

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6556404号
(P6556404)

(45) 発行日 令和1年8月7日(2019.8.7)

(24) 登録日 令和1年7月19日(2019.7.19)

(51) Int.Cl.		F I			
G06F 3/06	(2006.01)	G06F	3/06	301A	
G06F 13/14	(2006.01)	G06F	3/06	301Z	
G06F 13/38	(2006.01)	G06F	13/14	320H	
		G06F	13/38	350	

請求項の数 12 (全 21 頁)

(21) 出願番号	特願2019-506878 (P2019-506878)	(73) 特許権者	000005108
(86) (22) 出願日	平成29年3月24日 (2017.3.24)		株式会社日立製作所
(86) 国際出願番号	PCT/JP2017/011969		東京都千代田区丸の内一丁目6番6号
(87) 国際公開番号	W02018/173245	(74) 代理人	110000279
(87) 国際公開日	平成30年9月27日 (2018.9.27)		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31年1月18日 (2019.1.18)	(72) 発明者	黒川 碧
			東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		(72) 発明者	山崎 優太
			東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		審査官	田中 啓介

最終頁に続く

(54) 【発明の名称】 ストレージシステム及びストレージシステムのバックエンド構築方法

(57) 【特許請求の範囲】

【請求項1】

同一パーティションに存在可能なマスタデバイスの数が規定されている通信インターフェースに従い通信を中継するスイッチと、

前記スイッチに接続された複数の記憶デバイスと、

1以上のメモリを含んだメモリ部と、

前記メモリ部と前記スイッチとに接続された1以上のプロセッサであり複数のマスタデバイスを有するプロセッサ部と

を有し、

前記スイッチは、論理的に複数のパーティションに分割され、

前記複数のパーティションは、複数の第1のパーティションと、1以上の第2のパーティションとを含み、

前記複数の第1のパーティションには、前記プロセッサ部の前記複数のマスタデバイスが複数のパス経路で接続され、前記複数の記憶デバイスは接続されず、

前記1以上の第2のパーティションには、前記複数の記憶デバイスが接続され、前記プロセッサ部は接続されず、

前記スイッチは、異なるパーティション間での転送を可能にする機能であるアドレス変換機能を有し、

前記1以上の第2のパーティションの各々に、仮想的なマスタデバイスが設けられ、

前記1以上の第2のパーティションの各々について、その第2のパーティションにお

10

20

る仮想的なマスタデバイスが、その第2のパーティションに接続されている全ての記憶デバイスの各々に対して初期設定を実行する、ストレージシステム。

【請求項2】

前記スイッチは、1以上のプロセッサである内部プロセッサ部を有し、

前記1以上の第2のパーティションの各々について、その第2のパーティションにおける仮想的なマスタデバイスは、前記内部プロセッサ部である、請求項1記載のストレージシステム。

【請求項3】

前記1以上の第2のパーティションの各々について、その第2のパーティションに接続されている全ての記憶デバイスの各々に対しての初期設定では、前記内部プロセッサ部が、前記内部プロセッサ部が使用可能なバス番号を使用する、請求項2記載のストレージシステム。

10

【請求項4】

前記スイッチが、前記複数の記憶デバイスにそれぞれ対応した複数のメッセージ出力先アドレスを管理しており、

前記複数のメッセージ出力先アドレスの各々は、前記メモリ部のうちの、そのアドレスに対応した記憶デバイスからのメッセージの出力先の領域のアドレスであり、

前記1以上の第2のパーティションの各々について、その第2のパーティションにおける仮想的なマスタデバイスが、その第2のパーティションに接続されているいずれかの記憶デバイスからのメッセージを検出した場合、そのメッセージを、前記複数のメッセージ出力先アドレスのうちの、そのメッセージを出力した記憶デバイスに対応したメッセージ出力先アドレスへと転送し、

20

前記プロセッサ部は、前記メモリ部における、前記複数のメッセージ出力先アドレスにそれぞれ対応した複数の領域の各々を、その領域にメッセージが格納されているか否かを検出するために定期的にチェックする、請求項1記載のストレージシステム。

【請求項5】

前記メッセージは、前記記憶デバイスの障害に関する情報である障害情報である、請求項4記載のストレージシステム。

30

【請求項6】

前記スイッチが、前記複数の記憶デバイスの各々について、メッセージ出力先アドレスの他に、前記複数のパスのうちの、メッセージの出力のために使用されるパスを管理しており、

前記1以上の第2のパーティションの各々について、仮想的なマスタデバイスが、前記検出されたメッセージを、そのメッセージを出力した記憶デバイスに対応したパス経由で出力する、

請求項4記載のストレージシステム。

【請求項7】

前記プロセッサ部が、前記複数の記憶デバイスのうちのいずれかの記憶デバイスにデータのI/Oコマンドの送信する場合、そのI/Oコマンドが経由するパスであるコマンドパスを、前記複数のパスの各々のリンク状態と、前記複数のパスの負荷とのうちの少なくとも1つに基づいて前記複数のパスについて決定されたパス重みに従い選択する、請求項1記載のストレージシステム。

40

【請求項8】

前記プロセッサ部が、

前記複数のパスの少なくとも1つに非接続状態のパスがある場合、前記コマンドパスを、前記複数のパスの各々のリンク状態に基づいて前記複数のパスについて決定されたパス重みに従い選択し、

前記複数のパスの負荷に偏りがある場合、前記コマンドパスを、前記複数のパスの負

50

荷に基づいて前記複数のパスについて決定されたパス重みに従い選択する、
請求項 7 記載のストレージシステム。

【請求項 9】

前記スイッチ、前記メモリ部及び前記プロセッサ部の各々は二重化されており、
前記複数の記憶デバイスは、二重化されたスイッチである第 1 及び第 2 のスイッチの両方に接続されており、

前記第 1 及び第 2 のスイッチに、二重化されたプロセッサ部である第 1 及び第 2 のプロセッサ部がそれぞれ接続されており、

前記第 1 及び第 2 のプロセッサ部に、二重化されたメモリ部である第 1 及び第 2 のメモリ部がそれぞれ接続されており、

前記第 1 のメモリ部における、前記第 1 のプロセッサ部が使用するアドレス範囲と、前記第 2 のメモリ部における、前記第 2 のプロセッサ部が使用するアドレス範囲との少なくとも一部が重複しており、

前記第 1 のメモリ部における、前記第 1 のスイッチがサポートするアドレス範囲と、前記第 2 のメモリ部における、前記第 2 のスイッチがサポートするアドレス範囲とが異なっている、

請求項 1 記載のストレージシステム。

【請求項 10】

前記通信インターフェースは、P C I e (PCI-Express) であり、

前記スイッチは、マルチルートに対応可能な P C I e スイッチであり、

前記マスタデバイスは、ルートコンプレックスであり、

前記仮想的なマスタデバイスは、仮想的なルートコンプレックスであり、

前記アドレス変換機能は、N T B (Non Transparent Bridge) である、

前記初期設定は、コンフィギュレーションアクセスに従う設定である、

請求項 1 記載のストレージシステム。

【請求項 11】

前記複数の記憶デバイスは、複数の N V M e - S S D である、

請求項 10 記載のストレージシステム。

【請求項 12】

ストレージシステムのバックエンドの構築方法であって、

前記ストレージシステムは、

同一パーティションに存在可能なマスタデバイスの数が規定されている通信インターフェースに従い通信を中継するスイッチと、

前記スイッチに接続された複数の記憶デバイスと、

1 以上のメモリを含んだメモリ部と、

前記メモリ部と前記スイッチとに接続された 1 以上のプロセッサであり複数のマスタデバイスを有するプロセッサ部と

を有し、

前記構築方法は、

前記スイッチを、複数の第 1 のパーティションと、1 以上の第 2 のパーティションとを含む複数のパーティションに論理的に分割し、

前記複数の第 1 のパーティションには、前記プロセッサ部の前記複数のマスタデバイスが複数のパス経路で接続され、前記複数の記憶デバイスは接続されず、

前記 1 以上の第 2 のパーティションには、前記複数の記憶デバイスが接続され、前記プロセッサ部は接続されず、

前記スイッチは、異なるパーティション間での転送を可能にする機能であるアドレス変換機能を有し、

前記 1 以上の第 2 のパーティションの各々に、仮想的なマスタデバイスを設け、

前記 1 以上の第 2 のパーティションの各々について、その第 2 のパーティションにおける仮想的なマスタデバイスにより、その第 2 のパーティションに接続されている全ての記

