

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7032207号

(P7032207)

(45)発行日 令和4年3月8日(2022.3.8)

(24)登録日 令和4年2月28日(2022.2.28)

(51)国際特許分類

F I

G 0 6 F 3/06 (2006.01)

G 0 6 F 3/06 3 0 1 X

G 0 6 F 3/08 (2006.01)

G 0 6 F 3/06 3 0 1 Z

G 0 6 F 13/10 (2006.01)

G 0 6 F 3/08 H

G 0 6 F 13/38 (2006.01)

G 0 6 F 13/10 3 4 0 A

G 0 6 F 13/38 3 5 0

請求項の数 15 (全23頁)

(21)出願番号 特願2018-68067(P2018-68067)

(22)出願日 平成30年3月30日(2018.3.30)

(65)公開番号 特開2018-173959(P2018-173959
A)

(43)公開日 平成30年11月8日(2018.11.8)

審査請求日 令和3年3月4日(2021.3.4)

(31)優先権主張番号 62/480113

(32)優先日 平成29年3月31日(2017.3.31)

(33)優先権主張国・地域又は機関

米国(US)

(31)優先権主張番号 62/483913

(32)優先日 平成29年4月10日(2017.4.10)

(33)優先権主張国・地域又は機関

米国(US)

(31)優先権主張番号 15/618081

最終頁に続く

(73)特許権者 390019839

三星電子株式会社

Samsung Electronics
Co., Ltd.大韓民国京畿道水原市靈通区三星路12
9129, Samsung-ro, Yeon
gtong-gu, Suwon-si
, Gyeonggi-do, Repub
lic of Korea

(74)代理人 110000051

特許業務法人共生国際特許事務所

(72)発明者

ラムダス ビー. カチャーレ

アメリカ合衆国 95014 カリフォル
ニア州 クバーチノ ノルマンディ ウェイ

最終頁に続く

(54)【発明の名称】 NVMe - oF 装置用ストレージ集積方法、NVMe - oF イーサネットSSDのグルー
プにNVMe - oF SSD容量を集積する方法、及び集積されたイーサネットSSDグ

(57)【特許請求の範囲】

【請求項1】

NVMe オーバーファブリック (NVMe - oF) 装置用ストレージ集積方法において、
前記方法は、集積グループを複数のNVMe - oF SSDを含む集積されたイーサネット (登録商標
) SSDとして確認するステップと、前記集積グループの前記NVMe - oF SSDのいずれかを1次NVMe - oF SSD
として選択するステップと、前記集積グループの前記NVMe - oF SSDのうち、残りのものを2次NVMe - oF
SSDとして選択するステップと、前記NVMe - oF SSDを管理するため、プロセッサが前記1次NVMe - oF S
SD内のマップ割り当てテーブルを初期化するステップと、を含み、

前記1次NVMe - oF SSDのみがホストで見ることができ、

前記1次NVMe - oF SSDは、前記ホストから受信した命令を、マップ割り当てテ
ーブルを参照して複数のサブ命令に分割して、前記2次NVMe - oF SSDへ送るス
テップと、1つ以上の前記2次NVMe - oF SSDは、サブ命令を実行し、サブ命令完了エント
リーを前記1次NVMe - oF SSDへ送るステップと、全てのサブ命令完了エントリーを受信すると、前記1次NVMe - oF SSDが完了エ
ントリーを生成して前記ホストに伝送するステップと、をさらに含むことを特徴とするN

VMe オーバーファブリック装置用ストレージ集積方法。

【請求項 2】

前記プロセッサが、前記マップ割り当てテーブルを初期化するステップは、前記集積グループに連結されたストレージ管理者 (storage administrator) のガイド下で遂行されることを特徴とする、請求項 1 に記載の NVMe オーバーファブリック装置用ストレージ集積方法。

【請求項 3】

前記マップ割り当てテーブルは、前記集積グループの前記 NVMe - oF SSD のそれぞれに対する、NVMe - oF SSD の容量、NVMe - oF SSD のアドレス、及び NVMe - oF SSD の残留容量を含むことを特徴とする、請求項 1 に記載の NVMe オーバーファブリック装置用ストレージ集積方法。

10

【請求項 4】

前記 1 次 NVMe - oF SSD のアドレスをユーザー応用プログラムに提供して、前記集積グループと前記ユーザー応用プログラムとの間のデータ伝送を活性化させるステップをさらに含むことを特徴とする、請求項 3 に記載の NVMe オーバーファブリック装置用ストレージ集積方法。

【請求項 5】

前記 1 つ以上の 2 次 NVMe - oF SSD のうち、対応する 2 次 NVMe - oF SSD が 1 つ以上の前記サブ命令のいずれか 1 つのサブ命令を受信するステップと、前記対応する 2 次 NVMe - oF SSD からのデータを前記 1 次 NVMe - oF SSD に伝送するかどうかを前記 1 つのサブ命令に応じて判断するステップと、をさらに含むことを特徴とする、請求項 1 に記載の NVMe オーバーファブリック装置用ストレージ集積方法。

20

【請求項 6】

前記 1 次 NVMe - oF SSD が、ネームスペース生成命令又はネームスペース削除命令を受信するステップと、前記 1 次 NVMe - oF SSD が前記マップ割り当てテーブルを参照するステップと、命令が前記ネームスペース生成命令であれば、前記 1 次 NVMe - oF SSD 及び / 又は前記 1 つ以上の前記 2 次 NVMe - oF SSD に容量を割り当て、前記命令が前記ネームスペース削除命令であれば、前記 1 次 NVMe - oF SSD 及び / 又は前記 1 つ以上の 2 次 NVMe - oF SSD の内の対応する 1 つを検索するステップと、前記マップ割り当てテーブルをアップデートするステップと、をさらに含むことを特徴とする、請求項 1 に記載の NVMe オーバーファブリック装置用ストレージ集積方法。

30

【請求項 7】

前記 1 次 NVMe - oF SSD がリード / ライト命令を受信するステップと、前記 1 次 NVMe - oF SSD が前記マップ割り当てテーブルを検索するステップと、1 つ以上のリード / ライトサブ命令を生成するステップと、前記 1 つ以上の前記 2 次 NVMe - oF SSD に前記 1 つ以上のリード / ライトサブ命令をそれぞれ伝送するステップと、前記リード / ライト命令に応じて前記ホストと前記 1 次 NVMe - oF SSD との間、及び / 又は前記 1 つ以上の 2 次 NVMe - oF SSD との間でデータを伝送するステップと、前記データを伝送した以後、前記ホストに完了メッセージを伝送するステップと、をさらに含むことを特徴とする、請求項 1 に記載の NVMe オーバーファブリック装置用ストレージ集積方法。

40

【請求項 8】

前記 1 つ以上の 2 次 NVMe - oF SSD のうち、対応する 2 次 NVMe - oF SSD が、前記 1 つ以上のリード / ライトサブ命令のいずれか 1 つのリード / ライトサブ命令を受信するステップと、前記 1 つのリード / ライトサブ命令に対応する伝送情報を抽出するステップと、

50

前記対応する2次NVM e - o F S S Dからのリード/ライト要請を前記ホストに発行するステップと、

をさらに含むことを特徴とする、請求項7に記載のNVM e オーバーファブリック装置用ストレージ集積方法。

【請求項9】

NVM e - o F (NVM e オーバーファブリック)イーサネットS S DのグループにNVM e - o F S S D容量を集積する方法において、

集積グループの複数のNVM e - o F S S Dを確認するステップと、

前記NVM e - o F S S Dのいずれかを1次NVM e - o F S S Dとして指定するステップと、

前記NVM e - o F S S Dのうち、残りのものを2次NVM e - o F S S Dとして指定するステップと、を含み、

ホストのホストドライバが見ることができるただ1つのNVM e - o F S S Dが、前記1次NVM e - o F S S Dであり、

前記1次NVM e - o F S S Dのみが前記ホストで見ることができ、

前記1次NVM e - o F S S Dが、前記ホストからの命令を受信するステップと、

前記1次NVM e - o F S S Dは、前記ホストから受信した命令を、マップ割り当てテーブルを参照して複数のサブ命令に分離するステップと、

対応する前記2次NVM e - o F S S Dのそれぞれに前記1次NVM e - o F S S Dからの前記サブ命令を伝送するステップと、を含み、

前記対応する2次NVM e - o F S S Dが、前記1次NVM e - o F S S Dからの前記サブ命令を受信するステップと、

前記サブ命令に対応するタスクを遂行するステップと、

前記タスクの完了によって前記対応する2次NVM e - o F S S Dが各サブ命令完了エントリーを前記1次NVM e - o F S S Dに伝送するステップと、

前記1次NVM e - o F S S Dを通じてサブ命令コンテキストテーブルを維持するステップと、

前記1次NVM e - o F S S Dが前記2次NVM e - o F S S Dからの前記サブ命令完了エントリーを受信するステップと、

前記受信されたサブ命令完了エントリーに応じて前記1次NVM e - o F S S Dが前記サブ命令の実行を追跡するステップと、

前記1次NVM e - o F S S Dは、全てのサブ命令完了エントリーを受信すると前記ホストに完了エントリーを送るステップと、

をさらに含むことを特徴とする、NVM e - o FイーサネットS S DのグループにNVM e - o F S S D容量を集積する方法。

【請求項10】

前記1次NVM e - o F S S Dが、前記ホストから前記1次NVM e - o F S S Dによって受信された命令に応じて前記マップ割り当てテーブルを維持するステップをさらに含み、

前記マップ割り当てテーブルは、前記集積グループの前記1次NVM e - o F S S D及び1つ以上の前記2次NVM e - o F S S Dの間で分割されたロジックブロックアドレス(LBA)スペースを示すことを特徴とする、請求項9に記載のNVM e - o FイーサネットS S DのグループにNVM e - o F S S D容量を集積する方法。

【請求項11】

前記NVM e - o F S S Dのいずれかを前記1次NVM e - o F S S Dとして指定することにより、プロセッサは、前記集積グループを構成するために前記マップ割り当てテーブルを初期化するステップをさらに含むことを特徴とする、請求項9に記載のNVM e - o FイーサネットS S DのグループにNVM e - o F S S D容量を集積する方法。

【請求項12】

前記1次NVM e - o F S S Dは、前記ホストから前記1次NVM e - o F S S Dによ

10

20

30

40

50

って受信された命令に応じて前記2次NVMe-oF SSDの容量を集積して、前記集積グループの前記複数のNVMe-oF SSDを前記ホストが1つの集積された論理容量として認識するステップをさらに含むことを特徴とする、請求項9に記載のNVMe-oFイーサネットSSDのグループにNVMe-oF SSD容量を集積する方法。

