

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6955466号  
(P6955466)

(45) 発行日 令和3年10月27日(2021.10.27)

(24) 登録日 令和3年10月5日(2021.10.5)

(51) Int.Cl.	F I
<b>G06F 3/06 (2006.01)</b>	G06F 3/06 304B
<b>G06F 3/08 (2006.01)</b>	G06F 3/06 301B
<b>G06F 13/10 (2006.01)</b>	G06F 3/06 301C
<b>G06F 13/38 (2006.01)</b>	G06F 3/06 302B
<b>H04L 12/771 (2013.01)</b>	G06F 3/06 305A

請求項の数 20 (全 25 頁) 最終頁に続く

(21) 出願番号	特願2018-52308 (P2018-52308)	(73) 特許権者	390019839
(22) 出願日	平成30年3月20日 (2018.3.20)		三星電子株式会社
(65) 公開番号	特開2018-173949 (P2018-173949A)		Samsung Electronics Co., Ltd.
(43) 公開日	平成30年11月8日 (2018.11.8)		大韓民国京畿道水原市靈通区三星路129
審査請求日	令和3年1月29日 (2021.1.29)		129, Samsung-ro, Yeongtong-gu, Suwon-si, Gyeonggi-do, Republic of Korea
(31) 優先権主張番号	62/480, 113	(74) 代理人	110000051
(32) 優先日	平成29年3月31日 (2017.3.31)		特許業務法人共生国際特許事務所
(33) 優先権主張国・地域又は機関	米国 (US)	(72) 発明者	カチュア, ラムダス
(31) 優先権主張番号	15/617, 925		アメリカ合衆国, 95014, カリフォルニア州, クパチーノ, ノルマンディー ウエイ 7665
(32) 優先日	平成29年6月8日 (2017.6.8)		最終頁に続く
(33) 優先権主張国・地域又は機関	米国 (US)		
早期審査対象出願			

(54) 【発明の名称】 SSDのデータ複製システム及び方法

(57) 【特許請求の範囲】

【請求項1】

複数のソリッドステートドライブ (eSSD)、前記複数のeSSDの各々にダウンリンクを提供するように構成されたファブリックスイッチ、並びに前記ファブリックスイッチ及び前記複数のeSSDを制御するように構成されたベースボード管理コントローラ (BMC) をシャーシ内に含むデータ複製システムのデータ複製方法であって、

前記BMCを用いて、

前記複数のeSSDの中の1つをアクティブeSSDとして構成する段階と、

前記複数のeSSDの中の1つ以上を1つ以上のパッシブeSSDとして構成する段階と、

前記ファブリックスイッチをプログラムして、前記アクティブeSSDにアドレス指定されたネットワークパケットを第1のダウンリンクポートを介して前記アクティブeSSD及び前記第1のダウンリンクポートとは異なる第2のダウンリンクポートを介して前記1つ以上のパッシブeSSDの両方に伝送する段階と、を有し、

前記アクティブeSSDで、前記ファブリックスイッチを通じて前記第1のダウンリンクポートを介してホストからのホストデータの書き込み命令を含むネットワークパケットを受信する段階と、

前記ホストデータの書き込み命令に関連するホストデータを前記アクティブeSSDのメモリ位置に格納する段階と、

前記アクティブeSSDから前記ホストデータの書き込み命令に対応する前記アクティブ

e S S Dのメモリ位置及び命令を前記1つ以上のパッシブ e S S Dに送信する段階と、

前記アクティブ e S S Dのメモリ位置及び前記アクティブ e S S Dから受信された命令並びに前記第2のダウンリンクポートを介して前記ファブリックスイッチによって伝送された前記ネットワークパケット内で受信された前記ホストデータを用いて前記1つ以上のパッシブ e S S Dに前記ホストデータのコピーを格納する段階と、を有することを特徴とする方法。

【請求項2】

前記アクティブ e S S Dは、前記ホストによって発見可能であり、

前記1つ以上のパッシブ e S S Dは、前記ホストによって発見不可能であることを特徴とする請求項1に記載の方法。

10

【請求項3】

前記BMCを用いて、前記アクティブ e S S Dと前記1つ以上のパッシブ e S S Dとの間の専用通信チャンネルを設定して、前記アクティブ e S S Dから前記メモリ位置及び前記命令を前記1つ以上のパッシブ e S S Dに送信する段階を更に含むことを特徴とする請求項1に記載の方法。

【請求項4】

前記専用通信チャンネルは、前記シャーシ内の前記ファブリックスイッチ又はP C I eスイッチを通じて設定されることを特徴とする請求項3に記載の方法。

【請求項5】

前記アクティブ e S S Dを用いて、前記ホストデータの書込み命令にตอบสนองしてR D M A ( r e m o t e d i r e c t m e m o r y a c c e s s ) 読出し要請を前記ホストに発行する段階と、

20

前記ホストからの前記ホストデータを1つ以上のR D M A 読出し応答パケット内で受信する段階と、を更に含むことを特徴とする請求項1に記載の方法。

【請求項6】

前記R D M A 読出し応答パケットは、前記ホストデータのデータチャンクを含むことを特徴とする請求項5に記載の方法。

【請求項7】

前記1つ以上のパッシブ e S S Dの各々は、前記R D M A 読出し応答パケットを除外して、前記ファブリックスイッチから受信されたパケットを破棄することを特徴とする請求項5に記載の方法。

30

【請求項8】

前記アクティブ e S S Dで、一貫性モードに基づいて、前記1つ以上のパッシブ e S S Dから応答を受信した後、又は前記1つ以上のパッシブ e S S Dから応答を受信せずに前記ホストデータを前記アクティブ e S S D内に格納した後、完了キューエントリを前記ホストに送信する段階を更に含むことを特徴とする請求項1に記載の方法。

【請求項9】

前記1つ以上のパッシブ e S S Dの各々は、前記メモリ位置及び前記命令を格納するための第1バッファ並びに前記ホストデータのデータチャンクを格納するための第2バッファを含むことを特徴とする請求項1に記載の方法。

40

【請求項10】

前記BMCを用いて、前記アクティブ e S S Dに関連するエラーを検出する段階と、

前記シャーシ内の複数の e S S Dの中の1つの e S S Dを選択する段階と、

前記 e S S Dを新たなアクティブ e S S Dとして構成する段階と、

新たなパッシブ e S S Dが必要であるか否かを判別する段階と、

前記新たなパッシブ e S S Dを構成する段階と、

前記新たなアクティブ e S S Dと前記1つ以上のパッシブ e S S Dとを関連付けるように前記ファブリックスイッチをプログラムする段階と、を更に含むことを特徴とする請求項1に記載の方法。

【請求項11】

50

複数のソリッドステートドライブ ( e S S D )、前記複数の e S S D の各々にダウンリンクを提供するように構成されたファブリックスイッチ、並びに前記ファブリックスイッチ及び前記複数の e S S D を制御するように構成されたベースボード管理コントローラ ( B M C ) を含むシャーシを備え、

前記 B M C は、前記複数の e S S D の中の 1 つをアクティブ e S S D として構成し、前記複数の e S S D の中の 1 つ以上を 1 つ以上のパッシブ e S S D として構成し、前記ファブリックスイッチをプログラムして、前記アクティブ e S S D にアドレス指定されたネットワークパケットを第 1 のダウンリンクポートを介して前記アクティブ e S S D 及び前記第 1 のダウンリンクポートとは異なる第 2 のダウンリンクポートを介して前記 1 つ以上のパッシブ e S S D の両方に伝送し、

10

前記アクティブ e S S D は、前記ファブリックスイッチを通じて前記第 1 のダウンリンクポートを介してホストからのホストデータの書き込み命令を含むネットワークパケットを受信し、前記ホストデータの書き込み命令に関連するホストデータをメモリ位置に格納するように構成され、

前記アクティブ e S S D は、前記ホストデータに対応する前記メモリ位置及び命令を前記 1 つ以上のパッシブ e S S D に送信するように更に構成され、

前記 1 つ以上のパッシブ e S S D の各々は、前記アクティブ e S S D から受信された前記アクティブ e S S D のメモリ位置及び命令並びに前記第 2 のダウンリンクポートを介して前記ファブリックスイッチによって伝送された前記ネットワークパケット内で受信された前記ホストデータを用いて前記ホストデータのコピーを格納するように構成されることを特徴とするデータ複製システム。

20

【請求項 1 2】

前記アクティブ e S S D は、前記ホストによって発見可能であり、前記 1 つ以上のパッシブ e S S D は、前記ホストによって発見不可能であることを特徴とする請求項 1 1 に記載のデータ複製システム。

【請求項 1 3】

前記アクティブ e S S D と前記 1 つ以上のパッシブ e S S D との間に専用通信チャンネルが設定されて、前記メモリ位置及び前記命令を送信することを特徴とする請求項 1 1 に記載のデータ複製システム。

【請求項 1 4】

前記専用通信チャンネルは、前記シャーシ内の前記ファブリックスイッチ又は P C I e スイッチを通じて設定されることを特徴とする請求項 1 3 に記載のデータ複製システム。

30

【請求項 1 5】

前記アクティブ e S S D は、前記ホストデータの書き込み命令にตอบสนองして R D M A ( r e m o t e d i r e c t m e m o r y a c c e s s ) 読出し要請を前記ホストに発行し、前記ホストからの前記ホストデータを 1 つ以上の R D M A 読出し応答パケット内で受信するように構成されることを特徴とする請求項 1 1 に記載のデータ複製システム。

【請求項 1 6】

前記 R D M A 読出し応答パケットは、前記ホストデータのデータチャンクを含むことを特徴とする請求項 1 5 に記載のデータ複製システム。

40

【請求項 1 7】

前記 1 つ以上のパッシブ e S S D の各々は、前記 R D M A 読出し応答パケットを除外して、前記ファブリックスイッチから受信されたパケットを破棄することを特徴とする請求項 1 5 に記載のデータ複製システム。

【請求項 1 8】

前記アクティブ e S S D は、一貫性モードに基づいて、前記 1 つ以上のパッシブ e S S D から応答を受信した後、又は前記 1 つ以上のパッシブ e S S D から応答を受信せずに前記ホストデータを前記アクティブ e S S D 内に格納した後、完了キューエントリを前記ホストに送信することを特徴とする請求項 1 1 に記載のデータ複製システム。

【請求項 1 9】

50

前記1つ以上のパッシブeSSDの各々は、前記メモリ位置及び前記命令を格納するための第1バッファ並びに前記ホストデータのデータチャンクを格納するための第2バッファを含むことを特徴とする請求項11に記載のデータ複製システム。

【請求項20】

前記BMCは、

前記アクティブeSSDに関連するエラーを検出し、

前記シャーシ内の複数のeSSDの中の1つのeSSDを選択し、

前記eSSDを新たなアクティブeSSDとして構成し、

新たなパッシブeSSDが必要であるか否かを判別し、

前記新たなパッシブeSSDを構成し、

前記新たなアクティブeSSD、前記1つ以上のパッシブeSSD、及び前記新たなパッシブeSSDを関連付けるように前記ファブリックスイッチをプログラムするように更に構成されることを特徴とする請求項11に記載のデータ複製システム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ複製に関し、より詳細には、NVMe-oF eSSD (solid-state drive) のデータ複製システム及び方法に関する。

