVOL 111及びVOL 112に対応付けられている。

【0168】

10

20

30

40

50

図6に示す構成においては、仮想デバイスDEV1及びDEV2から論理ボリュームVOL1及びVOL2に至るI/Oが全てSYS1を經由していた。その結果、SYS1に負荷が集中していた。しかし、図15に示す構成では、SYS1の役割をSYS3に行わせることによって、SYS1の負荷を軽減することができる。

【0169】

しかし、図15に示すように機器構成が変化した場合、管理者が交代パスソフトによって設定した各経路のI/O割り振り比率、及び、デバイスに対して設定した閾値等が、機器構成と整合しなくなる恐れがある。その結果、設定した閾値及びI/O割り振り比率等が適切に効果を発揮しなくなるという問題が発生する。第2の実施の形態では、このように機器構成が変化した場合、変化した後の機器構成に整合するように、閾値及びI/O割り振り比率等が再設定される。

10

【0170】

図16は、本発明の第2の実施の形態の構成管理プログラム33110が実行する閾値見直し処理のフローチャートである。

【0171】

管理サーバ30000は、管理サーバ30000によって管理されるストレージ装置25000、仮想ストレージ装置20000及びホストコンピュータ10000から構成情報を定期的に取得する(ステップ62000)。その際、構成管理プログラム33110は、取得した構成情報を装置構成管理表33200の内容と比較し、先述した負荷バランス再調整によって機器構成が変更されたか否かを判定する(ステップ62010)。

20

【0172】

ステップ62010において、構成変更がないと判定された場合、閾値を見直す必要がないため、構成管理プログラム33110は処理を終了する。

【0173】

一方、ステップ62010において、構成が変更されたと判定された場合、構成管理プログラム33110は、装置構成管理表33200、ストレージクラスタ管理表33400及びデバイスグループ管理表33500の内容を、ステップ62000にて取得した構成情報と整合するように更新する(ステップ62020)。その際、構成管理プログラム33110は、構成情報に基づいて、仮想デバイス16000から論理ボリューム29100へ至る経路のうち、どの経路が主経路であるかを判定する。

30

【0174】

例えば、第1の実施の形態と同様、仮想デバイス16000から論理ボリューム29100へ至るまでに經由する仮想ストレージ装置20000(又は仮想ボリューム24000)の数が少ない経路が主経路であると判定されてもよい。

【0175】

主経路と副経路が変更された場合、構成管理プログラム33110は、交代パス管理表33600のI/O割り振り比率14460を再計算する。再計算の方法の一つとしては、主経路と副経路のI/O割り振り比率を入れ替える方法が考えられる。具体的には、構成が変更される前の主経路に設定されていたI/O割り振り比率を、構成が変更された後の主経路に設定し、構成が変更される前の副経路に設定されていたI/O割り振り比率を、構成が変更された後の副経路に設定してもよい。

40

【0176】

そして、構成管理プログラム33110は、更新されたI/O割り振り比率をホストコンピュータ10000上の交代パスソフト14300に送信する(ステップ62030)。交代パスソフト14300は、受信したI/O割り振り比率を交代パス管理表14400のI/O割り振り比率14460に登録することによって、交代パス管理表14400を更新する。

【0177】

ステップ62030におけるI/O割り振り比率は、論理ボリューム29100と仮想ボリューム24000との対応付けに基づいて、上記以外の方法によって定められてもよ

50

い。例えば、主経路と副経路に対して設定される I / O 割り振り比率が予め管理者によって定められていてもよい。ただし、副経路が経由する仮想ボリュームの数が、主経路が経由する仮想ボリュームの数より多い場合、副経路には、主経路に設定されるものより低い I / O 割り振り比率が設定される。

【 0 1 7 8 】

次に、構成管理プログラム 3 3 1 1 0 は、装置性能管理表 3 3 3 0 0 を参照し、I / O 割り振り比率が変更された経路上のデバイスに対して閾値が設定されているか否かを調べる (ステップ 6 2 0 4 0)。

【 0 1 7 9 】

次に、ステップ 6 2 0 4 0 の結果に基づいて、処理が分岐する (ステップ 6 2 0 5 0)

10

【 0 1 8 0 】

I / O 割り振り比率が変更された経路上のデバイスに閾値が設定されていない場合 (ステップ 6 2 0 5 0)、構成管理プログラム 3 3 1 1 0 は、そのデバイスに対して閾値を設定せずに処理を終了する。

