

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6231685号
(P6231685)

(45) 発行日 平成29年11月15日(2017.11.15)

(24) 登録日 平成29年10月27日(2017.10.27)

(51) Int.Cl.		F I			
G06F 13/14	(2006.01)	G06F	13/14	310H	
G06F 13/10	(2006.01)	G06F	13/10	340A	
G06F 3/06	(2006.01)	G06F	3/06	301C	
		G06F	13/14	310K	

請求項の数 15 (全 35 頁)

(21) 出願番号	特願2016-534025 (P2016-534025)	(73) 特許権者	000005108
(86) (22) 出願日	平成26年7月16日 (2014.7.16)		株式会社日立製作所
(86) 国際出願番号	PCT/JP2014/068873		東京都千代田区丸の内一丁目6番6号
(87) 国際公開番号	W02016/009504	(74) 代理人	110000279
(87) 国際公開日	平成28年1月21日 (2016.1.21)		特許業務法人ウィルフォート国際特許事務所
審査請求日	平成28年7月22日 (2016.7.22)		所
		(72) 発明者	末次 通夫
			東京都千代田区丸の内一丁目6番6号 株
			株式会社日立製作所内
		(72) 発明者	川口 智大
			東京都千代田区丸の内一丁目6番6号 株
			株式会社日立製作所内
		(72) 発明者	斎藤 秀雄
			東京都千代田区丸の内一丁目6番6号 株
			株式会社日立製作所内

最終頁に続く

(54) 【発明の名称】 ストレージシステム及び通知制御方法

(57) 【特許請求の範囲】

【請求項1】

ホストシステムに接続されるストレージシステムであって、
複数の論理ボリュームを含む複数種類の複数のリソースを管理し、前記複数の論理ボリュームが1つに仮想化された論理ボリュームである仮想ボリュームを前記ホストシステムに提供する複数のストレージ装置を有し、

前記複数のストレージ装置のうちのいずれかのストレージ装置である第1ストレージ装置が、第1コマンドを受信した場合、第1リソースの状態を変更し、前記第1リソースは、前記仮想ボリュームの基である前記複数の論理ボリュームのうち前記第1ストレージ装置が有する論理ボリュームである第1論理ボリュームと前記第1論理ボリュームに関連し前記第1ストレージ装置が管理するリソースとのうちの少なくとも1つであり、

前記第1ストレージ装置が、前記仮想ボリュームを指定したRTPGコマンド (Report TargetPortGroupsコマンド) 及び前記仮想ボリュームを指定したI/Oコマンドのいずれかである第2コマンドを受信した場合、前記第1論理ボリュームに関わる状態変更の通知である状態変更通知を第2ストレージ装置に送信し、前記第2ストレージ装置は、第2論理ボリュームを有するストレージ装置であり、前記第2論理ボリュームは、前記仮想ボリュームの基であり前記第1論理ボリュームに関連付けられた論理ボリュームであり、

前記第2ストレージ装置は、前記状態変更通知を受信し、前記受信した状態変更通知に基づき状態変更を設定する、
ストレージシステム。

【請求項 2】

前記複数のリソースは、複数のポートを含み、

前記第 1 リソースは、前記第 1 論理ボリュームへのパスであり前記第 1 論理ボリュームが関連付けられたパスを経由する第 1 パスであり、

前記複数のストレージ装置の各々は、そのストレージ装置が管理するパスの状態を表す情報を保持しているが、そのストレージ装置以外のストレージ装置が管理するパスの状態を表す情報を保持しておらず、

前記第 2 ストレージ装置において、前記状態変更通知に基づく状態変更は、前記第 2 論理ボリュームに関して設定される、

請求項 1 記載のストレージシステム。

10

【請求項 3】

前記第 1 ストレージ装置は、前記第 1 コマンドを受信した場合、前記第 1 論理ボリュームに関して状態変更を設定し、

前記第 2 コマンドの受信は、前記仮想ボリュームを指定した I/O コマンドの受信であり、

前記第 2 ストレージ装置は、前記受信した状態変更通知に基づき状態変更を設定した場合、完了応答を前記第 1 ストレージ装置に送信し、

前記第 1 ストレージ装置は、前記完了応答を受信した場合、前記第 1 論理ボリュームに関して設定されている状態変更を含んだ、前記 I/O コマンドの応答を、前記ホストシステムに送信する、

20

請求項 2 記載のストレージシステム。

【請求項 4】

前記第 1 ストレージ装置は、

前記第 1 論理ボリュームに関して設定されている状態変更を含んだ前記応答を受信した前記ホストシステムから、前記仮想ボリュームを指定した R T P G コマンドを前記第 2 コマンドとして受信し、

前記第 1 論理ボリュームとペアを構成する前記第 2 論理ボリュームを有する前記第 2 ストレージ装置に、前記第 2 論理ボリュームに関連付けられたパスの状態の要求を送信し、前記要求を受信した前記第 2 ストレージ装置から、前記第 2 論理ボリュームに関連付けられた 1 以上のパスの状態を表す情報を含んだパス状態情報を受信し、

30

前記第 1 論理ボリュームに関連付けられ前記第 1 パスを含んだ 1 以上のパスの状態を表す第 1 パス状態と、前記第 2 論理ボリュームに関連付けられた 1 以上のパスの状態である第 2 パス状態とを含んだ、前記 R T P G コマンドの応答を、前記ホストシステムに送信する、

請求項 3 記載のストレージシステム。

【請求項 5】

前記第 1 ストレージ装置は、前記第 1 論理ボリュームと前記第 2 論理ボリュームについてのペア状態及び I/O 状態に応じて、前記 R T P G コマンドの応答に含まれる前記第 1 パス状態及び前記第 2 パス状態を制御する、

請求項 4 記載のストレージシステム。

40

【請求項 6】

前記第 1 パス状態は、前記第 1 論理ボリュームに関連付けられた前記 1 以上のパスの優先度である第 1 パス優先度を表し、前記第 2 パス状態は、前記第 2 論理ボリュームに関連付けられた前記 1 以上のパスの優先度である第 2 パス優先度を表し、

前記第 1 ストレージ装置は、

前記ペア状態及び前記 I/O 状態が、前記第 1 論理ボリュームと前記第 2 論理ボリュームが同期状態であることを表していれば、前記第 1 パス優先度を、前記第 1 ストレージ装置が保持し前記第 1 論理ボリュームに関連付けられた前記 1 以上のパスの優先度を表す情報が表す優先度とし、前記第 2 パス優先度を、前記受信したパス状態情報が表す優先度とし、

50

前記ペア状態及び前記 I / O 状態が、前記第 1 論理ボリュームと前記第 2 論理ボリュームがサスペンド状態であり前記第 1 論理ボリューム内のデータよりも前記第 2 論理ボリューム内のデータの方が新しいことを表していれば、前記第 1 パス優先度よりも前記第 2 パス優先度を高くし、

前記ペア状態及び前記 I / O 状態が、前記第 1 論理ボリュームと前記第 2 論理ボリュームがサスペンド状態であり前記第 2 論理ボリューム内のデータよりも前記第 1 論理ボリューム内のデータの方が新しいことを表していれば、前記第 2 パス優先度よりも前記第 1 パス優先度を高くする、

請求項 5 記載のストレージシステム。

【請求項 7】

前記第 1 ストレージ装置は、前記第 2 ストレージ装置との通信が不可能である場合、

前記ペア状態及び前記 I / O 状態が、前記第 1 論理ボリュームと前記第 2 論理ボリュームがサスペンド状態であり前記第 1 論理ボリューム内のデータよりも前記第 2 論理ボリューム内のデータの方が新しいことを表していれば、エラー応答を前記ホストシステムに送信し、

前記ペア状態及び前記 I / O 状態が、前記第 1 論理ボリュームと前記第 2 論理ボリュームが同期状態であること、又は、前記第 1 論理ボリュームと前記第 2 論理ボリュームがサスペンド状態であり前記第 2 論理ボリューム内のデータよりも前記第 1 論理ボリューム内のデータの方が新しいことを表していれば、前記第 1 パス優先度を前記第 2 パス優先度よりも高くすることを含んだ、前記 R T P G コマンド の応答を、前記ホストシステムに送信する、

請求項 6 記載のストレージシステム。

【請求項 8】

前記状態変更は、A L U A (Asymmetric Logical Unit Access) に従う U A (Unit Attention) であり、

前記複数のストレージ装置の各々に、前記仮想ボリュームに関連付けられた複数のターゲットポートグループが定義されており、

前記複数のターゲットポートグループは、1 以上のポートの集合である、

請求項 2 記載のストレージシステム。

【請求項 9】

前記 第 1 コマンド は、前記第 1 ストレージ装置が有する複数の論理ボリュームのうち少なくとも前記第 1 論理ボリュームを含んだ 2 以上の論理ボリュームについての コマンド であり、

前記 第 2 コマンド は、前記第 1 論理ボリュームについての コマンド であり、

前記第 1 ストレージ装置は、前記 第 1 コマンド に応答して前記第 1 パスを含む複数のパスの各々の状態を変更しても、前記 第 2 コマンド を受信した場合に送信する状態変更通知を、前記 2 以上の論理ボリュームのうちの前記第 1 論理ボリュームに関わる状態変更の通知とする、

請求項 2 記載のストレージシステム。

【請求項 10】

前記第 1 ストレージ装置は、前記第 1 コマンドを受信した場合、前記第 1 論理ボリュームに関して状態変更を設定し、

前記第 2 コマンドの受信は、前記仮想ボリュームを指定した R T P G コマンド の受信であり、

前記 R T P G コマンド を受信した前記第 1 ストレージ装置は、

前記第 1 論理ボリュームとペアを構成する前記第 2 論理ボリュームを有する前記第 2 ストレージ装置に、前記状態変更通知と、前記第 2 論理ボリュームに関連付けられたパスの状態の要求とを送信し、前記要求を受信した前記第 2 ストレージ装置から、前記第 2 論理ボリュームに関連付けられた 1 以上のパスの状態を表す情報を含んだパス状態情報を受信し、

10

20

30

40

50

前記第 1 論理ボリュームに関連付けられ前記第 1 パスを含んだ 1 以上のパスの状態を表す第 1 パス状態と、前記第 2 論理ボリュームに関連付けられた 1 以上のパスの状態である第 2 パス状態とを含んだ、前記 R T P G コマンド の応答を、前記ホストシステムに送信する、

請求項 2 記載のストレージシステム。

【請求項 1 1】

1 以上のホストシステムが、複数のホスト装置で構成され、

前記複数のホスト装置の各々が、そのホスト装置が属するグループの ID を保持し、

前記複数のストレージ装置の各々が、そのストレージ装置が属するグループの ID を保持し、

各グループには、少なくとも 1 つのホスト装置と少なくとも 1 つのストレージ装置が属し、

前記第 1 ストレージ装置は、グループ ID の問合せを、前記第 1 ストレージ装置に接続されたいずれかのホスト装置である第 1 ホスト装置から受信した場合、前記第 1 ストレージ装置が保持するグループ ID を含んだ応答を前記第 1 ホスト装置に送信し、

前記第 1 ホスト装置は、

前記応答内のグループ ID が、前記第 1 ホスト装置が保持するグループ ID と同一であれば、前記第 1 ストレージ装置が有し前記第 1 論理ボリュームを含んだ 1 以上の論理ボリュームに関連付けられており前記第 1 パスを含んだ 1 以上のパスの各々の優先度を、第 1 の優先度とし、

前記応答内のグループ ID が、前記第 1 ホスト装置が保持するグループ ID と異なっていれば、前記 1 以上のパスの各々の優先度を、前記第 1 の優先度より低い第 2 の優先度とする、

請求項 2 記載のストレージシステム。

【請求項 1 2】

前記第 2 コマンドの受信は、前記仮想ボリュームを指定した I / O コマンドの受信、及び、前記仮想ボリュームに関連した所定種類の問合せのうち少なくとも 1 つであり、

前記第 1 ストレージ装置は、前記第 1 論理ボリュームに関して状態変更を管理している間に、前記仮想ボリュームを指定した I / O コマンドを前記ホストシステムから受信した場合、状態変更を含んだ、前記 I / O コマンドの応答を、前記ホストシステムに送信する

請求項 1 記載のストレージシステム。

【請求項 1 3】

前記複数のストレージ装置の各々は、そのストレージ装置が管理するリソースの状態を表す情報を保持しているが、そのストレージ装置以外のストレージ装置が管理するリソースの状態を表す情報を保持しておらず、

前記第 1 ストレージ装置は、前記仮想ボリュームを指定した R T P G コマンド を前記第 2 コマンドとして受信した場合、

前記第 1 論理ボリュームとペアを構成する前記第 2 論理ボリュームを有する前記第 2 ストレージ装置に、第 2 リソースの状態の要求を送信し、前記要求を受信した前記第 2 ストレージ装置から、前記第 2 リソースの状態を表す情報を含んだパス状態情報を受信し、前記第 2 リソースは、前記第 1 リソースに関連付けられたリソースであり、前記第 2 論理ボリュームと前記第 2 論理ボリュームに関連し前記第 2 ストレージ装置が管理するリソースとのうちの少なくとも 1 つであり、

前記第 1 リソースの状態と、前記第 2 リソースの状態とを含んだ、前記 R T P G コマンド の応答を、前記ホストシステムに送信する、

請求項 1 記載のストレージシステム。

【請求項 1 4】

前記第 1 コマンドは、前記第 1 ストレージ装置が有する複数の論理ボリュームのうち少なくとも前記第 1 論理ボリュームを含んだ 2 以上の論理ボリュームについてのコマンドで

10

20

30

40

50

あり、

前記第 2 コマンドは、前記第 1 論理ボリュームについてのコマンドであり、

前記第 1 ストレージ装置は、前記第 1 コマンドに応答して前記第 1 リソースを含む複数のリソースの各々の状態を変更しても、前記第 2 コマンドを受信した場合に送信する状態変更通知を、前記 2 以上の論理ボリュームのうちの前記第 1 論理ボリュームに関わる状態変更の通知とする、

請求項 1 記載のストレージシステム。

【請求項 1 5】

複数の論理ボリュームを含む複数種類の複数のリソースを管理し前記複数の論理ボリュームが 1 つに仮想化された論理ボリュームである仮想ボリュームをホストシステムに提供する複数のストレージ装置を有したストレージシステムにおける通知制御方法であって、

前記複数のストレージ装置のうちいずれかのストレージ装置である第 1 ストレージ装置により、第 1 コマンドを受信した場合、第 1 リソースの状態を変更し、前記第 1 リソースは、前記仮想ボリュームの基である前記複数の論理ボリュームのうち前記第 1 ストレージ装置が有する論理ボリュームである第 1 論理ボリュームと前記第 1 論理ボリュームに関連し前記第 1 ストレージ装置が管理するリソースとのうちの少なくとも 1 つであり、

前記第 1 ストレージ装置により、前記仮想ボリュームを指定した R T P G コマンド (ReportTargetPortGroups コマンド) 及び前記仮想ボリュームを指定した I / O コマンドのいずれかである第 2 コマンドを受信した場合、前記第 1 論理ボリュームに関わる状態変更の通知である状態変更通知を第 2 ストレージ装置に送信し、前記第 2 ストレージ装置は、第 2 論理ボリュームを有するストレージ装置であり、前記第 2 論理ボリュームは、前記仮想ボリュームの基であり前記第 1 論理ボリュームに関連付けられた論理ボリュームであり、前記第 2 ストレージ装置により、前記状態変更通知を受信し、前記受信した状態変更通知に基づき状態変更を設定する、

