

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6975095号
(P6975095)

(45) 発行日 令和3年12月1日(2021.12.1)

(24) 登録日 令和3年11月9日(2021.11.9)

(51) Int.Cl.		F I			
G06F 3/06	(2006.01)	G06F 3/06	301N		
G06F 3/08	(2006.01)	G06F 3/08	H		
G06F 13/10	(2006.01)	G06F 13/10	340A		
G06N 20/00	(2019.01)	G06N 20/00			

請求項の数 19 (全 17 頁)

(21) 出願番号	特願2018-88873 (P2018-88873)	(73) 特許権者	390019839
(22) 出願日	平成30年5月2日(2018.5.2)		三星電子株式会社
(65) 公開番号	特開2018-198054 (P2018-198054A)		Samsung Electronics Co., Ltd.
(43) 公開日	平成30年12月13日(2018.12.13)		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4月26日(2021.4.26)		129, Samsung-ro, Yeongtong-gu, Suwon-si, Gyeonggi-do, Republic of Korea
(31) 優先権主張番号	62/510, 728	(74) 代理人	110000051
(32) 優先日	平成29年5月24日(2017.5.24)		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 ピー.
(31) 優先権主張番号	15/672, 223		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8月8日(2017.8.8)		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】 マシンラーニングを実行するデータストレージ及び処理システムとその動作方法

(57) 【特許請求の範囲】

【請求項 1】

ホストサーバ及びストレージ部を備えるデータストレージ及び処理システムであって、
前記ストレージ部は、ドライブメモリー及びドライブプロセッサを含むドライブと、
前記ホストサーバと前記ドライブメモリーとの間でデータを伝送し受信するために前記
ホストサーバを前記ドライブに連結させるように構成された外部スイッチと、

グラフィック処理装置と、を含み、

前記ドライブプロセッサは、処理命令及びデータを前記ドライブメモリーから前記グラフィック処理装置に伝送するように構成され、

前記グラフィック処理装置は、前記処理命令に従ってデータを処理して結果データを生成するように構成され、

前記ドライブは、前記グラフィック処理装置が変換を遂行可能なことを示すテーブルを含み、

前記ドライブプロセッサは、前記ドライブメモリーに格納されたデータチャンク及び前記データチャンクに適用される前記変換を識別し、前記データチャンクを前記グラフィック処理装置に伝送するように構成されることを特徴とするデータストレージ及び処理システム。

【請求項 2】

前記グラフィック処理装置は、U・2コネクタを含み、前記U・2コネクタを介して前記ドライブに連結され、前記U・2コネクタを介して前記処理命令及びデータを受

10

20

信することを特徴とする請求項 1 に記載のデータストレージ及び処理システム。

【請求項 3】

前記グラフィック処理装置は、前記結果データを前記ドライブプロセッサに伝送するように構成され、

前記ドライブプロセッサは、前記ドライブメモリに前記結果データを格納するように構成されることを特徴とする請求項 1 に記載のデータストレージ及び処理システム。

【請求項 4】

前記グラフィック処理装置は、前記外部スイッチに連結され、前記外部スイッチを利用して前記結果データを前記ホストサーバに伝送するように構成されることを特徴とする請求項 1 に記載のデータストレージ及び処理システム。

10

【請求項 5】

前記グラフィック処理装置は、前記処理命令に従ってデータの処理を完了した後、状態メッセージを前記ドライブプロセッサに伝送するように構成されることを特徴とする請求項 4 に記載のデータストレージ及び処理システム。

【請求項 6】

前記ドライブプロセッサは、前記テーブルから前記グラフィック処理装置のアドレスを取得し、前記データチャンクを前記グラフィック処理装置のアドレスに伝送するように構成されることを特徴とする請求項 1 に記載のデータストレージ及び処理システム。

【請求項 7】

前記ストレージ部は、前記グラフィック処理装置の能力を判別して前記グラフィック処理装置の能力に基づいて前記テーブルをアップデートするように構成されたベースボード管理制御器、を更に含むことを特徴とする請求項 1 に記載のデータストレージ及び処理システム。

20

【請求項 8】

前記ストレージ部は、各ドライブが別途のテーブルを有する複数のドライブと、複数のグラフィック処理装置と、を含み、

前記ベースボード管理制御器は、前記複数のドライブの各ドライブの各テーブルをアップデートするように構成されることを特徴とする請求項 7 に記載のデータストレージ及び処理システム。

【請求項 9】

前記ベースボード管理制御器及び前記グラフィック処理装置は、NVMe - MI (Non Volatile Memory express - Management Interface) のプロトコルを用いて通信するように構成され、