【 0 1 8 1 】

一方、I / O 割り振り比率が変更された経路上のデバイスに閾値が設定されている場合 (ステップ 6 2 0 5 0)、構成管理プログラム 3 3 1 1 0 は、交代パス管理表 3 3 5 0 0 を参照し、I / O 割り振り比率が変更された経路上のデバイスに対して設定されるべき閾値を算出する (ステップ 6 2 0 6 0)。

20

【 0 1 8 2 】

次に、構成管理プログラム 3 3 1 1 0 は、算出した閾値を装置性能管理表 3 3 3 0 0 のアラート実行閾値 3 3 3 4 0 に登録する (ステップ 6 2 0 7 0)。

【 0 1 8 3 】

ステップ 6 2 0 6 0 における閾値の算出方法の例を具体的に示す。図 6 に示す構成において、DEV 1 1 1、DEV 1 2 1 及び DEV 1 3 1 に対する I / O 割り振り比率が、それぞれ、70%、20% 及び 10% と設定され、かつ、VOL 1 1 1、VOL 1 2 1 及び VOL 1 3 1 の閾値として、それぞれ、100、20 及び 10 が設定されていたとする。

【 0 1 8 4 】

その後、計算機システムの構成が、図 1 5 に示すように変更されたとき、DEV 1 1 1、DEV 1 2 1 及び DEV 1 3 1 に対する I / O 割り振り比率が、それぞれ、10%、20% 及び 70% に再設定される。このとき、I / O 割り振り比率に合わせて、VOL 1 1 1 の閾値は 10、VOL 1 3 1 の閾値は 100 に再設定される。このとき、構成管理プログラム 3 3 1 1 0 は、VOL 1 1 1 を含む仮想ストレージ装置 SYS 1 上のポート 2 1 0 0 0 のうち、VOL 1 1 1 への I / O が経由するポート PORT 1 に設定された閾値を、VOL 1 1 1 の閾値の減少分 90 だけ減算する。さらに、構成管理プログラム 3 3 1 1 0 は、VOL 1 3 1 への I / O が経由するポート PORT 3 に設定された閾値を、VOL 1 3 1 の閾値の増加分 90 だけ加算する。

30

【 0 1 8 5 】

なお、上述した閾値算出方式は本発明の実現手段を限定するものではない。閾値は、上記以外の方法によって算出されてもよい。

40

【 0 1 8 6 】

以上が、本実施の形態における閾値見直し処理である。

【 0 1 8 7 】

以上、本発明の第 2 の実施の形態によれば、複数の仮想ストレージ装置 2 0 0 0 0 間でストレージ処理能力と入出力負荷とのバランス再調整が行われた結果、機器構成が変更された場合も、変更後の機器構成に整合するように各デバイスの閾値を算出し、設定することが可能となる。

【 0 1 8 8 】

次に、本発明の第 3 の実施の形態について説明する。

50

【0189】

第3の実施の形態では、ストレージ管理者が機器構成変更後のデータ割り振り比率をストレージ管理ソフトにあらかじめ登録しておく。そして、装置障害等によって、特定の仮想ストレージ装置20000の入出力処理が停止した結果、機器構成が変化したときに、データ割り振り比率と、関連するデバイスの閾値とが算出され、設定される。

【0190】

第3の実施の形態の計算機システムの構成、及び、各装置が保持する管理情報は、以下に説明するものを除き、第1の実施の形態と同様である。以下、第3の実施の形態が第1の実施の形態と異なる部分について説明する。

【0191】

仮想ストレージ装置20000は、いずれかの仮想ストレージ装置20000においてデータ入出力の停止が発生したとき、複数の仮想ストレージ装置20000間で装置構成を再調整する機能を備える。このようなデータ入出力の停止は、例えば、障害発生を原因として仮想ストレージ装置20000がダウンしたときに発生する。ここでは例として、図6において仮想ストレージ装置SYS1がダウンした場合を想定する。

【0192】

図17は、本発明の第3の実施の形態において、負荷バランスが再調整された後の機器構成を示す説明図である。

【0193】

SYS4は、ストレージエリアネットワーク40000を介してSYS3と接続されている。SYS4上の論理ボリュームVOL1及びVOL2は、それぞれ、SYS3上の仮想ボリュームVOL131及びVOL132に対応付けられている。一方、SYS3は、ストレージエリアネットワーク40000を介してSYS2と接続されている。SYS3上の仮想ボリュームVOL131及びVOL132は、それぞれ、SYS2上の仮想ボリュームVOL121及びVOL122に対応付けられている。