【背景技術】

20

【0002】

SSD (solid-state drives) は従来のハードディスクドライブ (hard disk drives: HDDs) を代替する現代情報技術 (information technology: IT) インフラの主要格納要素になっている。SSDは、低いレイテンシ (latency)、高いデータ読出し/書込み処理量、及び信頼できる使用者データ格納を提供する。NVMe-oF (non-volatile memory express over fabrics) は数百及び数千個のNVMe-互換SSDをイーサネット (登録商標) (Ethernet (登録商標)) (以降、「イーサネット」) の (登録商標) は省略する。) のようなネットワークファブリック (network fabric) を通じて連結する新技術である。イーサネット連結を通じてNVMe-oF標準と互換されるSSDは、イーサネット-連結 (Ethernet-attached) SSD、又は簡単にeSSDという。

30

【0003】

NVMe-oFプロトコルは多数のeSSDがネットワークファブリックを通じて遠隔ホストに連結されるrDAS (remote direct-attached storage) プロトコルを支援する。NVMe-oFプロトコルは、ネットワークファブリックを通じて遠隔ホストとeSSDとの間にNVMe命令、データ、及び応答を伝達することができる信頼性ある伝送サービスを提供するためにRDMA (remote direct memory access) プロトコルを支援する。RDMAサービスを提供する伝送プロトコルの例としてInfiniBand、iWARP、RoCE v1、及びRoCE v2がある。

40

【0004】

使用者データに対する信頼性あるアクセスはデータストレージシステムの最も重要な要求事項の中の1つである。信頼性及びデータ可用性を必要な程度に達成するためにデータストレージシステムの多様な地点で多様な技術及び方法が使用される。例えば、ストレージ装置に書き込まれたデータは、ストレージ装置が使用できなくなった時のバックアップストレージ装置として1つ以上の他のストレージ装置に複製される。このデータ複製スキームはしばしばミラーリング (mirroring) 又はバックアップ (back-up) という。

【0005】

50

複数のeSSDが接続されたシャーシ(chassis)内で、各eSSDはイーサネットを通じてホストに直接連結される。ホストはデータミラーリングのために2つ以上のeSSDに対してデータ複製を遂行する。しかし、ホストによるデータ複製機能の具現は、ホスト及びシステムソフトウェアに負担を加重させ、結果的にデータ入/出力(I/O)動作に付加的なレイテンシ(latency)を追加し、データストレージシステムに対する全体費用を増加させる。

【0006】

或いは、インライン(inline)RAID(redundant array of independent disks)コントローラがホストに対して透過的(トランスペアレント)にデータ複製を遂行するためにSSDに埋め込まれる。しかし、RAIDコントローラは一般的に高く付加的なレイテンシを追加してI/O性能を低下させる。また、RAIDコントローラが埋め込まれたSSDは、多くの負荷又は最大の秒当たり入出力動作(input/output operations per second: IOPS)で最大電力コスト又はそれと類似な電力を消費し、NVMe-oF標準に応じたeSSDの電力消費要求事項を充足させない。従って、eSSDに埋め込まれたRAIDコントローラによるデータ複製はNVMe-oF-互換データストレージシステムに対して適合するソリューションではない可能性がある。

【先行技術文献】

【特許文献】

【0007】

【特許文献1】米国特許第9,015,525号明細書

【特許文献2】米国特許第9,471,259号明細書

【特許文献3】米国特許出願公開第2016/0352831号明細書

【特許文献4】米国特許出願公開第2010/0082792号明細書

【特許文献5】米国特許出願公開第2016/0299704号明細書

【特許文献6】米国特許出願公開第2016/0004877号明細書

【発明の概要】

【発明が解決しようとする課題】

【0008】

本発明は、上記従来の問題点に鑑みてなされたものであって、本発明の目的は、SSDs(solid-state drives)のデータ複製システム及び方法を提供することにある。

【課題を解決するための手段】

【0009】

上記目的を達成するためになされた本発明の一態様によるデータ複製方法は、複数のソリッドステートドライブ(eSSD)、前記複数のeSSDの各々にダウンリンクを提供するためのファブリックスイッチ、並びに前記ファブリックスイッチ及び前記複数のeSSDを制御するためのベースボード管理コントローラ(BMC)をシャーシ内に含むデータ複製システムのデータ複製方法であって、前記BMCを用いて、前記複数のeSSDの中の1つをアクティブeSSDとして構成する段階と、前記複数のeSSDの中の1つ以上を1つ以上のパッシブeSSDとして構成する段階と、前記ファブリックスイッチをプログラムして、前記アクティブeSSDに向かうパケットを前記アクティブeSSD及び前記1つ以上のパッシブeSSDの両方に伝送する段階と、を有し、前記アクティブeSSDで、ホストからホストデータ書き込み命令を受信する段階と、前記アクティブeSSDから前記ホストデータに対応するアドレス及び命令を前記1つ以上のパッシブeSSDに送信する段階と、前記アクティブeSSDに前記ホストデータを格納する段階と、前記アクティブeSSDから受信されたアドレス及び命令並びに前記ファブリックスイッチによって伝送されたパケット内で受信された前記ホストデータを用いて前記1つ以上のパッシブeSSDの各々に前記ホストデータのコピーを格納する段階と、を含む。

【0010】

上記目的を達成するためになされた本発明の一態様によるデータ複製システムは、複数のソリッドステートドライブ（eSSD）、前記複数のeSSDの各々にダウンリンクを提供するためのファブリックスイッチ、並びに前記ファブリックスイッチ及び前記複数のeSSDを制御するためのベースボード管理コントローラ（BMC）をシャーシ内に備え、前記BMCは、前記複数のeSSDの中の1つをアクティブeSSDとして構成し、前記複数のeSSDの中の1つ以上を1つ以上のパッシブeSSDとして構成し、前記ファブリックスイッチをプログラムして、前記アクティブeSSDに向かうパケットを前記アクティブeSSD及び前記1つ以上のパッシブeSSDの両方に伝送し、前記アクティブeSSDは、ホストから受信されたホストデータの書き込み命令に応答してホストデータを格納し、前記アクティブeSSDは、前記ホストデータに対応するアドレス及び命令を前記1つ以上のパッシブeSSDに伝送し、前記1つ以上のパッシブeSSDの各々は、前記アクティブeSSDから受信されたアドレス及び命令並びに前記ファブリックスイッチによって伝送されたパケット内で受信された前記ホストデータを用いて前記ホストデータのコピーを格納する。

10

#### 【発明の効果】

#### 【0011】

本発明のデータ複製システム及び方法は、データ複製機能が主にeSSD（即ち、アクティブ（active）eSSD及びパッシブ（passive）eSSD）及びシャーシ（chassis）内のイーサネットスイッチで具現されるため、データ複製を提供するための費用の側面で効果的であり、効率的なソリューションである。本発明のデータ複製システム及び方法によれば、複雑なシステムソフトウェア及びハードウェアに対する必要性を無くし、追加演算、格納装置、及び電力のためのホスト側の負担及び費用を減少させる。本発明のデータ複製はホスト側の最小限の変更のみでeSSD及びイーサネットスイッチを含むシャーシ内に具現される。

20

本発明のデータ複製システム及び方法によれば、アクティブeSSDと関連するパッシブeSSDとの間に完全な一貫性を保証することができ、データ損失がない。一部のデータ損失を許容する場合、使用者アプリケーションは更に良いレイテンシを達成することができる。

#### 【図面の簡単な説明】

#### 【0012】

30

【図1】一実施形態によるNVMe-oFシステムの一例のブロック図である。

【図2】一実施形態によるNVMe-oFシステムの他の例のブロック図である。

【図3】一実施形態によるNVMe-oFシステムのデータ複製のプロセスを示す図である。

【図4】一実施形態によるシャーシ内のeSSD及びイーサネットスイッチを初期化及びプログラミングしてデータ複製を具現するためのフローチャートである。

【図5】一実施形態によるアクティブeSSDによってホスト命令を処理するためのフローチャートである。

【図6】一実施形態によるデータ複製プロセスの一例を示す図である。

【図7】一実施形態によるパッシブeSSDによってホスト命令を処理するためのフローチャートである。

40

【図8】一実施形態によるアクティブeSSDからパッシブeSSDにLBA及び命令を伝達するためのフローチャートである。

【図9】一実施形態によるパッシブeSSDによってデータ複製を処理するためフローチャートである。

【図10】一実施形態による誤謬解決動作のためのフローチャートである。

#### 【発明を実施するための形態】

#### 【0013】

以下、本発明を実施するための形態の具体例を、図面を参照しながら詳細に説明する。

#### 【0014】

50

具現及びイベントの組合せの多様な新規な細部事項を含む上記及び他の望ましい特徴は、図面を参照してより具体的に説明し、特許請求の範囲に記載する。ここに記述した特定システム及び方法は、単なる例示として図示したものであり、これに制限されないことを理解すべきである。当業者が理解することができるように、本明細書で説明する原理及び特徴は本発明の範囲を逸脱せずに多様であり数多くの実施形態に適用される。

【 0 0 1 5 】

本明細書で説明する一部に含まれる図面は、現在の望ましい実施形態を図示し、上述した一般的な説明及び下記に説明する望ましい実施形態の詳細な説明と共に本明細書にその原理を説明し、教示する役割をする。

【 0 0 1 6 】

図面は必ずしも寸法通りに描かれたものではなく、類似な構造又は機能の要素は図面の全体に亘って例示的な目的のために類似な参照符号で一般的に表示する。図面は本明細書で説明する多様な実施形態の説明を容易にするためのものである。図面は、本願に開示する教示の全ての様態を記述するものではなく、特許請求の範囲を制限しない。

【 0 0 1 7 】

本明細書に開示する各々の特徴及び教示は個別的に又は他の特徴及び教示と共に利用されてイーサネットＳＳＤでデータ複製を提供する。個別的に又は組合せとして、多くの追加の特徴及び教示を利用する代表的な例を、図面を参照しながらより詳細に説明する。詳細な説明は、本発明の様態を実施するための詳細な説明を当業者に教示するためのものであり、特許請求の範囲を制限しようとするものはない。従って、詳細な説明で開示する特徴の組合せは、最も広い意味での教示を遂行するために必須的ではなく、単なる本発明の教示の特に代表的な例を説明するために使用する。

【 0 0 1 8 】

以下の説明で、単なる説明の目的のために、特定名称を本開示の完全な理解を提供するために提示する。しかし、当業者にはこのような特定細部事項が本開示の教示を実行するために要求されないことが明確である。

【 0 0 1 9 】

本明細書の詳細な説明の一部分はコンピュータメモリ内のデータビットに対する演算のアルゴリズム及びシンボル表現で提供される。このようなアルゴリズムの説明及び表現はデータ処理技術分野の当業者がそれらの作業内容を当業者に効果的に伝達するために使用される。本明細書のアルゴリズムは一般的に所望する結果を誘導する段階の一貫性あるシーケンスと看做す。各段階は物理的量を物理的に操作しなければならない段階である。一般的に必ずしもそうではないが、このような量は格納、伝送、結合、比較、及びその他の操作が可能な電気又は磁気信号の形態を有する。殆どの信号をビット、値、エレメント、記号、文字、用語、数字等を参照して一般的に使用することが便利であることが立証される。