前記ベースボード管理制御器は、前記 NVMe - MI の識別コマンドを利用することで、前記グラフィック処理装置の能力を判別することを特徴とする請求項 7 に記載のデータストレージ及び処理システム。

30

【請求項 10】

前記ストレージ部は、複数のグラフィック処理装置を含み、

前記ベースボード管理制御器は、前記複数のグラフィック処理装置の各グラフィック処理装置に対する負荷を判別し、前記複数のグラフィック処理装置の各グラフィック処理装置に対する負荷に基づいて前記テーブルをアップデートするように構成されることを特徴とする請求項 7 に記載のデータストレージ及び処理システム。

40

【請求項 11】

前記外部スイッチは、イーサネット（登録商標）スイッチであり、

前記ドライブは、イーサネット（登録商標）ソリッドステートドライブであることを特徴とする請求項 1 に記載のデータストレージ及び処理システム。

【請求項 12】

前記ホストサーバは、遠隔直接アクセスストレージプロトコルを用いて前記ストレージ部と通信することを特徴とする請求項 1 に記載のデータストレージ及び処理システム。

【請求項 13】

50

前記ストレージ部は、シャシー及びベースボード管理制御器を更に含み、
前記グラフィック処理装置は、現場交換可能な装置であり、
前記ベースボード管理制御器は、前記グラフィック処理装置が前記シャシーに挿入されることに対応してプラグイン (plug-in) イベントを感知するように構成されることを特徴とする請求項 1 に記載のデータストレージ及び処理システム。

【請求項 14】

ドライブプロセッサ及びドライブメモリーを含むドライブとグラフィック処理装置とを備えるデータストレージ及び処理システムの動作方法であって、
前記ドライブプロセッサが、ホストサーバからデータを受信するステップと、
前記データを前記ドライブメモリーに格納するステップと、
トリガーを感知するステップと、
前記トリガーに応答し、前記ドライブメモリーに格納されたデータチャンク及び前記データチャンクに適用する変換を識別するアルゴリズムを実行するステップと、
前記変換を遂行する装置に対応するアドレスを識別するために能力テーブルを検索するステップと、
前記データチャンク及び前記データチャンクを処理する処理命令を前記グラフィック処理装置のアドレスに伝送するステップと、を含むことを特徴とする方法。

【請求項 15】

前記グラフィック処理装置が、前記データチャンク及び前記処理命令を受信するステップと、
前記処理命令に従って前記データチャンクを処理して結果データを生成するステップと、
を更に含むことを特徴とする請求項 14 に記載の方法。

【請求項 16】

前記グラフィック処理装置が、前記結果データを前記ドライブプロセッサに伝送するステップと、
前記ドライブプロセッサが、前記結果データを前記ドライブメモリーに格納するステップと、
を更に含むことを特徴とする請求項 15 に記載の方法。

【請求項 17】

ベースボード管理制御器が、前記グラフィック処理装置の能力を感知するステップと、
前記グラフィック処理装置の能力に基づいて前記能力テーブルをアップデートするステップと、
を更に含むことを特徴とする請求項 14 に記載の方法。

【請求項 18】

前記データストレージ及び処理システムは、各ドライブが能力テーブルを含む複数のドライブを含み、
前記ベースボード管理制御器が、前記グラフィック処理装置の能力に基づいて前記複数のドライブの各ドライブの各能力テーブルをアップデートするステップを含むことを特徴とする請求項 17 に記載の方法。

【請求項 19】

ドライブプロセッサ及びドライブメモリーを含むドライブとグラフィック処理装置とを備えるデータストレージ及び処理システムであって、
ホストサーバからデータを受信する手段と、
前記データを前記ドライブメモリーに格納する手段と、
トリガーを感知する手段と、
前記トリガーに応答し、前記ドライブメモリーに格納されたデータチャンク及び前記データチャンクに適用する変換を識別するアルゴリズムを実行する手段と、
前記変換を遂行する装置に対応するアドレスを識別するために能力テーブルを検索する手段と、
前記データチャンク及び前記データチャンクを処理する処理命令を前記グラフィック処理装置のアドレスに伝送する手段と、
前記処理命令に従って前記データチャンクを処理して結果データを生成し、前記結果デ