10

20

30

40

50

憶デバイスの各々に対して初期設定を実行する、構築方法。

【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のシステム構築に関する。

【背景技術】

【0002】

複数の記憶デバイス（以下、ドライブ）を搭載したストレージシステムが知られている。ドライブとして、例えばSSD（Solid State Drive）がある。

10

【0003】

また、ドライブを接続するための通信インターフェース（プロトコル）として、NVMe（NVM（Non-Volatile Memory）Express）が知られている。

【0004】

そこで、ストレージシステムに搭載されるSSDとして、NVMe-SSDが採用されることが予測される。NVMe-SSDが採用された場合、プロセッサとSSD間の通信インターフェースとして、PCIe（PCI-Express）が採用される。具体的には、例えば、プロセッサとNVMe-SSDを、PCIeスイッチ経由で接続することが考えられる。PCIeスイッチに関する技術として、例えば、特許文献1～特許文献3が知られている。

20

【先行技術文献】

【特許文献】

【0005】

【特許文献1】US2006/0242330

【特許文献2】US8756360

【特許文献3】US9141571

【発明の概要】

【発明が解決しようとする課題】

【0006】

PCIeは、同一空間内に存在可能なマスタデバイス数が規定されている通信インターフェースの一例である。PCIeでは、一般に、ルートコンプレックスがマスタデバイスの一例であり、エンドポイントがスレイブデバイスの一例であり、ドメインが空間の一例である。PCIeでは、同一ドメイン内に存在可能なルートコンプレックスは1つでありその1つのルートコンプレックスに1以上のエンドポイントを接続可能である。

30

【0007】

マルチルート（複数のルートコンプレックス）に接続可能なPCIeスイッチであるマルチルートPCIeスイッチが知られている。

【0008】

ストレージシステムにおいて、PCIeスイッチとしてマルチルートPCIeスイッチが採用された場合、マルチルートPCIeスイッチに、プロセッサ部（1以上のプロセッサ）が有する複数のルートコンプレックス（例えばポート）と、複数のNVMe-SSD（エンドポイント）が接続されることになる。この場合、バックエンド（プロセッサ部とNVMe-SSD間）の構成として、下記の構成A及びBを含んだ構成を採用することが考えられる。

40

（構成A）マルチルートPCIeスイッチが、論理的に複数のパーティション（ドメイン）に区切られる。各パーティションに、プロセッサ部のルートコンプレックスと、1以上のNVMe-SSDとが接続される。1ドメインにつき1ルートコンプレックスであることを維持するためである。

（構成B）NTB（Non Transparent Bridge）が、マルチルートPCIeスイッチに搭載される。各ルートコンプレックスを、そのルートコンプレックスが接続されているパーテ

50

ィションとは異なるパーティションに接続されている NVMe - SSD にアクセス可能にするためである。

【0009】

しかし、構成 A 及び B を含んだバックエンド構成では、下記の課題 X 乃至 Z があると考えられる。

(課題 X) マルチルート PCIe スイッチに接続されている各 NVMe - SSD について、その NVMe - SSD の初期設定を、その NVMe - SSD が接続されているパーティションに接続されているパス経路で、プロセッサ部が行う必要がある。このため、そのパスが未接続状態又は障害状態であれば、その NVMe - SSD の初期設定を行うことができない。

10

(課題 Y) プロセッサ部が使用可能なデバイス数 (例えば、BDF (バス番号、デバイス番号及びファンクション番号) の数) に制限がある。そして、一般に、プロセッサ部には、マルチルート PCIe スイッチ以外の PCI デバイス (例えば、通信インターフェースデバイス) も接続される。このため、プロセッサ部が使用可能なデバイス数と同数の NVMe - SSD をマルチルート PCIe スイッチに接続したとしても、プロセッサ部は、それら全ての NVMe - SSD にアクセスすることはできない。

(課題 Z) バックエンド構成の設定 (例えば、NVMe - SSD に関する設定) が、煩雑である (例えば、プロセッサ部とマルチルート PCIe スイッチとを結ぶパスの増減に依存する)。具体的には、例えば、NVMe - SSD をどのパーティションに接続するか、いずれの障害が生じたらいずれのパーティションへの切り替えを行うか等が必要である。

20

【0010】

課題 X 乃至 Z のうちの少なくとも 1 つのような課題は、PCIe に限らず、同一空間内に存在可能なマスタデバイスの数が規定されている他種の通信インターフェースが、プロセッサ部とドライブ間 (バックエンド) の通信インターフェースとして採用される場合にも、あり得る。

【課題を解決するための手段】

【0011】

ストレージシステムは、同一パーティションに存在可能なマスタデバイスの数が規定されている通信インターフェースに従い通信を中継するスイッチと、スイッチに接続された複数の記憶デバイスと、1 以上のメモリを含んだメモリ部と、メモリ部とスイッチとに接続された 1 以上のプロセッサであり複数のマスタデバイスを有するプロセッサ部とを有する。スイッチは、論理的に複数のパーティションに分割される。複数のパーティションは、複数の第 1 のパーティションと、1 以上の第 2 のパーティションとを含む。複数の第 1 のパーティションには、プロセッサ部の複数のマスタデバイスが複数のパス経路で接続され、複数の記憶デバイスは接続されない。1 以上の第 2 のパーティションには、複数の記憶デバイスが接続され、プロセッサ部は接続されない。スイッチは、異なるパーティション間での転送を可能にする機能であるアドレス変換機能を有する。1 以上の第 2 のパーティションの各々に、仮想的なマスタデバイスが設けられる。1 以上の第 2 のパーティションの各々について、その第 2 のパーティションにおける仮想的なマスタデバイスが、その第 2 のパーティションに接続されている全ての記憶デバイスの各々に対して初期設定を実行する。

30

40

【発明の効果】

【0012】

プロセッサ部とスイッチにおける一部の第 1 のパーティション間のパスが未接続状態又は障害状態であっても、各記憶デバイスに対する初期設定を行うことができる。また、接続可能な記憶デバイスの数が、プロセッサ部が使用可能なデバイス数に制限されない。また、記憶デバイスは第 2 のパーティションに接続されればよいので、プロセッサ部とスイッチ間のパスの増減に非依存である。

【図面の簡単な説明】

【0013】

50

- 【図 1】実施例 1 に係るストレージシステムのバックエンド構成の概要を示す。
- 【図 2】一比較例に係るバックエンド構成の概要を示す。
- 【図 3】実施例 1 に係るストレージシステムを含んだ情報システムの構成を示す。
- 【図 4】ドライブ状態管理テーブルの構成を示す。
- 【図 5】ドライブ初期設定管理テーブルの構成を示す。
- 【図 6】ドライブコンフィギュレーション管理テーブルの構成を示す。
- 【図 7】パス管理テーブル 1 4 6 の構成を示す。
- 【図 8】アドレス変換テーブルの構成を示す。
- 【図 9】コマンド管理テーブルの構成を示す。
- 【図 10】ドライブ初期設定処理のフローを示す。 10
- 【図 11】障害情報通知の概要を示す。
- 【図 12】障害情報通知に関し CPU (ファームウェア) が行う処理のフローを示す。
- 【図 13】障害情報通知に関し内部 CPU (ファームウェア) が行う処理のフローを示す。
- 【図 14】パス選択処理のフローを示す。
- 【図 15】実施例 2 に係るコントローボックスの構成の一部を示す。
- 【図 16】全体アドレス管理テーブルの構成を示す。
- 【図 17】第 1 のスイッチアドレス管理テーブルの構成を示す。
- 【図 18】第 2 のスイッチアドレス管理テーブルの構成を示す。
- 【発明を実施するための形態】 20
- 【0014】
以下、幾つかの実施例を説明する。
- 【0015】
なお、以下の説明では、「xxx テーブル」といった表現にて情報を説明することがあるが、情報は、どのようなデータ構造で表現されていてもよい。すなわち、情報がデータ構造に依存しないことを示すために、「xxx テーブル」を「xxx 情報」と言うことができる。また、以下の説明において、各テーブルの構成は一例であり、1 つのテーブルは、2 以上のテーブルに分割されてもよいし、2 以上のテーブルの全部又は一部が 1 つのテーブルであってもよい。
- 【0016】 30
また、以下の説明では、「インターフェース部」は、1 以上の通信インターフェースデバイスを含む。1 以上の通信インターフェースデバイスは、1 以上の同種のインターフェースデバイス (例えば 1 以上の NIC (Network Interface Card)) であってもよいし 2 以上の異種のインターフェースデバイス (例えば NIC と HBA (Host Bus Adapter)) であってもよい。
- 【0017】
また、以下の説明では、「メモリ部」は、1 以上のメモリを含む。メモリ部に関して少なくとも 1 つのメモリは、揮発性メモリでよい。メモリ部は、主に、プロセッサ部による処理の際に使用される。
- 【0018】 40
また、以下の説明では、「プロセッサ部」は、1 以上のプロセッサを含む。少なくとも 1 つのプロセッサは、典型的には、CPU (Central Processing Unit) のようなマイクロプロセッサである。1 以上のプロセッサの各々は、シングルコアでもよいしマルチコアでもよい。プロセッサは、処理の一部または全部を行うハードウェア回路を含んでもよい。
- 【0019】
また、以下の説明では、「ホストシステム」は、1 以上の物理的なホスト計算機 (例えばホスト計算機のクラスタ) であってもよいし、少なくとも 1 つの仮想的なホスト計算機 (例えば VM (Virtual Machine)) を含んでもよい。以下、ホストシステムを、単に「ホスト」と呼ぶ。ホストは、ストレージシステムにおいてホストとして動作する VM (例