【0194】

しかし、図17に示すように機器構成が変化した結果、管理者が交代パスソフトによって設定した各経路のI/O割り振り比率、及び、デバイスに対して設定した閾値等が、機器構成と整合しなくなる恐れがある。その結果、設定した閾値及びI/O割り振り比率等が適切に効果を発揮しなくなるという問題が発生する。第3の実施の形態は、このように機器構成が変化したときに、変化した後の機器構成に整合するように、閾値及びI/O割り振り比率等が再設定される。

【0195】

次に、第3の実施の形態の計算機システムの構成について説明する。図18は管理サーバ30000の構成を示し、図19は管理サーバ30000に保持される管理情報を示す。なお、ホストコンピュータ10000、仮想ストレージ装置20000及びストレージ装置25000の構成は、第1の実施の形態と同様である(図2から図4参照)。

【0196】

図18は、本発明の第3の実施の形態の管理サーバ30000の構成を示すブロック図である。

【0197】

第3の実施の形態の管理サーバ30000の構成は、メモリ33000に図19に示すI/O比率管理表33700が追加され、図20で示す閾値見直し処理が加わった構成管理プログラム33120を備えることを除き、第1の実施の形態の管理サーバ30000(図5参照)と同じである。

【0198】

図19は、本発明の第3の実施の形態の管理サーバ30000が保持するI/O比率管理表33700の説明図である。

【0199】

I/O比率管理表33700は、フィールド33710から33730までの3フィー

10

20

30

40

50

ルドからなる。フィールド 33710 には、ストレージクラスタ構成を実現するために計算機システムが備える仮想ストレージ装置 20000 の台数が登録される。フィールド 33720 には、フィールド 33710 に登録された台数の仮想ストレージ装置 20000 によって構成されるストレージクラスタ構成において、主経路 (Master ルート) に対して設定されるべき I/O 割り振り比率が登録される。フィールド 33730 には、フィールド 33710 に登録された台数の仮想ストレージ装置 20000 によって構成されるストレージクラスタ構成において、副経路 (Slave ルート) に対して設定されるべき I/O 割り振り比率が登録される。

【0200】

図 19 には、管理サーバ 30000 が保持する I/O 比率管理表 33700 に登録された具体的な値の一例を示している。この例は、ストレージクラスタ構成を組む仮想ストレージ装置の台数が 2 台である場合、主経路 (Master ルート) に対して設定されるべき I/O 割り振り比率は 90%、副経路 (Slave ルート) に対して設定されるべき I/O 割り振り比率は 10% であることを示している。

10

【0201】

図 20 は、本発明の第 3 の実施の形態の構成管理プログラム 33120 が実行する閾値見直し処理のフローチャートである。

【0202】

管理サーバ 30000 は、管理サーバ 30000 によって管理されるストレージ装置 25000、仮想ストレージ装置 20000 及びホストコンピュータ 10000 から構成情報を定期的に取得する (ステップ 63000)。その際、構成管理プログラム 33120 は、取得した構成情報を装置構成管理表 33200 の内容と比較し、先述した障害発生によって機器構成が変更されたか否かを判定する (ステップ 63010)。

20

【0203】

ステップ 63010 において、構成変更がないと判定された場合、閾値を見直す必要がないため、構成管理プログラム 33120 は処理を終了する。

【0204】

一方、ステップ 63010 において、構成が変更されたと判定された場合、構成管理プログラム 33120 は、装置構成管理表 33200、ストレージクラスタ管理表 33400 及びデバイスグループ管理表 33500 の内容を、ステップ 63000 にて取得した構成情報を整合するように更新する (ステップ 63020)。その際、構成管理プログラム 33120 は、構成情報に基づいて、仮想デバイス 16000 から論理ボリューム 29100 へ至る経路のうち、どの経路が主経路であるかを判定する。この判定は、図 16 のステップ 62020 と同様の方法によって実行されてもよい。