10

20

30

40

50

ータを前記ホストサーバに伝送する手段と、

現場交換可能なグラフィック処理装置のプラグインイベントを感知し、前記現場交換可能なグラフィック処理装置の能力を判別し、前記現場交換可能なグラフィック処理装置の能力に基づき前記プラグインイベントの感知に応答して前記能力テーブルをアップデートする手段と、を含むことを特徴とするデータストレージ及び処理システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、マシンラーニングのためのシステムに関し、より詳しくは、SSDフォームファクターで内蔵されたGPUを有するNVMe-oF SSDシャシーに実装されたストレージシステムにおいて、複数のグラフィック処理装置を利用してマシンラーニングアルゴリズムを実行するデータストレージ及び処理システムとその動作方法に関する。

10

【背景技術】

【0002】

収集及び格納されるデータ量が増加している。モノのインターネット(Internet of Things)からソーシャルネットワーク、デジタルヘルスの領域に至るまで数多くのアプリケーションは膨大な量のデータを生成する。このようなデータはデータセンターのようなデータシステムでホスト(host)される。データは、しばしば格納、処理、及び分析が要求される。マシンラーニングアルゴリズムのようなアルゴリズムはデータから特徴又は有用な情報を抽出するためにビッグデータ(big data)セットに適用される。このデータを、迅速、効率的、経済的、弾力的に、格納、処理、及び分析するための技法及びアーキテクチャが必要である。

20

【0003】

背景技術で開示する上記の情報は、単に本発明の背景に対する理解を増進させるためのものであり、従って通常の技術者にとって既に知られた先行技術ではない情報を含む。

【先行技術文献】

【特許文献】

【0004】

特許文献1：米国特許第8819335号明細書

特許文献2：米国特許第9317204号明細書

特許文献3：米国特許第9430412号明細書

特許文献4：米国特許第9483431号明細書

特許文献5：米国特許出願公開第2017/0210076号明細書

特許文献6：米国特許出願公開第2017/0010995号明細書

特許文献7：米国特許出願公開第2017/0019312号明細書

特許文献8：米国特許出願公開第2017/0060442号明細書

30

【発明の概要】

【発明が解決しようとする課題】

【0005】

本発明は、上記従来技術に鑑みてなされたものであって、本発明の目的は、ストレージシステムでマシンラーニングアルゴリズムを実行するデータストレージ及び処理システムとその動作方法を提供することにある。

40

【課題を解決するための手段】

【0006】

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とデータストレージ及び処理システムは、ホストサーバ及びストレージ部を備え、前記ストレージ部は、ドライブメモリ及びドライブプロセッサを含むドライブと、前記ホストサーバと前記ドライブのメモリとの間でデータを伝送し受信するために前記ホストサーバを前記ドライブに連結させる外部スイッチと、グラフィック処理装置と、を含み、前記ドライブプロセッサは、処理命令及びデータを前記ドライブメモリから前記グラフィック処理装置に伝送し、前記グ

50

ラフィック処理装置は、前記処理命令に従ってデータを処理して結果データを生成する。
【0007】

前記グラフィック処理装置は、U . 2コネクタを含み、前記U . 2コネクタを介して前記ドライブに連結され、前記U . 2コネクタを介して前記処理命令及びデータを受信し得る。

前記グラフィック処理装置は、前記結果データを前記ドライブプロセッサに伝送し、前記ドライブプロセッサは、前記ドライブメモリーに前記結果データを格納し得る。

前記グラフィック処理装置は、前記外部スイッチに連結され、前記外部スイッチを利用して前記結果データを前記ホストサーバに伝送し得る。

前記グラフィック処理装置は、前記処理命令に従ってデータの処理を完了した後、状態メッセージを前記ドライブプロセッサに伝送し得る。

前記ホストサーバは、トリガー命令を前記ドライブプロセッサに伝送し、前記ドライブプロセッサは、前記トリガー命令に応答し、実行時に前記処理命令及びデータを前記グラフィック処理装置に伝送する過程を含むデータ処理アルゴリズムを実行し得る。

前記ドライブは、能力テーブルを含み、前記ドライブプロセッサは、前記ドライブメモリーに格納されたデータチャンク及び前記データチャンクに適用される変換を識別し、前記グラフィック処理装置が前記変換を遂行可能なことを示す前記能力テーブルから前記グラフィック処理装置のアドレスを検索し、前記データチャンクを前記グラフィック処理装置のアドレスに伝送し得る。

前記ストレージ部は、前記グラフィック処理装置の能力を判別し、前記グラフィック処理装置の能力に基づいて前記能力テーブルをアップデートするベースボード管理制御器 (base board management controller) を、更に含み得る。