えばサーバVM)でもよい。

【0020】

また、以下の説明では、「ストレージシステム」は、1以上の物理的なストレージ装置であってもよいし、少なくとも1つの仮想的なストレージ装置(例えばSDS(Software Defined Storage))を含んでもよい。例えば、ストレージシステムは、サーバVMと、ストレージコントローラとして動作するVMであるストレージVMとを実行してもよい。ストレージVMは、サーバVMからのI/O(Input/Output)要求に応答してI/Oを実行してよい。

【0021】

また、以下の説明では、同種の要素を区別しないで説明する場合には、参照符号(又は参照符号における共通部分)を使用し、同種の要素を区別して説明する場合は、要素のID(又は要素の参照符号)を使用することがある。

【実施例1】

【0022】

図1は、実施例1に係るストレージシステムのバックエンド構成の概要を示す。

【0023】

CPU(プロセッサ部の一例)115が、2個のルートコンプレックス(複数のマスタデバイスの一例)116A及び116Bを有する。各ルートコンプレックス116は、例えば、ポート(以下、CPUポート)である。

【0024】

マルチルートPCIeスイッチ(以下、PCIe-SW)57は、複数のポート(以下、SWポート)321を有する。複数のSWポート321には、下位デバイスに接続されるDSP(ダウンストリームポート)と、上位デバイスに接続されるUSP(アップストリームポート)とが含まれる。DSPは、属性がダウンストリームであるポートである。USPは、属性がアップストリームであるポートである。「下位デバイス」は、例えば、別のPCIeスイッチ(例えば別のPCIe-SW)又はエンドポイント(典型的にはドライブ117)である。「上位デバイス」は、例えば、CPU115である。

【0025】

また、PCIe-SW57は、複数のパーティション(ドメイン)51に論理的に分割される。

【0026】

本実施例では、少なくとも下記の構成が採用される。

(構成1)PCIe-SW57に、CPU115が有する2個のルートコンプレックス116A及び116Bがそれぞれ接続される2個のパーティション(以下、CPUパーティション)P0及びP1に加えて、ドライブ用のパーティション(以下、ドライブパーティション)P2が設けられる。ドライブパーティションP2には、ドライブ117が接続され(例えばドライブ117のみが接続され)、少なくともCPU115は接続されない。本実施例では、CPUパーティションP0及びP1の各々には、ドライブ117が接続されず、ドライブパーティションP2には、ドライブ117のみが接続されるものとする。なお、CPUパーティションP0及びP1の各々について、USPに、パス52経由で、CPU115のルートコンプレックス116が接続される。ドライブパーティションP2が有する複数のDSPに、それぞれ、複数のドライブ117がそれぞれ有する複数のポート(以下、ドライブポート)53が接続される。

(構成2)PCIe-SW57に、異なるパーティション51間の転送を実現するアドレス変換機能の一例であるNTB(Non Transparent Bridge)61が設けられる。NTB61が、CPUパーティションP0又はP1とドライブパーティションP2間の転送を実現する。

(構成3)ドライブパーティションP2に、仮想的なルートコンプレックスが設けられる。本実施例では、PCIe-SW57のCPU114(以下、内部CPU114)が、ルートコンプレックスの役割を持つ。内部CPU114が、各ドライブ(ドライブパーティ

10

20

30

40

50

ションP2でのエンドポイント)の初期設定を実行する。具体的には、例えば、ドライブパーティションP2に、仮想的なUSPが設けられ、そのUSPに、仮想的なルートコンプレックス(内部CPU114)が接続される。

【0027】

図2に示す一比較例によれば、上述した構成A及びB相当の構成が採用されている。このため、上述した課題X乃至Z相当の課題がある。すなわち、一比較例によれば、CPU201とPCIe-SW257間の或るパス252が未接続状態又は障害状態であれば、そのパス252が接続されているパーティション202に接続されているドライブ217の初期設定を行うことができない。また、一比較例によれば、CPU201が使用可能なデバイス数と同数のドライブ217をPCIe-SW257に接続したとしても、CPU201は、それら全てのドライブ217にアクセスすることはできない。また、一比較例によれば、バックエンド構成の設定(例えば、ドライブ217に関する設定)が、CPU201とPCIe-SW257間のパス252の数に依存する。

10

【0028】

そこで、構成1によれば、CPU115が接続されるパーティション51にはドライブ117が接続されない(CPU115から見えるバス番号のドメイン内にはドライブ117が配置されない)。このため、CPU115は、ドライブ117のためにBDF(バス番号、デバイス番号及びファンクション番号の組)を使用(消費)しないで済む。結果として、CPU115が使用するBDFの数を削減できる。また、バックエンドの構成の設定を、CPU115とPCIe-SW57間のパス52の数に非依存とすることができる。

20

【0029】

また、構成2によれば、CPU115は、CPUパーティションP0又はP1経由で、ドライブパーティションP2に接続されているドライブ117にアクセスできる。

【0030】

なお、構成1(及び構成2)によれば、CPU115はドライブパーティションP2に接続されていないため、ドライブ117への初期設定(ドライブ117のコンフィギュレーションレジスタにアクセスすること)をCPU115から直接行うことができない。

【0031】

そこで、構成3によれば、PCIe-SW57の内部CPU114が、ルートコンプレックスとして、ドライブ117に対する初期設定(ドライブ117のコンフィギュレーションレジスタにアクセスすること)をできる。

30

【0032】

以下、本実施例を詳細に説明する。なお、2個のCPUパーティションP0及びP1は、複数の第1のパーティションの一例である。1個のドライブパーティションP2は、1以上の第2のパーティションの一例である。CPUパーティションの数は、CPU115が有するルートコンプレックス116の数と同数でよい。

【0033】

図3は、実施例1に係るストレージシステムを含んだ情報システムの構成を示す。

40

【0034】

情報システムは、1又は複数のホストシステム(以下、ホスト)101と、ストレージシステム105とを有する。ホスト101とストレージシステム105は、通信ネットワーク71(例えば、SAN(Storage Area Network)又はLAN(Local Area Network))に接続される。

【0035】

ホスト101は、ストレージシステム105にI/O(Input/Output)要求を送信する。I/O要求は、I/O先の場所を表すI/O先情報を含む。I/O先情報は、例えば、I/O先のLU(Logical Unit)のLUN(Logical Unit Number)と、そのLUにおける領域のLBA(Logical Block Address)とを含む。LUは、ストレージシステム10

50

5 から提供される論理ボリューム（論理的な記憶デバイス）である。I/O先情報を基に、I/O先の論理領域が特定され、その論理領域に基づくドライブ117が特定される。