【 0 0 2 0 】

しかし、これら及び類似な用語の全ては、適切な物理的量に関連しており、単なるこれらの量に適用される便利なラベルであることに留意しなければならない。以下の説明から明確なように、特別に言及しない限り、説明の全体に亘って“処理”、“コンピューティング”、“計算”、“決定”、“ディスプレイ”のような用語を使用する論議はコンピュータシステム又は類似な電子コンピューティング装置の動作及びプロセスを言及し、コンピュータシステムのレジスター及びメモリ内の物理（電子）量に表現されたデータをコンピュータシステムメモリ若しくはレジスター又はその他の情報ストレージ装置、伝送装置若しくはディスプレイ装置内で物理量として類似に表現される他のデータに操作及び変換するコンピュータシステム又は類似な電子コンピューティング装置の動作及びプロセスを称する。

【 0 0 2 1 】

ここに提示するアルゴリズムは本質的に任意の特定コンピュータ又は他の装置に関連しない。多様な汎用システム、コンピュータサーバー、又はパーソナルコンピュータは、本

10

20

30

40

50

明細書の教示に応じるプログラムと共に使用されるか、或いは要求された方法の段階を遂行するためにより専門化された装置を構成することが便利なることもある。このような多様なシステムに必要な構造は下記の説明で示される。本明細書で説明するように本発明の開示内容を具現するために多様なプログラミング言語が使用されることを理解すべきである。

#### 【0022】

また、代表例及び従属請求項の多様な特徴は本発明の教示の追加的な有用な実施形態を提供するために具体的に及び明示的に列挙されない方式で結合される。また、全ての値の範囲又は個体のグループは請求された主題を制限するための目的のみならず、独創的な公開目的に全ての可能な中間値又は中間個体を開示する点に明示的に留意する。また、図面に示す構成要素の寸法及び形状は本発明の教示がいかに遂行されるかを理解することを助けるために設計されたものとして、実施形態に図示する寸法及び形状を制限しようとするものではない点を明示的に言及する。

#### 【0023】

本開示はNVMe-oF標準（本明細書でイーサネットSSD又はeSSDと称する）と互換可能な1つ以上のSSDを使用したデータ複製（data replication）システム及び方法を提供する。本発明のデータ複製システム及び方法によって用いられるデータ複製は、ホスト参加、又はCPU（central processing unit）やRAIDコントローラのような高価で性能に影響を及ぼす外部構成要素無しに達成される。

#### 【0024】

一実施形態によると、シャーシ160内のeSSD170の各々は、アクティブ、パッシブ、及びノーマルの3つのモードに構成されて動作する。シャーシのベースボード管理コントローラ（base board management controller：BMC）は、シャーシ内のeSSDの中の1つをアクティブeSSDとしてプログラムし、残りのeSSDの中の1つ以上をパッシブeSSDとしてプログラムする。アクティブeSSD及びパッシブeSSDのどちらにも構成されないeSSDは、ノーマルモードで動作し、標準NVMe-oF SSDとして動作する。アクティブ及びパッシブeSSDの中で、アクティブeSSDのみが遠隔ホスト（remote host）に可視化され（発見され）、遠隔ホストと共にNVMe-oFプロトコルの終了を遂行する。反面、パッシブeSSDは遠隔ホストに可視化され（発見され）ない。ネットワークファブリック（例えば、イーサネット）に位置する遠隔ホストは説明の便宜のために単にホスト（host）と称する。

#### 【0025】

アクティブeSSDは、ホストによって発行された命令を処理し、ホストにそしてホストからのデータ伝送を遂行し、ホスト命令の成功的な完了又は実行失敗を示すために完了キュー（completion queue：CQ）エントリ（entries）をホストに送信する。シャーシ内のパッシブeSSDは、ホスト命令を実行するか、或いはホストプロトコルに参加することもない。代わりにパッシブeSSDはアクティブeSSDに向かう進入トラフィックパケット（ingress traffic packets）を受動的に受信する。例えば、シャーシ内のイーサネットスイッチ（ファブリックスイッチ）はアクティブeSSDに向かう全ての進入パケットをアクティブeSSDに関連するようにプログラムされたパッシブeSSDに複製するようにプログラムされる。

#### 【0026】

一実施形態によると、パッシブeSSDはRDMA読出し応答（RDMA READ RESPONSE）パケットを除外した全ての複製された進入パケットを破棄する。RDMA読出し応答パケットはホスト書込みデータを提供し、ホスト書込みデータはホストのFIFOバッファに維持される。FIFOバッファ内のホスト書込みデータはアクティブeSSDによってフェッチ（fetch）される。ホスト書込みデータをフェッチした後、アクティブeSSDは命令（例えば、書込み命令）に関連するLBA（logical



block address) 及び / 又はホスト書込みデータを書き込むためのネームスペースを関連するパッシブ eSSD に提供する。アクティブ eSSD によって提供された命令及び LBA、そして複製された RDMA 読出し応答パケットから受信されたデータを用いて、パッシブ eSSD は複製されたデータを自体の不揮発性メモリに書き込む。パッシブ eSSD の不揮発性メモリはフラッシュメモリ又は永久メモリ (persistent memory) である。構成によって、1 つ以上のパッシブ eSSD がホスト書込みデータを複製する。

【0027】

一実施形態によると、アクティブ eSSD が関連するパッシブ eSSD に送信した LBA / ネームスペース及び書込み命令は送信キュー及び受信キューを含む RDMA キュー対 (queue pair: QP) 内に細部事項を含む。アクティブ eSSD とパッシブ eSSD との間に設定された経路内にホストによってイネーブルされた QP を通じて追加情報交換を要求する多数のアクティブ QP がある。一実施形態によると、パッシブ eSSD は全ての RDMA 読出し応答で QP 及びデータシーケンシングのようなホスト書込み関連に関するデータチャンク上の追加的な情報のための RDMA ヘッダーを覗き見する。

【0028】

一実施形態によると、アクティブ eSSD は、ホストからホスト書込み要請を受信した後、1 つ以上の RDMA 読出し要請をホストに送信する前に書込み命令を 1 回送信する。アクティブ eSSD は、ホスト書込みがどのように 1 つ以上の個別 RDMA 読出し要請に分類されるか、そして関連するデータがどのようにセグメント分散収集リスト (segment scatter gather list: SGL) から受信されるかを予め計算する。これは PCIe バス又はネットワークファブリックのみならず、専用チャネルのトラフィックを減少させ、これにより多数の QP 連結が多数のデータ複製セットに亘ってなされる時に I/O 性能を向上させる。また、データがエンクロージャー (enclosures)、ラック (racks)、データセンター、及び様々な地理的位置に亘って複製される時にデータ複製のオーバーヘッドを減少させる。

【0029】

パッシブ eSSD は、ホスト書込みデータが成功的に維持されたか、又はホスト書込みデータの格納が失敗したことを示す確認をアクティブ eSSD に送信するように構成される。アクティブ eSSD は完了キューエントリをホストに送信する前に関連する全てのパッシブ eSSD からの確認を待機する。このような確認メカニズムは、常にホスト書込みデータのデータ持続性を保護し、保証する。

【0030】

本発明のデータ複製システム及び方法はデータ複製機能を支援するために RAID-on-Chip (ROC) のためのターゲット側の x86 サーバー又は外部 ASIC (application-specific integrated circuit) を必要としない。本発明のデータ複製システム及び方法は、ホストの付加的なハードウェア及びソフトウェアオーバーヘッドを除去又は減少させることによってデータストレージ性能を向上させながら、総所有コスト (total cost of ownership: TCO) を減少させる。本発明のデータ複製システム及び方法はデータ複製がデータ I/O 動作と同時に遂行される従来のインライン (in-line) データ複製ソリューションに比べて向上した I/O レイテンシ及び帯域幅性能を有する。

【0031】

本発明のデータ複製システム及び方法は、データ複製機能が主に eSSD (即ち、アクティブ (active) eSSD 及びパッシブ (passive) eSSD) 及びシャーシ (chassis) 内のイーサネットスイッチで具現されるため、データ複製を提供するための費用の側面で効果的であり、効率的なソリューションである。本発明のデータ複製システム及び方法は、複雑なシステムソフトウェア及びハードウェアに対する必要性を無くし、追加演算、格納装置、及び電力のためのホスト側の負担及び費用を減少させる。本発明のデータ複製機能は、ホスト側の最小限の変更で、又は変更無しで eSSD 及び

10

20

30

40

50

イーサネットスイッチを含むシャーシ内に具現される。

【0032】

図1は、一実施形態によるNVMe-oFシステムの一例のブロック図である。NVMe-oFシステム100はホスト110及び1つ以上のNVMe-oF-互換イーサネットSSD(eSSD)を含むシャーシ160(本明細書でeSSDシャーシとも称する)を含む。例えば、シャーシ160は24又は48個のeSSDを含む。シャーシ160内のeSSDは各々eSSD(170a~170n)(本明細書で集合的にeSSD170と称する)で表示される。ホスト110は、アプリケーション111、オペレーティングシステム(operating system:OS)とファイルシステム(file system:FS)112、及びNVMe-oFドライバー113を含む。ホスト110のイニシエータ(initiator、例えばアプリケーション111)はNVMe-oFドライバー113を用いてイーサネット150を通じてeSSD170との連結を設定する。シャーシ160は、イーサネットスイッチ161、ベースボード管理コントローラ(base board management controller:BMC)162、及びPCIeスイッチ163を含む。ファブリックスイッチとしてのイーサネットスイッチ161はミッドプレーン(midplane)165を通じてeSSD170にイーサネット連結を提供し、PCIeスイッチ163はミッドプレーン165を通じてeSSD170に管理インターフェイス164を提供する。この例で、eSSD170の中の1つはアクティブeSSDとして構成され、他のeSSDはパッシブeSSDとして構成される。BMC162はシステム管理者によって与えられた命令に応じてeSSD170をプログラムする。

10

20

【0033】

イーサネットスイッチ161はホスト110とeSSD170との間にネットワーク連結を提供する。イーサネットスイッチ161は1つ以上のホストに連結するための大容量(例えば、100Gbps)アップリンク(uplinks)を有する。また、イーサネットスイッチ161はeSSD170に連結するための多数の低容量(例えば、25Gbps)ダウンリンク(downlinks)を有する。例えば、イーサネットスイッチ161は100Gbpsの12個のアップリンク及び25Gbpsの24又は48個のダウンリンクを含む。イーサネットスイッチ161はBMC162に対する専用(special)構成/管理ポート(図示せず)を有する。