【請求項13】

前記1次NVMe-oF SSDを用いて1つ以上の前記2次NVMe-oF SSDに容量を割り当てるステップと、

前記1次NVMe-oF SSDを用いて前記割り当てられた容量及び関連するマッピングされたロジックブロックアドレス(LBA)の範囲を前記マップ割り当てテーブルに記録するステップと、をさらに含むことを特徴とする、請求項9に記載のNVMe-oFイーサネットSSDのグループにNVMe-oF SSD容量を集積する方法。

10

【請求項14】

前記1次NVMe-oF SSDが、前記2次NVMe-oF SSD及び前記1次NVMe-oF SSDの集積容量をオーバープロビジョニング(over provisioning)するステップをさらに含むことを特徴とする、請求項9に記載のNVMe-oFイーサネットSSDのグループにNVMe-oF SSD容量を集積する方法。

【請求項15】

前記サブ命令に基づいて前記対応する2次NVMe-oF SSDからのデータを前記ホストに直接伝送するステップをさらに含むことを特徴とする、請求項9に記載のNVMe-oFイーサネットSSDのグループにNVMe-oF SSD容量を集積する方法。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一つの大きな論理容量としてホストに感知されるよう、複数のメモリードライブ(例えば、eSSD)を集積するためのシステム及び方法に関する。

【背景技術】

【0002】

ソリッドステートドライブ(SSD)は、現代のITインフラの基本ストレージ装置として急速に定着しており、既存のハードディスクドライブ(HDD)を代替しつつある。SSDは、非常に短い待機時間、多くのデータリード/ライト処理量及び安定的なデータ格納などを提供する。

30

【0003】

NVMeオーバーファブリック(NVMe over fabrics;以下NVMe-oFという)は、数千、数百個のNVMe-oF装置(例えば、非揮発性メモリー(NVMe)SSD)がファイバーチャネル(Fiber channel:FC)、インフィニバンド(Infini Band:IB)、及びイーサネット(登録商標)(Ethernet、登録商標)のようなネットワークファブリックを介して連結される新しい技術である。NVMe-oFプロトコルは、遠隔ダイレクトアタッチストレージ(remote Direct Attach Storage:rDAS)の実行を活性化させる。これは、多くのSSDを遠隔ホストに連結することを許容する。NVMe-oFプロトコルは、NVMe命令、データ及び応答の信頼できる通信を提供するため、遠隔ダイレクトメモリアクセス(Remote Direct Memory Access:RDMA)を使用する。RDMAサービスを提供するための伝送プロトコルはiWARP、RoCE v1、及びRoCE v2を含む。

40

【0004】

NVMe-oFインターフェースは、多くのSSDが遠隔ホストに連結されるようにする。一般に、各NVMe-oF SSDに対するドライバインスタンスは遠隔ホストにて実行される。一部の応用プログラムの場合、一つのSSDから提供されるストレージ容量では充分でないからである。

【発明の概要】

50

【発明が解決しようとする課題】

【0005】

本発明の目的は、一つの大容量の論理ボリュームとしてホストに認識されるよう、複数のSSDを集積する方法及び前記方法を達成するためのネットワーク構造を提供することにある。

【課題を解決するための手段】

【0006】

本発明の実施例によれば、NVMe オーバーファブリック (NVMe - oF) 装置用ストレージ集積方法において、前記方法は、集積グループを複数のNVMe - oF SSDを含む集積されたイーサネットSSDとして確認するステップと、前記集積グループの前記NVMe - oF SSDのいずれかを1次NVMe - oF SSDとして選択するステップと、前記集積グループの前記NVMe - oF SSDのうち、残りのものを2次NVMe - oF SSDとして選択するステップと、前記NVMe - oF SSDを管理するため、プロセッサが前記1次NVMe - oF SSD内のマップ割り当てテーブルを初期化するステップとを含む。

【0007】

本発明の実施例によれば、前記プロセッサが前記マップ割り当てテーブルを初期化するステップは、前記集積グループに連結されたストレージ管理者 (storage administrator) のガイド下で遂行される。

【0008】

本発明の実施例によれば、前記マップ割り当てテーブルは、前記集積グループの前記NVMe - oF SSDのそれぞれに対する、NVMe - oF SSDの容量、NVMe - oF SSDのアドレス、及びMe - oF SSDの残留容量を含む。

【0009】

本発明の実施例によれば、前記方法は、前記1次NVMe - oF SSDのアドレスをユーザー応用プログラムに提供して、前記集積グループと前記ユーザー応用プログラムとの間のデータ伝送を活性化させるステップをさらに含む。

【0010】

本発明の実施例によれば、前記方法は、前記1次NVMe - oF SSDが、前記集積グループに連結されたホストからの命令を受信するステップと、前記命令に対応するデータが、前記1次NVMe - oF SSDに格納されるかどうか、又は1つ以上の前記2次NVMe - oF SSDに格納できるかどうかを判断するステップと、前記データが前記1つ以上の2次NVMe - oF SSDに格納された場合、前記命令を前記1つ以上の2次NVMe - oF SSDにそれぞれ対応する1つ以上のサブ命令に分割するステップと、前記ホストに前記データを伝送するステップと、前記1つ以上の2次NVMe - oF SSDからサブ命令完了エントリーを受信するステップと、前記1次NVMe - oF SSDが完了エントリーを生成して前記ホストに伝送するステップと、をさらに含む。

【0011】

本発明の実施例によれば、前記方法は、前記1つ以上の2次NVMe - oF SSDのうち、対応する2次NVMe - oF SSDが、1つ以上のサブ命令のいずれかのサブ命令を受信するステップと、前記対応する2次NVMe - oF SSDからのデータを前記1次NVMe - oF SSDに伝送するかどうかを前記一つのサブ命令に応じて判断するステップと、完了エントリーを生成するステップと、前記1次NVMe - oF SSDに前記完了エントリーを伝送するステップと、をさらに含む。

【0012】

本発明の実施例によれば、前記方法は、前記1次NVMe - oF SSDが、ネームスペース生成命令又はネームスペース削除命令を受信するステップと、前記1次NVMe - oF SSDが前記マップ割り当てテーブルを参照するステップと、命令が前記ネームスペース生成命令であれば、前記1次NVMe - oF SSD及び/又は1つ以上の前記2次NVMe - oF SSDに容量を割り当て、前記命令が前記ネームスペース削除命令であ

10

20

30

40

50

れば、前記1次NVMe - oF SSD及び/又は前記1つ以上の2次NVMe - oF SSDの内の対応する一つを検索するステップと、前記マップ割り当てテーブルをアップデートするステップと、をさらに含む。

【0013】

本発明の実施例によれば、前記方法は、前記1次NVMe - oF SSDがリード/ライト命令を受信するステップと、前記1次NVMe - oF SSDが前記マップ割り当てテーブルを検索するステップと、1つ以上のリード/ライトサブ命令を生成するステップと、1つ以上の前記2次NVMe - oF SSDに前記1つ以上のリード/ライトサブ命令をそれぞれ伝送するステップと、前記リード/ライト命令に応じてホストと前記1次NVMe - oF SSDとの間、及び/又は前記1つ以上の2次NVMe - oF SSDとの間でデータを伝送するステップと、前記データを伝送した以後、前記ホストに完了メッセージを伝送するステップと、をさらに含む。

10

【0014】

本発明の実施例によれば、前記方法は、前記1つ以上の2次NVMe - oF SSDのうち、対応する2次NVMe - oF SSDが、前記1つ以上のリード/ライトサブ命令のいずれかのリード/ライトサブ命令を受信するステップと、1つのリード/ライトサブ命令に対応する伝送情報を抽出するステップと、前記対応する2次NVMe - oF SSDからのリード/ライト要請を前記ホストに発行するステップと、前記リード/ライトサブ命令に応答してデータ伝送を完了した以後、前記1次NVMe - oF SSDに完了エントリーを伝送するステップと、をさらに含む。

20

【0015】

本発明の実施例によれば、NVMe - oF (NVMe オーバーファブリック) イーサネットSSDのグループにNVMe - oF SSD容量を集積する方法において、前記方法は、集積グループの複数のNVMe - oF SSDを確認するステップと、前記NVMe - oF SSDのいずれかを1次NVMe - oF SSDとして指定するステップと、前記NVMe - oF SSDのうち、残りのものを2次NVMe - oF SSDとして指定するステップと、を含み、ここでホストのホストドライバが見ることができるただ1つのNVMe - oF SSDが前記1次NVMe - oF SSDである。

【0016】

本発明の実施例によれば、前記方法は、前記1次NVMe - oF SSDが、前記ホストから前記1次NVMe - oF SSDによって受信された命令に応じてマップ割り当てテーブルを維持するステップをさらに含み、前記マップ割り当てテーブルは、前記集積グループの前記1次NVMe - oF SSD及び1つ以上の前記2次NVMe - oF SSDの間で分割されたロジックブロックアドレス(LBA)スペースを示す。

30

【0017】

本発明の実施例によれば、前記方法は、前記NVMe - oF SSDのいずれかを前記1次NVMe - oF SSDとして指定することにより、プロセッサは、前記集積グループを構成するためにマップ割り当てテーブルを初期化するステップをさらに含む。

【0018】

本発明の実施例によれば、前記方法は、前記1次NVMe - oF SSDは、前記ホストから前記1次NVMe - oF SSDによって受信された命令に応じて前記2次NVMe - oF SSDの容量を集積して、前記集積グループの前記複数のNVMe - oF SSDを前記ホストが1つの集積された論理容量として認識するステップをさらに含む。

40

【0019】

本発明の実施例によれば、前記方法は、前記1次NVMe - oF SSDを用いて1つ以上の前記2次NVMe - oF SSDに容量を割り当てるステップと、前記1次NVMe - oF SSDを用いて前記割り当てられた容量及び関連するマッピングされたロジックブロックアドレス(LBA)の範囲をマップ割り当てテーブルに記録するステップと、をさらに含む。

【0020】

50

本発明の実施例によれば、前記方法は、前記１次NVMe - oF SSDが、前記２次NVMe - oF SSD及び前記１次NVMe - oF SSDの集積容量をオーバプロビジョニング(over provisioning)するステップをさらに含む。

【0021】

本発明の実施例によれば、前記方法は、前記１次NVMe - oF SSDが前記ホストからの命令を受信するステップと、前記命令をそれぞれ、前記２次NVMe - oF SSDのうち、対応する２次NVMe - oF SSDに対応する複数のサブ命令に分離するステップと、前記対応する２次NVMe - oF SSDのそれぞれに前記１次NVMe - oF SSDからの前記サブ命令を伝送するステップと、を含む。