30

【0205】

主経路と副経路が変更された場合、構成管理プログラム 33120 は、I/O 比率管理表 33700 に登録された比率に基づいて、交代パス管理表 33600 の I/O 割り振り比率 14460 を再設定する。そして、構成管理プログラム 33120 は、更新された I/O 割り振り比率をホストコンピュータ 10000 上の交代パスソフト 14300 に送信する (ステップ 63030)。交代パスソフト 14300 は、受信した I/O 割り振り比率を交代パス管理表 14400 の I/O 割り振り比率 14460 に登録することによって、交代パス管理表 14400 を更新する。

40

【0206】

なお、図 16 のステップ 62030 と同様、I/O 割り振り比率は、上記以外の方法によって定められてもよい。

【0207】

次に、構成管理プログラム 33120 は、装置性能管理表 33300 を参照し、I/O 割り振り比率が変更された経路上のデバイスに対して閾値が設定されているか否かを調べる (ステップ 63040)。

【0208】

50

次に、ステップ63040の結果に基づいて、処理が分岐する(ステップ63050)。

【0209】

I/O割り振り比率が変更された経路上のデバイスに閾値が設定されていない場合(ステップ63050)、構成管理プログラム33120は、そのデバイスに対して閾値を設定せずに処理を終了する。

【0210】

一方、I/O割り振り比率が変更された経路上のデバイスに閾値が設定されている場合(ステップ63050)、構成管理プログラム33120は、交代パス管理表33500を参照し、I/O割り振り比率が変更された経路上のデバイスに対して設定されるべき閾値を算出する(ステップ63060)。

10

【0211】

次に、構成管理プログラム33120は、算出した閾値を装置性能管理表33300のアラート実行閾値33340に登録する(ステップ63070)。

【0212】

ステップ63060における閾値の算出方法の例を具体的に示す。図6に示す構成において、DEV111、DEV121及びDEV131に対するI/O割り振り比率が、それぞれ、70%、20%及び10%と設定され、かつ、VOL111、VOL121及びVOL131の閾値として、それぞれ、100、20及び10が設定されていたとする。

【0213】

20

その後、計算機システムの構成が、図17に示すように変更されたとき、I/O比率管理表33700に登録された比率に基づいて、DEV121及びDEV131に対するI/O割り振り比率が、それぞれ、10%及び90%に再設定される。このとき、I/O割り振り比率に合わせて、VOL111の閾値は0、VOL121の閾値は10、VOL131の閾値は100に再設定される。このとき、構成管理プログラム33120は、VOL111を含む仮想ストレージ装置SYS1上のポート21000のうち、VOL111へのI/Oが経由するポートPORT1に設定された閾値を、VOL111の閾値の減少分100だけ減算する。構成管理プログラム33120は、VOL121へのI/Oが経由するポートPORT2に設定された閾値を、VOL121の閾値の減少分10だけ減算する。さらに、構成管理プログラム33120は、VOL131へのI/Oが経由するポ

30

【0214】

なお、上述した閾値算出方式は本発明の実現手段を限定するものではない。閾値は、上記以外の方法によって算出されてもよい。

【0215】

以上が、本実施の形態における閾値見直し処理である。

【0216】

上記の第3の実施の形態は、いずれかの仮想ストレージ装置SYS1に障害が発生した場合を例として説明した。しかし、障害発生以外の理由によって機器の構成が変更された場合にも、本実施の形態を適用することができる。例えば、負荷バランスの再調整ために機器の構成が変更された場合にも、第3の実施の形態を適用することができる。

40

【0217】

以上、本発明の第3の実施の形態によれば、装置障害によって特定の仮想ストレージ装置20000の入出力処理が停止した結果、機器構成が変更された場合も、変更後の機器構成に整合するように各デバイスの閾値を算出し、設定することが可能となる。

【0218】

以上の本発明の第1から第3の実施の形態をまとめると、管理ソフトウェアを用いる管理者は、管理下のデバイスに対して、交代パスソフトに設定したI/O割り振り比率に応じて算出した閾値を設定することが可能となる。そのため、ホストコンピュータ上で仮想ストレージ装置内のデバイスに対するI/O割り振り比率が定まっている場合も、I/O

50

割り振り比率に応じた適切な閾値をデバイスに対し設定することができる。

【0219】

さらに、複数の仮想ストレージ装置間でストレージ処理能力と入出力負荷とのバランスが再調整された結果、主経路と副経路の入れ替えが発生したとき、変更された機器構成に応じたI/O割り振り比率を交代パスソフトに設定することが可能となる。このため、I/O割り振り比率に応じた適切な閾値をデバイスに対して設定することができる。