前記ストレージ部は、各ドライブが別途の能力テーブルを有する複数のドライブと、複数のグラフィック処理装置と、を含み、前記ベースボード管理制御器は、前記複数のドライブの各ドライブの各能力テーブルをアップデートし得る。

前記ベースボード管理制御器及び前記グラフィック処理装置は、NVMe - MI (Non Volatile Memory express - Management Interface) のプロトコルを用いて通信し、前記ベースボード管理制御器は、前記NVMe - MIの識別 (Identify) コマンドを利用することで、前記グラフィック処理装置の能力を判別し得る。

前記ストレージ部は、複数のグラフィック処理装置を含み、前記ベースボード管理制御器は、前記複数のグラフィック処理装置の各グラフィック処理装置に対する負荷を判別し、前記複数のグラフィック処理装置の各グラフィック処理装置に対する負荷に基づいて前記能力テーブルをアップデートし得る。

前記外部スイッチは、イーサネット (登録商標) スイッチであり、前記ドライブは、イーサネット (登録商標) ソリッドステートドライブ (Ethernet (登録商標) solid state drive) であり得る。

前記ホストサーバは、遠隔直接アクセスストレージ (remote direct access storage) プロトコルを用いて前記ストレージ部と通信し得る。

前記ストレージ部は、シャシー (chassis) 及びベースボード管理制御器を更に含み、前記グラフィック処理装置は、現場交換可能な装置であり、前記ベースボード管理制御器は、前記グラフィック処理装置が前記シャシーに挿入されることに対応してプラグイン (plug - in) イベントを感知し得る。

【0008】

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ドライブプロセッサ及びドライブメモリーを含むドライブとグラフィック処理装置とを備えるデータストレージ及び処理システムの動作方法は、前記ドライブプロセッサが、ホストサーバからデータを受信するステップと、前記データを前記ドライブメモリーに格納するステップと、トリガーを感知するステップと、前記トリガーに応答し、前記ドライブメモリーに格納されたデータチ

10

20

30

40

50

ヤンク及び前記データチャンクに適用する変換を識別するアルゴリズムを実行するステップと、前記変換を遂行する装置に対応するアドレスを識別するために能力テーブルを検索するステップと、前記データチャンク及び前記データチャンクを処理する処理命令を前記グラフィック処理装置のアドレスに伝送するステップと、を含むことを特徴とする。

【0009】

前記方法は、前記グラフィック処理装置が、前記データチャンク及び前記処理命令を受信するステップと、前記処理命令に従って前記データチャンクを処理して結果データを生成するステップと、を更に含み得る。

前記方法は、前記グラフィック処理装置が、前記結果データを前記ドライブプロセッサに伝送するステップと、前記ドライブプロセッサが、前記結果データを前記ドライブメモリーに格納するステップと、を更に含み得る。

10

前記データストレージ及び処理システムは、ベースボード管理制御器を更に備え、前記方法は、前記ベースボード管理制御器が、前記グラフィック処理装置の能力を感知するステップと、前記グラフィック処理装置の能力に基づいて前記能力テーブルをアップデートするステップと、を更に含み得る。

前記データストレージ及び処理システムは、各ドライブが能力テーブルを含む複数のドライブを含み、前記方法は、前記ベースボード管理制御器が、前記グラフィック装置の能力に基づいて前記複数のドライブの各ドライブの各能力テーブルをアップデートするステップと、を含み得る。

【0010】

20

上記目的を達成するためになされた本発明の他の態様によるデータストレージ及び処理システムは、ドライブプロセッサ及びドライブメモリーを含むドライブとグラフィック処理装置とを備え、ホストサーバからデータを受信する手段と、前記データを前記ドライブメモリーに格納する手段と、トリガーを感知する手段と、前記トリガーにตอบสนองし、前記ドライブメモリーに格納されたデータチャンク及び前記データチャンクに適用する変換を識別するアルゴリズムを実行する手段と、前記変換を遂行する装置に対応するアドレスを識別するために能力テーブルを検索する手段と、前記データチャンク及び前記データチャンクを処理する処理命令を前記グラフィック処理装置のアドレスに伝送する手段と、前記処理命令に従って前記データチャンクを処理して結果データを生成し、前記結果データを前記ホストサーバに伝送する手段と、現場交換可能なグラフィック処理装置のプラグインイベントを感知し、前記現場交換可能なグラフィック処理装置の能力を判別し、前記現場交換可能なグラフィック処理装置の能力に基づき前記プラグインイベントの感知にตอบสนองして前記能力テーブルをアップデートする手段と、を含むことを特徴とする。