【0022】

本発明の実施例によれば、前記方法は、前記サブ命令に基づいて前記対応する２次NVMe - oF SSDからのデータを前記ホストに直接伝送するステップをさらに含む。

【0023】

本発明の実施例によれば、前記方法は、前記対応する２次NVMe - oF SSDが前記１次NVMe - oF SSDからの前記各サブ命令を受信するステップと、前記サブ命令に対応するタスクを遂行するステップと、前記タスクの完了によって前記対応する２次NVMe - oF SSDが各サブ命令完了エントリーを前記１次NVMe - oF SSDに伝送するステップと、をさらに含む。

【0024】

本発明の実施例によれば、前記方法は、前記１次NVMe - oF SSDを通じてサブ命令コンテキストテーブルを維持するステップと、前記１次NVMe - oF SSDが前記２次NVMe - oF SSDからのサブ命令完了エントリーを受信するステップと、前記受信されたサブ命令完了エントリーに応じて前記１次NVMe - oF SSDが前記サブ命令の実行を追跡するステップと、をさらに含む。

【0025】

本発明の実施例によれば、集積されたイーサネットSSDグループにおいて、イーサネットSSDシャーシと、ホストドライバと通信を可能なようにするための前記イーサネットSSDシャーシ上のイーサネットスイッチと、前記イーサネットスイッチと結合されたプロセッサと、ボード管理コントローラーに結合されたPCleスイッチと、複数のNVMe - oF (NVMe オーバーファブリック) SSDと、を含み、前記複数のNVMe - oF SSDは、１次NVMe - oF SSDと、前記イーサネットスイッチ及び前記PCleスイッチを含む専用通信チャンネルを通じて前記１次NVMe - oF SSDに連結された複数の２次NVMe - oF SSDと、を含み、ここで、前記１次NVMe - oF SSDのみが前記ホストで見ることができ、前記ボード管理コントローラーは、前記NVMe - oF SSDのいずれかに前記１次NVMe - oF SSDが含まれるかどうかを初期に判断するように構成される。

【0026】

従って、１つの１次eSSDが、全てのNVMe - oF プロトコル処理を遂行しながら２次eSSDによる全ての関連するサブ命令の完了を追跡し、２次eSSDがホストで見えないように維持するため、eSSDの集積グループのホストに１つの大きな論理容量として認識される。

【発明の効果】

【0027】

本発明によれば、ただ１つの１次eSSDだけをホストが見、全ての関連するサブ命令の完了を２次eSSDで追跡しながらNVMe - oF プロトコルの処理を全て遂行するため、eSSDの集積グループは、１つの大きな論理容量としてホストから見えるようにできる。

【図面の簡単な説明】

【0028】

各実施例は、添付の図面と結合して以下の説明からより詳細に理解することができる。

10

20

30

40

50

【 0 0 2 9 】

【図 1】本発明の一実施例に係る 1 つの e S S D シャーシに複数の集積 e S S D を含む N V M e - o F イーサネット S S D (e S S D) ストレージに使用されるシステムの構成を示すブロック図である。

【図 2】本発明の一実施例に係るイーサネット S S D ラック内に共に連結された図 1 の実施例に示す複数の e S S D シャーシのブロック図である。

【図 3】本発明の一実施例に係る 1 次 e S S D によって維持される「マップ割り当てテーブル」の例を示す図である。

【図 4】本発明の一実施例に係るマップ割り当てテーブルの初期化を説明するためのフローチャートである。

10

【図 5】本発明の一実施例に係る複数のラックにそれぞれ位置する複数の e S S D シャーシ内で、複数の集積された e S S D へのデータの流れを含む N V M e - o F イーサネット S S D (e S S D) ストレージに使用されるシステム構成を示すブロック図である。

【図 6】例示的命令コンテキストを説明するテーブルである。

【図 7】本発明の一実施例に係る 1 次 e S S D による命令の処理を説明するためのフローチャートである。

【図 8】本発明の一実施例に係るネームスペース生成及び削除命令の実行を説明するためのフローチャートである。

【図 9】本発明の一実施例に係る P - e S S D の制御下で、リード/ライト命令の実行を説明するためのフローチャートである。

20

【図 10】本発明の一実施例に係る S - e S S D によるサブ命令の実行を説明するためのフローチャートである。

【図 11】本発明の一実施例に係る S - e S S D によるリード/ライトサブ命令の実行を説明するためのフローチャートである。

【図 12】本発明の一実施例に係る S - e S S D に同期したデータ伝送及びサブ命令完了を説明する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】

【 0 0 3 0 】

本発明の概念特徴及びこれを達成する方法は、以下の実施例に関する詳細な説明及び添付の図面を参照すれば、より容易に理解できる。以下、添付の図面を参照して本発明の実施例を詳細に説明し、ここで同一の参照符号は同一の構成要素を示す。しかし、本発明は多様な形態で具体化でき、ここに示した実施例のみに限定されるものと解釈されてはならない。これらの実施例は、本発明が徹底且つ完全になるように当業者に本発明の側面及び特徴を十分に伝達できるように例として提供されたものである。従って、本発明の側面及び特徴の完全な理解のために当業者にとって不必要なプロセス、要素及び技術は説明しないこともある。特に言及しない限り、添付の図面及び詳細な説明の全般にわたって同一の参照符号は同一の構成要素を示し、それに対する説明は反復しない。図面において、要素、層及び領域の相対的な大きさは、明確性のために誇張することもある。

30

【 0 0 3 1 】

以下の記述においては、説明を目的として、多くの特定の細部事項を多様な実施例の完全な理解のために提供する。しかし、多様な実施例が、これらの特定の細部事項又は 1 つ以上の等価の構成なしに実施されることもあるのは自明である。他の例で、公知の構造及び装置は、不必要に多様な実施例を曖昧にすることを避けるため、ブロック図の形態で示す。

40

【 0 0 3 2 】

「第 1」、「第 2」、「第 3」などの用語を、本明細書において多様な構成要素、成分、領域、層及び/又はセクションを説明するために使用するが、これらの構成要素、成分、領域、層及び/又はセクションは、これらの用語によって制限されてはいけない。これらの用語は、1 つの要素、成分、領域、層又はセクションを他の要素、成分、領域、層又はセクションと区別するために使用する。よって、以下に説明する第 1 要素、成分、領域、層又はセクションは、本発明の思想及び範囲を逸脱することなく、第 2 の要素は、成分、

50

領域、層又はセクションと称することもできる。

【0033】

「底に」、「下に」、「下部に」、「上に」、「上部に」などのような空間的に相対的な用語は、図面に示すように1つの構成要素又は他の要素又は特徴との特徴関係を説明するため便宜のために使用する。また、空間的に相対的な用語は、図面に示した方位に加えて、使用又は作動時に装置の相異なる方位を含むものと理解できる。例えば、図面の装置が裏返されると、他の要素又は特徴の「下」又は「底」又は「下部」と記述された要素は、他の要素又は特徴の「上」で方位される。よって、「下」及び「下部」の例示的な用語は、上と下の方向の両方を含む。前記装置は、他の方向に方位され（例えば、90度又は他の方位に回転できる）、本明細書で使用する空間的に相対的な記述用語は、状況に合わせて解釈すべきである。

10

【0034】

構成要素、層、領域又は成分が、他の構成要素、層、領域又は成分の「上に」、「連結された」又は「結合された」と言及する場合、これは、他の構成要素、層、領域又は成分の直接上にあるか、連結されるか、結合され、1つ以上の介在した構成要素、層、領域又は成分などが存在することもある。また、1つの構成要素又は層が、2つの構成要素又は層の「間に」とあると言及する時、2つの構成要素又は層の間に1つの構成要素又は層が存在するか、又はその間に1つ以上の構成要素又は層が存在することもある。

【0035】

本明細書の目的上、「X、Y及びZうちの少なくとも1つ」及び「X、Y及びZからなるグループから選ばれた少なくとも1つ」は、Xのみ、Yのみ、Zのみと解釈でき、又はXYZ、XYX、YZ及びZZのようなX、Y及びZのうち、2つ以上の任意の組み合わせを含む。同一の符号は同一の構成要素を示す。本明細書で使用する「及び/又は」という用語は、1つ以上の列挙された関連項目の任意及び全ての組み合わせを含む。

20

【0036】

以下の実施例において、x軸、y軸及びz軸は、直角座標系の3つの軸に限定されず、より広義に解釈できる。例えば、x軸、y軸及びz軸は、互いに垂直であるか、互いに垂直でなく相異なる方向を示す。

【0037】

本明細書において使用する用語は、特定の実施例を説明するためのものであり、本発明を限定するためのものではない。ここで使用する単数形態は、文脈上で特に指示しない限り、複数形態を含む。本明細書で使用する「含む。」、「含む〜」、「備える。」及び「備える〜」という用語は、明示した特徴、定数、段階、動作、構成要素及び/又は成分の存在を特定するが、1つ以上の他の特徴、定数、段階、動作、構成要素、成分及び/又はこれらのグループの存在又は追加を排除しない。本明細書で使用する「及び/又は」という用語は、1つ以上の列挙した関連項目の任意及び全ての組み合わせを含む。「少なくとも1つ」のような表現が要素目録の前にある時、要素の全体目録を修正し、目録の個別要素を修正しない。

30

【0038】

本明細書で使用する用語「実質的に」、「略」及びこれと類似の用語は、近似語として使用され、程度用語として使用されるものではなく、当業者にとって認識できる測定された又は計算された数値の固有の偏差を説明するためである。また、本発明の実施例を説明する時、「できる」という用語を使用するのは、「本発明が1つ以上の実施例」からなることを意味する。本明細で使用する「使用する。」、「使用する〜」及び「使用される〜」という用語は「活用する。」、「活用する〜」及び「活用される〜」とそれぞれ同義語として見做すことができる。なお、「例示的な」という用語は、例又は説明を意味する。

40

【0039】

特定の実施例が異なるように実施される場合、特定処理順序が説明した順序と異なるように遂行できる。例えば、2つが連続的に記述されたプロセスは、実質的に同時に遂行されるか、又は説明した順序と反対の順序で遂行できる。