30

【0034】

BMC162は、イーサネットスイッチ161、PCIeスイッチ163、及びeSSD170を含むシャーシ160内の内部構成要素を管理する。BMC162はシステム管理のためにPCIe及び/又はSMBusインターフェイスを支援する。BMC162はアクティブeSSD及びパッシブeSSDを構成し、イーサネットスイッチ161をプログラムする。アクティブeSSDが失敗すると、BMC162は誤謬解決スイッチング(failover switching)を遂行し、新たなアクティブeSSDを指定し、必要に応じて追加的なパッシブeSSDを構成する。

【0035】

一実施形態によると、シャーシ160内のeSSD170は、アクティブ、パッシブ、及びノーマルの3つのモードの中の1つで動作する。アクティブeSSD及び1つ以上のパッシブeSSDとして構成されたeSSD170は所望のデータ複製を提供する。一部の実施形態で、シャーシ160は多数のアクティブ/パッシブeSSDセットを含む。所定のeSSDに対してデータ複製を必要としない場合、BMC162はeSSDをノーマルモードに構成する。ノーマルeSSDで構成されたeSSDは標準NVMe-oFSSDとして動作する。

40

【0036】

図2は、一実施形態によるNVMe-oFシステムの他の例のブロック図である。NVMe-oFシステム200はホスト210及びイーサネット250を通じてホスト210に連結される複数のeSSDシャーシ(260a~260m)(集合的にeSSDシャ

50

ーシ260と称する)を含む。ホスト210は、アプリケーション211、オペレーティングシステム(OS)とファイルシステム(FS)212、及びイーサネット250を通じてラック(rack)270内のeSSDシャーシ260の各々と連結するためのNVMe-oFドライバー213を含む。ラック270はラック270内の多数のeSSDシャーシ260の間に連結性を提供するTOR(top-of-rack)スイッチ271を含む。同様に、NVMe-oFシステム200は異なる位置に配置される多数のラック270を含む。ラック270はTORスイッチ271を通じて互いに連結される。多数のラック270のTORスイッチ271はイーサネット250を通じて専用線又は外部スイッチを通じて直接互いに連結される。

#### 【0037】

図3は、一実施形態によるNVMe-oFシステムのデータ複製のプロセスを示す図である。NVMe-oFシステム300はホスト310及び多数のeSSD(370a~370n)を含むシャーシ360を含む。例えば、シャーシ360は24又は48個のeSSDを含む。ホスト310は、アプリケーション311、オペレーティングシステム(OS)とファイルシステム(FS)312、及びNVMe-oFドライバー313を含む。ホスト310のイニシエータ(initiator、例えばアプリケーション311)はNVMe-oFドライバー313を用いてイーサネット350を通じてアクティブeSSD370aとのNVMe-oF連結を設定する。シャーシ360は、イーサネットスイッチ361、BMC362、及びPCIeスイッチ363を含む。イーサネットスイッチ361はミッドプレーン365を通じてeSSD370にイーサネット連結を提供し、PCIeスイッチ363はミッドプレーン365を通じてeSSD370に管理インターフェイス364を提供する。BMC362は多数のeSSD370の中の1つをアクティブeSSD370aとしてプログラムし、他のeSSDをパッシブeSSDとしてプログラムする。例示の便宜上、本例示は1つのアクティブeSSD370a及び1つのパッシブeSSD370bを図示する。しかし、eSSDの複製グループは、1つのアクティブeSSD及び1つ以上のパッシブeSSDを含み、本開示の範囲を逸脱しない限り、シャーシ内に1つ以上の複製グループが存在する。パッシブeSSD370bはアクティブeSSD370aに関連するように構成された多数のパッシブeSSDの中の1つを示す。パッシブeSSD370bによって遂行される活動はシャーシ360内の他のパッシブeSSDによって同時に遂行される。例えば、NVMe-oFシステム300がこのように構成された場合、同一のデータが2つ以上のパッシブeSSDに複製される。一部の実施形態で、シャーシ360は1つ以上の複製グループを含み、各々の複製グループは1つのアクティブeSSD及び複製グループ内のアクティブeSSDと共にデータ複製を提供するようにプログラムされた1つ以上のパッシブeSSDを含む。

#### 【0038】

アクティブeSSD370aはホスト310との全てのNVMe-oFプロトコル処理及び終了処理を行う。ホスト310とアクティブeSSD370aとの間のNVMe-oFシステム300の伝送プロトコルは、提出キュー(submission queue:SQ)及び完了キュー(completion queue:CQ)を用いて具現される。ホスト310で実行中のアプリケーション311によって呼出されるNVMe-oFドライバー313にアクティブeSSD370aが見えるように、アクティブeSSD370aの位置アドレスはホスト310に通知される。アクティブeSSD370aがホスト書込みデータを受信すると、アクティブeSSD370aは、ホスト書込みデータ及び必要な命令に関連する1つ以上のLBAを、専用通信チャンネルを通じて関連するパッシブeSSD370bに送信する。一実施形態で、アクティブeSSD370a及びパッシブeSSD370bは、イーサネットスイッチ361を通じて、又はシャーシ360内のPCIeスイッチ363を通じて専用の低い帯域幅通信チャンネルを形成する。

#### 【0039】

パッシブeSSD370bはアクティブeSSDに書き込まれたホストデータのコピーを格納する。BMC362は、パッシブeSSD370bがアクティブeSSD370a

10

20

30

40

50

に向かう全てのイングレス（進入）パケット（`ingress packets`）のコピーを受信するために、アクティブ `eSSD 370a` に指定されたアクティブイングレストラフィック（`ingress traffic`）をパッシブ `eSSD 370b` に複製するようにイーサネットスイッチ 361 をプログラムする。パッシブ `eSSD 370b` は、受信されたパケットをフィルタリングし、受信されたパケットの中の `RDMA` 読出し応答（`RDMA READ RESPONSE`）パケットのみを維持する。パッシブ `eSSD 370b` は、受信された `RDMA` 読出し応答パケットを分析し、`RDMA` 読出し応答パケットに含まれるホスト書込みデータチャンク（`chunks`）を取出す。また、パッシブ `eSSD 370b` はアクティブ `eSSD` からデータ複製命令及びデータチャンクに対応する `LBA` を受信する。パッシブ `eSSD 370b` は `LBA` 情報及び `RDMA` 読出し応答パケットから受信されたデータを用いてホストデータを自分の格納媒体に維持する。一部の実施形態で、イーサネットスイッチ 361 はイングレスパケットをフィルタリングして `RDMA` 読出し応答パケットのみをパッシブ `eSSD` に提供する。ホストデータが維持された後、パッシブ `eSSD 370b` は設定された専用通信チャンネルを通じてアクティブ `eSSD 370a` に確認（`confirmation`）を送信する。

10

#### 【0040】

`BMC 362` はストレージ管理者の指導の下にデータ複製機能を設定する役割を担う。一実施形態で、`BMC 362` は `eSSD 370` を構成するために `SMBus` 又は `PCIe` バスを通じて `NVMe-MI` プロトコルを使用する。まず、`BMC 362` は、`eSSD` を識別し、`eSSD` をアクティブモードにプログラムする。また、`BMC 362` は、1つ以上の `eSSD` を選択し、`eSSD` をパッシブ `eSSD` としてプログラムする。`BMC 362` によって構成されたパッシブ `eSSD` の数はストレージ管理者によって指定されたデータコピーの数に依存する。`BMC 362` はこのような `eSSD` をパッシブモードにプログラムする。

20

#### 【0041】

アクティブ及びパッシブ `eSSD 370` が適切に構成されると、`BMC 362` はイーサネットスイッチ 361 を構成する。イーサネットスイッチ 361 はアクティブ `eSSD 370a` に対応するダウンリンクポートに向かうイングレスパケットをパッシブ `eSSD 370b` に対応するダウンリンクポートに複製して伝送する。`eSSD 370` 及びイーサネットスイッチ 361 が適切に構成された後に、格納管理者及び / 又は `BMC 362` は `eSSD 370` がデータ格納及び複製のために準備されたことをホスト 310 のアプリケーション 311 及び `OS / ファイルシステム 312` に通知する。

30

#### 【0042】

図 4 は、一実施形態によるシャーシ内の `eSSD` 及びイーサネットスイッチを初期化及びプログラミングしてデータ複製を具現するためのフローチャートである。初期化プロセスはストレージ管理者によって発行された初期化命令によって開始（`trigger`）される（400 段階）。初期化命令にตอบสนองして、シャーシの `BMC` は `eSSD` を選択し（401 段階）、アクティブ `eSSD` として構成されるように選択された `eSSD` をプログラムする（402 段階）。一実施形態によると、初期化命令はアクティブ `eSSD` に対する固有識別子（例えば、`IP` アドレス、`MAC` アドレス）を含み、`BMC` は `SMBus` 又は `PCIe` バスを通じて `NVMe-MI` プロトコルを使用してアクティブ `eSSD` に構成命令を送信してアクティブモードに構成する。

40

#### 【0043】

アクティブ `eSSD` がプログラムされた後、シャーシ内の `BMC` は 1 つ以上の `eSSD` を選択し（403 段階）、選択された `eSSD` をパッシブ `eSSD` としてプログラムする（404 段階）。一実施形態によると、シャーシの `BMC` はストレージ管理者から受信された初期化命令で指示されたようにパッシブ `eSSD` を選択してプログラミングする。例えば、`BMC` に対する初期化命令はパッシブ `eSSD` の識別子（例えば、`IP` アドレス、`MAC` アドレス）を含み、`BMC` は `NVMe-MI` プロトコルを使用して `SMBus` 又は `PCIe` バスを通じて各々のパッシブ `eSSD` に構成命令を送信してパッシブ `eSSD` を

50

パッシブモードに構成する。

【 0 0 4 4 】

一実施形態によると、アクティブ e S S D 及びパッシブ e S S D の選択及びプログラミングは同時に又は反対の順序で遂行される。即ち、B M C は、先ずパッシブ e S S D を選択してプログラムし、アクティブ e S S D を選択及びプログラムするか、又はアクティブ e S S D 及びパッシブ e S S D を同時に構成する。

【 0 0 4 5 】

e S S D がプログラムされた後に、B M C はシャーシのイーサネットスイッチをプログラムする ( 4 0 5 段階 )。例えば、イーサネットスイッチはアクティブ e S S D とパッシブ e S S D との間に関連性を提供するマッピングテーブルを生成する。プログラムされた後、イーサネットスイッチはアクティブ e S S D に向かうインGRESパケットを関連するパッシブ e S S D に複製する。一実施形態で、イーサネットスイッチは、インGRESパケットをフィルタリングし、R D M A 読出し応答パケットのみをパッシブ e S S D に送信する。

【 0 0 4 6 】

アクティブ e S S D は、ホストによって発見可能であり、ホストで実行中の N V M e - o F ドライバーと通信する。ホストは N V M e - o F プロトコルを使用してアクティブ e S S D にホスト命令を送信する。アクティブ e S S D は、全てのホスト命令を実行し、必要な全てのデータ伝送を遂行し、そして命令完了キューエントリをホストに送信する。即ち、アクティブ e S S D は、N V M e - o F プロトコル処理を遂行し、ノーマル e S S D によって遂行されるようにホスト命令を実行する。また、アクティブ e S S D は 1 つ以上のパッシブ e S S D と併せてデータ複製を支援する。データ複製は主にホスト書込み命令に関連する。ホスト読出し命令の場合、アクティブ e S S D は特定機能を支援する必要がない。