【0036】

ストレージシステム105は、コントローラボックス78と、コントローラボックス78に接続された1以上のドライブボックス103とを有する。

【0037】

コントローラボックス78は、ストレージコントローラ79を有する。ストレージコントローラ79は、複数の通信インターフェースデバイス（インターフェース部の一例）と、メモリ111（メモリ部の一例）と、PCIe-SW57（スイッチの一例）と、それらに接続されたCPU115（プロセッサ部の一例）とを有する。

10

【0038】

複数の通信インターフェースデバイスは、複数のホスト101と通信するための1以上の通信インターフェースデバイスである1以上のH-I/F112（例えば、Fibre Channel、iSCSI、FCoE又はPCIeのデバイス）を含む。複数の通信インターフェースデバイスは、ドライブボックス103A内のドライブ87と通信するための1以上の通信インターフェースデバイスである1以上のD-I/F113（例えば、SASコントローラ又はPCIeスイッチ）を含んでもよい。ドライブボックス103A内の各ドライブ87は、NVMe-SSDでもよいし、他種のPCIeドライブでもよいし、他種のSSD（例えばSAS-SSD）でもよいし、他種のドライブ（例えばHDD（Hard Disk Drive））でもよい。

20

【0039】

メモリ111は、CPU115により実行される1以上のコンピュータプログラムと、CPU115により参照又は更新される情報とを格納する。1以上のコンピュータプログラムは、例えば、ファームウェア142のようなマイクロプログラムを含む。情報は、例えば、複数のテーブルである。複数のテーブルは、例えば、ドライブ実装状態管理テーブル143、ドライブ設定状態管理テーブル144、ドライブ初期設定管理テーブル145、パス管理テーブル146、アドレス変換テーブル147及びコマンド管理テーブル148を含む。また、メモリ111には、ホスト101からのI/O要求に応答してドライブ117（又は87）に入出力されるデータであるユーザデータが一時的に格納されるキャッシュ領域が設けられる。

30

【0040】

PCIe-SW57は、複数のSWポート321と、ブリッジ121と、内部CPU114（内部プロセッサ部の一例）と、内部メモリ188（内部メモリ部の一例）とを有する。

【0041】

ブリッジ121に、複数のSWポート321が接続されている。内部CPU114に、内部メモリ188とブリッジ121とが接続されている。内部CPU114と、SWポート321に接続されているデバイス間の通信は、ブリッジ121経由である。

【0042】

内部メモリ188は、内部CPU114により実行される1以上のコンピュータプログラムと、内部CPU114により参照又は更新される情報とを格納する。1以上のコンピュータプログラムは、例えば、ファームウェア189のようなマイクロプログラムを含む。情報は、例えば、1以上のテーブル（例えば、後述のメッセージ通知管理テーブル42）である。

40

【0043】

1以上のドライブボックス103は、少なくとも、PCIe-SW57に接続される複数のドライブ117を有するドライブボックス103Bを含む。ドライブ117は、NVMe-SSDである。ドライブ117は、他種のPCIeドライブでもよい。1以上のドライブボックス103は、更に、上述したドライブボックス103Aを含んでもよい。

【0044】

50

図4は、ドライブ実装状態管理テーブル143の構成を示す。

【0045】

ドライブ実装状態管理テーブル143は、各ドライブ117の実装状態に関する情報を保持するテーブルである。ドライブ実装状態管理テーブル143は、ドライブパーティションP2におけるDSP毎にエントリを有する。各エントリが、ドライブ#401、実装状態402及びリンク状態403といった情報を保持する。

【0046】

ドライブ#401は、ストレージシステム105のプログラムが管理しているドライブ番号である。

【0047】

実装状態402は、ドライブ117がDSPに対して実装（物理的に接続）されているか否かを示す。実装状態402の値として、“1”は、実装を意味し、“0”は、未実装を意味する。

【0048】

リンク状態403は、ドライブ117がDSPに対してリンクアップ（通信可能な状態）か否かを示す。リンク状態403の値として、“1”は、リンクアップを意味し、“0”は、リンクダウン（通信不可能な状態）を意味する。

【0049】

図5は、ドライブ設定状態管理テーブル144の構成を示す。

【0050】

ドライブ設定状態管理テーブル144は、各ドライブ117の設定状態に関する情報を保持するテーブルである。ドライブ設定状態管理テーブル144は、ドライブパーティションP2におけるDSP毎にエントリを有する。各エントリが、ドライブ#501、完了状態502及び結果状態503といった情報を保持する。

【0051】

ドライブ#501は、ストレージシステム105のプログラムが管理しているドライブ番号である。

【0052】

完了状態502は、ドライブ117に対する初期設定が完了したか否かを示す。完了状態502の値として、“1”は、完了を意味し、“0”は、未完了を意味する。

【0053】

結果状態503は、ドライブ117に対する初期設定が完了した結果として成功か否かを示す。結果状態503の値として、“1”は、成功を意味し、“0”は、失敗を意味する。

【0054】

図6は、ドライブ初期設定管理テーブル145の構成を示す。

【0055】

ドライブ初期設定管理テーブル145は、各ドライブ117のコンフィグレーションレジスタに設定する情報を保持するテーブルである。ドライブ初期設定管理テーブル145は、ドライブパーティションP2におけるDSP毎にエントリを有する。各エントリが、ドライブ#601、バス#602、デバイス#603、ファンクション#604、ベースアドレスレジスタ605及びMSI（Message Signaled Interrupt）テーブルレジスタ606といった情報を保持する。

【0056】

ドライブ#601は、ストレージシステム105のプログラムが管理しているドライブ番号である。

【0057】

バス#602は、BDFのうちのバス番号である。デバイス#603は、BDFのうちのデバイス番号である。ファンクション#604は、BDFのうちのファンクション番号である。このBDFは、CPU115が使用可能なBDFのうちのBDFではなく、内部

10

20

30

40

50

C P U 1 1 4 が使用可能な B D F のうちの B D F であり、ドライブ 1 1 7 のコンフィギュレーションレジスタに設定される。

【 0 0 5 8 】

ベースアドレスレジスタ 6 0 5 は、ドライブ 1 1 7 に対応したアドレス (C P U 1 1 5 のメモリマップ空間におけるアドレス) を示す。ドライブ 1 1 7 のコンフィギュレーションレジスタは、例えば、ドライブ 1 1 7 が有する記憶領域 (例えば、ドライブポート 5 3 (図 1 参照) のレジスタ) でよい。

【 0 0 5 9 】

M S I テーブルレジスタ 6 0 6 は、ドライブ 1 1 7 に対応したメッセージ (例えば、ドライブ 1 1 7 の障害に関する障害情報) の出力先のアドレス (例えば、メモリ 1 1 1 における領域のアドレス) を示す。

10

【 0 0 6 0 】

図 7 は、パス管理テーブル 1 4 6 の構成を示す。

【 0 0 6 1 】

パス管理テーブル 1 4 6 は、C P U 1 1 5 と P C I e - S W 5 7 間のパスに関する情報を保持するテーブルである。パス管理テーブル 1 4 6 は、パス 5 2 毎にエントリを有する。各エントリが、パス # 7 0 1、データ転送量 7 0 2、リンク状態 7 0 3 及び重み 7 0 4 といった情報を保持する。

【 0 0 6 2 】

パス # 7 0 1 は、パス 5 2 の識別番号である。データ転送量 7 0 2 は、パス 5 2 の負荷の一例であり、そのパス 5 2 を経由のデータ転送量 (例えば単位時間当りの転送量) を示す。リンク状態 7 0 3 は、パス 5 2 のリンク状態を示す。リンク状態 7 0 3 として、“ 1 ” は、正常 (接続状態) を意味し、“ 0 ” は、異常 (未接続状態又は障害状態) を意味する。重み 7 0 4 は、パス 5 2 の重みを示す。

20

【 0 0 6 3 】

図 8 は、アドレス変換テーブル 1 4 7 の構成を示す。

【 0 0 6 4 】

アドレス変換テーブル 1 4 7 は、変換前後のアドレス関係に関する情報を保持する。アドレス変換テーブル 1 4 7 は、ドライブ 1 1 7 初期設定処理後に C P U 1 1 5 により発行された I / O コマンド毎にエントリを有する。各エントリが、タグ # 8 0 1、パス # 8 0 2、実メモリアドレス 8 0 3 及びコマンド指定メモリアドレス 8 0 4 といった情報を保持する。なお、ドライブ初期設定処理後に C P U 1 1 5 により発行された管理対象のコマンドとして、本実施例では、I / O コマンドが採用されているが、I / O コマンドに加えて他種のコマンドも管理対象とされてよい。