る通知制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、概して、ストレージシステムを構成する複数のストレージ装置間の情報の通知を制御する技術に関する。

【背景技術】

【0002】

例えば、特許文献 1 によれば、1 つのストレージ装置が 2 つのコントローラを有し、2 つのコントローラのうちどちらを利用するかが判断され、その判断結果を基にパス優先度が決定される。

【先行技術文献】

【特許文献】

【0003】

【特許文献 1】US2012/0254657

【発明の概要】

【発明が解決しようとする課題】

【0004】

ところで、ホストシステムと論理ボリューム (ストレージ装置がホストシステムに提供する論理的な記憶デバイス) との間のパスの特定に関し、A L U A (Asymmetric Logical Unit Access) と呼ばれるプロトコルが知られている。ホストシステム (イニシエータ) は、A L U A を利用して、論理ボリューム (ターゲット) についてパスの状態 (例えば優先度) を照会できる。A L U A が利用される環境では、ストレージ装置が有する複数のポートのうち 1 以上のポートを「ターゲットポートグループ」(以下、T P G) と定義することができる。

【0005】

10

20

30

40

50

A L U A が利用される環境では、1つのT P Gの状態が変更されると、そのT P G（以下、状態変更T P G）に関連するT P G（以下、関連T P G）に対して、状態変更が通知されなければならない。状態変更T P Gと全ての関連T P Gが1つのストレージ装置にあれば、複数のT P Gについてそのストレージ装置内で共有される構成情報を更新することで（例えば全ての関連T P Gの各々に関する情報を更新することで）、状態変更の通知が可能である。

【0006】

一方、複数のストレージ装置がそれぞれ有する複数の論理ボリュームを仮想的に1つの論理ボリュームとしてホスト装置に提供するストレージシステムでは、状態変更T P Gと少なくとも1つの関連T P Gが異なる2以上のストレージ装置に分かれていることがある。この場合、状態変更T P Gを有するストレージ装置から関連T P Gを有するストレージ装置へ状態変更が通知されなければならない。そのような状態変更の通知を実現する1つの仕組みとして、全てのストレージ装置が、それぞれ、各ストレージ装置の構成情報を保持し、全てのストレージ装置間で構成情報を同期することが考えられる。しかし、そのような仕組みでは、各ストレージ装置が保持する構成情報の総量が膨大となり、さらに、構成情報の同期をとるためにストレージシステム全体の性能が低下してしまう。

【0007】

この種の問題は、論理ボリューム以外のオブジェクトが仮想化されるケースについてもあり得る。また、この種の問題は、T P G及びバス以外のリソースの状態変更（例えば優先度変更）についてもあり得る。また、この種の問題は、A L U Aを利用するストレージシステムに限らず、複数のストレージ装置を有する他種のストレージシステム、例えば、いずれかのストレージ装置でのリソースの状態変更をホストシステムからのアクセスに回答してそのホストシステムに通知したり（気付かせたり）、状態変更気付いたホストシステムから所定種類の問合せを受けた場合にリソースの変更後の状態をホストシステムに通知するようになっているストレージシステム、についてもあり得る。

【課題を解決するための手段】

【0008】

ストレージシステムは、複数のストレージ装置を有する。複数のストレージ装置は、複数の論理ボリュームを含む複数種類の複数のリソースを管理し、複数の論理ボリュームが1つに仮想化された論理ボリュームである仮想ボリュームをホストシステムに提供する。第1ストレージ装置（いずれかのストレージ装置）が、第1イベントを検出した場合、第1リソースの状態を変更する。その後、第1リソースの状態をホストシステムに将来通知し得ることを意味する第2イベントを検出した場合、第1ストレージ装置は、第1論理ボリュームに関わる状態変更の通知である状態変更通知を、第1論理ボリュームに関連付けられ仮想ボリュームの基である第2論理ボリュームを有する第2ストレージ装置に送信する。第2ストレージ装置は、状態変更通知を受信し、受信した状態変更通知に基づき状態変更を設定する。なお、第1リソースは、仮想ボリューム（複数の論理ボリュームが1つに仮想化された論理ボリューム）の基である複数の論理ボリュームのうち第1ストレージ装置が有する論理ボリュームである第1論理ボリュームと、第1論理ボリュームに関連し第1ストレージ装置が管理するリソース（例えばバス）とのうちの少なくとも1つである。また、ホストシステムは、1以上のホスト装置（例えば、クラスタを構成する2以上のホスト装置）である。1以上のホスト装置のうちの少なくとも1つは、物理的なホスト装置でよく、1以上のホスト装置は、物理的なホスト装置で実行される仮想的なホスト装置を含んでもよい。

【発明の効果】

【0009】

全てのストレージ装置が、それぞれ、各ストレージ装置の構成情報を保持し、且つ、全てのストレージ装置間で構成情報を同期すること無しに、第1論理ボリューム又はそれに関連したリソースであり第1ストレージ装置が管理する第1リソースの状態の変更を、ホストシステムにその状態変更を知られる前に、第1論理ボリュームに関連付けられた第2

10

20

30

40

50

論理ボリュームを有する第2ストレージ装置に通知することができる。

【図面の簡単な説明】

【0010】

【図1】実施例1に係る計算機システムの構成を示す。

【図2】各ストレージ装置が有する管理情報の構成を示す。

【図3】ALUA管理テーブルの構成を示す。

【図4】ペア管理テーブルの構成を示す。

【図5】UA管理テーブルの構成を示す。

【図6】UA同期フラグビットリストの構成を示す。

【図7】パス管理テーブルの構成を示す。

10

【図8】ポート管理テーブルの構成を示す。

【図9】仮想ボックス管理テーブルの構成を示す。

【図10】相対ポートID管理テーブルの構成を示す。

【図11】有効ID範囲管理テーブルの構成を示す。

【図12】ホスト装置がストレージ装置のパス優先度を把握するまでの流れを示す。

【図13】異なる2つのストレージ装置のパス優先度をマージすることを模式的に示す。

【図14】ALUA適用状態変更処理の流れを示す。

【図15】ALUA適用状態確認処理の流れを示す。

【図16】パス優先度変更処理の流れを示す。

【図17】I/O処理の際のUA伝播の流れを示す。

20

【図18A】ペア状態「PAIR」の場合の優先パスとアクセス経路の一例を示す。

【図18B】ペア状態「SUSPEND」の場合の優先パスとアクセス経路の一例を示す。

【図19】パス優先度変更処理の流れを示す。

【図20】パス優先度レポート処理の流れの第1の部分を示す。

【図21】パス優先度レポート処理の流れの第2の部分を示す。

【図22】パス優先度レポート処理の流れの第3の部分を示す。

【図23】パス優先度レポート処理の流れの第4の部分を示す。

【図24】実施例2に係る計算機システムの構成、及び、各ストレージ装置が有する管理情報の構成を示す。

30

【図25】データセンタID問合せコマンドの構成を示す。

【図26】パス優先度管理テーブルの構成を示す。

【図27A】第1データセンタに属するホスト装置及びストレージ装置がそれぞれ有するデータセンタID管理テーブルを示す。

【図27B】第2データセンタに属するホスト装置及びストレージ装置がそれぞれ有するデータセンタID管理テーブルを示す。

【図28】パス優先度自動決定処理の流れを示す。

【図29】I/O制御処理の流れを示す。

【図30】実施例1に係る計算機システムの概要を示す。

【発明を実施するための形態】

40

【0011】

以下の説明では「aaaテーブル」等の表現にて情報を説明する場合があるが、これら情報は、テーブル等のデータ構造以外で表現されていてもよい。そのため、データ構造に依存しないことを示すために「aaaテーブル」等について「aaa情報」と呼ぶことがある。さらに、各情報の内容を説明する際に、「ID」、「番号」という表現を用いるが、これらに代えて又は加えて他種の識別情報が使用されてよい。

【0012】

また、以下の説明では、プログラムを主語として処理を説明する場合があるが、プログラムは、プロセッサ(例えばCPU(Central Processing Unit))によって実行されることで、定められた処理を、適宜に記憶資源(例えばメモリ)及び/又は通信インター

50

フェイスを用いながら行うため、処理の主語がプロセッサとされてもよい。プログラムを主語として説明された処理は、プロセッサを含むストレージコントローラ、ストレージ装置又はホスト装置が行う処理としてもよい。また、プロセッサが、処理の一部又は全部を行うハードウェア回路を含んでもよい。コンピュータプログラムは、プログラムソースからストレージ装置又はホスト装置にインストールされてもよい。プログラムソースは、例えば、プログラム配布サーバ、又は、計算機が読み取り可能な記憶メディアであってもよい。

【0013】

また、以下の説明では、同種の要素を区別して説明する場合は、その要素の参照符号を使用し（例えばストレージ装置20A、ストレージ装置20B）、同種の要素を区別しないで説明する場合は、その要素の参照符号のうちの共通符号のみ使用する（例えばストレージ装置20）ことがある。

10

【0014】

以下、図面を参照して、幾つかの実施例を説明する。なお、以下の実施例では、ALUA (Asymmetric Logical Unit Access) を例に取る。

【実施例1】

【0015】

図30は、実施例1に係る計算機システムの概要を示す。

【0016】

ホスト装置10と、ホスト装置10が接続されたストレージシステム303とを有する。ホスト装置10は、ホストシステムの一例である。ストレージシステム303は、ストレージ装置20A及び20Bを有する。ストレージ装置20A及び20Bは、複数のストレージ装置の一例であり、ストレージ装置20の数は3以上でもよい。以下の説明では、ストレージ装置20Aを「第1ストレージ装置20A」と言い、ストレージ装置20Bを「第2ストレージ装置20B」と言うことがある。

20

【0017】

ストレージ装置20A及び20Bは、複数のLDEV（論理記憶デバイス）301を含む複数種類の複数のリソースを管理する。LDEVは、論理ボリュームである。リソースとしては、LDEVの他に、ターゲットポートグループ（TPG）307及びパス309等がある。TPG307A及び307Bの各々は、本実施例では1つのポートであるが、2以上のポートの集合でもよい。また、LDEV301は、論理的な記憶デバイスであり、実体的なLDEVでも仮想的なLDEVでもよい。実体的なLDEVは、物理的な記憶資源（例えば1以上の物理記憶デバイス）に基づくLDEVである。仮想的なLDEVは、外部のストレージ装置（図示せず）の記憶資源（例えばLDEV）に基づいておりストレージ仮想化技術に従うLDEVである外部接続LDEVであってもよいし、複数の仮想ページ（仮想的な記憶領域）で構成されており容量仮想化技術（典型的にはThin Provisioning）に従うLDEVであってもよい。本実施例では、説明を分かり易くするため、各LDEVは、いずれかのVDEVの基になるとする。

30

【0018】

ストレージ装置20A及び20Bは、LDEV301A及び301Bが1つに仮想化されたLDEVであるVDEV（仮想ボリューム）305をホスト装置10に提供する。ホスト装置10が、VDEV305を認識（例えばマウント）する。

40

【0019】

ホスト装置10は、物理的又は仮想的な計算機であり、パス管理プログラム302を実行する。パス管理プログラム302は、ホスト装置10が認識したVDEV305に関連付けられている複数のパス309を管理する。パス309は、TPG307を経由し、そのTPG307に関連付けられているLDEV301へと繋がる。図30の例では、1つのVDEV305の基になっているLDEV301毎に1つのパス309であるが、1つのLDEV301につきパス309は1つ以上でよく、故に、1つのVDEV305につきパス309は2つ以上でよい。パス管理プログラム302は、各パスの優先度を管理し

50

ており、優先度の高いパス 309 を優先度の低いパス 309 よりも優先的に使用する。例えば、パス 309 B の優先度よりパス 309 A の優先度が高ければ、パス管理プログラム 302 は、パス 309 A を優先的に使用して I/O (Input/Output) コマンドをストレージシステム 303 に送信する。パス優先度が、パス状態の一例である。

【0020】

パス 309 A の優先度は、ホスト装置 10 及び第 1 ストレージ装置 20 A の各々において管理されており、パス 309 B の優先度は、ホスト装置 10 及び第 2 ストレージ装置 20 B の各々において管理されている。各ストレージ装置は、他のストレージ装置が管理するパスを管理しない。すなわち、第 1 ストレージ装置 20 A は、第 2 ストレージ装置 20 B が管理するパス 309 B について優先度を管理しないし、第 2 ストレージ装置 20 B は、第 1 ストレージ装置 20 A が管理するパス 309 A について優先度を管理しない。従って、本実施例では、全てのストレージ装置 20 が、それぞれ、各ストレージ装置のパス管理情報 (パス優先度を含んだ情報) を保持し、且つ、全てのストレージ装置間でパス管理情報を同期する必要が無い。

10

【0021】

この環境において、例えば、第 1 ストレージ装置 20 A において、パス 309 A の優先度に変更されたとする。第 1 ストレージ装置 20 A は、パス 309 A の優先度に変更された場合、状態変更を、例えば、パス 309 A が関連付けられている第 1 LDEV 301 A について設定する。状態変更の一例が、ALUA では UA (Unit Attention) である。

【0022】

この段階では、未だ、ホスト装置 10 は、パス 309 A の優先度に変更されたことを知らない。ホスト装置 10 は、ストレージシステム 303 の詳細な構成を把握しておらず、ALUA によれば、ストレージシステム 303 が、ホスト装置 10 に、ホスト装置 10 が使用するパスを決めさせる。具体的には、ストレージシステム 303 からホスト装置 10 にパス優先度を通知することで、ホスト装置 10 が、パス優先度を管理でき、以って、適切なパス 309 を選択できるようになる。より具体的には、ホスト装置 10 が、パス優先度の問合せを第 1 ストレージ装置 20 A に送信し、第 1 ストレージ装置 20 A からその問合せの応答を受信し、その応答から、パス 309 A の変更後の優先度を知ることができる。ALUA では、そのような問合せは、SCSI でサポートされている ReportTargetPort Groups コマンド (以下、RTPG コマンド) である。

20

30

【0023】

第 1 ストレージ装置 20 A は、VDEV 305 に関連した RTPG コマンドをホスト装置 10 から受信した場合 (ステップ 1)、第 1 LDEV 301 A とペアを構成する第 2 LDEV 301 B を有する第 2 ストレージ装置 20 B に、状態変更通知と、第 2 LDEV に関連付けられたパスの状態の要求であるパス状態要求とを送信する (ステップ 2)。状態変更通知の一例が、後述の UA 同期フラグビットリストである。UA 同期フラグビットリストからは、第 1 ストレージ装置 20 A が有するどの LDEV について UA が変更されたかがわかる。パス状態要求の一例が、第 2 ストレージ装置 20 B が管理するパスの優先度の要求であるパス優先度要求である。

【0024】

第 2 ストレージ装置 20 B は、UA 同期フラグビットリスト及びパス優先度要求を受信し、UA 同期フラグビットリストを基に、UA を、例えば第 1 LDEV 301 A とペアを構成する第 2 LDEV 301 B について設定し、且つ、第 2 ストレージ装置 20 B が保持するパス管理情報を基に、第 2 LDEV 301 B に関連付いたパス 309 B の優先度を表す情報を含んだパス優先度情報を第 1 ストレージ装置 20 A に応答する (ステップ 3)。