【 0 0 4 7 】

ホスト書込み命令に対して、アクティブ e S S D はホストからデータをフェッチ ( f e t c h ) する役割をする。一実施形態によると、アクティブ e S S D は R D M A 読出し ( R D M A R E A D ) 要請をホストに発行してホストデータをフェッチする。呼応して、ホストは 1 つ以上の R D M A 読出し応答 ( R D M A R E A D R E S P O N S E ) パケットを用いてデータチャンクを送信する。アクティブ e S S D は、R D M A 読出し応答パケット内で受信されたホストデータチャンクを内部フラッシュメモリに格納するか、或いはホストデータを電力損失防止バッファ ( p o w e r - l o s s - p r o t e c t e d b u f f e r ) に維持する。

【 0 0 4 8 】

R D M A 読出し応答パケット内で受信されたホストデータの全てのデータチャンクに対して、アクティブ e S S D は関連するネームスペース ( N A M E S P A C E ) 及び L B A 情報を書込み命令と共に関連する 1 つ以上のパッシブ e S S D の各々に送信する。このような通信はアクティブ及びパッシブ e S S D の間に設定された専用通信チャンネルを通じて遂行される。専用通信チャンネルはシャーシのイーサネットスイッチ又は P C I e スイッチを通じて形成される。アクティブ e S S D が動作しなくなるか、或いはアクティブ e S S D が受信されたデータチャンクの手書きに失敗した場合、アクティブ e S S D はパッシブ e S S D に破棄命令を送信してデータチャンクを破棄する。

【 0 0 4 9 】

ホストデータが局所的に維持され、N A M E S P A C E . L B A 情報がパッシブ e S S D に送信されると、アクティブ e S S D は一貫性モードに基づいてホストに命令完了 ( c o m m a n d c o m p l e t i o n ) を送信する。一実施形態によると、アクティブ e S S D は、データ持続性に関するパッシブ e S S D からの応答を待機せずに、直ちにホストに命令完了を送信する。本明細書で、略完全な一貫性モード ( a l m o s t - p e r f e c t c o n s i s t e n c y m o d e ) と称される即時命令完了応答 ( i m m e d i a t e c o m m a n d - c o m p l e t i o n - r e s p o n s e ) は、バックアッ

10

20

30

40

50

ブコピーが忠実に持続されないリスクで、わずかに優れたレイテンシ性能を提供する。このようなエラーケースが発生する確率は小さく、少量のデータ損失を許容するアプリケーションが存在する可能性がある。このようなエラーケースに対する詳細な分析は以下で詳細に説明する。

#### 【 0 0 5 0 】

他の実施形態によると、アクティブ e S S D はパッシブ e S S D からデータ持続性の確認を待つ。1つ以上のコピーが複製された場合、アクティブ e S S D は命令完了エントリをホストに送信する前に関連する1つから全てのパッシブ e S S D までの確認を待つ。データ複製はメインデータ経路（例えば、アクティブ e S S D への）内のホストデータのストレージと並列に発生するため、パッシブ e S S D からの確認を待つことは、ホストで実行中のアプリケーションが許容するか否かに関係なく適度なレイテンシを追加することになる。即時命令完了応答を使用する略完全な一貫性モードとは異なり、完全な一貫性モード（perfect consistency mode）と称される関連する全てのパッシブ e S S D からの応答を待つこの動作モードはデータ損失無しに忠実なデータ一貫性を保証する。

10

#### 【 0 0 5 1 】

一実施形態によると、N V M e - o F システムの一貫性モードは多様な使用者 Q o S （quality of service）ポリシーを使用して B M C によって設定される。一貫性とレイテンシとの間のトレードオフはアプリケーションに依存的であり、B M C はストレージ管理者の指示に応じて多様なポリシーを具現する。アクティブ及びパッシブ e S S D がデータ一貫性の観点で同期化されない場合、復旧されるか、或いは同一な同期化レベルになるようにする必要がある。このような復旧機能は B M C によって管理される。復旧作業を容易にするために、アクティブ e S S D は最後の“ n ”回の成功的な書込み動作ログを維持する。例えば、ログは N A M E S P A C E . L B A と特定書込み動作を示すマーカー（marker）を含む。

20

#### 【 0 0 5 2 】

また、アクティブ e S S D はデータ複製プロセスを容易にするために関連するパッシブ e S S D に一部の管理命令を送信する。このような管理命令の例としては、生成 / 削除ネームスペース（N A M E S P A C E ）及びトリム（Trim）命令等があるが、これらに限定されない。管理命令を実行した後、パッシブ e S S D は確認を再びアクティブ e S S D に送信する。

30

#### 【 0 0 5 3 】

図 5 は、一実施形態によるアクティブ e S S D によってホスト命令を処理するためのフローチャートである。パッシブ e S S D で構成された全ての e S S D はアクティブ e S S D との関連性を有する。アクティブ e S S D とパッシブ e S S D との間の関連性は図 4 に示した初期化の一部として B M C によって設定される。パッシブ e S S D の主要機能の中の1つはホストデータを複製することである。一実施形態で、パッシブ e S S D は2つの並列 F I F O キューを維持することによってホストデータのデータ複製を遂行し、2つの並列 F I F O キューの1つはホストデータを格納するためであり（以下、“受信メッセージ F I F O ”と称する）、他の1つは書込みアドレス及び命令を格納するためである（以下、“L B A 及び命令 F I F O ”と称する）。パッシブ e S S D は、ホストデータがアクティブ e S S D に書き込まれる時に、2つの F I F O を用いてホストデータを維持する。

40

#### 【 0 0 5 4 】

アクティブ e S S D は、ホストからホスト命令を受信し、ホスト命令が受信された時に受信されたホスト命令を提出キュー（submission queue : S Q ）に配置する（500段階）。アクティブ e S S D は提出キュー内のエントリを仲裁し（501段階）、提出キューから実行のためのホスト命令を選択する（502段階）。ホスト命令は、ホスト書込み命令（503段階）、管理命令（510段階）、又はノーマル命令（516段階）の中のいずれか1つである。

#### 【 0 0 5 5 】

50

選択されたホスト命令がホスト書込み命令である場合（５０３段階）、ホスト書込み命令はアクティブeSSDの内部ストレージ装置に書き込むためのホストデータを含む。アクティブeSSDは、例えばストレージ装置のストレージユニット（例えば、ページのサイズ又はフラッシュメモリのブロックのサイズ）に応じてホストデータをデータチャンクに分割（segment）する（５０４段階）。アクティブeSSDはホスト書込み命令に関連する全てのデータチャンクがフェッチされたか否かを検査する（５０５段階）。分割されたデータチャンクの中のフェッチする少なくとも１つのデータチャンクが残っている場合、アクティブeSSDはホストにRDMA読出し要請を発行して残りのデータチャンクをフェッチする（５０６段階）。各データチャンクに対してアクティブeSSDはNAME SPACE・LBA及び命令に関連するパッシブeSSDの各々に送信する（５０７段階）。RDMA読出し要請をホストに発行してNAME SPACE・LBA及び命令をパッシブeSSDに発行することは同時に遂行される。データチャンクに対するRDMA読出し要請パケットが受信された後に、アクティブeSSDは内部ストレージ装置（例えば、フラッシュメモリ）にデータチャンクを維持させる（５０９段階）。フェッチされたデータチャンクはアクティブeSSDのデータバッファに一時的に格納され、単一ホスト書込み命令に関連する全てのデータチャンクは同時に維持される。

#### 【００５６】

全てのデータチャンクがフェッチされて維持された場合、アクティブeSSDは一貫性モードをチェックする（５１３段階）。一実施形態によると、一貫性モードは“完全（perfect）”又は“略完全（almost perfect）”の中の１つに設定される。アクティブeSSDの一貫性モードが“完全”に設定された場合、アクティブeSSDはアクティブeSSDに関連する全てのパッシブeSSDからの確認を待機し（５１４段階）、完了キューエントリをホストに送信する（５１５段階）。アクティブeSSDの一貫性モードが“略完全”に設定された場合、アクティブeSSDは関連するパッシブeSSDからの確認を待機せずに完了キューエントリをホストに送信する（５１５段階）。

#### 【００５７】

受信されたホスト命令がネームスペース生成命令を含む管理命令又は関連する１つ以上のパッシブeSSDに関連するトリム命令である場合（５１０段階）、アクティブeSSDは、受信された命令を実行し（５１１段階）、命令を該当するパッシブeSSDに送信する（５１２段階）。一貫性モードに応じて（５１３段階）、アクティブeSSDは、関連するパッシブeSSDからの確認を待機するか、又は待機せずに完了キューエントリをホストに送信する（５１５段階）。受信されたホスト命令がホスト書込み命令でも管理命令でもない場合、アクティブeSSDは受信された命令がデータ複製に関連しないノーマル命令であるとして取扱い、それに応じてノーマル命令を実行する（５１６段階）。

#### 【００５８】

図６は、一実施形態によるデータ複製プロセスの一例を示す図である。シャース６６０内のイーサネットスイッチ６６１は、イーサネット連結を通じてホストからホスト命令を受信し、受信されたホスト命令をアクティブeSSD ６７０aに伝送する。受信されたホスト命令はアクティブeSSD ６７０aの提出キュー（SQs）に配置される。提出キューに配置されたホスト命令は、ホスト書込み命令、管理者命令（例えば、ネームスペース生成）、及びノーマル命令（例えば、ログページ取得）を含む。

#### 【００５９】

提出キュー内の受信されたホスト命令の中のホスト書込み命令を処理する時、アクティブeSSD ６７０aはホストにRDMA読出し要請を発行してホスト書込み命令に関連するデータをフェッチする。アクティブeSSD ６７０aはホスト書込み命令に関連するホスト書込みデータのデータチャンクの各々に対する一連のRDMA読出し要請を発行する。各データチャンクに対して、アクティブeSSD ６７０aはデータをパッシブeSSD ６７０bに格納するためにNAME SPACE・LBA及び命令を送信する。一方、イーサネットスイッチ６６１はアクティブeSSD ６７０aに向かうイングレストラフィッ

10

20

30

40

50

クをパッシブ e S S D 6 7 0 b に複製するようにプログラムされる。

#### 【 0 0 6 0 】

図 7 は、一実施形態によるパッシブ e S S D によってホスト命令を処理するためのフローチャートである。パッシブ e S S D はアクティブ e S S D に向かう全てのイングレストラフィックを受信する ( 7 0 1 段階 )。イングレストラフィックは主にホスト命令を伝達する R D M A 送信 ( R D M A S E N D ) パケット及びホスト書込み命令のためのホストデータを伝達する R D M A 読出し応答 ( R D M A R E A D R E S P O N S E ) パケットを含む。R D M A パケット以外にもイングレストラフィックは他のネットワークトラフィックを更に含む。パッシブ e S S D は受信されたイングレスパケットを分析する ( 7 0 2 段階 )。パッシブ e S S D は R D M A 読出し応答パケットのみを維持するためにパケットフィルタロジックを使用して受信された全てのパケットをフィルタリングする ( 7 0 3 段階 )。パッシブ e S S D は R D M A 読出し応答パケット以外のイングレスパケットを破棄する ( 7 0 6 段階 )。パッシブ e S S D は、R D M A 読出し応答パケットを R D M A データメッセージに組合せ ( a s s e m b l e ) ( 7 0 4 段階 )、R D M A 読出し応答パケットに含まれるホスト書込みデータを取出す ( 7 0 5 段階 )。