30

【 0 0 6 5 】

タグ # 8 0 1 は、I / O コマンドに関連付けられたタグの識別番号 (実質的に I / O コマンドの識別番号) である。パス # 8 0 2 は、I / O コマンドが経由するパス 5 2 の識別番号である。実メモリアドレス 8 0 3 は、アドレス変換前のアドレス、具体的には、I / O コマンドに従う I / O 対象のユーザデータが格納されている領域 (メモリにおける領域) のアドレスを示す。コマンド指定メモリアドレス 8 0 4 は、アドレス変換後のアドレス、具体的には、I / O コマンドで指定されるメモリアドレス (例えば、C P U 1 1 5 のメモリマップ空間における宛先ドライブ対応アドレス) を示す。

40

【 0 0 6 6 】

図 9 は、コマンド管理テーブル 1 4 8 の構成を示す。

【 0 0 6 7 】

コマンド管理テーブル 1 4 8 は、I / O コマンドに関する情報を保持するテーブルである。コマンド管理テーブル 1 4 8 は、ドライブ初期設定処理後に C P U 1 1 5 により発行された I / O コマンド毎にエントリを有する。各エントリが、タグ # 9 0 1、ドライブ # 9 0 2、レンジ 9 0 3、コマンド指定メモリアドレス 9 0 4 及びドライブアドレス 9 0 5 といった情報を保持する。

50

【 0 0 6 8 】

タグ# 9 0 1 は、I / O コマンドに関連付けられたタグの識別番号（実質的に I / O コマンドの識別番号）である。ドライブ# 9 0 2 は、I / O コマンドに従う I / O 先ドライブ 1 1 7 の識別番号である。レングス 9 0 3 は、I / O コマンドに従う I / O 対象のユーザデータのデータ長を示す。コマンド指定メモリアドレス 9 0 4 は、I / O コマンドで指定されるメモリアドレスを示す。ドライブアドレス 9 0 5 は、I / O 先ドライブ 1 1 7 における I / O 先領域のアドレス（例えば L B A (Logical Block Address)）を示す。

【 0 0 6 9 】

以下、本実施例で行われる幾つかの処理を説明する。

【 0 0 7 0 】

< ドライブ初期設定 >

【 0 0 7 1 】

図 1 0 は、ドライブ初期設定処理のフローを示す。

【 0 0 7 2 】

C P U 1 1 5 (ファームウェア 1 4 2) が、ドライブ 1 1 7 に関する状態の問合せであるドライブ問合せを、P C I e - S W 5 7 に送信する (S 5 0 1)。

【 0 0 7 3 】

内部 C P U 1 1 4 (ファームウェア 1 4 2) が、ドライブ問合せに回答して、ドライブパーティション P 2 における D S P 毎のドライブ情報 (ドライブ#、実装か未実装か、リンクアップかリンクダウンか、設定完了か否か、及び、設定成功か否かを示す情報) を C P U 1 1 5 に返す (S 5 0 2)。

【 0 0 7 4 】

C P U 1 1 5 (ファームウェア 1 4 2) が、内部 C P U 1 1 4 からの情報 (ドライブパーティション P 2 における D S P 毎のドライブ情報) を、ドライブ実装状態管理テーブル 1 4 3 及びドライブ設定状態管理テーブル 1 4 4 に登録する (S 5 0 3)。

【 0 0 7 5 】

C P U 1 1 5 (ファームウェア 1 4 2) が、テーブル 1 4 3 及び 1 4 4 を基に、1 以上の未設定ドライブ 1 1 7 があるか否かを判断する (S 5 0 4)。「未設定ドライブ 1 1 7」は、実装状態 4 0 2 “ 1 ”、リンク状態 4 0 3 “ 1 ” 及び完了状態 5 0 2 “ 0 ” に対応したドライブ 1 1 7 である。なお、「未設定ドライブ 1 1 7」は、更に、結果状態 5 0 3 “ 0 ” に対応したドライブ 1 1 7 によい。初期設定に失敗したドライブ 1 1 7 に対して初期設定をリトライするためである。

【 0 0 7 6 】

S 5 0 4 の判断結果が真の場合 (S 5 0 4 : Y e s)、C P U 1 1 5 (ファームウェア 1 4 2) が、1 以上の未設定ドライブ 1 1 7 から選択した 1 つのドライブ 1 1 7 である対象ドライブ 1 1 7 について、初期設定指示を、P C I e - S W 5 7 に送信する。初期設定指示には、対象ドライブ 1 1 7 のドライブ# が関連付けられる。

【 0 0 7 7 】

内部 C P U 1 1 4 (ファームウェア 1 8 9) が、その初期設定指示に回答して、対象ドライブ 1 1 7 に対して初期設定を実行する (S 5 0 6)。具体的には、例えば、内部 C P U 1 1 4 (ファームウェア 1 8 9) は、C P U 1 1 5 (ファームウェア 1 4 2) が定義したドライブ初期設定管理テーブル 1 4 5 で指示された B D F と、ベースアドレスレジスタと、M S I テーブルレジスタとを含んだ情報を、対象ドライブ 1 1 7 のコンフィギュレーションレジスタに設定する。

【 0 0 7 8 】

内部 C P U 1 1 4 (ファームウェア 1 8 9) が、その実行結果を C P U 1 1 5 に返す (S 5 0 7)。その実行結果は、その設定した情報 (B D F、ベースアドレスレジスタ及び M S I テーブルレジスタ) と、初期設定の成否を表す情報と、対象ドライブ 1 1 7 のドライブ# とを含む。

【 0 0 7 9 】

10

20

30

40

50

C P U 1 1 5 (ファームウェア 1 4 2) が、返却された実行結果に含まれている情報 (ドライブ #、B D F、ベースアドレスレジスタ及び M S I テーブルレジスタを含んだ情報) を、対象ドライブに対応したエントリ (ドライブ初期設定管理テーブル 1 4 5 におけるエントリ) に、登録する (S 5 0 8)。

【 0 0 8 0 】

C P U 1 1 5 (ファームウェア 1 4 2) が、返却された実行結果に、成功を表す情報が含まれているか否かを判断する (S 5 0 9)。

【 0 0 8 1 】

S 5 0 9 の判断結果が真の場合 (S 5 0 9 : Y e s)、C P U 1 1 5 (ファームウェア 1 4 2) が、対象ドライブ 1 1 7 に対応したエントリ (ドライブ設定状態管理テーブル 1 4 4 におけるエントリ) に、完了状態 5 0 2 “ 1 ” 及び結果状態 5 0 3 “ 1 ” を登録する (S 5 1 0)。その後、処理が、S 5 0 4 に戻る。

【 0 0 8 2 】

S 5 0 9 の判断結果が偽の場合 (S 5 0 9 : N o)、C P U 1 1 5 (ファームウェア 1 4 2) が、対象ドライブ 1 1 7 に対応したエントリ (ドライブ設定状態管理テーブル 1 4 4 におけるエントリ) に、完了状態 5 0 2 “ 1 ” 及び結果状態 5 0 3 “ 0 ” を登録する (S 5 1 1)。その後、処理が、S 5 0 4 に戻る。

【 0 0 8 3 】

以上のように、ドライブ初期設定処理では、C P U 1 1 5 が使用可能な B D F がドライブ 1 1 7 の初期設定について使用 (消費) されず、内部 C P U 1 1 4 が使用可能な B D F がドライブ 1 1 7 に対する初期設定で使用される。C P U 1 1 5 は、初期設定が完了 (且つ成功) したドライブ 1 1 7 のコンフィギュレーションレジスタに、そのドライブ 1 1 7 に対応したベースアドレス (内部 C P U 1 1 4 が使用したベースアドレス) を用いて、アクセスすることができる。