40

【0025】

第 1 ストレージ装置 20 A は、パス優先度情報を受信し、第 1 LDEV 301 A に関連付けられ第 1 パス 309 A を含んだ 1 以上のパスの優先度を表す第 1 パス優先度と、第 2 LDEV 301 B に関連付けられ第 2 パス 309 B を含んだ 1 以上のパスの状態である第 2 パス優先度とを含んだ応答 (RTPG コマンドの応答である RTPG 応答) を、ホスト

50

装置 10 に送信する (ステップ 4)。つまり、RTPG 応答において、2つのストレージ装置 20A 及び 20B のパス優先度がマージされている。ホスト装置 10 (パス管理プログラム 302) は、RTPG 応答に含まれる第 1 パス優先度及び第 2 パス優先度を、ホスト装置 10 が保持するパス管理情報に設定し、以後、DEV 305 に対する I/O コマンドを送信する場合、そのパス管理情報が表す第 1 パス優先度及び第 2 パス優先度を基に、パス優先度が高い方のパス 309A 又は 309B を選択する。

【0026】

以上の処理において、第 1 ストレージ装置 20A から第 2 ストレージ装置 20B への状態変更通知 (UA 同期フラグビットリスト) が、第 1 ストレージ装置 20A でのパス優先度変更を第 2 ストレージ装置 20B に通知することの一例である。以上の処理によれば、
10
全てのストレージ装置 20A 及び 20B が、それぞれ、各ストレージ装置のパス管理情報を保持し、且つ、全てのストレージ装置間でパス管理情報を同期すること無しに、第 1 LDEV 301A でのパス優先度変更を第 2 ストレージ装置 20B に通知することができる。

【0027】

また、以上の処理によれば、パス優先度変更の通知は、第 1 ストレージ装置 20A でのパス優先度変更の都度に行われるのではなく、RTPG コマンドを受信した場合に行われる。RTPG コマンドの受信は、変更後のパス優先度をホスト装置 10 に将来通知し得ることを意味する第 2 イベントの検出の一例である。パス優先度変更の都度にパス優先度変更を通知するとなると、第 1 ストレージ装置 20A の負荷が高くなり、結果として、
20
ストレージシステム 303 の性能が低下し得る。以上の処理によれば、そのような性能低下を避けることができる。なお、第 2 イベントの検出の別の例として、後述するように、I/O コマンドの受信 (パス優先度変更を表す UA をホスト装置 10 に送信するきっかけ) がある。

【0028】

なお、本実施例では、RTPG に設定される第 1 パス優先度及び第 2 パス優先度は、第 1 LDEV 301A と第 2 LDEV 301B とのペアの状態であるペア状態と、第 1 LDEV 301A と第 2 LDEV 301B とについての I/O モードとに応じて、第 1 ストレージ装置 20A (RTPG コマンドを受信した方のストレージ装置) により制御される (I/O モードが、I/O 状態の一例である)。これにより、第 1 ストレージ装置 20A は、
30
ペア状態及び I/O モードに応じた最適パスをホスト装置 10 に選択させることができるようになる。

【0029】

具体的には、ペア状態が「PAIR」であり I/O モードが「ミラー」、すなわち、第 1 LDEV 301A と第 2 LDEV 301B が同期状態であることが表されていれば、第 1 ストレージ装置 20A は、RTPG 応答に含められる第 1 パス優先度を、第 1 ストレージ装置が保持するパス管理情報から特定されるパス優先度とし、RTPG 応答に含められる第 2 パス優先度を、第 2 ストレージ装置 20B から受信したパス優先度情報が表すパス優先度とする。ペア状態が「SUSPEND」であり I/O モードが「リモート」(第 1 LDEV 301A 内のデータよりも第 2 LDEV 301B 内のデータの方が新しい) 場合、
40
第 1 ストレージ装置 20A は、第 1 パス優先度よりも第 2 パス優先度を高くする (例えば、第 1 パス優先度 = Active/non-optimized、第 2 パス優先度 = Active/optimized)。ペア状態が「SUSPEND」であり I/O モードが「ローカル」(第 2 LDEV 301B 内のデータよりも第 1 LDEV 301A 内のデータの方が新しい) 場合、第 1 ストレージ装置 20A は、第 2 パス優先度よりも第 1 パス優先度を高くする (例えば、第 2 パス優先度 = Active/non-optimized、第 1 パス優先度 = Active/optimized)。これにより、第 1 ストレージ装置 20A は、ペア状態及び I/O モードに応じた最適パスをホスト装置 10 に選択させることができるようになる。

【0030】

また、第 1 ストレージ装置 20A は、第 2 ストレージ装置 20B との通信が不可能であ

10

20

30

40

50

る場合、ペア状態が「SUSPEND」でありI/Oモードが「リモート」であれば、RTPG応答としてエラー応答をホスト装置10に送信し、一方、ペア状態が「PAIR」でありI/Oモードが「ミラー」、又は、ペア状態が「SUSPEND」でありI/Oモードが「ローカル」であれば、第2パス優先度を含まず第1パス優先度をActive/optimizedとしたRTPG応答を、ホスト装置10に送信する。これにより、第1ストレージ装置20Aは、ペア状態及びI/Oモードに加えて第2ストレージ装置20Bとの通信状態に応じた制御をホスト装置10に実行させることができるようになる。

【0031】

ホスト装置10は、定期的にRTPGコマンドを送信してもよいし、ストレージシステム303からUA (Unit Attention) を受信した場合にRTPGコマンドをストレージシステム303に送信してもよい。第1ストレージ装置20Aは、ホスト装置10からVDEV305 (又はVDEV305の基のLDEV301A又は301B) を指定したI/Oコマンドを受信した場合、そのI/Oコマンドに従いI/Oを実行しI/O応答をホスト装置10に送信するが、第1ストレージ装置20Aに設定されているUAが少なくともパス優先度変更を表す値であれば、I/O応答にそのUAを設定する。ホスト装置10は、I/O応答内のUAがパス優先度変更を表している場合に、I/Oコマンドで指定したVDEV305 (又はVDEV305に関連付けられているポート) を指定したRTPGコマンドを少なくとも第1ストレージ装置20Aに送信する。つまり、UAを含んだI/O応答の送信は、パス優先度変更をホスト装置10に気付かせるアクションの1つである。従って、第1ストレージ装置20Aは、I/Oコマンドを受信した場合に (UAを含んだI/O応答をホスト装置10に送信する前に)、パス優先度変更の通知を、ペア相手のストレージ装置20B (受信したI/Oコマンドで指定されているVDEVの基になっている第1LDEV301Aとペアを構成する第2LDEV301Bを有するストレージ装置20B) に送信する。これにより、パス優先度変更がペア相手のストレージ装置20Bに通知される。

【0032】

また、本実施例では、第1ストレージ装置20Aにおいて、パス優先度変更は、ストレージシステム303の管理システム (図示せず) からのパス優先度変更要求 (パス優先度の変更の要求) に応答して行われる。パス優先度変更要求は、第1イベントの一例である。パス優先度変更は、LDEV単位よりも大きな単位、例えば、TPG単位、又は、ホストグループ単位で行うことができる (ホストグループについては後述)。すなわち、1つのパス優先度変更要求により、そのパス優先度変更要求で指定されたりソースに関連付けられている複数のLDEVの各々について、パス優先度の変更が可能である。一方、RTPGコマンドもI/Oコマンドも、LDEV単位 (LUN単位) である。このため、第1ストレージ装置20Aは、パス優先度変更の通知 (状態変更通知の送信) を、LDEV単位で行ってよい。つまり、第1ストレージ装置20Aは、パス優先度変更の単位よりも小さい単位でパス優先度変更の通知を行ってよい。これにより、第1ストレージ装置20Aの負荷軽減が期待でき、以って、ストレージシステム303の性能向上が期待できる。

【0033】

以上が、実施例1の概要である。なお、以上の説明では、第1ストレージ装置20Aを代表的に例に取っているが、本実施例 (及び実施例2) で説明する第1ストレージ装置20Aの処理は、他の各ストレージ装置20も実行可能である。

【0034】

以下、実施例1を詳細に説明する。

【0035】

図1は、実施例1に係る計算機システムの構成を示す。

【0036】

1以上のホスト装置10と、ストレージ管理サーバ11が、ストレージシステムを構成する複数のストレージ装置20A及び20Bと、通信ネットワーク (例えばSAN (Storage Area Network) 12) を介して通信可能に接続されている。ストレージ管理サーバ1

10

20

30

40

50

1 が、管理システムの一例である。管理システムは、1 以上の計算機で構成されてよい。具体的には、例えば、管理計算機が情報を表示する場合（具体的には、管理計算機が自分の表示デバイスに情報を表示する、或いは、管理計算機が表示用情報を遠隔の表示用計算機に送信する場合）、管理計算機が管理システムである。

【0037】

ホスト装置 10 は、通信インターフェイスデバイスと、記憶デバイスと、それらに接続されたプロセッサとを有する。通信インターフェイスデバイスを介してストレージ装置 20 等と通信できる。記憶デバイスは、例えばメモリである。プロセッサは、例えば CPU であり、記憶デバイスに記憶されたプログラム（例えば、パス管理プログラム 302（図 30 参照）、アプリケーションプログラム及びオペレーティングシステム）を実行することができる。

10

【0038】

ストレージ装置 20A 及び 20B のうちストレージ装置 20A を例に取り、ストレージ装置の構成を説明する。

【0039】

ストレージ装置 20A は、1 以上の PDEV 群 28A と、1 以上の PDEV 群 28A に対する I/O を制御するストレージコントローラ 122A とを有する。

【0040】

PDEV 群 28A は、1 以上の PDEV の集合、例えば RAID (Redundant Array of Independent (or Inexpensive) Disks) グループである。PDEV は、不揮発性の物理記憶デバイスを意味し、例えば HDD (Hard Disk Drive) 又は SSD (Solid State Drive) である。PDEV 群 28A に基づき、実体的な LDEV を構築することができる。

20

【0041】

ストレージコントローラ 122A は、ホスト装置 10 及びストレージ管理サーバ 11 と通信するための FE I/F 部（フロントエンドのインターフェイス部）と、1 以上の PDEV 群 28A と通信するための BE I/F 部（バックエンドのインターフェイス部）と、メモリ部と、それらに接続された制御部とを有する。FE I/F 部が、1 以上の CHA (チャネルアダプタ) 22A により実現されている。BE I/F 部が、1 以上の DKA (ディスクアダプタ) 26A により実現されている。メモリ部が、SM (共有メモリ) 21A、CM (キャッシュメモリ) 24A、1 以上の CHA 22A 内の 1 以上の LM (ローカルメモリ) 222A、及び、1 以上の DKA 26A 内の 1 以上の LM 262A により実現されている。制御部が、1 以上の CHA 内の 1 以上の MP (マイクロプロセッサ) 221A、及び、1 以上の DKA 26A 内の 1 以上の MP 261A により実現されている。SM 21A 及び CM 24A は、それぞれ同一のメモリ上に設けられた領域であってもよいし、異なるメモリであってもよい。CHA 22A、SM 21A、CM 24A 及び DKA 26A が、接続部（例えばバス又はスイッチ）25 を介して通信可能である。

30

【0042】

1 以上の CHA 22A は、複数のポート 29A を有する。CHA 22A は、MP 221A 及び LM (例えば揮発メモリ) 222A を有する。CHA 22A の動作は MP 221A により制御される。CHA 22A は、I/O コマンドをホスト装置 10 から受信する。受信した I/O コマンドがライトコマンドの場合、ライトコマンドに従うライト対象のデータが CHA 22A により CM 24A に一時格納され、CM 24A 上のライト対象のデータが、DKA 26A により LDEV に書き込まれる。受信した I/O コマンドがリードコマンドの場合、DKA 26A により、リードコマンドに従い LDEV から読み出されたリード対象のデータ CM 24A に一時格納され、CHA 22A により、CM 24A 上のリード対象のデータがホスト装置 10 に送信される。I/O コマンドは、I/O 先の LDEV (例えば VDEV) の番号 (例えば LUN (Logical Unit Number)) と、その LDEV における領域のアドレス (例えば LBA (Logical Block Address)) とを含む。

40

【0043】

50

D K A 2 6 A は、M P 2 6 1 A 及び L M 2 6 2 A を有する。D K A 2 6 A の動作は M P 2 6 1 A により制御される。D K A 2 6 A は、L D E V に対する I / O (例えば、P D E V 群 2 8 A に対するデータの I / O) を制御し、その際、C M 2 4 A に対するデータの I / O も制御する。

【 0 0 4 4 】

C M 2 4 A は、揮発メモリ及び / 又は不揮発メモリであり、L D E V に対する I / O 対象データが一時的に格納される。

【 0 0 4 5 】

S M 2 1 A は、ポート管理情報等の情報が格納される。C H A 2 2 A 及び D K A 2 6 A の各々は、適宜、S M 2 1 A 内の情報を参照するか、或いは、S M 2 1 A 内の情報の少なくとも一部を、自分の L M (2 2 2 A 又は 2 6 2 A) に格納しその L M に格納された情報を参照する。

10

【 0 0 4 6 】

ストレージ装置 2 0 A が行う処理は、ストレージ装置 2 0 A のストレージコントローラ 1 2 2 A により行なわれる。

【 0 0 4 7 】

図 2 は、各ストレージ装置 2 0 が有する管理情報の構成を示す。なお、P D E V 群 2 8 は、図示の通り、1 以上の P D E V 2 8 1 により構成されている。以下、第 1 ストレージ装置 2 0 A を例に取るが、他の各ストレージ装置についても同様である。

【 0 0 4 8 】

20

第 1 ストレージ装置 2 0 A が保持する管理情報は、A L U A 管理テーブル 2 1 1 A、ペア管理テーブル 2 1 2 A、U A 管理テーブル 2 1 3 A、パス管理テーブル 2 1 4 A、ポート管理テーブル 2 1 5 A、仮想ボックス管理テーブル 2 1 6 A、相対ポート I D 管理テーブル 2 1 7 A、有効 I D 範囲管理テーブル 2 1 8 A 及び U A 同期フラグビットリスト 2 2 3 A を含む。これらのテーブル 2 1 1 A ~ 2 1 8 A 及びリスト 2 2 3 A のうちの 2 以上が 1 つにマージされてもよいし、これらのテーブル 2 1 1 A ~ 2 1 8 A 及びリスト 2 2 3 A のうちの少なくとも 1 つの各々が 2 以上に分割されていてもよい。また、本実施例では、テーブル 2 1 1 A ~ 2 1 8 A は、S M 2 1 A に格納され、U A 同期フラグビットリスト 2 2 3 A は、C H A 2 2 A の L M 2 2 2 A に格納されるが、リスト 2 2 3 A が L M 2 2 2 A に代えて又は加えて、S M 2 1 A に格納されてもよいし、テーブル 2 1 1 A ~ 2 1 8 A のうちの少なくとも 1 つが、S M 2 1 A に代えて又は加えて L M 2 2 2 A に格納されてもよい。以下、図 3 ~ 図 1 1 を参照して、テーブル 2 1 1 A ~ 2 1 8 A 及びリスト 2 2 3 A を説明する。