#### 【 0 0 6 1 】

再び図 6 を参照すると、パッシブ e S S D は、R D M A 読出し応答パケットを R D M A データメッセージに組合せ、これらを受信されたメッセージ F I F O 6 7 2 に配置する。受信されたメッセージ F I F O 6 7 2 はホストデータチャンクをパッシブ e S S D 6 7 0 b のストレージ媒体に格納する前にこれらを一時的に格納する。また、パッシブ e S S D 6 7 0 b はアクティブ e S S D 6 7 0 a から受信された N A M E S P A C E . L B A 及び命令を別のバッファ (ここでは L B A 及び命令 F I F O 6 7 1 と称する) に格納する。L B A 及び命令 F I F O 6 7 1 に格納された N A M E S P A C E . L B A 及び命令は、識別子 ( i d e n t i f i e r ) 又はマーカー ( m a r k e r ) を使用して受信されたメッセージ F I F O 6 7 2 に配置された対応するホストデータを示す。パッシブ e S S D 6 7 0 b は命令 (例えば、書込み) 及び対応するデータチャンクを使用するアドレスに基づいてホストデータを自体のストレージ媒体に維持する。ホストデータが維持された後、パッシブ e S S D 6 7 0 b は専用通信チャンネルを通じてアクティブ e S S D 6 7 0 a に確認を送信する。

#### 【 0 0 6 2 】

図 8 は、一実施形態によるアクティブ e S S D からパッシブ e S S D に L B A 及び命令を伝達するためのフローチャートである。受信された各々のデータチャンクに対して、アクティブ e S S D は関連する L B A 及び命令を関連するパッシブ e S S D の各々に送信し、パッシブ e S S D は L B A 及び命令を L B A 及び命令 F I F O に配置する ( 8 0 1 段階 )。パッシブ e S S D は受信されたイングレスパケットを分析して R D M A 読出し応答パケットに含まれるホスト書込みデータのみを受信されたメッセージ F I F O に配置する ( 8 0 2 段階 )。アクティブ e S S D から受信された L B A 及び命令を用いて、パッシブ e S S D はストレージ媒体にホストデータを維持させる ( 8 0 3 段階 )。パッシブ e S S D はアクティブ e S S D によって命令された通りにデータを複製するプロセスを繰り返す。

#### 【 0 0 6 3 】

アクティブ e S S D による命令は、一般的にパッシブ e S S D が対応するデータチャンクを提供された N A M E S P A C E . L B A で格納するように指示する “ 書き込み ( W r i t e ) ” 命令である。アクティブ e S S D にデータチャンクを書き込めない場合は稀である。この場合、アクティブ e S S D はパッシブ e S S D に “ 破棄 ( D i s c a r d ) ” 命令を送信する。パッシブ e S S D が “ 破棄 ” 命令に遭遇すると、パッシブ e S S D は単に該当データチャンクを投げる ( t h r o w )。また、アクティブ e S S D がパッシブ e S S D に送信する他の命令がある。このような命令は、ホスト又はアクティブ e S S D 及びパッシブ e S S D を含むシャーシ内のイーサネットスイッチ、B M C、又は P C I e スイッチからアクティブ e S S D が受信するネームスペース管理命令 (例えば、ネームスペース生成) 及びトリム命令のような管理命令である。管理命令に対するプロトコルは



、本発明の範囲を逸脱しない限り、他のホスト命令を含むように拡張される。パッシブ eSSD が LBA 及び命令 FIFO で非データ (non-data) 命令に遭遇すると、パッシブ eSSD は受信されたメッセージ FIFO に影響を与えずに非データ命令を実行する。

#### 【0064】

命令がパッシブ eSSD によって実行される時、パッシブ eSSD は選択的にアクティブ eSSD に完了確認 (completion confirmation) を送信する。このモードで、アクティブ eSSD は元のホスト命令に対する完了キューエントリをホストに送信する前に関連する 1 つ以上のパッシブ eSSD からの確認を待つ。この動作モードはアクティブ及びパッシブ eSSD が常に互いに一致していることを保証する。

10

#### 【0065】

データチャンクの書込みに失敗するか又は命令の実行が失敗する場合は稀であるため、アクティブ eSSD は関連するパッシブ eSSD からの確認を待機しないこともある。関連するパッシブ eSSD の中のいずれかがデータ書込みに失敗するか又は命令を実行できない場合、関連するパッシブ eSSD の中の 1 つで誤謬が発生する。一実施形態によると、パッシブ eSSD はアクティブ eSSD によって命令された通りにホストデータを複製することに対する誤謬をシャーシ内の BMC に報告する。BMC は例外処理及び一貫性状態復旧のために適切な措置を取る。

#### 【0066】

図 9 は、一実施形態によるパッシブ eSSD によってデータ複製を処理するためのフローチャートである。BMC は、アクティブ eSSD を構成し、アクティブ eSSD に 1 つ以上のパッシブ eSSD を関連させる。

20

#### 【0067】

パッシブ eSSD は LBA 及び命令 FIFO 内のエントリを検査し (900 段階)、LBA 及び命令 FIFO から次の命令を引き出す (901 段階)。命令が書込み命令である場合 (902 段階)、パッシブ eSSD は NAME SPACE . LBA に基づいて書込み命令に関連するデータチャンクを維持する (903 段階)。命令が書込み命令ではない場合 (902 段階)、パッシブ eSSD は命令が破棄命令であるか否かを判別する (906 段階)。パッシブ eSSD は破棄命令に応答してデータチャンクを破棄する (907 段階)。命令が破棄命令ではない場合 (906 段階)、パッシブ eSSD は命令が管理命令 (909 段階) であるか否かを更に判別し、管理命令を実行する (910 段階)。命令が管理命令ではない場合 (909 段階)、パッシブ eSSD は BMC (911 段階) に警告し、エラー状態を報告する (912 段階)。

30

#### 【0068】

パッシブ eSSD の一貫性モードが完全 (perfect) モードである場合 (904 段階)、パッシブ eSSD は確認をアクティブ eSSD に送信する (905 段階)。パッシブ eSSD の一貫性モードが略完全な (almost-perfect) モード (904 段階) である場合、パッシブ eSSD はエラーがあるか否かを確認する (908 段階)。エラーが発生した場合、パッシブ eSSD はエラーを BMC に報告する (911、912 段階)。エラーが無い場合、パッシブ eSSD は LBA 及び命令 FIFO を検査し (900 段階)、LBA 及び命令 FIFO で次の命令を実行する (901 段階)。

40

#### 【0069】

一実施形態によると、BMC はアクティブ eSSD 及びパッシブ eSSD の健康状態 (health status) を周期的にモニターリングする。例えば、BMC は、健康状態モニターリングのために NVMe - MI プロトコル、及び特に NVMe - MI "Health Status Poll" 命令を使用する。BMC は管理目的として PCIe インターフェイス又は SMBus インターフェイスを使用する。アクティブ eSSD 又はパッシブ eSSD の中の 1 つの eSSD は様々な理由でシャーシ内から引き出される。BMC は周期的にシャーシ内の各 eSSD から "現在 (Present)" 信号を受信する。eSSD が引き出された場合、BMC はこの場合に対する報告を受信する。

50

## 【0070】

BMCが現在のアクティブeSSDが失敗したか又はシャーシにそれ以上存在しないと判別した場合、BMCは誤謬解決(failover)動作を開始する。先ず、BMCはアクティブモードでシャーシ内のパッシブeSSD又はノーマルeSSDの中の1つをプログラムする。必要な場合、BMCはシャーシ内にある他のeSSDを新たなパッシブeSSDとして選択する。BMCはイーサネットスイッチを再びプログラムして新たなアクティブeSSDのイングレストラフィックを複製する。

## 【0071】

図10は、一実施形態による誤謬解決動作のためのフローチャートである。BMCはアクティブeSSDが失敗したか又はシャーシから引き出されたかを検出する(1000段階)。BMCはシャーシ内の新たなアクティブeSSDを選択する(1001段階)。一実施形態によると、シャーシ内のパッシブeSSDの中の1つ又はノーマルeSSDが新たなアクティブeSSDに変換される。その後、BMCは新たなパッシブeSSDが必要であるか否かを判別する(1002段階)。新たなパッシブeSSDが要求された場合(1002段階)、BMCは、シャーシ内でeSSDを選択し(1005段階)、eSSDをパッシブeSSDとしてプログラムし、eSSDを新たなアクティブeSSDに関連させる(1006段階)。パッシブeSSDが新たなアクティブeSSDに変換されると、ホストデータのコピー(copies)の数が減少する。従って、新たなパッシブeSSDはシャーシ内に存在するノーマルeSSDの中の1つから構成される。或いは、他のアクティブeSSDに関連するパッシブeSSDがアクティブeSSDに変換される。BMCは、シャーシ内のイーサネットスイッチをプログラムし(1003段階)、新たなアクティブeSSDをプログラムする(1004段階)。

## 【0072】

データ複製セットアップ(setup)で、複製されたデータコピーを一貫性あるように維持することが重要である。即ち、eSSDのセット(アクティブeSSD及び関連するパッシブeSSD)は互いに一貫している必要がある。一貫しない場合、どのデータコピーが正しいか否かを判断することができず、このような状況が発生すると、データが損失する可能性がある。一部の使用ケースは特定データ損失イベントを許容するが、大部分の使用ケースはデータ損失を予想しない。

## 【0073】

一実施形態によると、本発明のデータ複製システム及び方法は、アクティブeSSDと関連するパッシブeSSDとの間に完全な一貫性を保証し、データ損失はない。場合によって、稀なことであるが、一部のデータ損失を許容する場合、使用者アプリケーションはより良いレイテンシを達成することができるが、そのようなデータ損失の可能性は非常に低い。

## 【0074】

使用ケースに基づいて、BMCはデータ複製の一貫性モードを複製設定の一部として完全な一貫性モード又は略完全な一貫性モードに設定する。完全な一貫性モードでアクティブeSSDは完了キューエントリをホストに送信する前に関連する全てのパッシブeSSDからの確認を待つ。完全な一貫性モードで使用者データのコピーを維持した後、アクティブeSSDは全てのパッシブeSSDがホストデータを成功的に維持することを保証するためにホスト命令完了キューエントリを直ちに送信しない。パッシブeSSDからの確認を待つことはホスト書込みアクセスに若干のレイテンシを追加する。アクティブ及びパッシブeSSDは同時に使用者データを受信するため、パッシブeSSDからの確認による追加レイテンシは重要でない。遅延は主にパッシブeSSD確認がアクティブeSSDに到達するのに掛かる時間に対応する。アクティブeSSDとパッシブeSSDとの間の専用通信チャンネルはローカルであり、専用であるため、このようなeSSD間の通信時間は最小限になる。