50

【 0 0 4 0 】

また、本明細書に開示し、及び／又は引用する任意の数値範囲は、前記範囲内に含まれる同一の数値精度のすべての下位範囲を含む。例えば、「1.0乃至10.0」の範囲には、指定された1.0の値と10.0の値の間のすべてのサブ範囲、即ち、最少値を有する1.0より大きいかまたは同じであり、最大値の10.0より小さいかまたは同一の例を挙げて、2.4乃至7.6のようなサブ範囲が含まれる。本願に引用する最大数値限界は、ここに含まれた全てのより低い数値限界を含み、本明細書で引用する任意の最少数値限界は、ここに含まれた全てのより高い数値限界を含む。従って、出願人は、本明細書で明示的に列挙する範囲内に含まれた任意のサブ範囲を明示的に言及するため、請求の範囲を含む本明細書を補正する権利を保有する。

10

【 0 0 4 1 】

多様な実施例は、各実施例の概略的な説明の断面積説明及び／又は中間構造を参照して、ここに述べる。このように、例えば、製造技術及び／または許容誤差のような結果により説明の形態が変更されることは予想できる。従って、本願に開示する実施例は、具体的に説明する領域の形状に限定されるものと解釈されてはならず、例えば、製造のような形状の偏差を含まなければならない。例えば、長方形として説明した注入領域は、一般に注入領域から非注入領域への2値変化というよりは、ラウンド(r o u n d)又はカーブした(c u r v e d)特徴及び／又は角で注入濃度が変化される特性を有する。同様に、注入によって形成された埋め込み領域は、埋め込み領域と注入が起こる表面との間の領域に若干の注入を引き起こす。これによって、図面に示す領域は、本質的に概略的なものであり、その形状は、装置の実際形状を例示するものではなく、制限するものでもない。

20

【 0 0 4 2 】

ここに述べる本発明の実施例に係る電子または電気装置及び／又は任意の他の関連装置又は構成要素は、任意の適合したハードウェア、ファームウェア(例えば、注文型集積回路)、ソフトウェア、又はこれらの組み合わせを用いて実施される。例えば、これらの装置の多様な構成要素は、1つの集積回路(I C)チップ上に又は個別I Cチップ上に形成できる。また、これらの装置の多様な構成要素は、可撓性印刷回路フィルム、テープキャリアパッケージ(T C P)、印刷回路基板(P C B)上に具現されたり、1つの基板上に形成されることができる。また、これらの装置の多様な構成要素は、1つ以上のプロセッサ、1つ以上のコンピューティング装置で実行され、コンピュータプログラムの命令を遂行し、ここに説明する多様な機能を遂行するため、他のシステム構成要素と相互に作用するプロセッサ又はスレッド(t h r e a d)である。前記コンピュータプログラム命令は、例えばランダムアクセスメモリー(R A M)のような標準メモリー装置を利用するコンピューティング装置で具現されるメモリーに格納される。また、前記コンピュータプログラム命令は、例えばC D - R O M、フラッシュドライブなどのような非一時的コンピュータ判読可能媒体に格納される。また多様なコンピューティング装置の機能が、単一のコンピューティング装置に結合されるか、統合され、又は特定コンピューティング装置の機能が、本発明の思想と範囲を逸脱しない範囲内で1つ以上の他のコンピューティング装置にわたって分散される。

30

【 0 0 4 3 】

特に定義がない限り、本明細書で使用する全ての用語(技術用語及び科学用語を含む)は、本発明が属する技術の分野における当業者にとって一般的に理解できるものと同一の意味を持つ。また、一般的に使用される辞書で定義された用語のような用語は、関連技術及び／又は本明細書と関連してその意味と一致する意味を持つものと解釈すべきであり、理想的であるか、又は過度に形式的な意味として解釈されてはならない。

40

【 0 0 4 4 】

図1は、本発明の一実施例に係る1つのシャーシ120に複数の集積e S S Dを含むN V M e - o F イーサネットS S D (e S S D) ストレージに使用されるシステムの構成を示すブロック図である。図2は、本発明の一実施例に係るイーサネットS S D ラック内に共に連結された図1の実施例に示すe S S Dを含む複数のシャーシ120のブロック図で

50

ある。

【 0 0 4 5 】

上述したように、N V M e - o F インターフェースは、多くの e S S D 1 1 0 が遠隔ホスト 1 9 0 に連結されるようにする。一般的に、各 N V M e - o F e S S D に対するドライバインスタンスは遠隔ホスト 1 9 0 により実行される。しかし、一部の応用プログラムの場合、1つの e S S D 1 1 0 から提供されるストレージ容量では充分でないことがある。そのような応用プログラムには、数百テラバイトの容量を有する1つの論理ボリュームが必要である。したがって、そのような応用プログラムには、1つの「集積グループ」に全て集積される多くの個別 S S D を提供して、応用プログラムに1つの大きな論理ボリュームとして認識されるようにする本発明が有効である。

10

【 0 0 4 6 】

例えば、24個の16-テラバイト(16TB) e S S D は、1つの論理384TBドライブとして認識される。複数の集積された e S S D 1 1 0 を必要とするいくつかの応用プログラムは、ビッグデータマイニング及び分析、石油化学、ガス及びエネルギー探査、実験粒子物理学、及び医薬品の開発を含む。このような例は、大きな格納容量と高性能を必要とする高性能コンピューティング(HPC)を要求する。

【 0 0 4 7 】

基本となる e S S D 1 1 0 を集めて、1つの論理的で膨大なボリュームを提供するシステムソフトウェア層を有することができるが、そのようなシステムソフトウェアは、一般的に非常に複雑かつ精巧である。このようなソフトウェアは、遠隔ホスト 1 9 0 で実行される複数の N V M e - o F ドライバインスタンスを要求し、これによって、メモリー、CPU サイクル及び電力のようなシステム資源を消費する。このソリューションは、x86サーバー又は R O C (R A I D - o n - C h i p) システムを使用して、大容量を1つの論理ボリュームとして提供する。しかし、このようなソリューションは、一般的に複雑で高価であり、性能及びエネルギーペナルティを有する。例えば、CPUを使用してデータを受信し、伝送する時、本発明の実施例に係るDMAエンジン、ASICなどにより消費されるエネルギーよりも数倍多いエネルギーを消費する。

20

【 0 0 4 8 】

従って、本発明の各実施例は、効率的で低コストで複数の e S S D 1 1 0 を集積させるための N V M e - o F e S S D に使用される方法及び構造を提供する。

30

【 0 0 4 9 】

図1を参照すれば、e S S D 1 1 0 には、2つの役割(例えば、ストレージ管理者(storage administrator)の指示に従って)の1つが割り当てられて、各 e S S D 1 1 0 が、1次 e S S D (P - e S S D) 又は2次 e S S D (P - e S S D) の内の1つの役割を遂行する。1つのシャーシ 1 2 0 (又は、1つの与えられたラック内の複数のシャーシ 1 2 0 又は広い領域に分散された複数のラック 2 3 0 内の複数のシャーシ 1 2 0) 内の1つの P - e S S D 1 1 0 p、及び1セットの複数の S - e S S D 1 1 0 s は、1つの論理ドライブとして遠隔ホスト 1 9 0 によって使用される必須フラッシュメモリー容量を集散的に提供する。シャーシ 1 2 0 は、ボード管理コントローラ装置(BMC) 1 5 0 のようなプロセッサ及び外部接続のためのイーサネットスイッチ 1 6 0 と共に e S S D 1 1 0 を含む。シャーシ 1 2 0 は、以下の各実施例の説明において N V M e - o F 装置のグループを指すが、本発明の他の実施例は、物理的なハウジング(例えば、シャーシ、ラック又はコンテナ基盤ハウジング)に関係なく、任意の他の複数の N V M e - o F 装置に類似に適用される。また、e S S D 1 1 0 は、後述する実施例の N V M e - o F 装置を述べるために使用するが、他の N V M e - o F 装置は、本発明の他の実施例において同様に適用される。

40

【 0 0 5 0 】

従って、e S S D 1 1 0 の集積が対応するイーサネットスイッチ 1 6 0 を通じてラック 2 3 0 内の複数のシャーシ 1 2 0 にわたっているだけでなく、それぞれが複数のシャーシ 1 2 0 を含む複数のラック 2 3 0 にわたっている。

50

【0051】

P - e S S D 110 p は、遠隔ホスト N V M e - o F ドライバ 170 から見える (v i s i b l e t o) 唯一の e S S D 110 であり、N V M e - o F プロトコル終了を処理する。P - e S S D 110 p は、自身及び同じ集積グループにおける全ての残りの S - e S S D 110 s の代わりに、遠隔ホスト 190 に1つの大規模の集積論理容量を提供する。P - e S S D 110 p は、遠隔ホスト N V M e - o F ドライバ 170 から全ての入/出力 (I / O) の命令 180 を受信し、命令応答 (例えば、完了エントリー) 182 を遠隔ホスト 190 に提供する。

【0052】

また、P - e S S D 110 p は、e S S D 110 の同じ集積グループの一部又は全ての S - e S S D 110 s と共に、P - e S S D 110 p の間で分割された論理ブロックアドレス (L B A) 空間を示すマップ割り当てテーブル (M A T) を維持する。命令 180 が、P - e S S D 110 p により受信されれば、P - e S S D 110 p は、e S S D 110 のいずれか (例えば、P - e S S D 110 p、1つ以上の S - e S S D 110 s、又は2つ以上の全てのセット) が命令 180 を満足させるかどうかを決定するため、先に M A T (例えば、以下に述べる図3の M A T 300) を検索する。その後、M A T によって、P - e S S D 110 p は適切に修正された N V M e - o F I / O のサブ命令 132 を適切な S - e S S D 110 s のセットに伝送する。

【0053】

P - e S S D 110 p は、サブ命令 132 を伝送するために電源をオンした後、P C I e スイッチ 140 及び制御プレーン 135 を介してそれぞれの S - e S S D 110 s と、専用イーサネット R D M A 接続 (又は、独占的通信チャンネル) を構築する。この専用キューペア (Q P) 通信チャンネルは、P - e S S D 110 p が S - e S S D 110 s に I / O 命令 (例えば、サブ命令 132) を電送し、S - e S S D 110 s からの完了エントリーの受信に使用する。専用通信チャンネル 130 は、イーサネットであり得、イーサネットスイッチ 160 を通じてデータが伝送される。しかし、また専用通信チャンネル 130 は P C I e 基盤チャンネルであることもでき、データが P C I e スイッチ 140 を通じて伝送される。即ち、全ての e S S D 110 は、相互に通信するために2つ以上の通信モードを使用する。例えば、イーサネットチャンネルは、一般的なデータ伝送に使用され、P C I e 基盤チャンネルは、管理のために使用され、どのチャンネルも専用通信チャンネル 130 として使用される。