30

【 0 0 4 9 】

図 3 は、A L U A 管理テーブル 2 1 1 A の構成を示す。

【 0 0 5 0 】

A L U A 管理テーブル 2 1 1 A は、第 1 ストレージ装置 2 0 A が有する L D E V 毎に A L U A が適用されているか否かを表す情報を有する。具体的には、A L U A 管理テーブル 2 1 1 A は、第 1 ストレージ装置 2 0 A が有する L D E V 毎に、以下の情報、すなわち、L D E V 番号 (L D E V の L D E V 番号) 4 0 1 と、A L U A 適用状態 4 0 2 (A L U A が適用されているか否か) とを有する。例えば、A L U A 適用状態「 0 」は、A L U A 適用対象外を意味し、A L U A 適用状態「 1 」は、A L U A 適用対象を意味する。

40

【 0 0 5 1 】

なお、「L D E V 番号」とは、L D E V の識別番号でありホスト装置 1 0 には認識されずストレージ装置 2 0 において使用される番号である。L D E V には、L D E V 番号に加えて、後述の L U N (Logical Unit Number) が関連付けられる。L U N は、ホスト装置 1 0 に認識され、I / O コマンドに含まれる (I / O コマンドで指定される)。本実施例では、同一の L U N が、V D E V の基である L D E V 3 0 1 A 及び 3 0 1 B に関連付けられ、それ故、ホスト装置 1 0 は、V D E V の番号として L U N を認識できる。各ストレージ装置 2 0 は、I / O コマンドに含まれている L U N に対応した L D E V 番号を特定す

50

ることで、I/O先のLDEVを特定できる。

【0052】

図4は、ペア管理テーブル212Aの構成を示す。

【0053】

ペア管理テーブル212Aは、第1ストレージ装置20Aが有しVDEVの基になっているLDEV毎に、ペア相手のLDEV、ペア状態、I/Oモード、及び、ALUA適用状態をペア相手のストレージ装置に通知するか否かを表す情報、を有する。具体的には、ペア管理テーブル212Aは、第1ストレージ装置20Aが有しVDEVの基になっているLDEV毎に、以下の情報、なわち、LDEV番号(自ストレージ)501と、LDEV番号(ペア相手ストレージ)502と、ペア状態503と、I/Oモード504と、ALUA同期ビット505とを有する。

10

【0054】

LDEV番号(自ストレージ)501は、第1ストレージ装置20Aが有するLDEV、例えば第1LDEV301AのLDEV番号である。LDEV番号(ペア相手ストレージ)502は、第1ストレージ装置20Aが有するLDEVとペアを構成するLDEV(他ストレージ装置が有するLDEV)の番号、例えば第2ストレージ装置20Bが有する第2LDEV301BのLDEV番号である。

【0055】

ペア状態503は、LDEVペアの状態を表し、例えば「PAIR」(一方のLDEVにデータが格納されると他方のLDEVにそのデータがコピーされる状態)、「SUSPEND」(一方のLDEVにデータが格納されても他方のLDEVにそのデータがコピーされない状態)等がある。

20

【0056】

I/Oモード504は、LDEVペア(言い換えればVDEV)に対するI/Oのモード(状態)を表し、例えば「Mirror」(LDEVペアを構成するLDEV間が同期)、「Local」(自ストレージ内のLDEV内のデータがペア相手ストレージ内のLDEVよりも新しく故に自ストレージ内のLDEVへのI/Oを許可)、「Remote」(ペア相手ストレージ内のLDEV内のデータが自ストレージ内のLDEVよりも新しく故に自ストレージ内のLDEVへのI/Oを禁止しペア相手ストレージ内のLDEVへのI/Oを許可)、及び「Block」(エラーをホスト装置へ応答)等がある。

30

【0057】

ALUA同期ビット505は、ALUA適用状態をペア相手のストレージ装置に通知するか否かを表す。例えば、ALUA同期ビット「0」は、ペア相手LDEVのALUA適用状態の変更が不要を意味し、ALUA同期ビット「1」は、ペア相手LDEVのALUA適用状態の変更が必要を意味する。

【0058】

図5は、UA管理テーブル213Aの構成を示す。

【0059】

UA管理テーブル213Aは、UA管理単位毎にUAを有する。UA管理単位は、UAが表す状態変更についてのリソース種類によって異なってよい。本実施例では、UAはパス優先度変更の有無を表し、UA管理単位は、ポート及びホストグループ単位である。具体的には、UA管理テーブル213Aは、ポート29A毎に、以下の情報、すなわち、ポート番号601、ホストグループ番号602及びUA603を有する。

40

【0060】

ポート番号601は、ポート29Aの識別番号である。ホストグループ番号602は、ポート29Aに関連付けられたホストグループの番号である。なお、「ホストグループ」とは、1以上のホスト装置10の集合であり、ストレージ装置20のポートに関連付けられるリソースの一種である。ポート29Aを経由したI/Oは、I/O元のホスト装置10が、そのポート29Aに関連付けられたホストグループに含まれるホスト装置10の場合に、第1ストレージ装置20Aにより許可される。

50

【 0 0 6 1 】

U A 6 0 3 は、U A の値を表し、例えば、U A 「 0 」は、U A 設定無しを意味し、U A 「 1 」は、パス優先度変更を意味し、U A 「 2 」は、A L U A 適用状態変更を意味する。

【 0 0 6 2 】

図 6 は、U A 同期フラグビットリスト 2 2 3 A の構成を示す。

【 0 0 6 3 】

U A 同期フラグビットリスト 2 2 3 A は、第 1 ストレージ装置 2 0 A が有する L D E V 毎に、ペア相手ストレージにパス優先度変更の U A を通知するため自ストレージの L D E V についてパス優先度変更の U A (例えば U A 「 1 」) が設定されたか否かを表す。具体的には、U A 同期フラグビットリスト 2 2 3 A は、第 1 ストレージ装置 2 0 A が有する L D E V 毎に、以下の情報、すなわち、L D E V 番号 7 0 1 及び U A 同期フラグビット 7 0 2 を有する。L D E V 番号 7 0 1 は、第 1 ストレージ装置 2 0 A が有する L D E V の番号である。U A 同期フラグビット 7 0 2 は、L D E V についてパス優先度変更の U A が設定されたか否かを表す。例えば、U A 同期フラグ「 0 」は、パス優先度変更無しを意味し、U A 同期フラグ「 1 」は、パス優先度変更を意味する。

10

【 0 0 6 4 】

U A 同期フラグの変更は、ホスト装置 1 0 からの I / O コマンドの受付を止めることなく行われ、且つ、状態変更通知の送信 (U A 同期フラグビットリストの送信) は、U A 同期フラグの変更とは非同期に行われる。これにより、ストレージシステムの性能低下を回避できる。

20

【 0 0 6 5 】

図 7 は、パス管理テーブル 2 1 4 A の構成を示す。

【 0 0 6 6 】

パス管理テーブル 2 1 4 A は、パス管理情報の一例であり、パス管理単位毎のパス優先度を表す。パス管理単位は、本実施例ではホストグループ単位である。つまり、本実施例では、ホストグループ単位で、パス優先度が設定及び変更される。具体的には、パス管理テーブル 2 1 4 A は、ホストグループ毎に、以下の情報、すなわち、仮想ボックス番号 8 0 1、H G 番号 8 0 2、相対ポート I D 8 0 3 及びパス優先度 8 0 4 を有する。

【 0 0 6 7 】

仮想ボックス番号 8 0 1 は、仮想ボックスの番号であり、仮想ボックスは、仮想的なストレージシステムに相当する。H G 番号 8 0 2 は、ホストグループの番号である。

30

【 0 0 6 8 】

相対ポート I D 8 0 3 は、ポート 2 9 A のポート番号と異なる I D であり、1 つの仮想ボックス (仮想ストレージシステム) (異なるストレージ装置 2 0 A、2 0 B) においてポート番号が重複しないよう割り振られた I D である。相対ポートは、仮想ポートと呼ぶこともできる。

【 0 0 6 9 】

パス優先度 8 0 4 の具体的な値として、例えば、「Active/optimized」(パス優先度「高」であり I / O 実行可能)、「Active/non-optimized」(パス優先度「低」であり待機状態)、「Unavailable」(I / O 実行不可能)等がある。

40

【 0 0 7 0 】

なお、本実施例では、以下方針でパス優先度を管理することができる。すなわち、(1) 自ストレージが管理している、自ストレージ内のホストグループのパス優先度は、常に正しい情報である。(2) ペア相手ストレージが管理している、自ストレージ内のホストグループのパス優先度は、常に正しい情報である必要はない。

【 0 0 7 1 】

また、本実施例において、同一ストレージ装置内で、ポート番号は重複できない。しかし、仮想ボックス番号が異なれば、同一ストレージ装置内で、相対ポート I D は重複してよい。また、同一仮想ボックス (仮想ストレージ) 内で、ストレージ装置が異なれば、ポート番号は重複してよい。また、同一仮想ボックス (仮想ストレージ) 内で、相対ポート

50

IDは重複できない。

【0072】

図8は、ポート管理テーブル215Aの構成を示す。

【0073】

ポート管理テーブル215Aは、各ポート29Aに関連付けられている情報を表す。具体的には、ポート管理テーブル215Aは、ポート29A毎に、以下の情報、すなわち、ポート番号1001、インデックス番号1002、LUN1003、LDEV番号1004及びHG番号1005を有する。ポート番号1001は、ポート29Aの番号である。インデックス番号1002は、ポート29Aに関連付け可能なエントリ(LDEV)の通し番号である。LUN1003は、ポート29Aに関連付けられたLDEVに関連付けられて

10

【0074】

図9は、仮想ボックス管理テーブル216Aの構成を示す。

【0075】

仮想ボックス管理テーブル216Aは、仮想ボックス毎に、第1ストレージ装置20Aが管理するリソースを表す。具体的には、仮想ボックス管理テーブル216Aは、仮想ボックス毎に、以下の情報、すなわち、仮想ボックス番号1101、HG番号1102、ポート番号1003及びLDEV番号1104を有する。仮想ボックス番号1101は、仮想ボックスの番号である。HG番号1102は、仮想ボックス内のポート29Aに関連付けられたホストグループの番号である。ポート番号1003は、仮想ボックス内のポート29Aの番号である。LDEV番号1104は、仮想ボックス内のLDEV(仮想ボックス内のポート29Aに関連付けられたLDEV)の番号である。

20

【0076】

図10は、相対ポートID管理テーブル217Aの構成を示す。

【0077】

相対ポートID管理テーブル217Aは、ホストグループ毎に、ホストグループに関連付けられたポート29Aの相対ポートIDを表す。具体的には、相対ポートID管理テーブル217Aは、ホストグループ毎に、以下の情報、すなわち、HG番号1201及び相対ポートID1202を有する。HG番号1201は、ホストグループの番号である。相対ポートID1202は、ホストグループが関連付けられているポート29Aに割り振られた相対ポートIDである。

30

【0078】

図8～図10に示したテーブル215A～217Aを参照することで、ポート番号、LUN、LDEV番号及びHG番号の対応関係を把握することができる。

【0079】

図11は、有効ID範囲管理テーブル218Aの構成を示す。

【0080】

有効ID範囲管理テーブル218Aは、仮想ボックス毎に、付与可能な相対ポートIDの値(番号)の範囲を表す。具体的には、有効ID範囲管理テーブル218Aは、仮想ボックス毎に、仮想ボックス番号1301、開始ID1302及び終了ID1303を有する。仮想ボックス番号1301は、仮想ボックスの番号である。開始ID1302は、付与可能な相対ポートID範囲のうちの最初のID(値)である。終了ID1303は、付与可能な相対ポートID範囲のうちの最後のID(値)である。

40

【0081】

以上が、テーブル211A～218A及びリスト223Aの説明である。なお、テーブル211A～218A及びリスト223Aのうち少なくとも1つにおいて、少なくとも1種類の情報に代えて又は加えて、他種の情報が登録されてもよい。例えば、ペア管理テーブル212Aに、ペア相手ストレージ装置のID等が登録されてもよい。

50

【 0 0 8 2 】

図 1 2 は、ホスト装置 1 0 がストレージ装置 2 0 のパス優先度を把握するまでの流れを示す。

【 0 0 8 3 】

ホスト装置 1 0 が、第 1 L D E V 3 0 1 A の L U N を含んだ I / O コマンドを、パス 3 0 9 A 経由で第 1 ストレージ装置 2 0 A に送信する。第 1 ストレージ装置 2 0 A は、その I / O コマンドを受信した場合 (S 1 2 0 1)、その I / O コマンドの L U N から特定される第 1 L D E V 3 0 1 A についての U A 管理テーブル 2 1 3 A (図 5 参照) が「 0 」であれば、その I / O コマンドに従う I / O を行い、G o o d 応答をホスト装置 1 0 に送信する (S 1 2 0 2)。なぜなら、第 1 L D E V 3 0 1 A について U A 設定無しのためである。

10

【 0 0 8 4 】

その後、第 1 L D E V 3 0 1 A が関連付いているホストグループについてパス優先度の変更され、それに伴い、第 1 ストレージ装置 2 0 A により、そのホストグループについて U A 「 1 」が U A 管理テーブル 2 1 3 A (図 5 参照) に設定され、且つ、第 1 L D E V 3 0 1 A について U A 同期フラグ「 1 」が U A 同期フラグビットリスト 2 2 3 A (図 6 参照) に設定されたとする (S 1 2 0 3)。

【 0 0 8 5 】

その後、再び、第 1 ストレージ装置 2 0 A が、第 1 L D E V 3 0 1 A の L U N を含んだ I / O コマンドを、パス 3 0 9 A 経由で受信したとする (S 1 2 0 4)。この場合、第 1 L D E V 3 0 1 A が関連付いているホストグループの U A が「 1 」なので (正確には U A が「 0 」以外なので)、第 1 ストレージ装置 2 0 A が、その I / O コマンドに従う I / O を行った後、その U A 「 1 」を含んだ応答をホスト装置 1 0 に送信する (S 1 2 0 5)。

20

【 0 0 8 6 】

ホスト装置 1 0 は、U A 「 1 」を含んだ応答を受信した場合に、S 1 2 0 4 の I / O コマンド内の L U N についてパス優先度の変更があったことに気付く。ホスト装置 1 0 は、S 1 2 0 4 の I / O コマンド内の L U N と同一の L U N を含んだ A L U A 適用状態問合せを第 1 ストレージ装置 2 0 A に送信する。第 1 ストレージ装置 2 0 A は、L U N を含んだ A L U A 適用状態問合せを受信した場合 (S 1 2 0 6)、その L U N に対応した第 1 L D E V 3 0 1 A についての A L U A 適用状態 4 0 2 (図 3 参照) が表す A L U A 適用状態を含んだ応答をホスト装置 1 0 に送信する (S 1 2 0 7)。

30

【 0 0 8 7 】

ホスト装置 1 0 は、その応答から A L U A 適用であることを特定した場合、S 1 2 0 4 の I / O コマンド内の L U N と同一の L U N を含んだポート番号問合せを第 1 ストレージ装置 2 0 A に送信する。第 1 ストレージ装置 2 0 A は、L U N を含んだポート番号問合せを受信した場合 (S 1 2 0 8)、その L U N に対応したポート番号をポート管理テーブル 2 1 5 A から特定し、特定したポート番号を含んだ応答をホスト装置 1 0 に送信する (S 1 2 0 9)。