30

【発明の効果】

【0011】

本発明によると、データストレージ及び処理システムは、複数のグラフィック処理装置を利用してマシンラーニングアルゴリズムを効率的に実行することができ、複数のグラフィック処理装置の負荷の均衡をとることができる。

【図面の簡単な説明】

40

【0012】

【図1】図1は、関連技術によるストレージシステムのブロック図である。

【図2】図2は、本発明の一実施形態によるストレージシステムのブロック図である。

【図3】図3は、本発明の一実施形態によるグラフィック処理装置のブロック図である。

【図4】図4は、本発明の一実施形態によるグラフィック処理装置と通信するドライブに対する手順を示すフローチャートである。

【図5】図5は、本発明の一実施形態による能力テーブルの一例を示す図である。

【図6】図6は、本発明の一実施形態によるグラフィック処理装置を管理するベースボード管理制御器に対する手順を示すフローチャートである。

【発明を実施するための形態】

50

【 0 0 1 3 】

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。

【 0 0 1 4 】

次の詳細な説明では、本発明の特定の実施形態のみを、図面を通じて示して説明する。当業者（通常の技術者）が認識するように、本発明は多様な形態で具現され、本明細書に記載する実施形態に制限されるものと解釈すべきではない。各実施形態に有る特徴又は態様の説明は、典型的に他の実施形態で他の類似した特徴及び態様に利用されるものとして考慮されなければならない。本明細書で同一の参照番号は同一の要素を示す。

【 0 0 1 5 】

図 1 は、関連技術によるストレージシステムのブロック図である。ホストサーバ 1 1 0 は、ネットワーク 1 2 0 を通じて一つ以上のストレージ部 1 3 0 に連結される。ホストサーバ 1 1 0 は、一つ以上のアプリケーション 1 1 2 とオペレーティングシステム（OS）及びファイルシステム 1 1 4 を実行する。ホストサーバ 1 1 0 は、一つ以上のストレージ部と相互作用するために利用される遠隔直接アクセスストレージ（rDAS：remote Direct Access Storage）ドライバを含む。

【 0 0 1 6 】

ストレージ部 1 3 0 は、シャシー（chassis）1 3 1、イーサネット（登録商標）（Ethernet（登録商標）：以下、「イーサネット」の（登録商標）の記載は省略する。）スイッチ 1 3 2、ベースボード管理制御器（BMC：baseboard management controller）1 3 4、PCIe スイッチ 1 3 6、及び複数のイーサネットソリッドステートドライブ（eSSD：ethernet Solid State Drives）（1 4 2 A ~ 1 4 2 C）を含む。PCIe スイッチ 1 3 6 はシャシーミッドプレーン（midplane）1 4 0 を通じて BMC 1 3 4 を eSSD（1 4 2 A ~ 1 4 2 C）に連結する。BMC 1 3 4 は eSSD（1 4 2 A ~ 1 4 2 C）を管理するために PCIe スイッチ 1 3 6 を利用する。イーサネットスイッチ 1 3 2 は eSSD（1 4 2 A ~ 1 4 2 C）をネットワーク 1 2 0 に連結する。ホストサーバ 1 1 0 は、イーサネットスイッチ 1 3 2 を通じてデータを eSSD（1 4 2 A ~ 1 4 2 C）に送信し、eSSD（1 4 2 A ~ 1 4 2 C）からデータを受信するために rDAS ドライバを利用する。

【 0 0 1 7 】

マシンラーニングのようなプロセスは、大量のデータが処理されるように要求する。計算リソース（例えば、プロセッサ）はデータに対してアルゴリズム（例えば、マシンラーニングアルゴリズム）を実行する。ストレージ部 1 3 0 に格納されたデータにこのような処理を遂行する際、ホストサーバ 1 1 0 はストレージ部 1 3 0 から処理されるデータを要求しなければならない。ストレージ部 1 3 0 はネットワーク 1 2 0 を通じてデータをホストサーバ 1 1 0 に伝送しなければならない。ホストサーバ 1 1 0 は、例えばマシンラーニングアルゴリズムを実行することにより、ホストサーバ 1 1 0 に位置する CPU 又 GPU を利用してデータを処理し、格納のためにストレージ部 1 3 0 に処理結果を再び伝送する必要がある。