【0054】

S - e S S D 110 s は、N V M e - o F プロトコルを使用して遠隔ホスト 190 にデータを伝達し、そして遠隔ホスト 190 からデータを受送信するだけのノーマル N V M e - o F S S D である。このようなデータ伝送は、R D M A リード及びライトサービスを使用して遠隔ホスト 190 に直接行われる。S - e S S D 110 s は、P - e S S D 110 p から命令 (例えば、サブ命令 132) を受信し、遠隔ホスト N V M e - o F ドライバ 170 からは直接命令を受信しない。S - e S S D 110 s は、サブ命令 132 の完了を示すために、遠隔ホスト 190 ではなく P - e S S D 110 p にサブ命令完了エントリー 134 を伝送する。

【0055】

P - e S S D 110 p は、全ての N V M e - o F プロトコル終了を処理し、ホスト命令及び完了キューイング (例えば、提出キュー/完了キュー (S Q / C Q)) を全て処理し、遠隔ホストイニシエータ上で実行される遠隔ホスト N V M e - o F ドライバ 170 が P - e S S D 110 p を見ることができる。遠隔ホスト N V M e - o F ドライバ 170 が N V M e 管理命令又は命令 180 を発行すれば、P - e S S D 110 p に命令 180 が発行され、P - e S S D 110 p によって全ての命令 180 が実行される。しかし、命令 180 は、複数の e S S D 110 に拡散できる。

【0056】

また、P - e S S D 110 p は、命令 180 に応じて自身が占有したデータ伝送を遂行

10

20

30

40

50

できる。P - e S S D 1 1 0 p は、オリジナルの命令 1 8 0 に対応する命令応答 1 8 2 としての命令完了エンタリーを遠隔ホスト 1 9 0 に伝送する前、全てのサブ命令完了エンタリー 1 3 4 が専用通信チャンネル 1 3 0 から到着するの (S - e S S D 1 1 0 s から) を待つ。

【 0 0 5 7 】

なお、P - e S S D 1 1 0 p は、命令コンテキストテーブル (c o m m a n d c o n t e x t t a b l e) (例えば、図 6 参照) の実行時に各命令に対する「命令コンテキスト」を維持する。この命令コンテキストは、P - e S S D 1 1 0 p がサブ命令 1 3 2 の実行、データ伝送及び任意のエラー状態の追跡に使用される。全てのサブ命令 1 3 2 が完了すれば、命令応答 1 8 2 又は完了エンタリーが遠隔ホスト 1 9 0 に提供され、命令コン

10

【 0 0 5 8 】

図 2 を参照すれば、複数の e S S D のシャーシ 1 2 0 は、イーサネットのラック 2 3 0 内で共に連結され、T O R (T o p - o f - R a c k) スイッチ 2 4 0 は、共通のラック 2 3 0 内の複数のシャーシ 1 2 0 の間を連結させるために提供される。同様に、相互に異なる地理的位置に位置した複数のラック 2 3 0 は、相互に直接連結されるか、又は外部スイッチを介して相互に連結されたそれぞれの T O R スイッチ 2 4 0 を介して相互に連結される。イーサネットのラック 2 3 0 は、1 つのデータセンターの建物内にあるか、広範囲な地理的領域に分散される。

【 0 0 5 9 】

要約すれば、本発明の実施例は、1 つの大容量 N V M e - o F S S D として提供される複数のイーサネット N V M e - o F S S D (e S S D 1 1 0) を集積させるメカニズムを提供する。e S S D 1 1 0 は、1 つのシャーシ 1 2 0 に位置し、1 つのラック 2 3 0 内の複数のシャーシ 1 2 0 に位置し、又はそれぞれ複数のシャーシ 1 2 0 を有する複数のイーサネットのラック 2 3 0 にわたって分散される。e S S D 1 1 0 内の 1 つには、1 次 e S S D (P - e S S D 1 1 0 p) の役割が割り当てられる。他の e S S D 1 1 0 には 2 次 e S S D (S - e S S D 1 1 0 s) の役割が割り当てられる。S - e S S D 1 1 0 s は遠隔ホストイニシエータと直接データ伝送を遂行するが、S - e S S D 1 1 0 s は、P - e S S D 1 1 0 p からサブ命令 1 3 2 を受信し、サブ命令 1 3 2 を完了し、P - e S S D 1 1 0 p にこれらのサブ命令 1 3 2 に対するサブ命令完了エンタリー 1 3 4 を伝

20

30

【 0 0 6 0 】

図 3 は、本発明の一実施例に係る 1 次 e S S D によって維持される「マップ割り当てテーブル」の例を示す図である。図 4 は、本発明の一実施例に係るマップ割り当てテーブルの初期化を説明するためのフローチャートである。

【 0 0 6 1 】

図 3 及び図 4 を参照すれば、上述したように本実施例は、2 つのタイプの e S S D 1 1 0 (例えば、P - e S S D 1 1 0 p 及び S - e S S D 1 1 0 s) を活用する。P - e S S D 1 1 0 p 及び S - e S S D 1 1 0 s は、全て遠隔ホスト 1 9 0 にストレージサービスを提供するために N V M e - o F プロトコルを使用する。P - e S S D 1 1 0 p は、1 つの論理ボリュームとして遠隔ホスト 1 9 0 に認識される e S S D 1 1 0 の集積グループ内に S - e S S D 1 1 0 s の細部事項を含むテーブル (例えば、マップ割り当てテーブル (M A T)) を維持する。

40

【 0 0 6 2 】

M A T 3 0 0 は、P - e S S D 1 1 0 p と同一のシャーシ 1 2 0 内の B M C 1 5 0 により初期化される。B M C 1 5 0 は、イーサネットスイッチ 1 6 0 及び e S S D 1 1 0 と同一の構成要素、及びイーサネットのシャーシ 1 2 0 を管理する。B M C 1 5 0 は、システム管理目的として P C I e 及び S M B u s インターフェースを備える。なお、B M C 1 5 0 は、どのような e S S D 1 1 0 が集積されるかを決定し (S 4 1 0) (例えば、ストレ

50

ージ管理者の指示下で)、eSSD110が決定されれば、BMC150はイーサネットスイッチ160を構成する。

【0063】

MAT300の左側にある3つのコラムは、ストレージ管理者の案内に従ってBMC150により初期化される。BMC150及びストレージ管理者は、集積グループ/ストレージシステムに存在する全てのeSSD110を見られ、これらの情報を持つ。このような情報は、eSSD110それぞれの容量311及びアドレス位置312を含む。ストレージ管理者は、「集積イーサネットSSD」(例えば、集積グループ)の形成に必要なeSSD110を決定する。BMC150及びストレージ管理者は、遠隔ホストNVMe-oFドライバ170に対応するユーザー応用プログラムが集積されたイーサネットSSDを見つかる場所を知ることができるよう、P-eSSD 110pのネットワークアドレスをユーザーに知らせるか又は提供する。また、BMC150及びストレージ管理者は、eSSD110のいずれかをP-eSSD 110pとして選択又は指名し(S420)、初期指定後、様々な理由でP-eSSD 110pに指定されたeSSD110を変更できる。その後、BMC150は、集積グループの1次及び2次モードをプログラムする(S430)。

10

【0064】

また、P-eSSD 110pは、BMC150上にMAT300のコピーを維持し、BMC150と共に格納されたMAT300のコピーを定期的にアップデートする。一部の実施例において、P-eSSD 110pだけが公式的なMAT300を含み、「0」のeSSDインデックス313はP-eSSD 110pを示し、残りのeSSDインデックス値はS-eSSD 110sのそれぞれに対応する。

20

【0065】

P-eSSD 110pは、遠隔ホストNVMe-oFドライバ170に対するNVMe-oFプロトコルを終了し、遠隔ホストNVMe-oFドライバ170によって発行された全ての命令180を実行する。遠隔ホスト190の命令180が完了すれば、P-eSSD 110pは、命令応答182としての完了エントリーを遠隔ホストNVMe-oFドライバ170に戻す形態で送付する。遠隔ホスト190の命令180と関連して、遠隔ホストNVMe-oFドライバ170は、S-eSSD 110sの存在を完全に認識できない。なお、P-eSSD 110pは、提出キュー(SQ)を維持し、命令完了を完了キュー(CQ)に提出する。

30

【0066】

MAT300の右側にある3つのコラムは、P-eSSD 110pによってアップデートされ、維持される。遠隔ホストNVMe-oFドライバ170が「ネームスペース」を生成すれば、特定のフラッシュ容量がそのネームスペースに割り当てられる。ネームスペースLBA範囲314は、eSSD110のセットにマッピングされ、P-eSSD 110pによって維持されるMAT300に記録される。このプロセスの細部事項は、以下の図8を参照して説明する。

【0067】

また、P-eSSD 110pは特定初期化動作を遂行する。MAT300がP-eSSD 110pにより初期化されれば(S440)、P-eSSD 110pは、eSSD110のいずれが対応するS-eSSD 110sであることを認識する。その後、P-eSSD 110pは、集積グループ内のS-eSSD 110sのそれぞれと専用通信チャンネル130を設定する。専用通信チャンネル130は、シャーシ120のイーサネットPCTeスイッチを通過するPCIeインターフェースを通じてなされるか、又はシャーシ120のPCIeスイッチ140を通過するPCIeインターフェースを通じてなされる。S-eSSD 110sの内の1つが同一のラック230に位置した他のシャーシ120に位置すれば、専用通信チャンネル130は、TORSスイッチ240を通じて設定される。同様に、与えられたシャーシ120内の専用通信チャンネル130は、PCIeスイッチ140を通じても設定される。

40

50

【 0 0 6 8 】

図 5 は、本発明の一実施例に係る複数のラックにそれぞれ位置する複数の e S S D シャーシ内で、複数の集積された e S S D へのデータの流れを含む N V M e - o F イーサネット S S D (e S S D) ストレージに使用されるシステム構成を示すブロック図である。

【 0 0 6 9 】

図 5 を参照すれば、P - e S S D 1 1 0 p は、外部ネットワークスイッチ及びルータを介して広域ネットワーク (W A N) 内の S - e S S D 1 1 0 s とイーサネット通信チャンネル 5 3 0 を設立する。このような専用のイーサネット通信チャンネル 5 3 0 は、R D M A キューペア (Q P) になるか独占的な方法になる。イーサネット通信チャンネル 5 3 0 は、サブ命令 1 3 2 及び関連した完了メッセージの交換に使用される。