【 0 0 8 8 】

ホスト装置 1 0 は、S 1 2 0 4 の I / O コマンド内の L U N と同一の L U N を含んだ R T P G (ReportTargetPortGroups) コマンドを送信する。第 1 ストレージ装置 2 0 A は、R T P G コマンドを受信した場合 (S 1 2 1 0)、図 3 0 を参照して説明した処理を行うことで、例えば図 1 3 に示すように、異なるストレージ装置 2 0 A 及び 2 0 B の異なるパス管理テーブル 2 1 4 A 及び 2 1 4 B の該当部分 (図 1 3 の太枠内の情報の全て) が R T P G 応答にマージされる。具体的には、第 1 ストレージ装置 2 0 A は、パス管理テーブル 2 1 4 A 内の相対ポート I D 及びパス優先度の組と、パス管理テーブル 2 1 4 B 内の相対ポート I D 及びパス優先度の組とを、R T P G 応答に設定する。第 1 ストレージ装置 2 0 A は、R T P G コマンド内の L U N に対応したポート番号をポート管理テーブル 2 1 5 A、相対ポート I D 管理テーブル、仮想ボックス管理テーブルから特定し、特定したポート番号を、R P G 応答に設定してもよい。第 1 ストレージ装置 2 0 A は、その R T P G 応答を

40

50

ホスト装置 10 に送信する (S 1 2 1 1)。これにより、ホスト装置 10 は、変更後のパス優先度を知ることができる。

【 0 0 8 9 】

以下、本実施例で行われる処理を説明する。なお、以下、説明を分かり易くするため、コマンド受信ストレージ装置を、第 1 ストレージ装置 20 A とし、ペア相手ストレージを、第 2 ストレージ装置 20 B とする。コマンド受信ストレージ装置とは、ストレージ管理サーバ 11 又はホスト装置 10 からコマンド (I / O コマンド又は何らかの問合せ) を受信するストレージ装置である。ペア相手ストレージ装置とは、コマンド受信ストレージ装置が受信したコマンドに関わる L D E V とペアを構成する L D E V を有するストレージ装置である。

10

【 0 0 9 0 】

図 1 4 は、A L U A 適用状態変更処理の流れを示す。

【 0 0 9 1 】

第 1 ストレージ装置 20 A は、ストレージ管理サーバ 11 から A L U A 適用状態変更コマンドを受信する (S 1 4 0 1)。A L U A 適用状態変更コマンドは、H G 番号を含む。以下、その H G 番号を、図 1 4 の説明において「対象 H G 番号」と言う。

【 0 0 9 2 】

第 1 ストレージ装置 20 A は、対象 H G 番号に対応する L D E V を、ポート管理テーブル 2 1 5 A から特定する (S 1 4 0 2)。第 1 ストレージ装置 20 A は、特定した L D E V の A L U A 適用状態 4 0 2 (図 3 参照) を、S 1 4 0 1 で受信したコマンドに従い変更する (S 1 4 0 3)。また、第 1 ストレージ装置 20 A は、対象 H G 番号に対応する U A 6 0 3 (図 5 参照) を、「 2 」 (A L U A 適用状態の変更を意味する値) に変更する (S 1 4 0 4)。また、第 1 ストレージ装置 20 A は、S 1 4 0 2 で特定した L D E V に対応した A L U A 同期ビット 5 0 5 (図 4 参照) を「 1 」 (A L U A 適用状態の変更を意味する値) に変更する (S 1 4 0 5)。そして、第 1 ストレージ装置 20 A は、S 1 4 0 1 で受信したコマンドの応答として完了応答をストレージ管理サーバ 11 に送信する (S 1 4 0 6)。

20

【 0 0 9 3 】

S 1 4 0 5 と非同期に、第 1 ストレージ装置 20 A は、A L U A 同期ビット 5 0 5 が「 1 」の L D E V の番号を含んだ変更通知を、第 2 ストレージ装置 20 B に送信する。その変更通知は、第 1 ストレージ装置 20 A の L D E V の番号 (又は L U N)、又は、その L D E V のペア相手の L D E V の番号を含む。第 2 ストレージ装置 20 B は、その変更通知が表す L D E V に対応した A L U A 適用状態 (A L U A 管理テーブル 2 1 1 A の A L U A 適用状態) を「 1 」に変更する (S 1 4 0 7)。また、第 2 ストレージ装置 20 B は、その変更通知内の対象 H G 番号に対応した U A (U A 管理テーブル 2 1 3 B の U A 6 0 3) を「 2 」に変更する (S 1 4 0 8)。第 2 ストレージ装置 20 B は、S 1 4 0 8 と非同期に、その変更通知に対して完了通知を S 1 4 0 8 の完了後に第 1 ストレージ装置 10 A に送信する。完了通知は、A L U A 適用状態の変更が完了した L D E V の番号 (第 2 ストレージ装置 20 B の L D E V の番号)、又は、その L D E V とペア相手の L D E V の番号 (第 1 ストレージ装置 20 A の L D E V の番号) を含んでよい。

30

40

【 0 0 9 4 】

第 1 ストレージ装置 20 A は、完了通知を受信した場合、すなわち、第 2 ストレージ装置 20 B において A L U A 適用状態の変更が完了した場合、変更完了した L D E V に対応した A L U A 同期ビット 5 0 5 (図 4 参照) を「 0 」に変更する (S 1 4 0 9)。

【 0 0 9 5 】

なお、V D E V の基になる L D E V 間でペアが構成されていない場合、その V D E V の基になる各 L D E V について A L U A 同期ビットを「 1 」に変更する必要は無い。

【 0 0 9 6 】

図 1 5 は、A L U A 適用状態確認処理の流れを示す。

【 0 0 9 7 】

50

第1ストレージ装置20Aは、ストレージ管理サーバ11からALUA適用状態問合せを受信する(S1501)。ALUA適用状態問合せは、HG番号を含む。以下、そのHG番号を、図15の説明において「対象HG番号」と言う。

【0098】

第1ストレージ装置20Aは、対象HG番号に対応するLDEVを、ポート管理テーブル215Aから特定する(S1502)。第1ストレージ装置20Aは、特定したLDEVのALUA適用状態402(図3参照)を、S1501で受信した問合せに従い特定する(S1503)。

【0099】

第1ストレージ装置20Aは、特定したALUA適用状態を含んだ応答をストレージ管理サーバ11に送信する(S1504)。

10

【0100】

なお、S1504で送信された応答に含まれているALUA適用状態は、S1502で特定されたLDEVとペアを構成するLDEVのALUA適用状態でもある。なぜなら、図14の処理により、ペア相手LDEVのALUA適用状態も同じ状態に変更されているからである。従って、管理者(ストレージ管理サーバ11を使用するユーザ)は、LDEVペアを構成する2つのLDEVのうち一方のLDEVのALUA適用状態を確認すればよい。

【0101】

図16は、パス優先度変更処理の流れを示す。

20

【0102】

第1ストレージ装置20Aは、ストレージ管理サーバ11からパス優先度変更コマンドを受信する(S1601)。パス優先度変更コマンドは、HG番号を含む。以下、そのHG番号を、図16の説明において「対象HG番号」と言う。

【0103】

第1ストレージ装置20Aは、対象HG番号に対応するLDEVを、ポート管理テーブル215Aから特定する(S1602)。第1ストレージ装置20Aは、特定したLDEVのUA同期フラグビット702(図6参照)を、「1」(パス優先度変更を意味する値)に変更する(S1603)。第1ストレージ装置20Aは、対象HG番号に対応するパス優先度804(図7参照)を、S1601で受信したコマンドに従い変更する(S1604)。第1ストレージ装置20Aは、対象HG番号に対応したUA603(図5参照)を「1」(パス優先度変更を意味する値)に変更する(S1605)。

30

【0104】

第1ストレージ装置20Aは、完了を表す応答をストレージ管理サーバ11に送信する(S1606)。

【0105】

なお、上述したように、各ストレージ装置20は、パス優先度変更を意味する値のUAが存在する場合にそのUAに対応したLDEVをI/O先としたI/Oコマンドを受信すると、UAを含んだI/O応答(I/Oコマンドの応答)を送信することになる。しかし、パス優先度変更処理の最中に、第2ストレージ装置20B(ペア相手ストレージ装置)が、そのようなI/Oコマンドを受信しても、UAを含んだI/O応答を送信しない。なぜなら、第2ストレージ装置20Bは未だUA変更を知らないからである。しかし、それは、以下の理由により問題ではない。

40

(1) パス優先度変更後にUA変更が第2ストレージ装置20Bに通知されていないということは、パス優先度変更後に第1ストレージ装置20Aが、UA変更に関わるLDEVを指定したI/Oコマンドを受信していないということ。従って、その時点では、UA変更に関わるLDEVペアを構成する2つのLDEV内のデータは同じである。このため、どちらのLDEVからデータが読み出されてもよい。

(2) パス優先度変更後に第1ストレージ装置20Aは、I/Oコマンドを受信すると、UAを含んだI/O応答がホスト装置10に送信する。それにより、ホスト装置10は、

50

パス優先度の変更があったことを知ることでき、どこのパス優先度に変更されたかを R T P G コマンドで確認できる。

(なお、U A 変更通知 (伝播) 前に第 2 ストレージ装置 2 0 B がライトコマンドを受信しそれにより第 2 ストレージ装置 2 0 B 内の L D E V (ペア相手の L D E V) 内のデータが新しくなったとしても、それは、以下の理由により問題ではない。すなわち、同期コピーのため、U A 設定無しの L D E V を指定したライトコマンドを第 2 ストレージ装置 2 0 B が受信した場合、その U A 設定無しの L D E V とペアの L D E V (U A 設定有の L D E V) を有する第 1 ストレージ装置 2 0 A に対して、ストレージ装置 20B からライトデータが伝播される。ストレージ装置 20A はストレージ装置 20B からのライトデータ伝播を契機に、
図 1 7 の U A 伝播を実施する。それにより、U A 設定無の L D E V (第 2 ストレージ装置 2 0 B) にも U A が伝播され、第 2 ストレージ装置 2 0 B がホストに U A を応答することができる。この場合、図 1 7 における S 1 7 0 1 (ホスト装置からの I / O コマンド受信) を、「第 2 ストレージ装置からのコピーコマンド受信」に読み替えることができる。)

10

【 0 1 0 6 】

図 1 7 は、I / O 処理の際の U A 伝播の流れを示す。

【 0 1 0 7 】

第 1 ストレージ装置 2 0 A は、ホスト装置 1 0 から I / O コマンドを受信する (S 1 7 0 1)。I / O コマンドは L U N を含んでいる。

【 0 1 0 8 】

第 1 ストレージ装置 2 0 A は、U A 同期フラグビット 7 0 2 が「 1 」である L D E V が存在するか否かを、U A 同期フラグビットリスト 2 2 3 A (図 6 参照) を基に判断する (S 1 7 0 2)。

20

【 0 1 0 9 】

S 1 7 0 2 の判断結果が否定の場合 (S 1 7 0 2 : N o)、第 1 ストレージ装置 2 0 A は、ホスト装置 1 0 に G o o d 応答 (G o o d を表す I / O 応答) を送信する (S 1 7 0 3)。

【 0 1 1 0 】

一方、S 1 7 0 2 の判断結果が肯定の場合 (S 1 7 0 2 : Y e s)、第 1 ストレージ装置 2 0 A は、ペア管理テーブル 2 1 2 A を基に、U A 同期フラグビット「 1 」の L D E V のペア相手 L D E V を特定する (S 1 7 0 4)。なお、I / O 先 L D E V (I / O コマンド内の L U N に対応した L D E V) の U A 同期フラグビットが「 0 」であっても、U A 同期フラグビットリスト 2 2 3 A 上の少なくとも 1 つの U A 同期フラグビットが「 1 」であれば、S 1 7 0 2 の判断結果は肯定である。また、S 1 7 0 4 では、U A 同期フラグが「 1 」である 1 以上の L D E V について 1 以上のペア相手 L D E V が特定され、特定された 1 以上のペア相手 L D E V を有する 1 以上のストレージ装置の各々について、後述の S 1 7 0 5 ~ S 1 7 0 7 が行われる。ここでは、説明を分かり易くするため、U A 同期フラグが「 1 」である L D E V は 1 つであり、故に、ペア相手 L D E V も 1 つであり、ペア相手ストレージ装置は第 2 ストレージ装置 2 0 B であるとする。

30

【 0 1 1 1 】

第 1 ストレージ装置 2 0 A は、S 1 7 0 4 で特定したペア相手 L D E V を有する第 2 ストレージ装置 2 0 B に、U A 同期フラグビットリスト 2 2 3 A を送信する (S 1 7 0 5)。つまり、第 1 ストレージ装置 2 0 A は、第 1 ストレージ装置 2 0 A の U A 変更を第 2 ストレージ装置 2 0 B に通知する。なお、送信される U A 同期フラグビットリスト 2 2 3 A は、I / O コマンドを受信した C H A 2 2 A 内の L M 2 2 A 上のリスト 2 2 3 A である。リスト 2 2 3 A の中身は、全ての C H A 2 2 A で同じでもよい。或いは、アクセス可能な L D E V (受信する I / O コマンドで I / O 先として指定され得る L D E V) が C H A 2 2 A によって異なっていれば、C H A 2 2 A によって、リスト 2 2 3 A の中身が違っていてもよい。

40

【 0 1 1 2 】

第 2 ストレージ装置 2 0 B は、U A 同期フラグビットリスト 2 2 3 A を受信し、受信し

50

たリスト 2 2 3 A から U A 同期フラグビットが「1」の L D E V を特定し、特定した L D E V をペア相手 L D E V とする L D E V をペア管理テーブル 2 1 2 B を基に特定し、特定した L D E V に関連付いている H G 番号をポート管理テーブル 2 1 5 B を基に特定し、特定した H G 番号に対応する U A (U A 管理テーブル 2 1 3 B における U A) を「1」(パス優先度変更を意味する値)に変更する (S 1 7 0 6)。これにより、第 1 ストレージ装置 2 0 A の U A 変更が第 2 ストレージ装置 2 0 B に通知された (通知完了) となる。第 2 ストレージ装置 2 0 B は、完了を第 1 ストレージ装置 2 0 A に通知する。

【 0 1 1 3 】

第 1 ストレージ装置 2 0 A は、完了の通知を第 2 ストレージ装置 2 0 B から受信した場合、U A 同期フラグビットリスト 2 2 3 A における全ての U A 同期フラグビット「1」を「0」に変更する (S 1 7 0 7)。そして、第 1 ストレージ装置 2 0 A は、U A 応答 (U A = 1 を含んだ I / O 応答) をホスト装置 1 0 に送信する (S 1 7 0 8)。

10

【 0 1 1 4 】

図 1 7 の I / O 処理によれば、I / O 先 L D E V の U A 同期フラグビットが「0」であっても、少なくとも 1 つの他の L D E V の U A 同期フラグビットが「1」でありさえすれば、ホスト装置 1 0 に U A 応答が送信される。しかし、必ずしもそうである必要は無く、例えば、下記変形例が考えられる。