【 0 0 8 4 】

< 障害情報通知 >

【 0 0 8 5 】

ドライブパーティション P 2 には、C P U 1 1 5 は接続されていない。ドライブパーティション P 2 において、内部 C P U 1 1 4 が、仮想的なルートコンプレックス (例えばルートコンプレックス) である。このため、ドライブ 1 1 7 が、ドライブ 1 1 7 の障害を検出した場合、ルートコンプレックスである内部 C P U 1 1 4 宛に障害情報 (メッセージの一例) が発行される。しかし、ドライブパーティション P 2 に接続されていない C P U 1 1 5 には、その障害情報は届かない (課題 1)。C P U 1 1 5 から見えるバス番号のドメイン内にはドライブ 1 1 7 が配置されていないためである。

【 0 0 8 6 】

課題 1 を解決する方法として、内部 C P U 1 1 4 が、ドライブ 1 1 7 からの障害情報を検出した場合に、C P U 1 1 5 に割込みで障害を通知する方法が考えられる。しかし、その方法では、割込みを受けた C P U 1 1 5 の処理が中断するので、1 個のドライブ 1 1 7 の障害が、ストレージシステム 1 0 5 全体に影響を及ぼすことになり得る (課題 2)。C P U 1 1 5 には、バックエンドの P C I e - S W 5 7 だけでなく、H - I / F 1 1 2 のような他のデバイスも接続されているためである。

【 0 0 8 7 】

課題 2 を解決する方法として、C P U 1 1 5 が、P C I e - S W 5 7 の内部 C P U 1 1 4 に対して定期的に障害問合せ (いずれかのドライブ 1 1 7 で障害が発生したか否かの問合せ) を発行する方法が考えられる。しかし、障害問合せは、C P U 1 1 5 と P C I e - S W 5 7 間のパス 5 2 を経由するため、ドライブ 1 1 7 に対する I / O の性能が低下する可能性もある。

【 0 0 8 8 】

そこで、本実施例では、図 1 1 に示すように、ドライブ 1 1 7 からの障害情報が、内部メモリ 1 8 8 に格納され、その障害情報は、内部 C P U 1 1 4 (又は D M A (Direct Mem

10

20

30

40

50

ory Access))により、内部メモリ188から、メモリ111に転送され格納される。CPU115は、内部メモリ188に代えて、メモリ111に、障害情報の有無をチェックするために定期的にアクセスする。これにより、障害情報の有無のチェックのためにパス52が使用されないため、ドライブ117に対するI/Oの性能の低下を避けることができる。

【0089】

以下、障害情報通知に関し、CPU115及び内部CPU114の各々が行う処理のフローを説明する。なお、障害情報は、障害が発生した部位(例えばドライブ#)を示す情報と、障害の詳細を示す情報とを含む。

【0090】

図12は、障害情報通知に関しCPU115(ファームウェア142)が行う処理のフローを示す。

【0091】

内部メモリ188内に、メッセージ通知管理テーブル42が格納されている。メッセージ通知管理テーブル42は、障害情報のようなメッセージの出力先に関する情報をドライブ117毎(DSP毎)に保持するテーブルである。CPU115が、PCIe-SW57のメモリ内のメッセージ通知管理テーブル42に、ドライブ117毎の情報を設定する(S1201)。ドライブ117毎の情報は、ドライブ#(ドライブ117の識別番号)、MSIテーブルレジスタ(ドライブ117からの障害情報の転送先領域(メモリにおける領域)のアドレス)、及び、パス#(その障害情報の転送に使用されるパスのパス#)を含む。メッセージの通知のために使用されるパスが、複数のドライブ117に均等に分散されていれば、メッセージ通知のために特定のパスに負荷が集中することを避けることが期待できる。

【0092】

CPU115は、メモリ111の各領域(MSIテーブルレジスタ606が示す領域)を、その領域にメッセージ(例えば障害情報)が格納されているか否かを検出するために定期的にチェック(参照)する(S1202)。

【0093】

いずれかの領域に障害情報が格納されていることを検出した場合(S1202:Yes)、CPU115は、その領域から障害情報を取得し、その障害情報を基に障害対処処理を実行する(S1203)。例えば、CPU115は、その障害情報を基に特定されたドライブ117である障害ドライブ117に対するI/Oを停止し、その障害ドライブ117を閉塞する。

【0094】

図13は、障害情報通知に関し内部CPU114(ファームウェア189)が行う処理のフローを示す。

【0095】

内部CPU114は、いずれかのドライブ117の障害を検出した場合(S1101:Yes)、そのドライブ117からの障害情報を、内部メモリ188に格納する(S1102)。内部メモリ188内の障害情報は、メッセージ通知管理テーブル42を基に、内部メモリ188から、障害ドライブ117に対応したMSIテーブルレジスタが示す領域(メモリ111内の領域)に、障害ドライブ117に対応したパス52経由で、内部CPU114(又はDMA)により転送される(S1103)。内部CPU114は、転送された障害情報を内部メモリ188から削除する(S1104)。

【0096】

<パス選択>

【0097】

CPU115と各ドライブ117間に複数のパス52が存在する。ドライブ117へのI/Oコマンドで指定するメモリアドレスにより、そのI/Oコマンドが経由するパスを決定する必要がある。

10

20

30

40

50

【0098】

そこで、本実施例では、バス管理テーブル146（データ転送量702及びリンク状態703のうちの少なくとも1つに基づき）、I/Oコマンドで指定するバスとメモリアドレスの関係が決定される。

【0099】

以下、バス選択処理を説明する。なお、CPU115（ファームウェア142）は、バスのデータ転送量及びリンク状態を定期的にチェックしてバス管理テーブル146にチェック結果を登録するようになっている。

【0100】

図14は、バス選択処理のフローを示す。

10

【0101】

CPU115（ファームウェア142）が、I/O要求に従うユーザデータの配置先領域（メモリ111における領域）を決定する（S1701）。CPU115（ファームウェア142）が、アドレス変換テーブル147におけるエントリ（そのI/Oコマンドに対応したエントリ）に、その配置先領域を示す実メモリアドレス803と、そのユーザデータのI/OのためのI/Oコマンドのタグ#801とを登録する。また、CPU115（ファームウェア142）が、コマンド管理テーブル148におけるエントリ（そのI/Oコマンドに対応したエントリ）に、そのI/Oコマンドのタグ#901と、ユーザデータのI/O先のドライブ117のドライブ#902と、ユーザデータのデータ長を示すレングス903と、ユーザデータのI/O先のドライブアドレス905とを登録する。

20

【0102】

CPU115（ファームウェア142）が、各バスのリンク状態703が“1”（正常）か否かを判断する（S1702）。

【0103】

S1702の判断結果が偽の場合（少なくとも1つのバスのリンク状態703が“0”の場合）（S1702：No）、CPU115（ファームウェア142）が、リンク状態703“1”のバスについての重みとして、リンク状態703“0”よりも高い重みをバス管理テーブル146に設定する（S1703）。なお、その際、CPU115（ファームウェア142）が、リンク状態703“1”に対応したバスの重みを、それらのバスの各々のデータ転送量702に基づき、調整してよい。例えば、データ転送量702が高い程高い重みとしてよい。S1703の後、CPU115（ファームウェア142）が、重みづけラウンドロビンに従いバスを決定、すなわち、バスの重みの比率に基づきバスを決定する（S1707）。例えば、バス0とバス1の重みの比が1：2の場合、バス1が2回選択された後にバス0が1回選択される。S1707の後、CPU115（ファームウェア142）が、選択したバスに対応したメモリアドレスを決定し、アドレス変換テーブル147及びコマンド管理テーブル148を更新する（S1708）。すなわち、CPU115（ファームウェア142）が、アドレス変換テーブル147におけるエントリ（I/Oコマンドに対応するエントリ）に、選択したバスのバス#802と、選択したバスに対応するメモリアドレスであるコマンド指定メモリアドレス804とを登録する。また、CPU115（ファームウェア142）が、コマンド管理テーブル148におけるエントリ（I/Oコマンドに対応するエントリ）に、選択したバスに対応するメモリアドレスであるコマンド指定メモリアドレス904を登録する。

30

40

【0104】

S1702の判断結果が真の場合（全てのバスのリンク状態703が“1”の場合）（S1702：Yes）、CPU115（ファームウェア142）が、リンク状態703“1”のバスについてのデータ転送量702に偏りがあるか否かを判断する（S1704）。ここで言う「偏り」は、例えば、バス管理テーブル146のうちの最大のデータ転送量702と最小のデータ転送量702との差が所定値以上であることでよい。