## 【0075】

アクティブeSSDが引き出されるか又は急に非動作状態になる場合、提出キューに一

10

20

30

40

50

部のホスト書込み命令があることがある。この場合、アクティブ e S S D は部分書込み命令に対する命令完了キューエントリを送信しない。従って、ホストの観点でデータ無欠性を維持する。ホストは失敗した書込み命令を他のストレージユニットに再発行するか又は誤謬解決スイッチングを待つ。

【 0 0 7 6 】

略完全な一貫性モードで、アクティブ e S S D は関連するパッシブ e S S D からの確認を待たない。アクティブ e S S D が自体の書込みデータコピーを維持すると直ちにアクティブ e S S D は命令完了をホストに送信する。その結果、完全な一貫性モードと比較すると、ホストに対するレイテンシが向上する。このモードにある間にアクティブ e S S D が引き出されるか又は非動作状態になる場合、一部のホスト書込み命令が実行中であることがある。既に送信された命令完了がない全てのホスト書込み命令に対してデータ無欠性が維持される。ホストは代替ストレージユニットに書込み命令を再び発行する。1つの特殊な場合は、アクティブ e S S D がホストに書込み命令完了を成功的に送信したが、何らかの理由でパッシブ e S S D がそのデータを維持できないことである。この特殊な場合、ホストは、そのデータが成功的に書き込まれたが、データのアクティブコピーをそれ以上使用することができず、データのパッシブコピーの維持が失敗したと判別する。書込み命令完了を送信した直後にアクティブ e S S D が引き出されるか又は非動作状態になってパッシブ e S S D が該当時点で関連する書込みデータを維持させることができない確率は非常に少ない。略完全な一貫性モードは、若干より良いレイテンシを提供するが、小さいデータ損失の恐れがある。このような状況に耐えられるアプリケーションの場合、略完全な一貫性モードが実行可能なオプションになる。

【 0 0 7 7 】

本発明の一実施形態によるデータ複製方法は、複数のソリッドステートドライブ ( e S S D )、前記複数の e S S D の各々にダウンリンクを提供するためのファブリックスイッチ、並びに前記ファブリックスイッチ及び前記複数の e S S D を制御するためのベースボード管理コントローラ ( B M C ) をシャーシ内に含むデータ複製システムのデータ複製方法であって、前記 B M C を用いて、前記複数の e S S D の中の1つをアクティブ e S S D として構成する段階と、前記複数の e S S D の中の1つ以上をパッシブ e S S D として構成する段階と、前記ファブリックスイッチをプログラムして、前記アクティブ e S S D に向かうパケットを前記アクティブ e S S D 及び前記1つ以上のパッシブ e S S D の両方に伝送する段階と、を有し、前記アクティブ e S S D で、ホストからホストデータの書込み命令を受信する段階と、前記アクティブ e S S D から前記ホストデータに対応するアドレス及び命令を前記1つ以上のパッシブ e S S D に送信する段階と、前記アクティブ e S S D に前記ホストデータを格納する段階と、前記アクティブ e S S D から受信されたアドレス及び命令並びに前記ファブリックスイッチによって伝送されたパケット内で受信された前記ホストデータを用いて前記1つ以上のパッシブ e S S D の各々に前記ホストデータのコピーを格納する段階と、を含む。

【 0 0 7 8 】

前記アクティブ e S S D は、前記ホストによって発見可能であり、前記1つ以上のパッシブ e S S D は、前記ホストによって発見不可能であり得る。

【 0 0 7 9 】

前記方法は、前記 B M C を用いて、前記アクティブ e S S D と前記1つ以上のパッシブ e S S D との間の専用通信チャンネルを設定し、前記アクティブ e S S D から前記アドレス及び前記命令を前記1つ以上のパッシブ e S S D に送信する段階を更に含む得る。

【 0 0 8 0 】

前記専用通信チャンネルは、前記シャーシ内の前記ファブリックスイッチ又は P C I e スイッチを通じて設定され得る。

【 0 0 8 1 】

前記方法は、前記アクティブ e S S D を用いて、前記ホストデータの書込み命令に応答して R D M A ( r e m o t e d i r e c t m e m o r y a c c e s s ) 読出し要請

を前記ホストに発行する段階と、前記ホストから、前記ホストデータを1つ以上のR D M A 読出し応答パケット内で受信する段階と、を更に含み得る。

【0082】

前記R D M A 読出し応答パケットは、前記ホストデータのデータチャンクを含み得る。

【0083】

前記1つ以上のパッシブe S S Dの各々は、前記ファブリックスイッチから受信されたパケットの中の前記R D M A 読出し応答パケットを除外した残りのパケットを破棄し得る。

【0084】

前記方法は、前記アクティブe S S Dで、一貫性モードに基づいて、前記1つ以上のパッシブe S S Dから応答を受信した後、又は前記1つ以上のパッシブe S S Dから応答を受せずに前記ホストデータを前記アクティブe S S D内に格納した後、完了キューエントリを前記ホストに送信する段階を更に含み得る。

10

【0085】

前記1つ以上のパッシブe S S Dの各々は、前記アドレス及び前記命令を格納するための第1バッファ並びに前記ホストデータのデータチャンクを格納するための第2バッファを含み得る。

【0086】

前記方法は、前記B M Cを用いて、前記アクティブe S S Dに関連するエラーを検出する段階と、前記シャーシ内の複数のe S S Dの中の1つのe S S Dを選択する段階と、前記選択されたe S S Dを新たなアクティブe S S Dとして構成する段階と、新たなパッシブe S S Dが必要であるか否かを判別する段階と、前記新たなパッシブe S S Dを構成する段階と、前記新たなアクティブe S S D、前記少なくとも1つのパッシブe S S D、及び前記新たなパッシブe S S Dに関連させるように前記ファブリックスイッチをプログラムする段階と、を更に含み得る。

20

【0087】

本発明の一実施形態によるデータ複製システムは、複数のソリッドステートドライブ(e S S D)、前記複数のe S S Dの各々にダウンリンクを提供するためのファブリックスイッチ、並びに前記ファブリックスイッチ及び前記複数のe S S Dを制御するためのベースボード管理コントローラ(B M C)をシャーシ内に備え、前記B M Cは、前記複数のe S S Dの中の1つをアクティブe S S Dとして構成し、前記複数のe S S Dの中の1つ以上をパッシブe S S Dとして構成し、前記ファブリックスイッチをプログラムして、前記アクティブe S S Dに向かうパケットを前記アクティブe S S D及び前記1つ以上のパッシブe S S Dの両方に伝送し、前記アクティブe S S Dは、ホストから受信されたホストデータの書き込み命令に応答してホストデータを格納し、前記アクティブe S S Dは、前記ホストデータに対応するアドレス及び命令を前記1つ以上のパッシブe S S Dに送信し、前記1つ以上のパッシブe S S Dの各々は、前記アクティブe S S Dから受信されたアドレス及び命令並びに前記ファブリックスイッチによって伝送されたパケット内で受信された前記ホストデータを用いて前記ホストデータのコピーを格納する。

30

【0088】

前記アクティブe S S Dは、前記ホストによって発見可能であり、前記1つ以上のパッシブe S S Dは、前記ホストによって発見不可能であり得る。

40

【0089】

前記アドレス及び前記命令を送信するために前記アクティブe S S Dと前記1つ以上のパッシブe S S Dのとの間に専用通信チャンネルが設定され得る。

【0090】

前記専用通信チャンネルは、前記シャーシ内の前記ファブリックスイッチ又はP C I e スイッチを通じて設定され得る。

【0091】

前記アクティブe S S Dは、前記ホストデータの書き込み命令に応答してR D M A ( r e

50

memory access) 読出し要請を前記ホストに発行し、前記ホストから、前記ホストデータを1つ以上のRDMA読出し応答パケット内で受信し得る。

【0092】

前記RDMA読出し応答パケットは、前記ホストデータのデータチャンクを含み得る。

【0093】

前記1つ以上のパッシブeSSDの各々は、前記ファブリックスイッチから受信されたパケットの中の前記RDMA読出し応答パケットを除外した残りのパケットを破棄し得る。

【0094】

前記アクティブeSSDは、一貫性モードに基づいて、前記1つ以上のパッシブeSSDから応答を受信した後、又は前記1つ以上のパッシブeSSDから応答を受信せずに前記ホストデータを前記アクティブeSSD内に格納した後、完了キューエントリを前記ホストに送信し得る。

【0095】

前記1つ以上のパッシブeSSDの各々は、前記アドレス及び前記命令を格納するための第1バッファ並びに前記ホストデータのデータチャンクを格納するための第2バッファを含み得る。

【0096】

前記BMCは、前記アクティブeSSDに関連するエラーを検出し、前記シャーシ内の複数のeSSDの中の1つのeSSDを選択し、前記選択されたeSSDを新たなアクティブeSSDとして構成し、新たなパッシブeSSDが必要であるか否かを判別し、前記新たなパッシブeSSDを構成し、前記新たなアクティブeSSD、前記少なくとも1つ以上のパッシブeSSD、及び前記新たなパッシブeSSDに関連させるように前記ファブリックスイッチをプログラムし得る。

【0097】

以上、本発明の実施形態について図面を参照しながら詳細に説明したが、本発明は、上述の実施形態に限定されるものではなく、本発明の技術的範囲から逸脱しない範囲内で多様に変更実施することが可能である。

【符号の説明】

【0098】

100、200、300	NVMe-oFシステム
110、210、310	ホスト
111、211、311	アプリケーション
112、212、312	OS/ファイルシステム
113、213、313	NVMe-oFドライバー
150、250、350	イーサネット
160、360、660	シャーシ
161、361、661	イーサネットスイッチ
162、362	ベースボード管理コントローラ(BMC)
163、363	PCIeスイッチ
164、364	管理インターフェイス
165、365	ミッドプレーン
170、170a~170n、370a~370n	eSSD
260、260a~260m	eSSDシャーシ
270	ラック
271	TORスイッチ
370a、670a	アクティブeSSD
370b、670b	パッシブeSSD
671	LBA及び命令FIFO