(変形例 1) S 1 7 0 2 において、第 1 ストレージ装置 1 0 は、I / O 先 L D E V の U A 同期フラグビットが「1」であるか否かを判断する。I / O 先 L D E V の U A 同期フラグビットが「0」の場合、少なくとも 1 つの他の L D E V の U A 同期フラグビットが「1」であっても、S 1 7 0 2 の判断結果が否定であり、I / O 先 L D E V の U A 同期フラグビットが「1」の場合、S 1 7 0 2 の判断結果が肯定である。この変形例 1 では、S 1 7 0 2 の判断結果が肯定の場合、S 1 7 0 5 で送信される U A 同期フラグビットは、I / O 先 L D E V に対応した U A 同期フラグビットのみでよい。S 1 7 0 7 で「0」に変更される U A 同期フラグビットは、I / O 先 L D E V に対応したビットのみでよい。

20

(変形例 2) S 1 7 0 2 の判断結果が肯定の場合、S 1 7 0 5 で送信される U A 同期フラグビットは、U A 同期フラグビットリスト 2 2 3 A のうち、I / O コマンドの送信元ホスト装置を含んだ H G に関連付いた L D E V に対応した U A 同期フラグビットのみである。

【 0 1 1 5 】

図 1 8 A 及び図 1 8 B は、ペア状態とパス優先度に関連性があることの一例を示す。具体的には、図 1 8 A は、ペア状態「P A I R」の場合の適切パスの一例を示し、図 1 8 B は、ペア状態「S U S P E N D」の場合の適切パスの一例を示す。

30

【 0 1 1 6 】

図 1 8 A に示すように、ペア状態「P A I R」(且つ I / O モード「ミラー」) の場合、ホスト装置 1 0 は、パス管理テーブルのパス優先度に従い、パス 3 0 9 A を使用する、すなわち、第 1 ストレージ装置 2 0 A の L D E V に対して I / O を行うのが望ましい。

【 0 1 1 7 】

しかし、ペア状態が、例えば図 1 8 B に示すように、ペア状態「S U S P E N D」(且つ I / O モード「リモート」) に変更された場合、パス管理テーブルのパス優先度に関わらず、パス 3 0 9 B を使用する、すなわち、第 2 ストレージ装置 2 0 B の L D E V に対して I / O を行うのが望ましい。従って、このケースでは、パス優先度が、パス 3 0 9 A についてパス優先度「Active/non-optimized」に変更され、パス 3 0 9 B についてパス優先度「Active/optimized」に変更されることが望ましい。本実施例では、次の流れにより、ホスト装置は、パス優先度を知ることができる。

40

- (1) ペア状態により、ストレージ装置で管理されているパス優先度に変更される。
- (2) ペア状態変更によるパス優先度変更の際も U A が設定される。具体的には、図 1 9 の処理が行われた後に、I / O コマンドを受信する等の U A 変更イベント (図 1 7 の S 1 7 0 1) が行われると、図 1 7 の処理により、U A が伝播する。
- (3) U A 応答により、ホスト装置は、パス優先度変更を知る。
- (4) ホスト装置は、R T P G コマンドを送信することにより、パス優先度を知る。

50

(なお、図18Bに破線矢印で示すように、パス309A経由で第2ストレージ装置20BのLDEVに対してI/Oを行うことも可能ではあるが、パス309Bを使用することに比べて性能が低い。)

【0118】

そこで、本実施例では、ペア状態の変更を行ったストレージ装置(例えば20A)は、ペア状態の変更に伴い、UAの値を「1」(パス優先度変更を意味する値)に変更する。UA「1」の場合、後述するように、やがてホスト装置10からRTPGコマンドを第1ストレージ装置20Aが受信することになるが、RTPGコマンドの受信を含むパス優先度レポート処理では、ペア状態及びI/Oモードに応じたパス優先度を表すレポートを含んだRTPG応答が作成されホスト装置10に対して送信される。これにより、ホスト装置10が保持するパス管理情報が表すパス優先度が、ペア状態及びI/Oモードに応じたパス優先度に変更され、以後、ホスト装置10により、その変更後のパス管理情報を基に、パスが使用されることになる。

10

【0119】

図19は、パス優先度変更処理の流れを示す。

【0120】

第1ストレージ装置20Aは、ストレージ管理サーバ11からペア状態変更コマンドを受信する(S1901)。ペア状態変更コマンドは、ペア状態の変更対象のLDEVペアのうち第1ストレージ装置20Aが有するLDEVの番号を含む。以下、そのLDEVの番号を、図19の説明において「対象LDEV番号」と言う。

20

【0121】

第1ストレージ装置20Aは、対象LDEV番号に対応したHG番号をポート管理テーブル215Aから特定し、特定したHG番号に対応したUA603(図5参照)を「1」(パス優先度変更を意味する値)に変更する(S1902)。また、第1ストレージ装置20Aは、対象LDEV番号に対応したUA同期フラグビット702(図5参照)を「1」に変更する(S1903)。

【0122】

第1ストレージ装置20Aは、完了を表す応答をストレージ管理サーバ11に送信する(S1904)。

【0123】

なお、S1902でUAが「1」に変更されたが、パス優先度変更処理においては、そのUA変更は、ペア相手ストレージ装置(20B)に通知されない。第1ストレージ装置20AがI/Oコマンドを受信等のUA伝搬イベント発生した場合に(図17参照)、UA変更がペア相手ストレージ装置に通知される。

30

【0124】

図20~図23は、パス優先度レポート処理の流れを示す。

【0125】

図20に示すように、第1ストレージ装置20Aが、RTPGコマンドをホスト装置10から受信する(S2001)。RTPGコマンドは、LUNを含む。このLUNは、例えば、UA応答に対応したI/Oコマンド内のLUNと同一のLUNである。すなわち、ホスト装置10は、I/Oコマンドの応答としてUA応答を受信した場合、そのI/Oコマンド内のLUNと同一LUNを含んだRTPGコマンドを、そのI/Oコマンドの送信先ストレージ装置と同一のストレージ装置に送信できる。なお、ホスト装置10は、RTPGコマンドを、UA応答を受信したというイベントとは別のイベントに回答して(例えば定期的に)送信してよい。以下、図20~図23の説明において、RTPGコマンド内のLUNに対応したLDEV(第1ストレージ装置20Aが有するLDEV)を「対象LDEV」と言い、対象LDEVの番号を、「対象LDEV番号」と言う。

40

【0126】

第1ストレージ装置20Aは、対象LDEVがVDEVの基になるLDEVペアを構成しているか否かを、ペア管理テーブル212Aを基に判断する(S2002)。対象LD

50

E V 番号がペア管理テーブル 2 1 2 A にあれば、S 2 0 0 2 の判断結果が肯定であり、対象 L D E V 番号がペア管理テーブル 2 1 2 A に無ければ、S 2 0 0 2 の判断結果が否定である。

【 0 1 2 7 】

S 2 0 0 2 の判断結果が否定の場合 (S 2 0 0 2 : N o)、第 1 ストレージ装置 2 0 A は、自ストレージ装置内のパス優先度 (例えば、図示しないテーブルから特定されたパス優先度) を含んだ応答をホスト装置 1 0 に送信する (S 2 0 0 5)。

【 0 1 2 8 】

S 2 0 0 2 の判断結果が肯定の場合 (S 2 0 0 2 : Y e s)、第 1 ストレージ装置 2 0 A は、対象 L D E V 番号に対応した I / O モードを、ペア管理テーブル 2 1 2 A から特定する (S 2 0 0 3)。

10

【 0 1 2 9 】

S 2 0 0 4 で特定された I / O モードが「 B l o c k」の場合、第 1 ストレージ装置 2 0 A は、ホスト装置 1 0 にエラー応答を送信する (S 2 0 0 6)。S 2 0 0 4 で特定された I / O モードが「 R e m o t e」の場合、第 1 ストレージ装置 2 0 A は、図 2 1 に示す処理を行う。S 2 0 0 4 で特定された I / O モードが「 M i r r o r」の場合、第 1 ストレージ装置 2 0 A は、図 2 2 に示す処理を行う。S 2 0 0 4 で特定された I / O モードが「 L o c a l」の場合、第 1 ストレージ装置 2 0 A は、図 2 3 に示す処理を行う。

【 0 1 3 0 】

図 2 1 に示すように、I / O モードが「リモート」の場合、第 1 ストレージ装置 2 0 A は、対象 L D E V に関連付けられている全てのポートの情報 (例えばポート番号及び相対ポート I D) と、対象 L D E V のペア相手 L D E V (ペア相手ストレージ) の情報とを、対象 L D E V 番号を用いて、ペア管理テーブル 2 1 2 A、相対ポート I D 管理テーブル 2 1 7 A、仮想ボックス管理テーブル 2 1 6 A 及びポート管理テーブル 2 1 5 A から特定する (S 2 1 0 1)。特定される情報 (ポートの情報及びペア相手 L D E V の情報) は、例えば、対象 L D E V を含んだ仮想ボックスの範囲とされる。これは、図 2 2 (S 2 2 0 3) 及び図 2 3 (S 2 3 0 3) でも同様である。以下、図 2 1 ~ 図 2 3 の説明において、対象 L D E V のペア相手ストレージ装置は、第 2 ストレージ装置 2 0 B とする。

20

【 0 1 3 1 】

第 1 ストレージ装置 2 0 A は、特定した第 2 ストレージ装置 2 0 B (ペア相手ストレージ装置) と通信可能か否かを判断する (S 2 1 0 2)。例えば、第 1 ストレージ装置 2 0 A は、所定のコマンドを第 2 ストレージ装置 2 0 B に送信して所定時間内に応答が返ってくるか否かを判断する。

30

【 0 1 3 2 】

S 2 1 0 2 の判断結果が否定の場合 (S 2 1 0 2 : N o)、第 1 ストレージ装置 2 0 A は、ホスト装置 1 0 にエラー応答を送信する (S 2 1 0 3)。

【 0 1 3 3 】

S 2 1 0 2 の判断結果が肯定の場合 (S 2 1 0 2 : Y e s)、第 1 ストレージ装置 2 0 A は、U A 同期フラグビットリスト 2 2 3 A と、パス優先度要求とを、第 2 ストレージ装置 2 0 B に送信する (S 2 1 0 4)。

40

【 0 1 3 4 】

第 2 ストレージ装置 2 0 B は、対象 L D E V 番号と、U A 同期フラグビットリスト 2 2 3 A と、パス優先度要求とを受信する。第 2 ストレージ装置 2 0 B は、受信した U A 同期フラグビットリスト 2 2 3 A を基に、U A 同期フラグビットが「 1」の全ての L D E V の各々について、H G 番号を特定し、特定した H G 番号に対応した U A を「 1」に変更する (S 2 1 0 5)。また、第 2 ストレージ装置 2 0 B は、パス優先度要求に応答して、対象 L D E V に対応したペア相手 L D E V に関連付いた H G 番号のパス優先度 (パス管理テーブル 2 1 4 B から特定されたパス優先度) を表すパス優先度情報、第 1 ストレージ装置 2 0 A に通知する (S 2 1 0 6)。

【 0 1 3 5 】

50

第1ストレージ装置20Aは、パス優先度情報を第2ストレージ装置20Bから受信した場合、UA同期フラグビットリスト223Aにおける全てのUA同期フラグビット「1」を「0」に変更し(S2107)、第1パス優先度(対象LDEVを含む1以上のLDEV(第1ストレージ装置20A内のLDEV)についてのパス優先度)と第2パス優先度(ペア相手LDEVを含む1以上のLDEV(第2ストレージ装置20B内のLDEV)についてのパス優先度)とをマージしたレポートを含むRTPG応答をホスト装置10に送信する(S2108)。但し、そのレポートでは、第1ストレージ装置20Aが有する第1パス優先度情報(対象LDEVのHG番号に対応したパス優先度を表す情報)が表すパス優先度に関わらず、第1パス優先度は「Active/non-optimized」とされ、第2ストレージ装置20Bから受信した第2パス優先度情報(ペア相手LDEVのHG番号に対応したパス優先度を表す情報)が表すパス優先度に関わらず、第2パス優先度は「Active/optimized」とされる。つまり、以後、ペア相手LDEVへのパスが優先的に使用されるようにする。なぜなら、I/Oモードが「リモート」の場合、ペア相手LDEV内のデータの方が対象LDEV内のデータより新しい可能性があり信頼できるデータであるためである。

10

【0136】

なお、S2104では、UA同期フラグビットリスト223A全体ではなく、UA同期フラグビットリスト223Aのうち対象LDEVのUA同期フラグビットが送信されてもよい。このため、S2105において、UA変更は、対象LDEVのペア相手LDEVに関連付いたHG番号についてのみ行われてよい。また、S2104において、UA同期フラグビットリスト223Aに少なくとも1つのUA同期フラグビット「1」があるか否か(或いは、うち対象LDEVのUA同期フラグビットが「1」であるか否か)が判断され、その判断の結果が肯定の場合、UA同期フラグビットリスト(又は対象LDEVのUA同期フラグビット)が第2ストレージ装置20Bに送信され、その判断の結果が否定の場合、UA同期フラグビットリスト(又は対象LDEVのUA同期フラグビット)が第2ストレージ装置20Bに送信されず且つS2105がスキップされてよい。これは、図22(S2206及びS2207)及び図23(S2306及びS2307)についても同様である。

20

【0137】

図22に示すように、I/Oモードが「ミラー」の場合、第1ストレージ装置20Aは、対象LDEVが閉塞しているか否かを判断する(S2201)。

30

【0138】

S2201の判断結果が肯定の場合(S2101:Yes)、第1ストレージ装置20Aは、ホスト装置10にエラー応答を送信する(S2202)。

【0139】

S2201の判断結果が否定の場合(S2101:No)、第1ストレージ装置20Aは、図21のS2101と同じ処理を行う(S2203)。

【0140】

第1ストレージ装置20Aは、特定した第2ストレージ装置20B(ペア相手ストレージ装置)と通信可能か否かを判断する(S2204)。

40

【0141】

S2204の判断結果が否定の場合(S2102:No)、第1ストレージ装置20Aは、第1パス優先度情報が表すパス優先度に関わらず第1パス優先度を「Active/optimized」としたレポートを含むRTPG応答をホスト装置10に送信する(S2205)。なぜなら、第2ストレージ装置20Bと通信できないので、以後、対象LDEVへのパスが優先的に使用されるようにするためである。

【0142】

S2204の判断結果が肯定の場合(S2204:Yes)、図21の2104~S2107と同じ処理が行われる(S2206~S2209)。第1ストレージ装置20Aは、第1パス優先度及び第2パス優先度をマージしたレポートを含むRTPG応答をホスト

50

装置 10 に送信する (S 2 2 1 0)。そのレポートでは、第 1 パス優先度は、第 1 ストレージ装置 20 A が有する第 1 パス優先度情報が表すパス優先度通りであり、第 2 パス優先度は、第 2 ストレージ装置 20 B から受信した第 2 パス優先度情報が表すパス優先度通りである。