【0105】

S1704の判断結果が偽の場合（S1704：No）、CPU115（ファームウェア

50

ア 1 4 2) が、パス管理テーブル 1 4 6 に、データ転送量 7 0 2 (パスの負荷) が比較的
低いパスの重みとして比較的高い重みを設定し、データ転送量 7 0 2 が比較的高いパスの
重みとして比較的低い重みを設定する (S 1 7 0 5) 。 S 1 7 0 5 の後、上述した S 1 7
0 7 及び S 1 7 0 8 が行われる。

【 0 1 0 6 】

S 1 7 0 4 の判断結果が真の場合 (S 1 7 0 4 : Y e s) 、 C P U 1 1 5 (ファームウ
ェア 1 4 2) が、ラウンドロビンに従いパスを決定する (S 1 7 0 6) 。 S 1 7 0 6 の後
、上述した S 1 7 0 8 が行われる。

【 実施例 2 】

【 0 1 0 7 】

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
について説明を省略又は簡略する。

【 0 1 0 8 】

図 1 5 は、実施例 2 に係るコントローボックスの構成の一部を示す。

【 0 1 0 9 】

コントローラボックス 7 8 は、二重化されたストレージコントローラ 7 9 としての第 1
及び第 2 のストレージコントローラ 7 9 A 及び 7 9 B を有する。従って、 P C I e - S W
5 7 も二重化されている。第 1 及び第 2 の P C I e - S W 5 7 A 及び 5 7 B に、複数のド
ライブ 1 1 7 が接続されている。

【 0 1 1 0 】

第 1 のメモリ 1 1 1 A における、第 1 の C P U 1 1 5 A が使用するアドレス範囲と、第
2 のメモリ 1 1 1 B における、第 2 の C P U 1 1 5 B が使用するアドレス範囲との少なく
とも一部が重複している。具体的には、例えば、第 1 及び第 2 の C P U 1 1 5 A 及び 1 1
5 B は、それぞれ、ユーザデータやテーブルを、第 1 及び第 2 のメモリ 1 1 1 A 及び 1 1
1 B の同一アドレス範囲に格納することができる。

【 0 1 1 1 】

第 1 及び第 2 の P C I e - S W 5 7 A 及び 5 7 B の各々において、各ドライブ 1 1 7 は
C P U 1 1 5 が接続されている C P U パーティションに接続されてないものの、 S A S (
Serial Attached SCSI) 又は F C (Fibre Channel) と異なり、 C P U 1 1 5 、メモリ 1
1 1 、 P C I e - S W 5 7 及びドライブ 1 1 7 の各々が、 P C I e という同一の通信イン
ターフェースで通信できる。

【 0 1 1 2 】

ここで、もし、或るドライブ 1 1 7 が暴走すると、そのドライブ 1 1 7 の暴走の影響に
より、両方のメモリ 1 1 1 A 及び 1 1 1 B の同一のアドレスに対して不正なアクセスがさ
れ、その同一のアドレスが示す両方の領域内のデータ (例えば、ユーザデータの少なく
とも一部、又は、1 以上のテーブルのうちの少なくとも一部) が破壊され得る。結果として
、システムダウンが生じ得る。

【 0 1 1 3 】

各 P C I e - S W 5 7 には、 N T B が設けられている。 N T B は、アドレス変換を行う
。このため、メモリへのアクセス先アドレスが、 N T B が使用可能なアドレス範囲外のアド
レス (サポート範囲外のアドレス) であれば、 P C I e - S W 5 7 によって、ドライブ
1 1 7 の暴走によってメモリへ不正なアクセスが生じることを防ぐことができる。しかし
、メモリへのアクセス先アドレスが、サポート範囲内のアドレスの場合その不正なアクセ
スを防ぐことはできない。

【 0 1 1 4 】

そこで、本実施例では、第 1 の P C I e - S W 5 7 A がサポートする第 1 のアドレス範
囲と第 2 の P C I e - S W 5 7 B がサポートする第 2 のアドレス範囲と異なる。第 1 及び
第 2 のアドレス範囲は一部の重複も無いことが望ましい。これにより、第 1 のメモリ 1 1
1 A へのアクセス先アドレスが第 1 のアドレス範囲内であったとしても、第 2 のメモリ 1
1 1 B への同一アクセス先アドレスは第 2 のアドレス範囲外となる。結果として、システ

10

20

30

40

50

ムダウンが生じる可能性を軽減することができる。

【0115】

具体的には、例えば、下記の通りである。

【0116】

第1及び第2のCPU115A及び1115Bにそれぞれ接続された第1及び第2のメモリ111A及び111Bの各々が、図16に示す全体アドレス管理テーブル1600を格納する。全体アドレス管理テーブル1600は、バス毎にエントリを有する。各エントリが、CTL#1601、バス#1602、スタートアドレス1603及びサイズ1604といった情報を保持する。CTL#1601は、バスを含んだストレージコントローラ79A又は79Bの識別番号である。バス#1602は、バスの識別番号である。スタートアドレス1603は、バスに対応したアドレス範囲の開始アドレス(メモリアドレス)である。サイズ1604は、アドレス範囲のサイズである。スタートアドレス1603及びサイズ1604により、アドレス範囲が定義される。

10

【0117】

第1のPCIe-SW57Aの内部メモリ188Aが、図17に示す第1のスイッチアドレス管理テーブル1700Aを格納する。第1のスイッチアドレス管理テーブル1700Aは、第1のPCIe-SW57Aに接続されているバス毎にエントリを有する。各エントリが、バス#1701A、スタートアドレス1702A及びサイズ1703Aといった情報を保持する。バス#1701Aは、バスの識別番号である。スタートアドレス1702Aは、バスに対応したアドレス範囲の開始アドレス(メモリアドレス)である。サイズ1703Aは、アドレス範囲のサイズである。スタートアドレス1702A及びサイズ1703Aにより、第1のPCIe-SW57Aがサポートするアドレス範囲が定義される。

20

【0118】

第2のPCIe-SW57Bのメモリ188Bが、図18に示す第2のスイッチアドレス管理テーブル1700Bを格納する。第2のスイッチアドレス管理テーブル1700Bの構成は、第1のスイッチアドレス管理テーブル1700Aの構成と同様である。すなわち、第2のスイッチアドレス管理テーブル1700Bは、第2のPCIe-SW57Bに接続されているバス毎にエントリを有し、各エントリが、バス#1701B、スタートアドレス1702B及びサイズ1703Bといった情報を保持する。

30

【0119】

図16~図18のテーブル1600、1700A及び1700Bによれば、第1のPCIe-SW57Aがサポートする第1のアドレス範囲(メモリ111Aにおける領域のアドレス範囲)と、第2のPCIe-SW57Bがサポートする第2のアドレス範囲(メモリ111Bにおける領域のアドレス範囲)が異なる。これにより、上述したように、第1のメモリ111Aへのアクセス先アドレスが第1のアドレス範囲内であったとしても、第2のメモリ111Bへの同一アクセス先アドレスは第2のアドレス範囲外となる。結果として、システムダウンが生じる可能性を軽減することができる。

【0120】

以上、幾つかの実施例を説明したが、これらは本発明の説明のための例示であって、本発明の範囲をこれらの実施例にのみ限定する趣旨ではない。本発明は、他の種々の形態でも実行することが可能である。例えば、本発明は、PCIe-SW57に代えて、同一ドメインに存在可能なマスタデバイスの数が規定されている通信インターフェースに従い通信を中継する他種のスイッチが採用されても、適用することができる。