【0220】

また、装置障害によっていずれかの仮想ストレージ装置の入出力処理が停止した結果、機器構成が変更されたときも、変更された機器構成に応じたI/O割り振り比率を交代パスソフトに設定することが可能となる。このため、I/O割り振り比率に応じた適切な閾値をデバイスに対し設定することができる。

10

【図面の簡単な説明】

【0221】

【図1】本発明の第1の実施の形態の計算機システムの構成を示すブロック図である。

【図2】本発明の第1の実施の形態のホストコンピュータの詳細な構成を示すブロック図である。

【図3】本発明の第1の実施の形態の仮想ストレージ装置の詳細な構成例を示すブロック図である。

【図4】本発明の第1の実施の形態のストレージ装置の詳細な構成例を示すブロック図である。

20

【図5】本発明の第1の実施の形態の管理サーバの詳細な構成を示すブロック図である。

【図6】本発明の第1の実施の形態のストレージクラスタ構成の説明図である。

【図7】本発明の第1の実施の形態のホストコンピュータが保持する交代パス管理表の説明図である。

【図8】本発明の第1の実施の形態の仮想ストレージ装置が保持する仮想ボリューム管理表の説明図である。

【図9】本発明の第1の実施の形態の管理サーバが保持する装置性能管理表の説明図である。

【図10】本発明の第1の実施の形態において表示されるアラート実行閾値設定画面の説明図である。

30

【図11】本発明の第1の実施の形態の管理サーバが保持するストレージクラスタ管理表の説明図である。

【図12】本発明の第1の実施の形態の管理サーバが保持するデバイスグループ管理表の説明図である。

【図13】本発明の第1の実施の形態の構成管理プログラムが実行する閾値設定処理のフローチャートである。

【図14】本発明の第1の実施の形態の変形例のストレージクラスタ構成の説明図である。

【図15】本発明の第2の実施の形態において、負荷バランスが再調整された後の機器構成を示す説明図である。

40

【図16】本発明の第2の実施の形態の構成管理プログラムが実行する閾値見直し処理のフローチャートである。

【図17】本発明の第3の実施の形態において、負荷バランスが再調整された後の機器構成を示す説明図である。

【図18】本発明の第3の実施の形態の管理サーバの構成を示すブロック図である。

【図19】本発明の第3の実施の形態の管理サーバが保持するI/O比率管理表の説明図である。

【図20】本発明の第3の実施の形態の構成管理プログラム33120が実行する閾値見直し処理のフローチャートである。

【図21】本発明の第1の実施の形態の交代パスソフト14300が実行するI/O振り

50

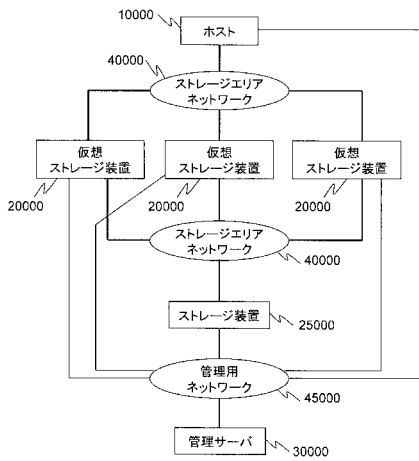
分け処理のフローチャートである。

【符号の説明】

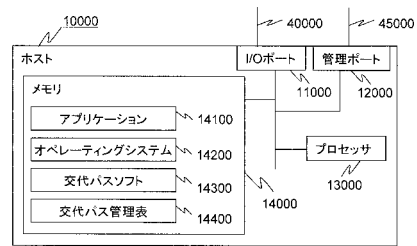
【0222】

- 10000 ホストコンピュータ
- 15000 デバイスファイル
- 16000 仮想デバイス
- 20000 仮想ストレージ装置
- 24000 仮想ボリューム
- 29100 論理ボリューム
- 25000 ストレージ装置
- 30000 管理サーバ
- 40000 ストレージエリアネットワーク
- 45000 管理用ネットワーク