10

20

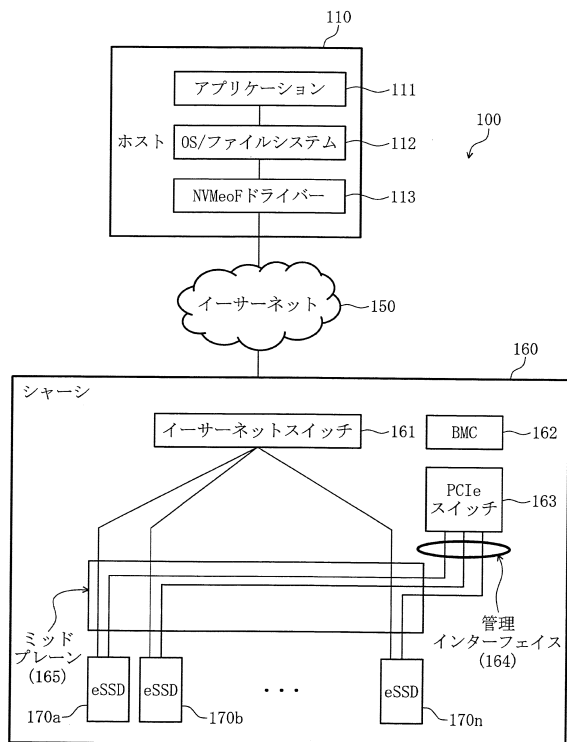
30

40

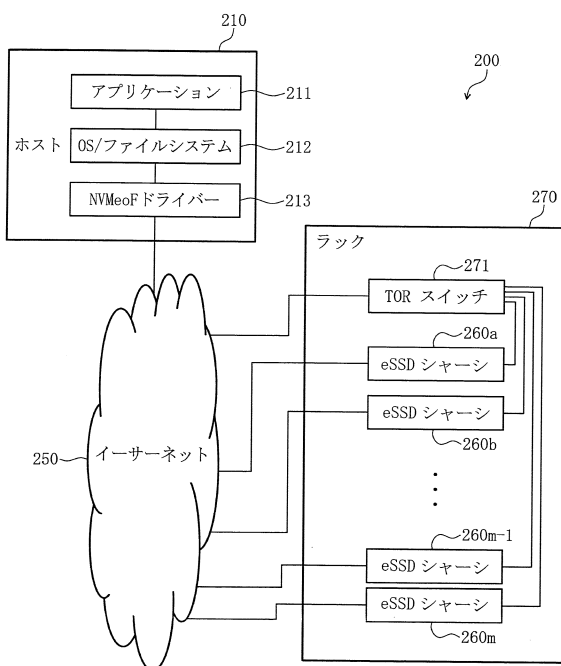
50

## 672 受信されたメッセージFIFO

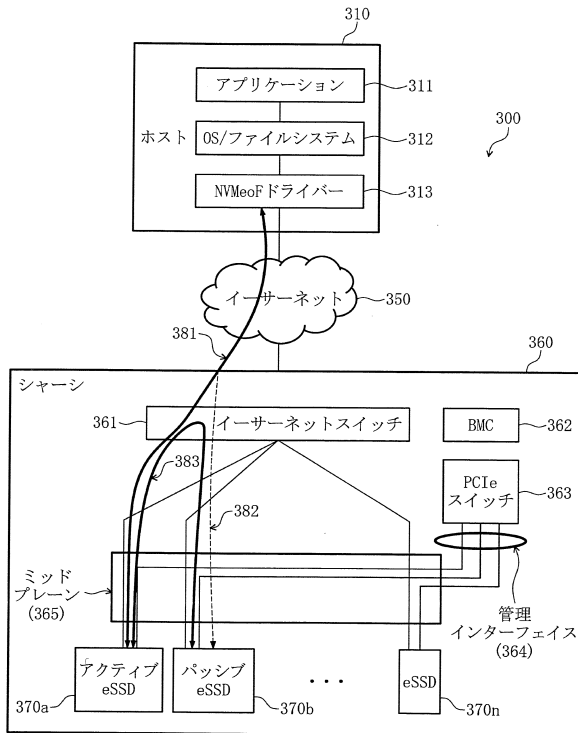
【図1】



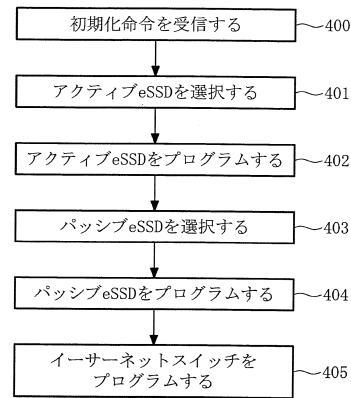
【図2】



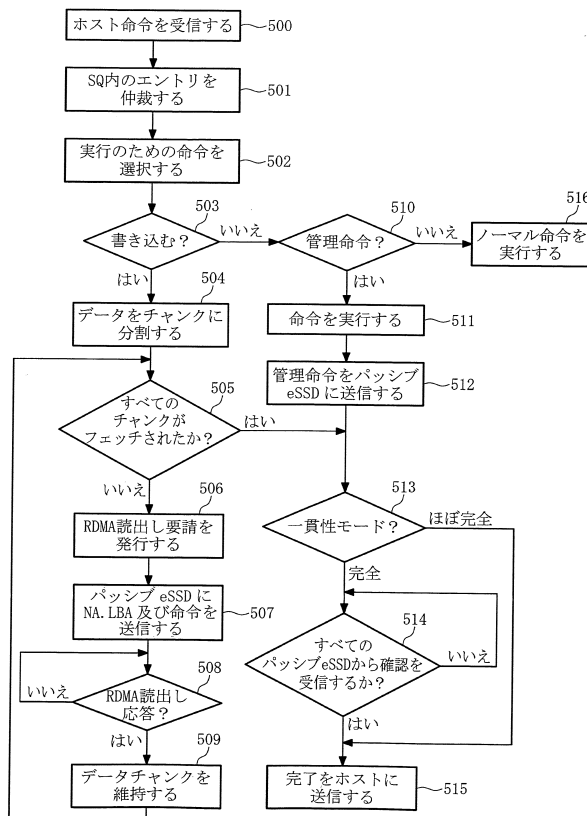
【図 3】



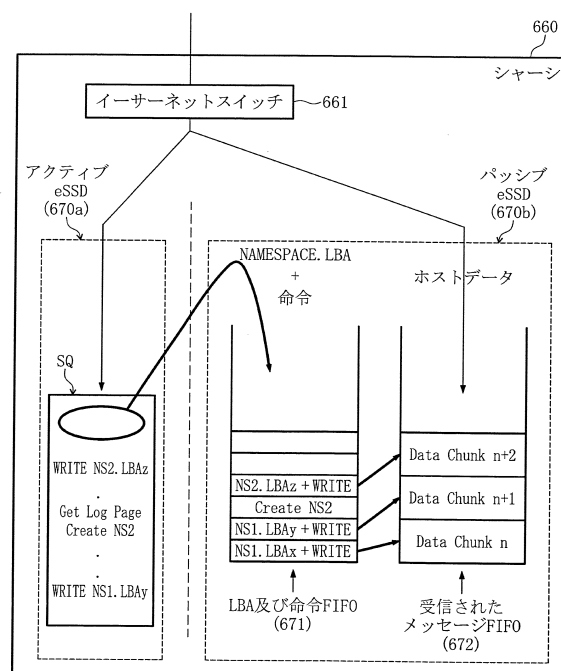
【図 4】



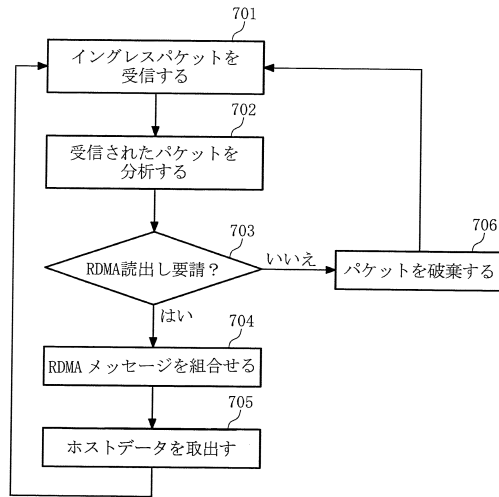
【図 5】



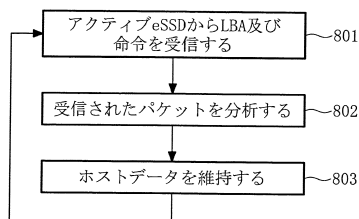
【図 6】



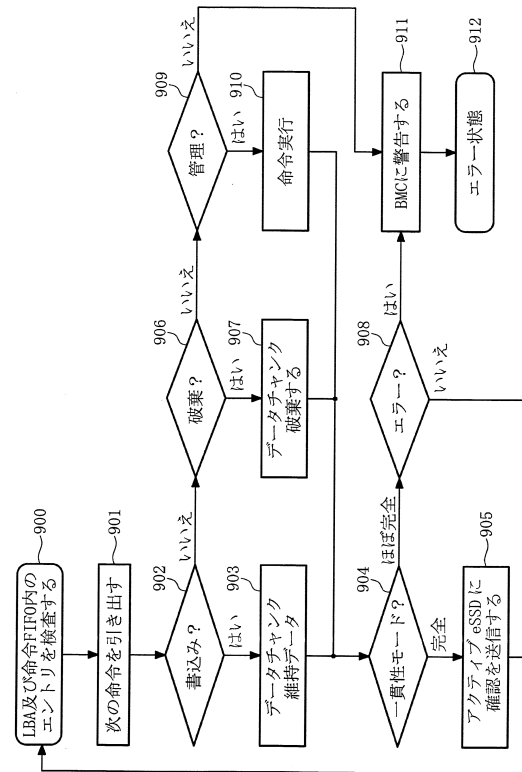
【図 7】



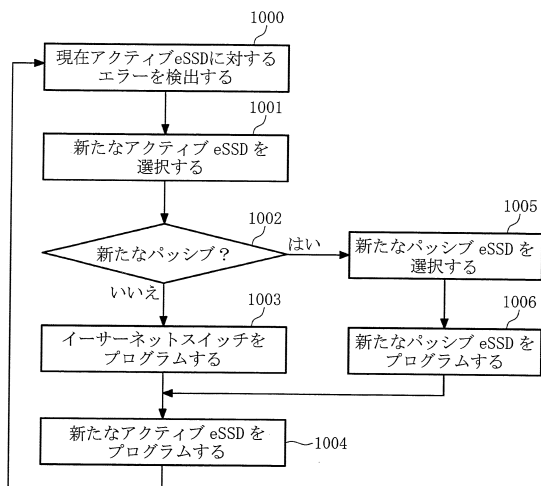
【図 8】



【図 9】



【図 10】





## フロントページの続き

(51)Int.Cl.

F I

G 0 6 F	3/08	H
G 0 6 F	13/10	3 4 0 A
G 0 6 F	13/38	3 5 0
H 0 4 L	12/771	

(72)発明者 ワーレイ, フレッド

アメリカ合衆国, 9 5 1 2 9, カリフォルニア州, サン ノゼ, グリーン ドライブ 1  
4 7 1

(72)発明者 オラリグ, ソンボン ポール

アメリカ合衆国, 9 4 5 6 6, カリフォルニア州, プレザントン, パセオ グラナダ 3  
0 5 0

(72)発明者 ピント, オスカー

アメリカ合衆国, 9 5 1 3 2, カリフォルニア州, サン ノゼ, ベサニー アベニュー  
1 7 7 4

審査官 松平 英

(56)参考文献 特開2001-337788(JP,A)

特開2011-209892(JP,A)

米国特許出願公開第2015/0012607(US,A1)

米国特許出願公開第2015/0261720(US,A1)

(58)調査した分野(Int.Cl., DB名)

G 0 6 F	3 / 0 6
G 0 6 F	3 / 0 8
G 0 6 F	1 3 / 1 0
G 0 6 F	1 3 / 3 8
H 0 4 L	1 2 / 7 7 1