【 0 1 4 3 】

図 2 3 に示すように、I / O モードが「ローカル」の場合、図 2 2 の S 2 2 0 1 ~ S 2 2 0 9 と同じ処理が行われる (S 2 3 0 1 ~ S 2 3 0 9)。第 1 ストレージ装置 20 A は、第 1 パス優先度と第 2 パス優先度とをマージしたレポートを含む R T P G 応答をホスト装置 10 に送信する (S 2 3 1 0)。但し、そのレポートでは、第 1 ストレージ装置 20 A が有する第 1 パス優先度情報が表すパス優先度に関わらず、第 1 パス優先度は「Active / optimized」とされ、第 2 ストレージ装置 20 B から受信した第 2 パス優先度情報が表すパス優先度に関わらず、第 2 パス優先度は「Active / non-optimized」とされる。つまり、以後、対象 L D E V へのパスが優先的に使用されるようにする。なぜなら、I / O モードが「ローカル」の場合、対象 L D E V 内のデータの方がペア相手 L D E V 内のデータより新しい可能性があり信頼できるデータであるためである。

10

【 実施例 2 】

【 0 1 4 4 】

実施例 2 を説明する。その際、実施例 1 との相違点を主に説明し、実施例 1 との共通点については説明を省略又は簡略する。

【 0 1 4 5 】

まず、実施例 2 の概要を述べる。

20

【 0 1 4 6 】

ホスト装置 10 の数が多いと、ホストグループが増える可能性があり、そうすると、ホストグループ単位にパス優先度を設定するには手間がかかる。

【 0 1 4 7 】

そこで、実施例 2 では、パス優先度を設定する手間が、実施例 1 より軽減される。具体的には、実施例 2 では、管理者はパス優先度を設定する必要は無い。1 以上のホスト装置と 1 以上のストレージ装置でグループが定義される。本実施例では、そのグループを「データセンタ」と呼ぶ。ホスト装置と同一のデータセンタに属するストレージ装置との間の物理的な距離は、そのホスト装置と異なるデータセンタに属するストレージ装置との間の物理的な距離よりも短いことが望ましい。

30

【 0 1 4 8 】

データセンタ毎に I D が割り振られる。データセンタ I D は、管理者又はユーザ (仮想ボックスのユーザ) 等により割り振られてもよい。データセンタ I D は、そのデータセンタ I D が割り振られたデータセンタに属するホスト装置及びストレージ装置がそれぞれ保持する。ホスト装置は、同一データセンタ内のストレージ装置 (そのホスト装置が保持するデータセンタ I D と同一のデータセンタ I D を保持するストレージ装置) とそのストレージ装置内の L D E V との間のパスの優先度を、異なるデータセンタ内のストレージ装置 (そのホスト装置が保持するデータセンタ I D と異なるデータセンタ I D を保持するストレージ装置) とそのストレージ装置内の L D E V との間のパスの優先度よりも高くする。具体的には、例えば、ホスト装置 2 4 1 0 は、同一データセンタに属するストレージ装置 2 4 2 0 内の L D E V との間のパスの優先度を「Active / optimized」に決定し、異なるデータセンタに属するストレージ装置 2 4 2 0 内の L D E V との間のパスの優先度を「Active / non-optimized」に決定する。以上のことを、例えば下記のように表現できる。

40

少なくとも 1 つのホスト装置と、

前記複数のホスト装置に接続された複数のストレージ装置と

を有し、

前記複数のホスト装置の各々が、そのホスト装置が属するグループの I D を保持し、

前記複数のストレージ装置の各々が、そのストレージ装置が属するグループの I D を保持し、

50

各グループには、少なくとも1つのホスト装置と少なくとも1つのストレージ装置が属し、
 前記少なくとも1つのホスト装置の各々は、
 グループIDの問合せをいずれかのストレージ装置に送信し、
 前記問合せの応答として、前記問合せの送信先のストレージ装置が保持するグループID
 を含んだ応答を受信し、
 前記応答内のグループIDが、当該ホスト装置が保持するグループIDと同一であれば、
 前記応答の送信元ストレージ装置が有する全ての論理ボリュームへの1以上のパスの各々
 の優先度を、第1の優先度とし、
 前記応答内のグループIDが、当該ホスト装置が保持するグループIDと異なっていれば、
 前記応答の送信元ストレージ装置が有する全ての論理ボリュームへの1以上のパスの各々
 の優先度を、第1の優先度よりも低い第2の優先度とする、
 計算機システム。

10

【0149】

図24は、実施例2に係る計算機システムの構成、及び、各ストレージ装置が有する管理情報の構成を示す。

【0150】

複数のホスト装置2410と複数のストレージ装置2420がある。図24の例では、
 第1及び第2ホスト装置2410A及び2410Bと、第1及び第2ストレージ装置24
 20A及び2420がある。以下、説明を分かり易くするために、第1ホスト装置241
 0Aと第1ストレージ装置2420Aが第1データセンタに属し、第2ホスト装置241
 0Bと第2ストレージ装置2420Bが第2データセンタに属しているとする。

20

【0151】

ホスト装置2410は、データセンタID問合せをストレージ装置2420に送信する
 ことで、そのストレージ装置2420が保持するデータセンタIDを含んだ応答をストレ
 ージ装置2420から受信する。データセンタID問合せは、例えば図25に示す構成で
 よい。すなわち、データセンタID問合せは、オペレーションコードフィールドとデー
 タセンタIDフィールドとを有してよい。オペレーションコードフィールドには、オペレ
 ーションコードが記述される。オペレーションコードは、この問合せを受信したストレ
 ージ装置に実施させたい命令の種類をあらわすコードである。また、データセンタID
 問合せにおけるデータセンタIDフィールドは、ブランクでよく、応答の際に、この問
 合せを受信したストレージ装置2420が保持するデータセンタIDがそのストレージ装
 置2420に記述され、図示の問合せの構造体を含んだ応答(データセンタIDが記述され
 た構造体を含んだ応答)が、ストレージ装置2420からホスト装置2410に返される。

30

【0152】

ホスト装置2410は、記憶デバイスの一例としてメモリ101を有する。メモリ10
 1は、パス優先度管理テーブル102を記憶する。パス優先度管理テーブル102は、例
 えば図26に示すように、ストレージ装置毎に、ストレージ製番(ストレージの製造番
 号)2601、パス状態2602及びパス優先度2603を有する。つまり、本実施例では
 、パス優先度の設定は、ストレージ装置単位で行われる。すなわち、ホスト装置24
 10がアクセス可能な1つのストレージ装置内の複数のLDEVに対応した複数のパスにつ
 いて、まとめて、パス優先度が設定される。なお、ストレージ製番2601は、ストレ
 ージ装置のIDの一例である。パス状態2602の値としては、「Normal」(ストレージ装
 置と通信可能)と、「Blocked」(ストレージ装置と通信不可能)がある。ホスト装置24
 10は、パス状態「Blocked」のストレージ装置についてのパスの優先度が「Active/opti
 mized」でも、そのパスを使用せず、パス状態「Normal」の他のストレージ装置につ
 いてのパスを使用する。

40

【0153】

また、ホスト装置2410のメモリ101は、データセンタID管理テーブル103を
 記憶し、ストレージ装置2420(例えばSM21)も、データセンタID管理テーブル

50

219を記憶する。テーブル103及び219はそれぞれ構成が同じであり、また、同一データセンタ内のテーブル103及び219は保持する値も同じである。具体的には、第1ホスト装置2410Aと第1ストレージ装置2420Aが第1データセンタに属するが、テーブル103A及び219Aは、図27Aに示すように、同一のデータセンタID「0x0001」を保持する。同様に、第2ホスト装置2410Bと第2ストレージ装置2420Bが第2データセンタに属するが、テーブル103B及び219Bは、図27Bに示すように、同一のデータセンタID「0x0002」を保持する。

【0154】

以下、実施例2で行われる処理を説明する。

【0155】

図28は、バス優先度自動決定処理の流れを示す。なお、この処理は、ホスト装置2410内のバス管理プログラム(図30参照)により行うことができる。また、この処理の主体として、ホスト装置2410Aを例に取る。

【0156】

ホスト装置2410Aは、データセンタID問合せを、ストレージ装置2420に送信し(S2801)、そのストレージ装置2420から、そのストレージ装置2420が保持するデータセンタIDを含んだ応答を受信する(S2802)。

【0157】

ホスト装置2410Aは、受信した応答内のデータセンタIDと、データセンタID管理テーブル103A内のデータセンタIDとが同一か否かを判断する(S2803)。

【0158】

S2803の判断結果が否定の場合(S2803:No)、S2802で受信した応答の送信元ストレージ装置がホスト装置2410Aと異なるデータセンタに属しているということである。この場合、ホスト装置2410Aは、S2802で受信した応答の送信元ストレージ装置について、バス優先度2603として「Active/non-optimized」を設定する(S2804)。これは、応答の送信元ストレージ装置がストレージ装置2420Bの場合に行われる処理である。

【0159】

一方、S2803の判断結果が肯定の場合(S2803:Yes)、S2802で受信した応答の送信元ストレージ装置がホスト装置2410Aと同一データセンタに属しているということである。この場合、ホスト装置2410Aは、S2802で受信した応答の送信元ストレージ装置について、バス優先度2603として「Active/optimized」を設定する(S2805)。これは、応答の送信元ストレージ装置がストレージ装置2420Aの場合に行われる処理である。

【0160】

図29は、I/O制御処理の流れを示す。この処理は、ホスト装置2410内のバス管理プログラムにより行うことができる。また、この処理の主体として、ホスト装置2410Aを例に取る。この処理は、VDEVに対するI/Oを、例えばバス管理プログラムがアプリケーションプログラムから指示された場合に行われる。

【0161】

ホスト装置2410Aは、バス優先度2603が「Active/optimized」のストレージ装置をバス優先度管理テーブル102Aから検索する(S2901)。ホスト装置2410Aは、その検索により見つかったストレージ装置のバス状態2602が「Blocked」か否かを判断する(S2902)。

【0162】

S2902の判断結果が否定の場合(S2902:No)、ホスト装置2410Aは、I/O先VDEVの基であり見つかったストレージ装置のLDEVへのバスを選択し、選択したバスを使用してI/Oコマンドを送信する(S2903)。

【0163】

S2902の判断結果が肯定の場合(S2902:Yes)、ホスト装置2410Aは

10

20

30

40

50

、パス状態 2 6 0 2 が「Blocked」ではないストレージ装置があるか否かを、パス優先度管理テーブル 1 0 2 A を基に判断する (S 2 9 0 4) 。

【 0 1 6 4 】

S 2 9 0 4 の判断結果が肯定の場合 (S 2 9 0 4 : Y e s)、ホスト装置 2 4 1 0 A は、I / O 先 V D E V の基でありパス状態 2 6 0 2 が「Blocked」ではないいずれかのストレージ装置の L D E V へのパスを選択し、選択したパスを使用して I / O コマンドを送信する (S 2 9 0 5)。パス状態 2 6 0 2 が「Blocked」ではないストレージ装置として、パス優先度 2 6 0 3 が「Active/optimized」のストレージ装置とパス優先度 2 6 0 3 が「Active/non-optimized」のストレージ装置の両方があれば、パス優先度 2 6 0 3 が「Active/optimized」のストレージ装置が選択される。

10

【 0 1 6 5 】

S 2 9 0 4 の判断結果が否定の場合 (S 2 9 0 4 : N o)、ホスト装置 2 4 1 0 A は、通信不可で異常終了 (例えばエラーを発行) する (S 2 9 0 6) 。

【 0 1 6 6 】

以上、幾つかの実施例及び変形例を説明したが、本発明は、これらの実施例及び変形例に限定されるものでなく、その要旨を逸脱しない範囲で種々変更可能であることはいうまでもない。例えば、本発明は、A L U A に限らず、ストレージ装置が状態変更をホストシステムに通知することでホストシステムが状態変更を知ることができるようになっている環境下のシステムにも適用できる。

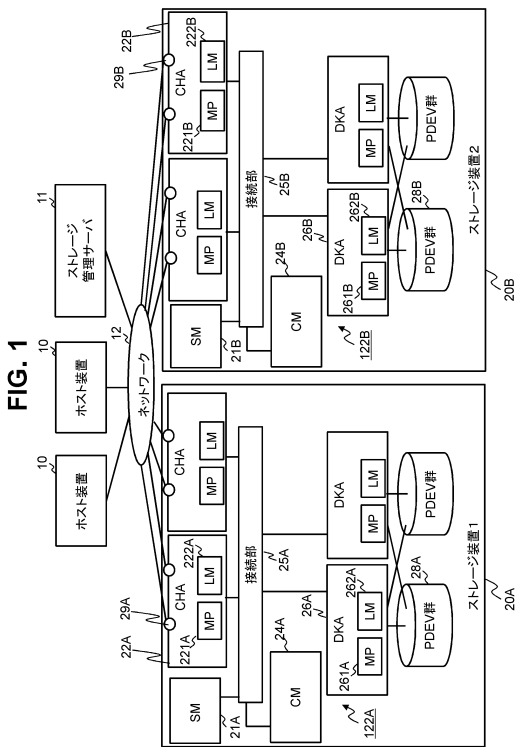
20

【 符号の説明 】

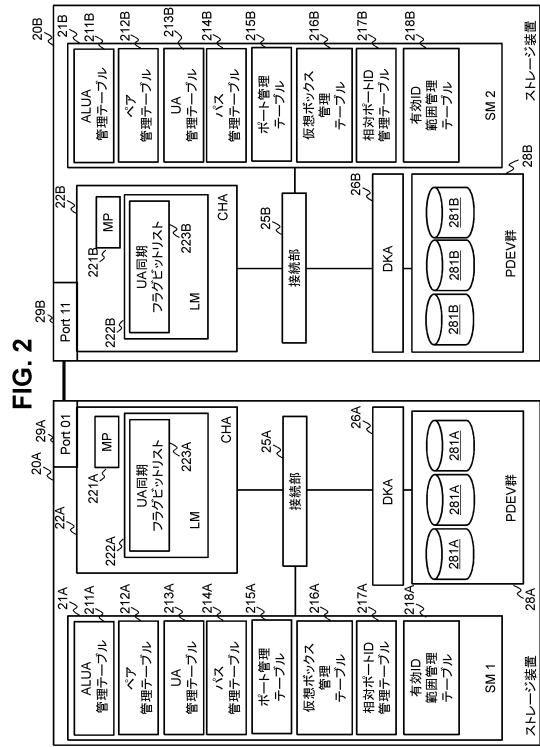
【 0 1 6 7 】

1 0 : ホスト装置 2 0 A : 第 1 ストレージ装置 2 0 B : 第 2 ストレージ装置

【 図 1 】



【 図 2 】



【 図 3 】

FIG. 3
ALUA管理テーブル
211A

LDEV番号	ALUA適用状態
01.01	0
01.02	1
....

【 図 4 】

FIG. 4
ペア管理テーブル
212A

LDEV番号 (自ストレージ)	LDEV番号 (ペア相手ストレージ)	ペア状態	I/Oモード	ALUA同期 ビット
01.01	02.01	PAIR	Mirror	0
01.02	02.02	SUSPEND	Local	0
01.03	02.03	SUSPEND	Remote	0
01.04	02.04	SUSPEND	Block	1
....

【 図 8 】

FIG. 8
ポート管理テーブル
215A

ポート番号	インデックス 番号	LUN	LDEV番号	HG番号
01	001	00	00.01	HG01
01	002	01	00.02	HG02
...
02	001	11	11.11	HG11
02	002	12	11.12	HG12
...