10

【 0 0 7 0 】

図 6 は、例示的命令コンテキストを説明するテーブルである。

【 0 0 7 1 】

図 6 を参照すれば、各サブ命令 1 3 2 は命令 I D 6 4 0 を有し、P - e S S D 1 1 0 p がサブ命令完了エントリー 1 3 4 を受信すれば、サブ命令完了エントリー 1 3 4 がオリジナルの命令 1 8 0 により再び追跡できるように、「命令タグ」6 1 0 を持つ。サブ命令 1 3 2 のサブ命令完了エントリー 1 3 4 を逆追跡すれば、サブ命令の番号フィールド 6 2 0 が減少し、受信されたエラー状態が現状態でラッチされる。サブ命令の番号フィールド 6 2 0 が「 0 」に到達すれば、サブ命令 1 3 2 に対応する命令 1 8 0 が完了し、P - e S S D 1 1 0 p は、命令応答 1 8 2 としての完了エントリーを遠隔ホスト 1 9 0 へ戻す。この時点で、P - e S S D 1 1 0 p は蓄積されたエラー状態 6 3 0 を持つ完了エントリーを生成し、それを関連する C Q に置く。

20

【 0 0 7 2 】

図 7 は、本発明の一実施例に係る P - e S S D 1 1 0 p による命令 1 8 0 の処理を説明するためのフローチャートである。

【 0 0 7 3 】

図 1 及び図 7 を参照すれば、上述したように、各 P - e S S D 1 1 0 p は、提出キュー S Q を維持する。P - e S S D 1 1 0 p によって受信され、実行のために利用可能な命令 1 8 0 が存在すれば (S 7 1 0)、P - e S S D 1 1 0 p は S Q を調整し、実行のための命令 1 8 0 を選択する。即ち、P - e S S D 1 1 0 p は、全ての N V M e 命令 (例えば、管理命令及び命令 1 8 0) を実行する、即ち、S - e S S D 1 1 0 s が遠隔ホスト 1 9 0 に直接データを伝送するが、S - e S S D 1 1 0 s は、遠隔ホスト 1 9 0 から命令を直接受信せず、遠隔ホスト 1 9 0 に直接完了エントリーを伝送しない。P - e S S D 1 1 0 p によって実行される命令 1 8 0 は、任意の S - e S S D 1 1 0 s とどのような通信も必要としないことがある。

30

【 0 0 7 4 】

命令 1 8 0 を受信した後、P - e S S D 1 1 0 p はデータがどこにあるか、そしてそれが全てのデータにアクセスできるかを判断する (S 7 2 0)。P - e S S D 1 1 0 p が全てのデータを持つ場合、P - e S S D 1 1 0 p は遠隔ホスト 1 9 0 にデータを伝送する (S 7 7 0)。

40

【 0 0 7 5 】

P - e S S D 1 1 0 p が全てのデータを持っていないと判断すれば (S 7 2 0)、P - e S S D 1 1 0 p は、M A T 3 0 0 を参照して要請されたデータがどこに位置するかを判断する (S 7 3 0)。関連した e S S D 1 1 0 セットが確認されれば、P - e S S D 1 1 0 p は、その命令 1 8 0 の実行を続ける。関連して要請された全てのデータが P - e S S D 1 1 0 p 自体で利用可能であれば、P - e S S D 1 1 0 p はデータ伝送 1 8 5 を遂行する。しかし、要請されたデータが P - e S S D 1 1 0 p 及び / 又は S - e S S D 1 1 0 s のセットに分散されていると、P - e S S D 1 1 0 p は、オリジナルの命令 1 8 0 を適切な個数のサブ命令 1 3 2 に分割する (S 7 4 0)。サブ命令 1 3 2 の個数は、要請されたデータが分散されている e S S D 1 1 0 の個数に対応する。それぞれのサブ命

50

令 1 3 2 は、それぞれの e S S D 1 1 0 が所有する要請されたデータの部分に対応する。

【 0 0 7 6 】

P - e S S D 1 1 0 p は、サブ命令 1 3 2 内に適切な開示 L B A (S L B A)、ブロック数 (N L B)、及び遠隔分散 / 収集目録 (S G L) を置く。S G L は、アドレス、キー及び遠隔ホスト 1 9 0 上の伝送バッファサイズを含む。次に、P - e S S D 1 1 0 p は、命令分割プロセッサで専用通信チャンネル 1 3 0 を通じてそれぞれの S - e S S D 1 1 0 s にサブ命令 1 3 2 を伝送 (S 7 5 0) し、それぞれの S - e S S D 1 1 0 s からサブ命令完了エントリー 1 3 4 を受信するために待機する (S 7 6 0)。従って、オリジナルの命令 1 8 0 は、適切な e S S D 1 1 0 が並列にデータ伝送を遂行できるようにサブ命令 1 3 2 に分割され、データーが遠隔ホスト 1 9 0 に伝送されるようにする (S 7 7 0)。

10

【 0 0 7 7 】

P - e S S D 1 1 0 p は、図 6 と関連する実行中の命令 1 8 0 に対する命令コンテキストを生成する。命令コンテキストは、サブ命令 1 3 2 の実行及びサブ命令 1 3 2 の中間エラー状態を追跡するために使用される (S 7 8 0)。P - e S S D 1 1 0 p が、命令 1 8 0 の完了を確認すれば、P - e S S D 1 1 0 p は遠隔ホスト 1 9 0 に完了エントリーを伝送する (S 7 9 0)。

【 0 0 7 8 】

従って、遠隔ホスト N V M e - o F ドライバ 1 7 0 によって発行された全ての命令 1 8 0 が、P - e S S D 1 1 0 p によって受信 (S 7 1 0) 及び実行される。P - e S S D 1 1 0 p は、複数の又は全ての命令 1 8 0 を単独で完了できる (例えば、P - e S S D 1 1 0 p が「S 7 2 0」ステップで単独で命令 1 8 0 を完了するのに必要な全ての情報を持っていると判断すれば)。一部の場合、P - e S S D 1 1 0 p は、命令 1 8 0 を完了する前に S - e S S D 1 1 0 s から特定情報を持ってくる。P - e S S D 1 1 0 p が S - e S S D 1 1 0 s から特定の非ユーザーデータ情報を探索すれば、P - e S S D 1 1 0 p はサブ命令 1 3 2 を生成して、それぞれの S - e S S D 1 1 0 s に伝送する。S - e S S D 1 1 0 s は、それらの間の専用通信チャンネル 1 3 0 を使用して、P - e S S D 1 1 0 p に任意の必要なデータ及びサブ命令完了エントリー 1 3 4 を戻す。

20

【 0 0 7 9 】

図 8 は、本発明の一実施例に係るネームスペース生成及び削除命令の実行を説明するためのフローチャートである。

30

【 0 0 8 0 】

図 8 を参照すれば、P - e S S D 1 1 0 p は、ネームスペース生成命令 (S 8 1 0) 及び / 又は削除命令 (S 8 1 1) を受信して実行する。P - e S S D 1 1 0 p がネームスペース生成命令を受信すれば (S 8 1 0)、P - e S S D 1 1 0 p は M A T 3 0 0 を検索し (S 8 2 0)、全体利用可能なプールから適切な容量を割り当てる (S 8 3 0)。新たに生成されたネームスペースは、P - e S S D 1 1 0 p のみにフラッシュ容量を有することもでき、または S - e S S D 1 1 0 s のうち、特定の一つのみにフラッシュ容量を有することもでき、または新たに生成されたネームスペースは、P - e S S D 1 1 0 p と S - e S S D 1 1 0 s の任意の組み合わせにフラッシュ容量を有する。以後、P - e S S D 1 1 0 p は、割り当てられた容量及び関連するマッピングされた L B A 範囲を M A T 3 0 0 に記録する (S 8 4 0)。

40

【 0 0 8 1 】

P - e S S D 1 1 0 p が、ネームスペースを削除するためのネームスペース削除命令を受信すれば (S 8 1 1)、P - e S S D 1 1 0 p は M A T 3 0 0 を検索し (S 8 2 1)、対応する e S S D を回収して (S 8 3 1)、関連容量の割り当てを解除し (S 8 4 1)、その後、M A T 3 0 0 をアップデートする。

【 0 0 8 2 】

S - e S S D 1 1 0 s によって実行されるネームスペース生成 / 削除命令と関連して、S - e S S D 1 1 0 s は、ネームスペース生成 / 削除命令を直接受信しない。一般に、

50

S - e S S D 110s は、全体容量を示す1つのネームスペースを含む。適切な場合、P - e S S D 110p は、ネームスペース生成命令又は削除命令をサブ命令132としてS - e S S D 110s に発行する。S - e S S D 110s は、それぞれこれらの命令を実行し、対応するサブ命令完了エントリー134をP - e S S D 110p に戻す。このような流れは、P - e S S D 110p から受信するサブ命令についても同様である。

【0083】

図9は、本発明の一実施例に係るP - e S S D 110pの制御下でリード/ライト命令の実行を説明するためのフローチャートである。

【0084】

図9を参照すれば、P - e S S D 110p は、リード/ライトの命令180を含む全ての命令180を受信して実行する。P - e S S D 110p は、リード/ライトの命令180を受信すれば(S910)、まずP - e S S D 110p はM A T 300を検索する(S920)。M A T 300から、P - e S S D 110p は関連するユーザーデータが位置するe S S D セットを確認する。

【0085】

図6を参照して説明したように、P - e S S D 110p は、P - e S S D 110p がサブ命令132の実行を追跡するよう、オリジナル命令に対する命令コンテキストを生成する(S930)。続いて、P - e S S D 110p は、対応するリード/ライトサブ命令を生成し(S940)、適切なサブ命令132を適切なS - e S S D 110s に伝送する(S950)。なお、P - e S S D 110p は、必要なすべての伝送ネットワーク関連情報(例えば、アドレス)をS - e S S D 110s に提供する。サブ命令132の一部として、S - e S S D 110s は、遠隔バッファアドレス、大きさ及び保安キーを含む遠隔ホスト190のS G Lを受信する。

【0086】

サブ命令132内のデータ伝送フィールドは適切に修正される。S - e S S D 110s は、遠隔ホスト190バッファに直接データを伝送し(S960)、完了すればS - e S S D 110s は、(遠隔ホスト190へ直接伝送せずに)P - e S S D 110p に完了メッセージを伝送する。必要であれば、P - e S S D 110p は、遠隔ホスト190へのデータ伝送を遂行する(S960)。