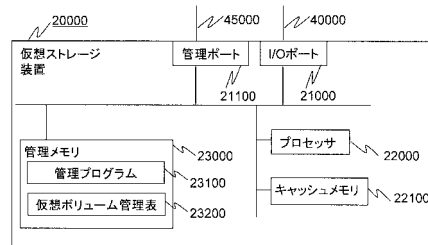
【図1】



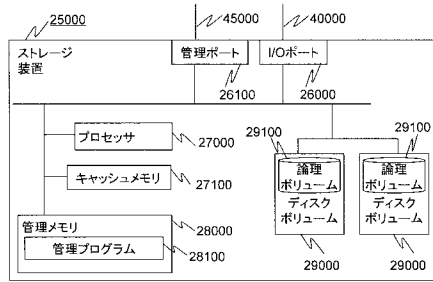
【図2】



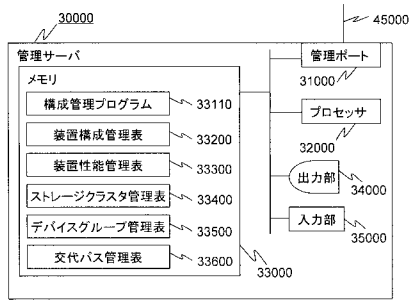
【図3】



【図4】



【図5】



【図7】

交代バス管理表

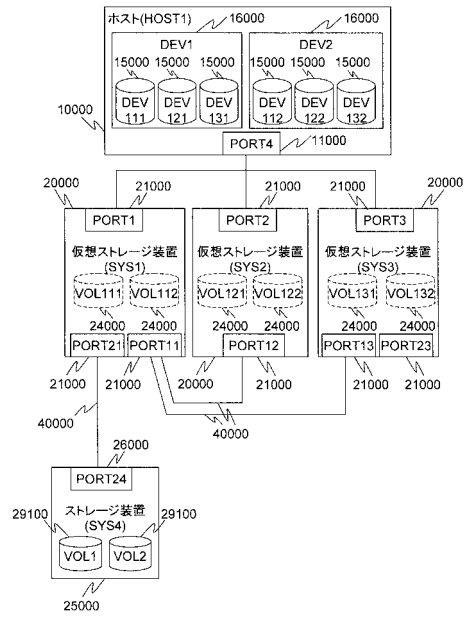
仮想デバイスID	デバイスファイルID	接続先装置ID	接続先ポートID	接続先ボリュームID	I/O割り振り比率
DEV1	DEV111	SYS1	PORT1	VOL111	70%
	DEV121	SYS2	PORT2	VOL121	20%
	DEV131	SYS3	PORT3	VOL131	10%
DEV2	DEV112	SYS1	PORT1	VOL112	80%
	DEV122	SYS2	PORT2	VOL122	10%
	DEV132	SYS3	PORT3	VOL132	10%
:	:	:	:	:	:

【図8】

仮想ボリューム管理表

仮想ボリュームID	ポートID	接続先装置ID	接続先ポートID	接続先ボリュームID
VOL111	PORT21	SYS4	PORT24	VOL1
VOL112	PORT21	SYS4	PORT24	VOL2
:	:	:	:	:

【図6】

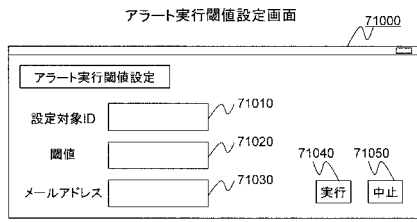


【図9】

装置性能管理表

装置ID	デバイスID	単位時間当りI/O量	アラート実行閾値
SYS1	VOL111	200	1000
SYS1	VOL112	40	200
SYS2	VOL121	20	100
SYS2	VOL122	240	800
SYS3	VOL131	30	100
SYS3	VOL132	30	100
HOST1	DEV111	140	-
HOST1	DEV112	40	-
HOST1	DEV121	20	-
HOST1	DEV122	240	-
HOST1	DEV131	30	-
HOST1	DEV132	30	-
SYS1	PORT1	300	1000
SYS2	PORT2	400	1000
SYS3	PORT3	360	1000
:	:	:	:

【図10】



【図12】

デバイスグループ管理表

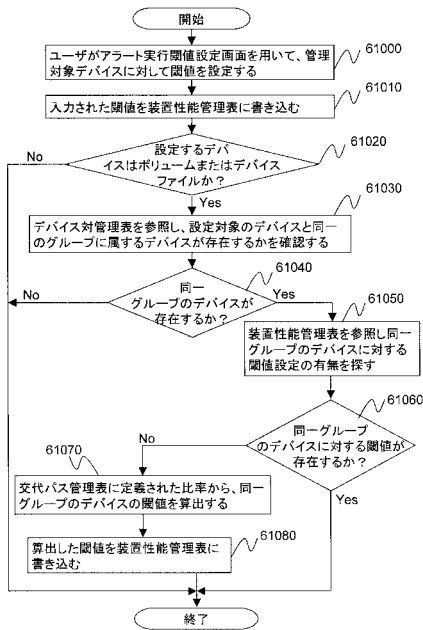
グループID	グループ構成要素
G1	DEV111(Master), DEV121, DEV131
G2	PORT1(Master), PORT2, PORT3
G3	VOL111(Master), VOL121, VOL131
G4	PORT12, PORT13
:	:

【図11】

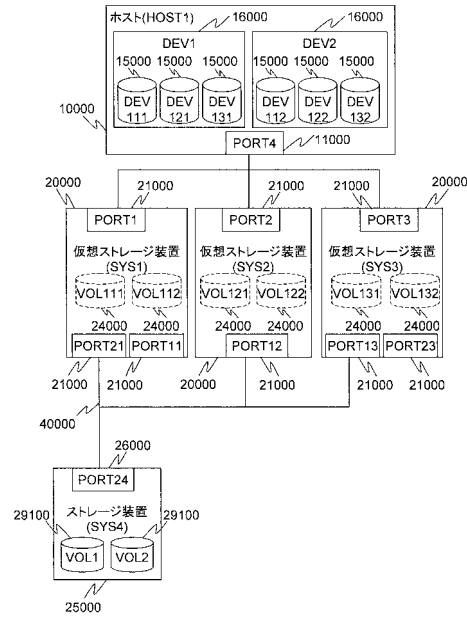
ストレージクラス管理表

仮想デバイスID	デバイスファイルID	仮想化装置ID	仮想ボリュームID	接続元装置ID	論理ボリュームID	Master/Slave種別
DEV1	DEV111	SYS1	VOL111	SYS4	VOL1	Master
	DEV121	SYS2, SYS1	VOL121, VOL111			Slave
	DEV131	SYS3, SYS1	VOL131, VOL111			Slave
DEV2	DEV112	SYS1	VOL112	SYS4	VOL2	Master
	DEV122	SYS2, SYS1	VOL122, VOL112			Slave
	DEV132	SYS3, SYS1	VOL132, VOL112			Slave
:	:	:	:	:	:	:

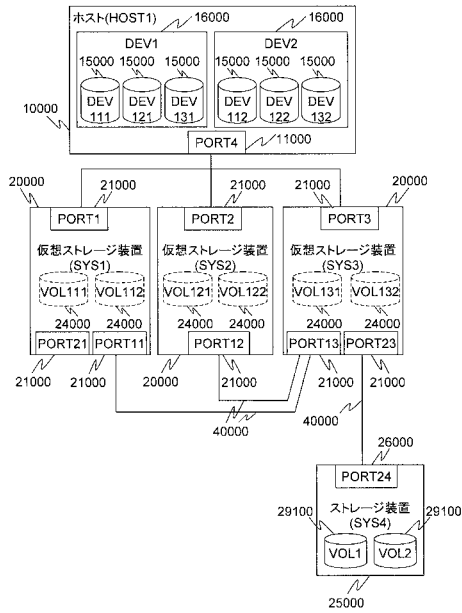
【図13】



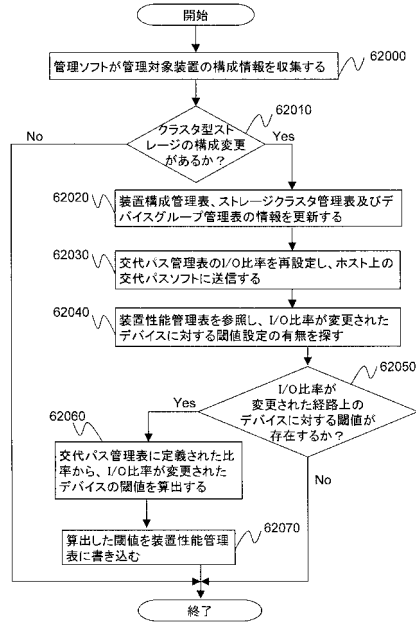
【図14】



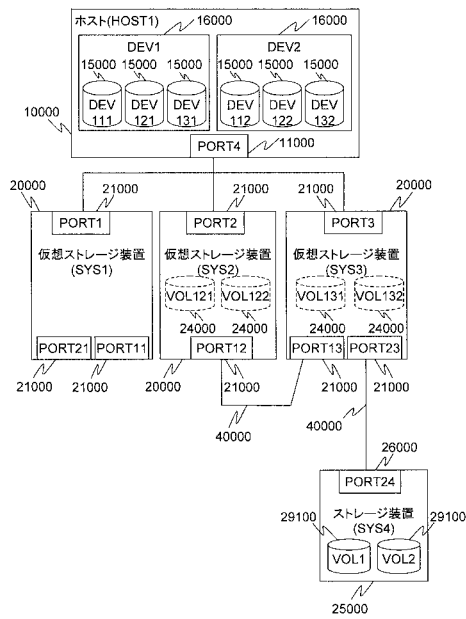
【図15】



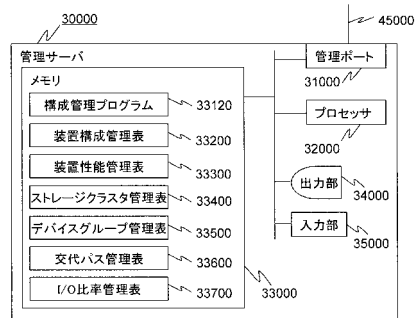
【図16】



【図17】



【図18】

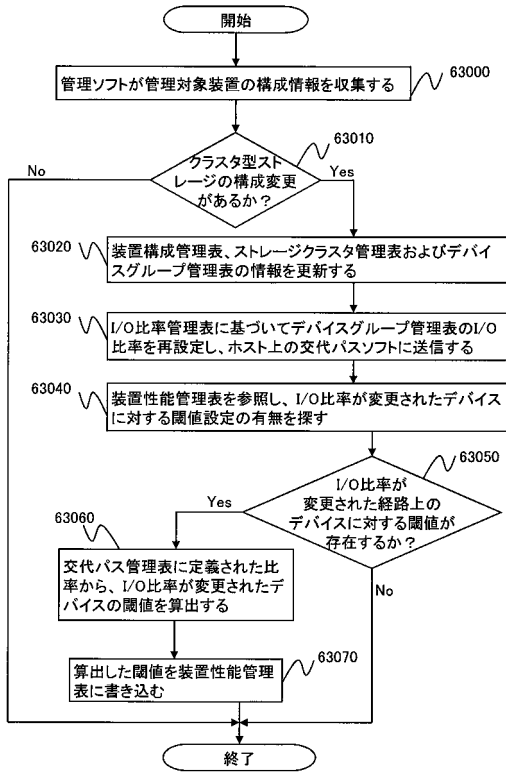


【図19】

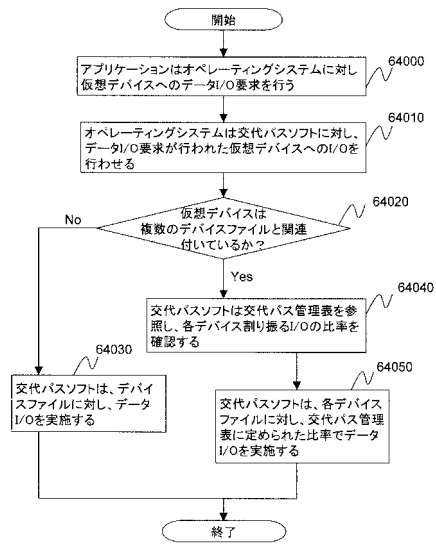
I/O比率管理表

装置台数	Master装置 I/O比率	Slave装置 I/O比率
2	90%	10%
3	80%	10%
4	70%	10%
⋮	⋮	⋮

【図20】



【図21】



フロントページの続き

(72)発明者 中島 淳

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 矢川 雄一

神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所 システム開発研究所内

審査官 坂東 博司

(56)参考文献 特開2004-086512(JP,A)

特開2004-302751(JP,A)

特開2005-242690(JP,A)

特開2004-334561(JP,A)

特開2004-206623(JP,A)

特開2006-154880(JP,A)

特開2006-092322(JP,A)

特開2003-316522(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06

G06F 13/14