【 0 0 1 8 】

処理のためにストレージ部 1 3 0 からホストサーバ 1 1 0 にデータを移動することは、相当な電氣的エネルギーを消費し、ネットワーク 1 2 0 を通じてデータ及び結果を前後に伝播することに関連する処理遅延を伴う。アーキテクチャは、ネットワーク 1 2 0 を通じて、データ及び結果を伝送するための適切な帯域幅とデータを処理するためのホストサーバ 1 1 0 における処理リソース及びシステムメモリーとを要求するため、アーキテクチャの費用が高くなる。更に、ホストサーバ 1 1 0 でデータ処理における処理リソースの速度はデータセットの処理にボトルネック現象を引き起こし、処理に用いられる処理リソースは追加、除去、又は交換が容易でない。

【 0 0 1 9 】

図 2 は、本発明の一実施形態によるストレージシステムのブロックである。図 2 を参照

10

20

30

40

50

すると、ホストサーバ 210 はネットワーク 220 を通じて一つ以上のストレージ部 230 に連結される。ホストサーバ 210 はアプリケーション 212 とオペレーティングシステム (OS) 及びファイルシステム 214 を実行する。また、ホストサーバ 210 は一つ以上のストレージ部 230 と相互に作用するために利用されるストレージドライバ 216 を含む。

【0020】

ストレージ部 230 は、シャーシ 231、外部スイッチ 232 (例えば、イーサネットスイッチ)、ベースボード管理制御器 (BMC) 234、内部スイッチ 236 (例えば、PCIe スイッチ)、一つ以上のドライブ (242A ~ 242B、単に 242 とも称する。)、及び一つ以上の U.2 グラフィック処理装置 (250A ~ 250B、単に 250 とも称する。)(U.2 GPU)を含む。U.2 GPU の用語を本明細書に亘って使用するが、GPU は U.2 連結以外の連結で作動し、このような連結は本発明の範囲に含まれる。この用語は明確性のため単純に使用される。内部スイッチ 236 は、BMC 234、ドライブ (242A ~ 242B)、及び U.2 GPU (250A ~ 250B) を、シャーシミッドプレーン 240 を通じて連結する。外部スイッチ 232 は、ドライブ (242A ~ 242B)、U.2 GPU (250A ~ 250B)、BMC 234、及びネットワーク 220 を連結する。本実施形態によると、シャーシ 231 はドライブ及び/又は U.2 GPU を受容するための複数のスロットを含む。

【0021】

例えば、いくつかの実施形態において、ストレージドライバ 216 は NVMe - オーバーファブリック (NVMe - oF : NVMe - over Fabric s) ドライバのような遠隔直接アクセスストレージ (rDAS) ドライバであり、ネットワーク 220 はイーサネットネットワークであり、外部スイッチ 232 はイーサネットスイッチであり、ドライブはイーサネットソリッドステートドライブ (eSSD) である。ホストサーバ 210 は、イーサネットネットワークを通じて、一つ以上のストレージ部 230 内の eSSD とデータ通信するために rDAS ドライバを利用する。いくつかの実施形態において、内部スイッチ 236 は PCIe スイッチである。

【0022】

図 3 は、本発明の一実施形態によるグラフィック処理装置 (U.2 GPU 350) のブロック図である。いくつかの実施形態において、図 2 の U.2 GPU (250A ~ 250B) は図 3 の U.2 GPU 350 で具現される。図 3 を参照すると、U.2 GPU 350 は、プロセッサ 360、DRAM 362、不揮発性メモリ (NVM) 364、電源供給器 / 電源調節器 366、及びコネクタ 352 を含む。コネクタ 352 はプロセッサ 360 を連結するためのインターフェースを提供する。コネクタ 352 は PCIe スイッチのような内部スイッチ及びイーサネットスイッチのような外部スイッチとインターフェースするためのインターフェースを提供する。例えば、コネクタは、イーサネットインターフェース 354、PCIe インターフェース 356、及びシステム管理バス (SMBus) インターフェース 358 を提供する。一実施形態において、コネクタ 352 は U.2 コネクタ / SFF - 8639 コネクタである。

【0023】

U.2 GPU 350 は、ドライブ (例えば、eSSD) と並んでストレージ部のシャーシに連結されてストレージ部の BMC と通信するように構成される。いくつかの実施形態において、コネクタ 352 は、eSSD のようなドライブと同様に、シャーシで同一のスロットに連結されるように構成される。いくつかの実施形態において、U.2 GPU 350 は、シャーシに連結された場合、イーサネットスイッチのような外部スイッチと通信する。いくつかの実施形態において、U.2 GPU 350 は、シャーシに挿入されるか又はシャーシから除去され、ストレージ部の動作に対して自動的に収容又は除去される現場交換装置として具現される。例示は、以下の図 6 で説明される。この方法において、U.2 GPU 350 の形態の処理リソースは、効率的に、ストレージ部に追加されるか、ストレージ部から除去されるか、又はストレージ部内で交換される。また、U.2 G