【0087】

さらに、S - e S S D 110s のそれぞれは、遠隔ホスト190に直接データ伝送を遂行するため、P - e S S D 110p から十分な情報を受信する(S960)。N V M e - o F プロトコルで、R D M A 伝送サービス(R D M A リード及びR D M A ライト)は、S - e S S D 110s から遠隔ホスト190にデータを伝送するために使用される。遠隔ホスト190は、複数のS - e S S D 110s が遠隔ホスト190にデータを伝送するようにR D M A プロトコルのS R C (S h a r e d R e c e i v e Q u e u e)機能を支援する必要がある(S960)。R D M A プロトコルは、イーサネット/IP/T C P (i W A R P)、イーサネット/I n f i n i B a n d (R o C E v 1)又はイーサネット/IP/U D P (R o C E v 2)で実行される。S - e S S D 110s と遠隔ホスト190と間の通信と関連して、S - e S S D 110s は、厳格に遠隔ホスト190とのデータ伝送だけを遂行する(S960)。つまり、遠隔ホスト190とS - e S S D 110s は、R D M A - リード及びR D M A - ライト動作だけを遂行する(S960)。サブ命令完了エントリー134及び任意の非ユーザーデータ伝送は、R D M A 伝送動作又は一部異なる独占プロトコルを使用してP - e S S D 110p により遂行される。

【0088】

全てのサブ命令が完了すれば(S970)、全てのサブ命令完了エントリーがP - e S S D 110p によって受信される時に示したように、P - e S S D 110p は、オリジナルの命令180に対する命令応答182としての完了エントリーを生成し、完了エントリーを遠隔ホスト190上の適切なC Q に伝送する(S980)。次に、P - e S S D 110p は、命令コンテキストの割り当てを解除する(S990)。

10

20

30

40

50

【 0 0 8 9 】

図 1 0 は、本発明の一実施例に係る S - e S S D によるサブ命令の実行を説明するためのフローチャートである。

【 0 0 9 0 】

図 1 0 を参照すれば、本実施例において、どのような S - e S S D 1 1 0 s も命令 1 8 0 又は任意の命令を遠隔ホスト N V M e - o F ドライバ 1 7 0 から直接受信しない。代わりに、P - e S S D 1 1 0 p は、必要時に応じてのみ S - e S S D 1 1 0 s にサブ命令 1 3 2 を伝送する (S 1 0 1 0)。S - e S S D 1 1 0 s は、任意のデータ伝送が必要であるかどうかを判断し (S 1 0 2 0)、専用通信チャンネル 1 3 0 を通じて P - e S S D 1 1 0 p と必要なデータ伝送を遂行する (S 1 0 3 0)。S - e S S D 1 1 0 s は、サブ命令完了エントリー 1 3 4 を生成して (S 1 0 4 0)、P - e S S D 1 1 0 p に伝送する (S 1 0 5 0)。他の実施例において、S - e S S D 1 1 0 s は、サブ命令完了エントリー 1 3 4 だけでなく、データを P - e S S D 1 1 0 p に伝送するために R D M A 伝送動作を使用する。

10

【 0 0 9 1 】

図 1 1 は、本発明の一実施例に係る S - e S S D 1 1 0 s によるリード/ライトサブ命令の実行を説明するためのフローチャートである。

【 0 0 9 2 】

図 1 1 を参照すれば、S - e S S D 1 1 0 s は、主にリード又はライトのサブ命令 1 3 2 (例えば、管理者サブ命令と対照的に) を処理する。すなわち、S - e S S D 1 1 0 s は、主にプロトコルプロセッシングの他の側面に侵入しなくても、遠隔ホスト 1 9 0 に、そして遠隔ホスト 1 9 0 からのデータ移動を遂行する。S - e S S D 1 1 0 s がリード/ライトのサブ命令 1 3 2 を受信すれば (S 1 1 1 0)、S - e S S D 1 1 0 s は受信伝送ネットワーク情報を利用して (S 1 1 2 0)、遠隔ホスト 1 9 0 に R D M A リード又は R D M A ライト要請を発行する (S 1 1 3 0)。サブ命令 1 3 2 の一部として、S - e S S D 1 1 0 s は、遠隔バッファアドレス/オフセット、大きさ及び保安キーの細部事項を受信する。必要なデータ伝送が完了すれば (S 1 1 4 0)、S - e S S D 1 1 0 s は、適切なエラー状態とともにサブ命令完了エントリー 1 3 4 を P - e S S D 1 1 0 p に伝送する (S 1 1 5 0)。

20

【 0 0 9 3 】

図 1 2 は、本発明の一実施例に係る S - e S S D 1 1 0 s に同期したデータ伝送及びサブ命令完了を説明するためのフローチャートである。

30

【 0 0 9 4 】

図 1 2 を参照すれば、与えられた命令 1 8 0 に対し、P - e S S D 1 1 0 p が命令応答 1 8 2 としての完了エントリーを遠隔ホスト 1 9 0 に伝送するが、S - e S S D 1 1 0 s のセットは遠隔ホスト 1 9 0 にデータ伝送を遂行する。与えられた命令 1 8 0 に対する完了エントリーが遠隔ホスト 1 9 0 に提供される前に、全てのデータは、関連したデータが定義されない行動/エラーをもたらす前に、遠隔ホスト 1 9 0 に到達する完了エントリーとして遠隔ホスト 1 9 0 に伝送されなければならない。

【 0 0 9 5 】

上述したように、P - e S S D 1 1 0 p は、遠隔ホスト 1 9 0 からリード/ライトの命令 1 8 0 を受信し (S 1 2 1 0)、命令 1 8 0 を複数のサブ命令 1 3 2 に分割する (S 1 2 2 0)。次に、P - e S S D 1 1 0 p からそれぞれのサブ命令 1 3 2 を受信した各 S - e S S D 1 1 0 s は、遠隔ホスト 1 9 0 にデータを伝送し (S 1 2 5 0)、データ伝送が完了したと判断すれば (S 1 2 6 0)、S - e S S D 1 1 0 s はサブ命令完了エントリー 1 3 4 を P - e S S D 1 1 0 p に伝送する (S 1 2 7 0)。次に、全ての関連した S - e S S D 1 1 0 s から全てのサブ命令完了エントリー 1 3 4 を受信すれば (S 1 2 3 0)、P - e S S D 1 1 0 p は、命令応答 1 8 2 としての完了エントリーを遠隔ホスト 1 9 0 に送る (S 1 2 4 0)。

40

【 0 0 9 6 】

50

統合イーサネットSSDデータが伝送され、完了エントリーポスティングがeSSD110セットを通じて分散されるため、データ伝送及び完了ポスティングの同期化が達成されなければならない。一般に、1つのeSSD110が命令実行（データ伝送＋完了ポスティング）の2つのステップを全て遂行する場合には、このような問題が発生しない。しかし、統合イーサネットSSDの場合には、そうではない。従って、P-eSSD110pは、命令180に対する命令応答182としての完了エントリーを送付する前に、それぞれのS-eSSD110sからの全てのサブ命令完了エントリー134を待たねばならない。なお、S-eSSD110sは、全てのデータ伝送が、サブ命令132のサブ命令完了エントリー134をP-eSSD110pに送る前に、完全で信頼性があるように終了されることを保障しなければならない。このような2つのステップの同期化プロセスは、統合イーサネットSSDシステムにおいて常にNVMe-oFプロトコルが完全になされることを保障する。

10

【0097】

上述によれば、ただ一つの1次eSSDだけをホストが見、すべての関連したサブ命令の完了を2次eSSDで追跡しながらNVMe-oFプロトコル処理を全て遂行するため、eSSDの集積グループは一つの大きな論理容量としてホストから見える。

【符号の説明】

【0098】

- 110 イーサネットSSD (eSSD)
- 110p P-eSSD
- 110s S-eSSD
- 120 シャーシ
- 130 専用通信チャンネル
- 132 サブ命令
- 134 サブ命令完了エントリー
- 135 制御プレーン
- 140 PCIeスイッチ
- 150 ボード管理コントローラ装置 (BMC)
- 160 イーサネットスイッチ
- 170 遠隔ホストNVMe-oFドライバ
- 180 命令
- 182 命令応答
- 185 データ伝送
- 190 遠隔ホスト
- 230 ラック
- 240 TORスイッチ
- 300 MAT
- 311 容量
- 312 アドレス位置
- 313 eSSDインデックス
- 314 ネームスペースLBA範囲
- 530 イーサネット通信チャンネル
- 610 命令タグ
- 620 サブ命令の番号フィールド
- 630 蓄積されたエラー状態
- 640 命令ID

20

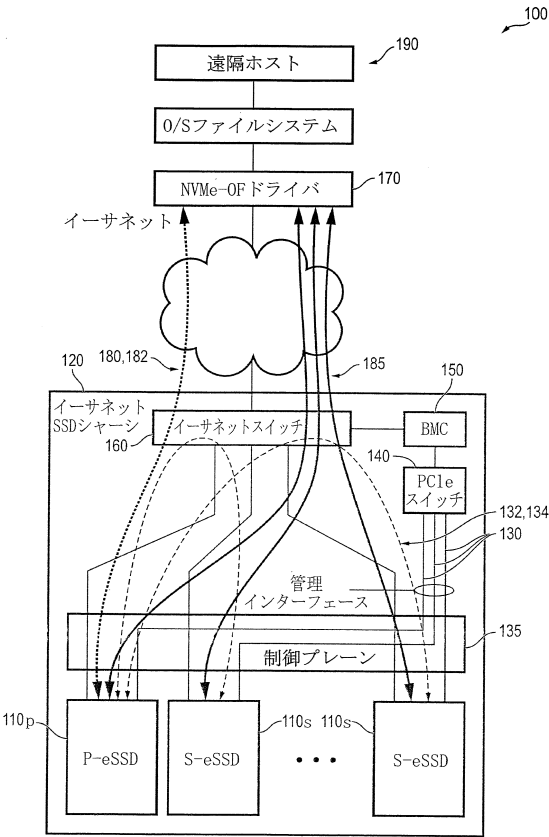
30

40

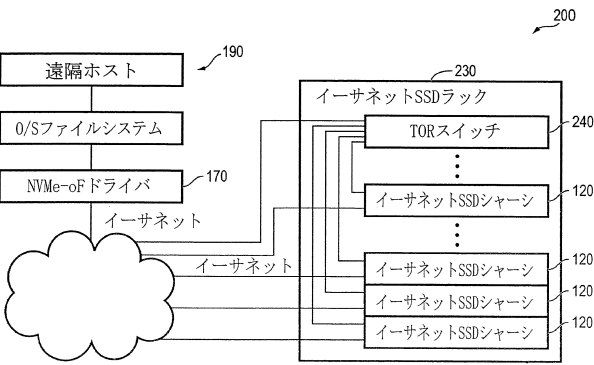
50

【図面】

【図 1】



【図 2】

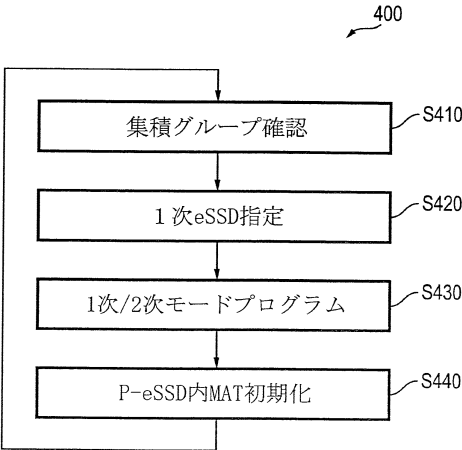


【図 3】

300

| 313 eSSD インデックス | 311 容量 (TB) | 312 伝送アドレス (MAC/ID) | 314 マッピングアドレス NS_LBA範囲 | フリール容量 (TB) | その他 |
|-----------------------|-------------------|---------------------------|------------------------------|----------------|-----|
| 0 | 16 | | | | |
| 1 | 8 | | | | |
| 2 | 32 | | | | |
| 3 | 16 | | | | |