40

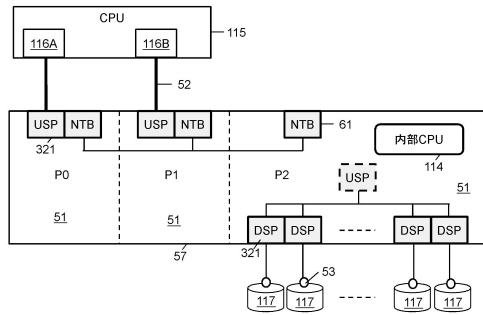
【符号の説明】

【0121】

105...ストレージシステム

【 図 1 】

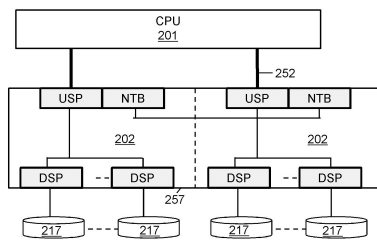
FIG. 1



【 図 2 】

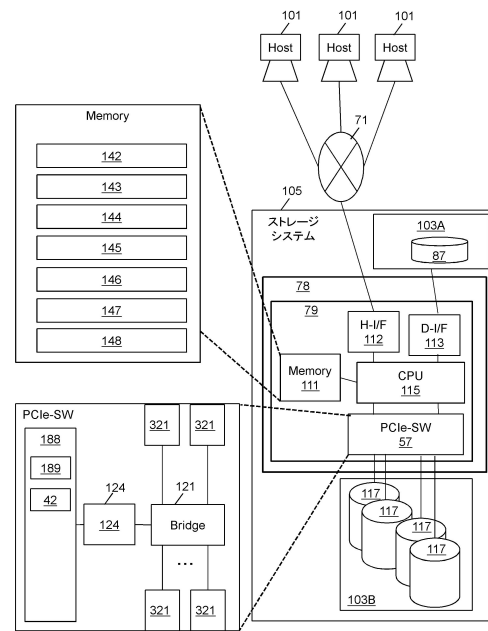
FIG. 2

比較例



【 図 3 】

FIG. 3



【 図 4 】

FIG. 4

ドライブ実装状態管理テーブル
143

ドライブ# 401	実装状態 402	リンク状態 403
0	1	1
1	1	1
2	0	0
...

【 図 7 】

FIG. 7

バス管理テーブル
146

Path# 701	データ転送量 [MB/s] 702	リンク状態 703	重み 704
0	6000	1	1
1	4000	1	2

【 図 5 】

FIG. 5

ドライブ設定状態管理テーブル
144

ドライブ# 501	完了状態 502	結果状態 503
0	0	0
1	0	0
2	1	1
...

【 図 8 】

FIG. 8

アドレス変換テーブル
147

Tag# 801	Path# 802	実メモリアドレス 803	コマンド指定メモリアドレス 804
1	1	0x00000000_00040000	0x0000A000_00040000
2	0	0x00000000_00080000	0x0000B000_00080000
...

【 図 6 】

FIG. 6

ドライブ初期設定管理テーブル
145

ドライブ# 601	Bus# 602	Device# 603	Function# 604	Base address register 605	MSI table register 606
0					
1					
2					
...

【 図 9 】

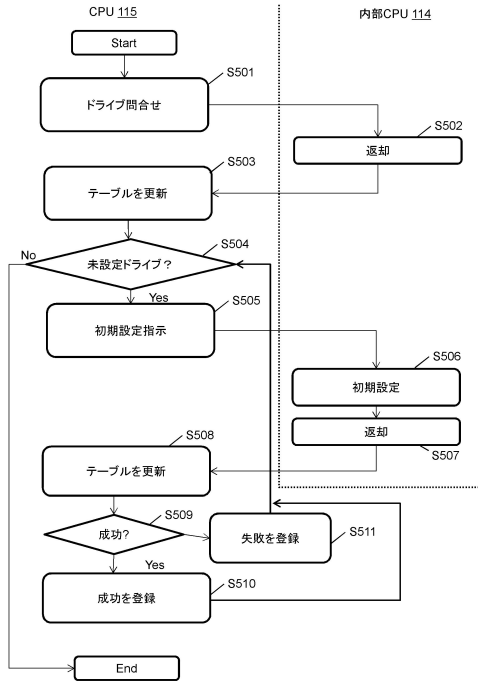
FIG. 9

コマンド管理テーブル
148

Tag# 901	ドライブ# 902	Length 903	コマンド指定メモリアドレス 904	ドライブアドレス 905
1	0	0x20000	0x0000A000_00040000	0x00005000
2	1	0x10000	0x0000B000_00080000	0x00001000
...

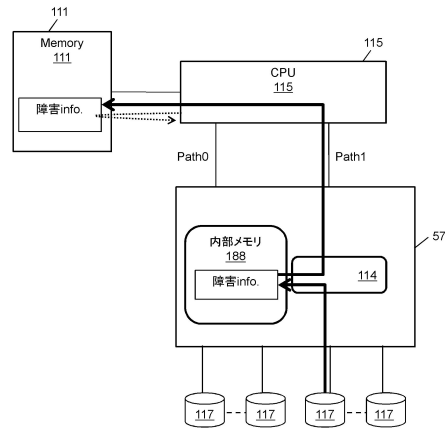
【図10】

FIG. 10



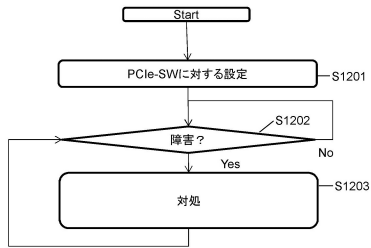
【図11】

FIG. 11



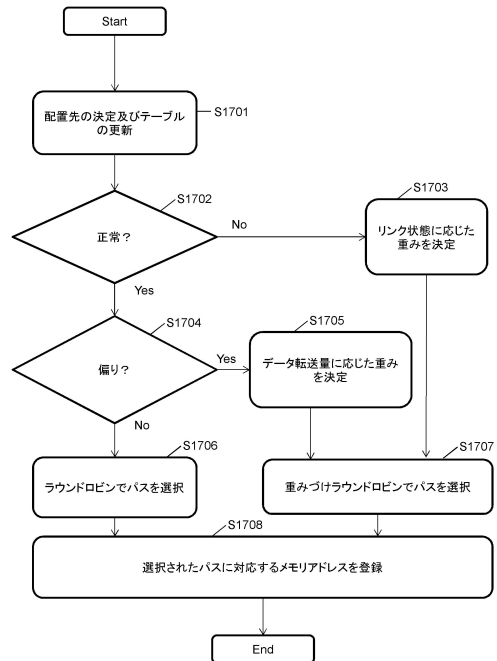
【図12】

FIG. 12



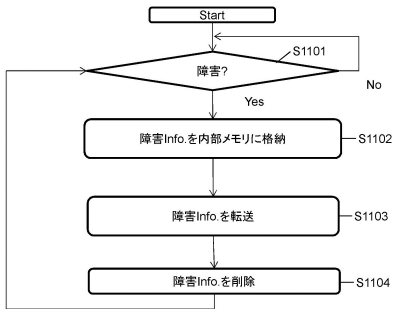
【図14】

FIG. 14



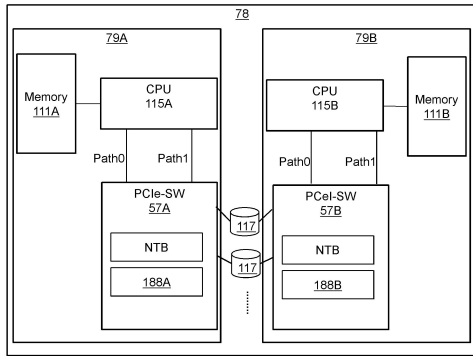
【図13】

FIG. 13



【 図 15 】

FIG. 15



【 図 16 】

FIG. 16

全体アドレス管理テーブル
1600

CTL# 1601	Path# 1602	Start address 1603	Size 1604
0	0	0x0000A000_00000000	0x00001000_00000000
0	1	0x0000B000_00000000	0x00001000_00000000
1	0	0x0000C000_00000000	0x00001000_00000000
1	1	0x0000D000_00000000	0x00001000_00000000

【 図 17 】

FIG. 17

第1のスイッチアドレス管理テーブル
1700A

Path# 1701A	Start address 1702A	Size 1703A
0	0x0000A000_00000000	0x00001000_00000000
1	0x0000B000_00000000	0x00001000_00000000

【 図 18 】

FIG. 18

第2のスイッチアドレス管理テーブル
1700B

Path# 1701B	Start address 1702B	Size 1703B
0	0x0000C000_00000000	0x00001000_00000000
1	0x0000D000_00000000	0x00001000_00000000

フロントページの続き

(56)参考文献 特開2011-248662(JP,A)
特開2014-002545(JP,A)
特表2015-501501(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F3/06-3/08
G06F13/10-13/14
G06F13/38-13/42