10

20

30

40

50

PU350は、与えられたスロットで利用可能な最大電力（コネクタが支援する最大電力）の提供を受けて、処理動作のための電力量を利用する。反面、スロット内のドライブは格納動作のために同一の電力量を利用する。例えば、一実施形態において、コネクタはU.2コネクタであり、U.2 GPUには25Wの電力が提供される。

【0024】

図2を再び参照すると、ドライブ(242A~242B)はプロセッサ244及びフラッシュメモリ又は他の不揮発性メモリのようなメモリ246を含む。いくつかの実施形態において、ドライブ(242A~242B)はプロセッサ244によって利用されるダイナミックランダムアクセスメモリ(DRAM)を含む。一つ以上のドライブ(242A~242B)において、プロセッサ244は、対応するドライブのメモリ246に格納されたデータに対し、マシンラーニングアルゴリズムのようなデータ処理アルゴリズムを実施する命令を実行する。このような命令を実行する過程として、プロセッサ244は、データを処理する命令だけでなく、メモリ246からのデータをU.2 GPU250に伝送する。いくつかの実施形態において、データ及び命令は、内部スイッチ236を経由してシャシーミッドプレーン240を通じ、ドライブ242とU.2 GPU250との間に伝送される。いくつかの実施形態において、データ及び命令は外部スイッチ232を通じてドライブ242とU.2 GPU250との間に伝送される。ドライブ(242A~242B)及びU.2 GPU(250A~250B)は両方共にU.2コネクタのようなコネクタを利用する。

【0025】

図4は、本発明の一実施形態によるグラフィック処理装置(U.2 GPU250)と通信するドライブ242(即ち、図2のドライブ242)に対する手順を示すフローチャートである。図4の手順はドライブ242に内蔵されたプロセッサ244により実行される。610ステップで、プロセッサ244はトリガー(trigger)を感知する。トリガーは、プロセッサ244がマシンラーニングアルゴリズムのようなアルゴリズムによりドライブ242のメモリ246に含まれるデータを処理しなければならないことを示す。例えば、いくつかの実施形態において、トリガーは特定の時間の経路である。即ち、アルゴリズムはバックグラウンドで周期的に自動的に実行されるようにトリガーされる。いくつかの実施形態において、トリガーはホストサーバ210又は他の外部のソースから受信された命令又はクエリー(query)である。いくつかの実施形態において、トリガーは、ドライブ242により受信されるデータの新しいブロック、又はホスト210から受信される他のデータの作業(例えば、リード(read)、ライト(write)、削除又はアップデート)である。いくつかの実施形態において、トリガーは、追加分析、処理、又は他の種類の処理を必要とするマシンラーニングアルゴリズムの結果である。いくつかの実施形態において、トリガーは、は臨界値を超過して格納されたデータ量のような、ドライブのいくつかの内部状態又は動作により生成される。いくつかの実施形態において、トリガーは他のドライブ242から受信され、ドライブ242は他のドライブ242にトリガーを伝達する。

【0026】

いくつかの実施形態において、アルゴリズムはドライブ242のメモリ246又はその外の位置(例えば、別途のDRAM)に格納され、トリガーは単純に格納されたアルゴリズムが実行されるべきであることを示す。いくつかの実施形態において、トリガーは実行されなければならないアルゴリズムを含む。いくつかの実施形態において、ドライブ242は格納された複数のアルゴリズムを有し、トリガーは実行するアルゴリズムを識別するタグ(tag)を含む。620ステップで、プロセッサ244はアルゴリズムの実行を開始する。いくつかの実施形態において、アルゴリズムはマシンラーニングアルゴリズムである。622ステップで、プロセッサ244は、例えばアルゴリズムにおける命令又は現在のプロセッサ244の利用に基づいて、メモリ246上のデータの一つ以上のデータチャンクで一つ以上の変換、関数、又は他の形態の処理が遂行されなければならないことを識別する。

【 0 0 2 7 】

6 2 4ステップで、プロセッサ2 4 2は、必要な変換、関数、又はその他の処理作業を遂行するU . 2 GPU 2 5 0のアドレスを識別するための能力テーブル (C A P t a b l e : c a p a b i l i t y t a b l e) を調査する。いつかの実施形態において、能力テーブルはドライブ2 4 2に格納される。図5は、本発明の一実施形態による能力テーブル5 0 0の一例を示す図である。能力テーブル5 0 0は、プロセッサ2 4 2がU . 2