【 図 9 】

FIG. 9
仮想ボックス管理テーブル
216A

仮想ボックス番号	HG番号	ポート番号	LDEV番号
01	HG01	01	00.01
01	HG02	02	00.02
...

【 図 10 】

FIG. 10
相対ポートID管理テーブル
217A

HG番号	相対ポートID
HG01	0x0001
HG02	0x0002
...	...

【 図 5 】

FIG. 5
UA管理テーブル
219A

ポート番号	ホストグループ番号	UA
01	HG01	0
02	HG02	1
....

【 図 6 】

FIG. 6
UA同期フラグビットリスト
262A

LDEV番号	UA同期フラグビット
01.01	0
01.02	1
....

【 図 7 】

FIG. 7
バス管理テーブル
214A

仮想ボックス番号	HG番号	相対ポートID	バス優先度
01	HG01	0x0000	Active/optimized
01	HG02	0x0001	Active/non-optimized
...
02	HG10	0x2000	Active/optimized
02	HG11	0x2001	Active/non-optimized
...

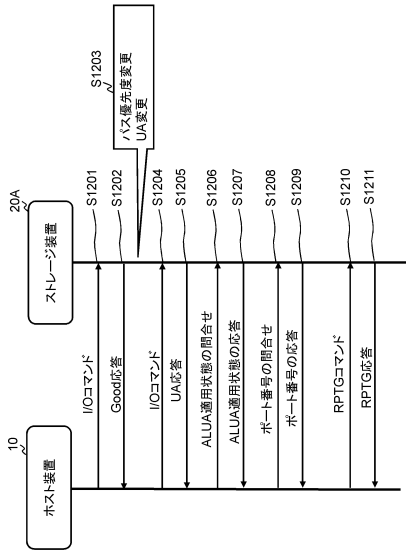
【 図 11 】

FIG. 11
有効ID範囲管理テーブル
218A

仮想ボックス番号	開始ID	終了ID
01	0x0001	0x1999
02	0x2000	0x2999
...

【 図 1 2 】

FIG. 12



【 図 1 3 】

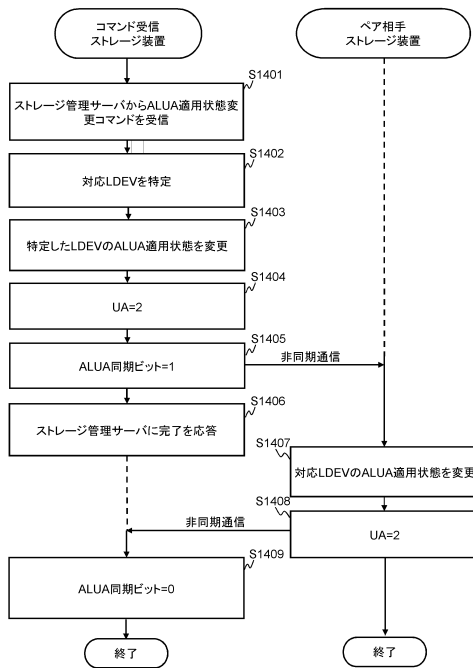
FIG. 13

FIG. 13 is a table with two sections, 214A and 214B, showing device status and optimization levels. The columns are: 装置のクラス番号 (Device Class Number), HG番号 (HG Number), 相対ポートID (Relative Port ID), ハス優先度 (Bus Priority), and ハス優先度 (Bus Priority).

装置のクラス番号	HG番号	相対ポートID	ハス優先度	ハス優先度
01	HG01	0x0000	Active/optimized	Active/non-optimized
01	HG02	0x0001	Active/optimized	Active/non-optimized
01
01	HG09	0x0999	Active/optimized	Active/non-optimized
02	HG10	0x2000	Active/optimized	Active/non-optimized
02	HG11	0x2001	Active/optimized	Active/non-optimized
...

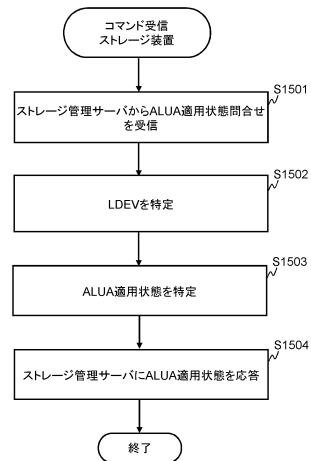
【 図 1 4 】

FIG. 14



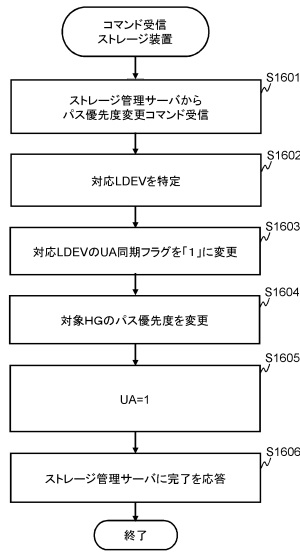
【 図 1 5 】

FIG. 15



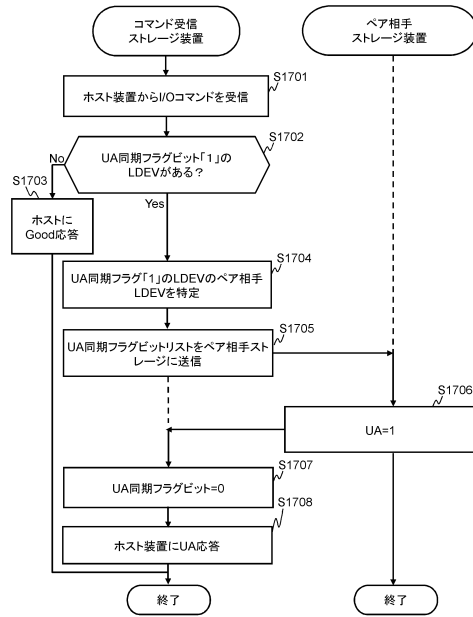
【 図 16 】

FIG. 16



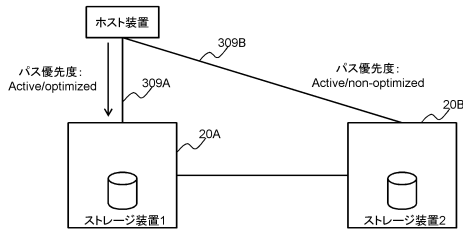
【 図 17 】

FIG. 17



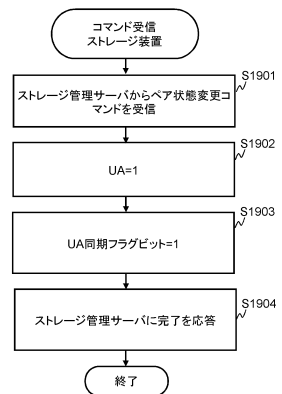
【 図 18 A 】

FIG. 18A



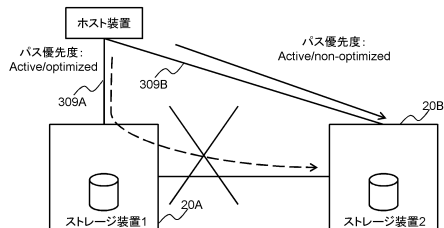
【 図 19 】

FIG. 19



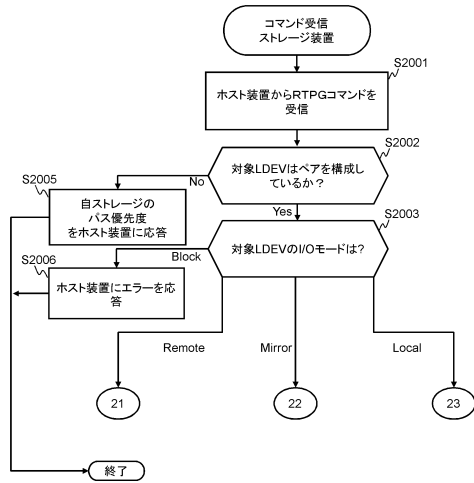
【 図 18 B 】

FIG. 18B



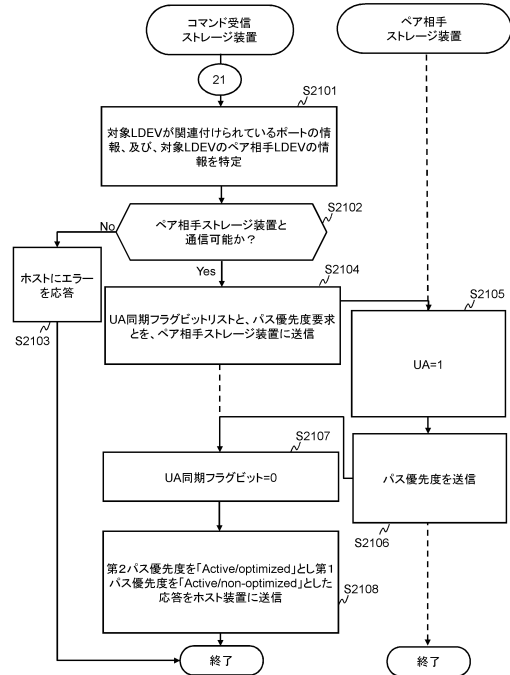
【図20】

FIG. 20



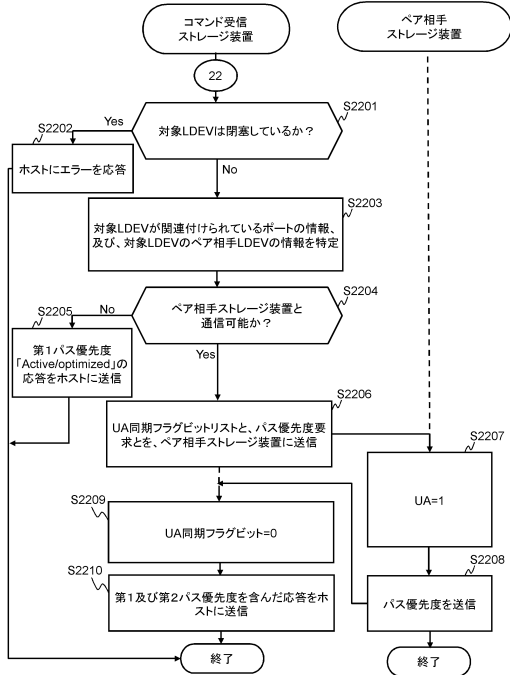
【図21】

FIG. 21



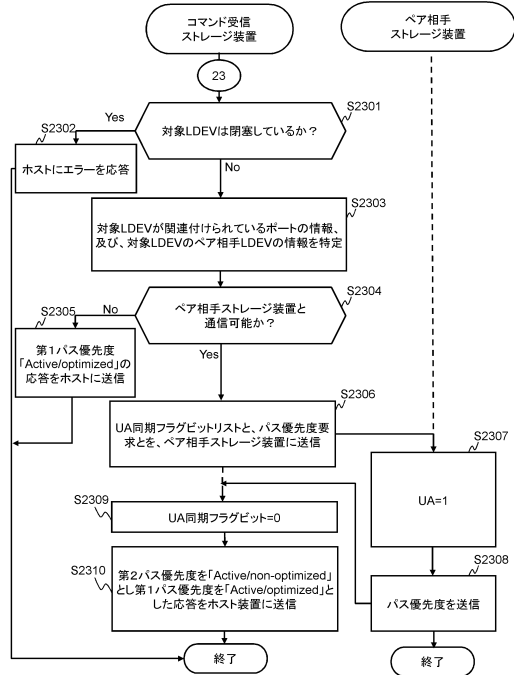
【図22】

FIG. 22

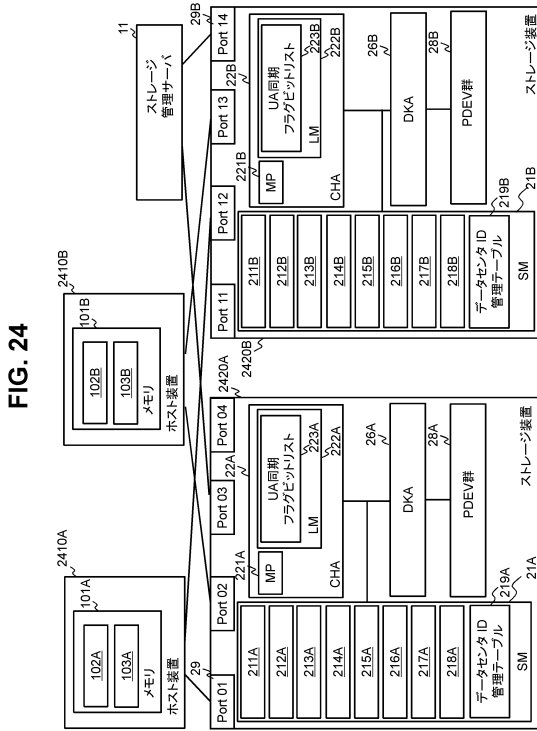


【図23】

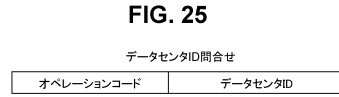
FIG. 23



【 図 2 4 】



【 図 2 5 】



【 図 2 6 】

FIG. 26 is a table titled 'バス優先度管理テーブル' (Bus Priority Management Table). It has three columns: 'ストレージ製番' (Storage Model Number), 'バス状態' (Bus Status), and 'バス優先度' (Bus Priority).

ストレージ製番	バス状態	バス優先度
11111	Normal	Active/optimized
22222	Normal	Active/non-optimized
...

【 図 2 7 A 】

FIG. 27A is a table titled 'データセンタID管理テーブル' (Data Center ID Management Table) for 103A, 219A. It has two columns: 'データセンタID' (Data Center ID) and '0x0001'.

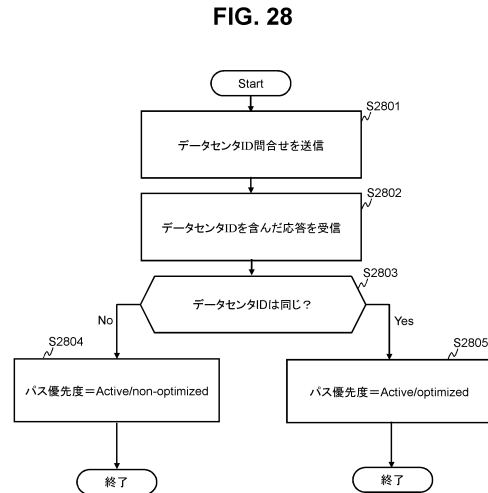
データセンタID	0x0001

【 図 2 7 B 】

FIG. 27B is a table titled 'データセンタID管理テーブル' (Data Center ID Management Table) for 103B, 219B. It has two columns: 'データセンタID' (Data Center ID) and '0x0002'.

データセンタID	0x0002

【 図 2 8 】



フロントページの続き

(72)発明者 南波 優

東京都品川区東品川四丁目12番7号 株式会社日立ソリューションズ内

審査官 田中 啓介

(56)参考文献 特表2012-504793(JP,A)

特開2014-048710(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F3/06-3/08

G06F13/00-13/14

G06F13/20-13/42

H04L12/00-12/28

H04L12/44-12/955

H04W8/26、24/00、28/02

H04W72/04、74/04、74/08

H04W84/12、88/08