【図 4】



10

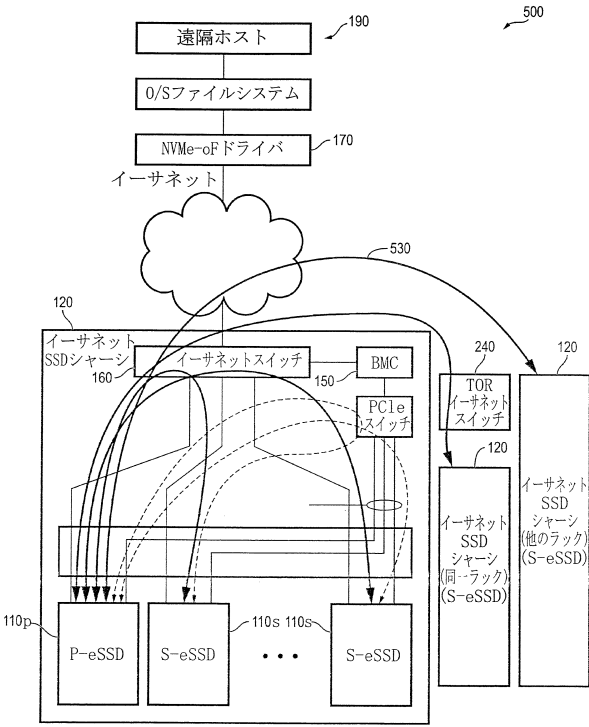
20

30

40

50

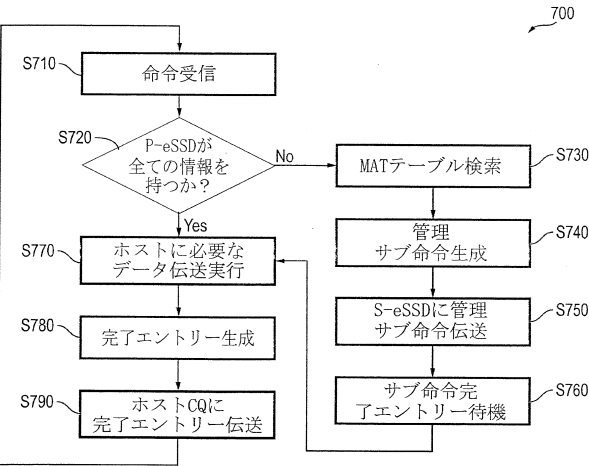
【図 5】



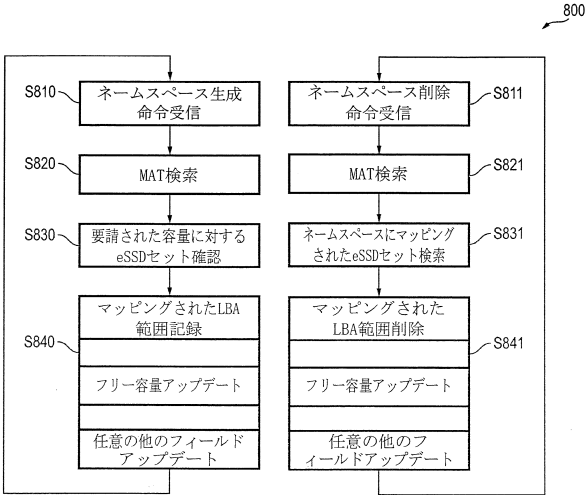
【図 6】

| 命令ID | 命令Tag | サブ命令の順序 | 蓄積されたエラー状態 |
|------|-------|---------|------------|
| 123 | | | |
| 15 | | | |
| 39 | | | |
| 3 | | | |

【図 7】



【図 8】



10

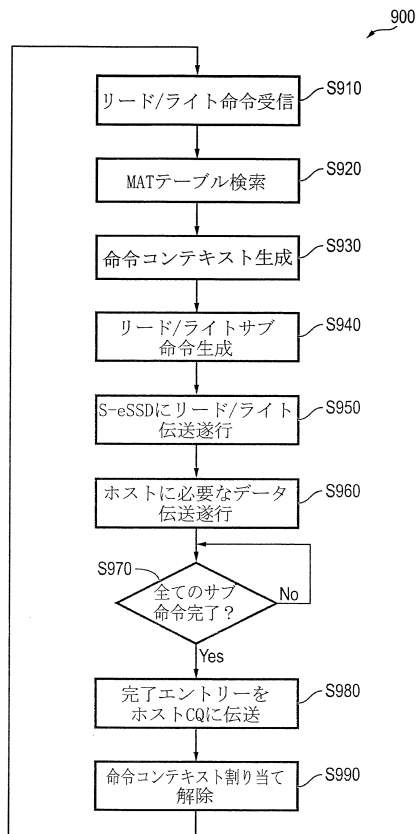
20

30

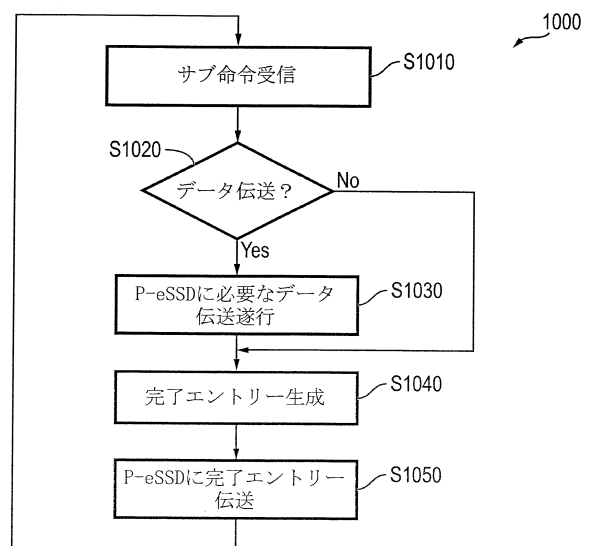
40

50

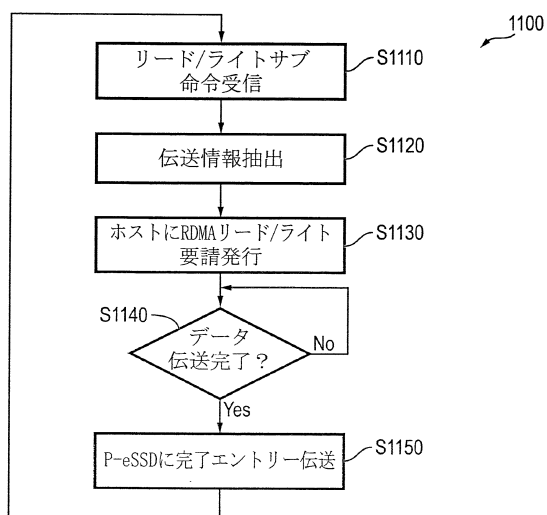
【図 9】



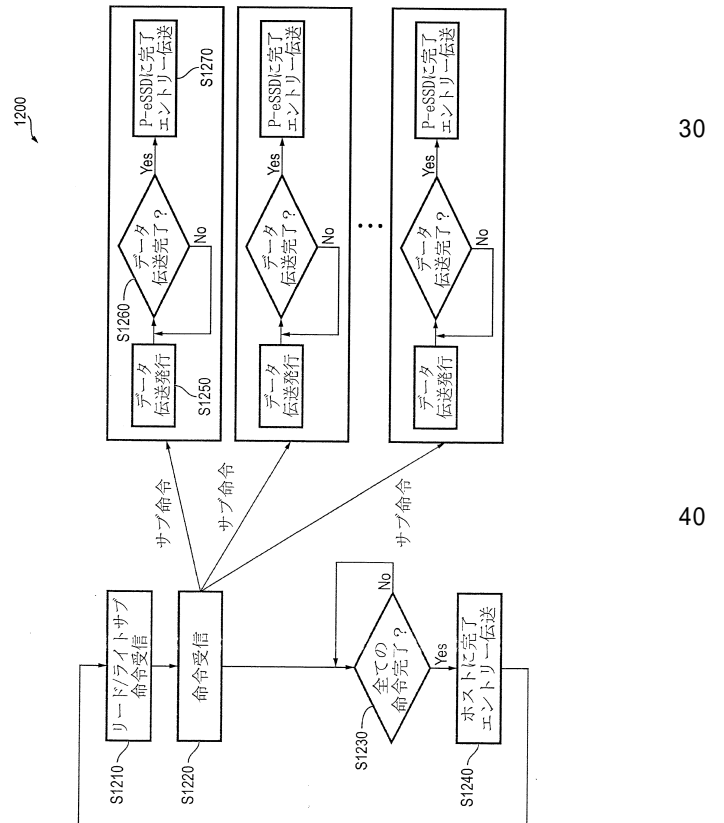
【図 10】



【図 11】



【図 12】



10

20

30

40

50

フロントページの続き

(54)【発明の名称】 ループ

(32)優先日 平成29年6月8日(2017.6.8)

(33)優先権主張国・地域又は機関
米国(US)

早期審査対象出願

前置審査

7 6 6 5

(72)発明者 ソンボン ポール オラリグ

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

(72)発明者 フレッド ウォーリー

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

審査官 打出 義尚

(56)参考文献 米国特許出願公開第 2 0 1 5 / 0 0 1 2 6 0 7 (U S , A 1)

米国特許出願公開第 2 0 1 5 / 0 2 6 1 7 2 0 (U S , A 1)

中国特許出願公開第 1 0 5 9 1 2 2 7 5 (C N , A)

(58)調査した分野 (Int.Cl. , D B 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