GPU 2 5 0を利用して遂行される変換、関数、又はその他の処理作業に対応する一連の変換1 ~ 変換nに対する項目を含む。各項目はプロセッサ2 4 2が与えられた変換を遂行するように利用するU . 2 GPU 2 5 0のアドレスで満たされる。例えば、能力テーブル5 0 0を格納するドライブに対するプロセッサは、変換1を遂行するためにスロット - 3内のU . 2 GPU又はスロット - 4内のU . 2 GPUを利用する。プロセッサ2 4 2は、アルゴリズムにより識別された変換を遂行するため、能力テーブルから識別されたU . 2 GPUを選択する。

10

【 0 0 2 8 】

6 3 0ステップで、プロセッサ2 4 4はU . 2 GPU 2 5 0により処理されるデータチャンクを検索する。データチャンクはメモリー2 4 6から取り出されるか、又はドライブ2 4 2のリード/ライトのキャッシュバッファから取り出される。6 3 2ステップで、プロセッサ2 4 4は、例えばメッセージのヘッダーに命令を置き、本文にデータチャンクを置くことで、データチャンク及びデータチャンクを処理するための命令を含むメッセージを生成する。6 3 4ステップで、プロセッサ2 4 4は、例えばP C I eプロトコルを利用する内部スイッチ2 3 6を通じて通信することで、U . 2 GPU 2 5 0にメッセージを伝達する。

20

【 0 0 2 9 】

6 2 8ステップで、プロセッサ2 4 4は、U . 2 GPU 2 5 0により処理される全てのデータチャンクがU . 2 GPUに伝送されたか否かを判別する。全てのデータチャンクが伝送されていない場合、プロセッサ2 4 4は、残っているデータチャンクに対して6 3 0、6 3 2、及び6 3 4ステップを反復する。プロセッサ2 4 4が、処理される全てのデータチャンクがU . 2 GPU 2 5 0に伝送されたと判別した場合、手順は6 3 6ステップに進行する。

【 0 0 3 0 】

ドライブ2 4 2のプロセッサ2 4 4からメッセージを受信すると、直ちにU . 2 GPU 2 5 0は、結果データを生成するため、メッセージから受信されたデータチャンクに対して、メッセージで識別される変換を遂行する。U . 2 GPU 2 5 0が結果データを生成すると、U . 2 GPU 2 5 0はプロセッサ2 4 4に変換が完了したことを示す処理応答のメッセージを伝送する。いくつかの実施形態において、処理応答のメッセージは結果データを含み、プロセッサ2 4 4はドライブ2 4 2のメモリー2 4 6に結果データを格納する。いくつかの実施形態において、U . 2 GPU 2 5 0は追加的に又は代替的に結果データを、U . 2 GPU 2 5 0を含むストレージ部1 3 0内の他の位置又はストレージ部1 3 0の外部に伝送する。例えば、U . 2 GPU 2 5 0は追加的に又は代替的に結果データをホストサーバ2 1 0に伝送するか、或いはプロセッサ2 4 4がU . 2 GPU 2 5 0に対するメッセージにアドレスを含ませ、U . 2 GPU 2 5 0は結果データを指定されたアドレスに伝達する。

30

40

【 0 0 3 1 】

6 3 6ステップで、プロセッサ2 4 4は6 3 4ステップで伝送されたメッセージを受信した各U . 2 GPU 2 5 0から処理応答のメッセージを受信したか否かを判別する。全ての処理応答のメッセージを受信していない場合、6 3 8ステップで、プロセッサ2 4 4は、残っている処理応答のメッセージを受信するために待機する。全ての処理応答のメッセージが受信された場合、手順は6 2 0ステップに戻り、プロセッサ2 4 4はアルゴリズムの実行を継続する。或いは、例えばアルゴリズムが完了した場合、手順は6 1 0ステップに戻り、プロセッサ2 4 4は続けて進行するために他のトリガーが感知されるのを待つ。

50

【0032】

図6は、本発明の一実施形態によるグラフィック処理装置(U.2 GPU250A~250B)を管理するベースボード管理制御器234(BMC)に対する手順を示すフローチャートである。510ステップで、BMC234は、パワーオン(power-on)、プラグイン(plug-in)、又はプラグアウト(plug-out)のイベントがあったことを感知する。例えば、BMC234は、シャーシ231のポート(ミッドプレーンスロット)で“現在の”ピン(pin)をモニターし、新しいU.2 GPU250が現在のピンに連結(ミッドプレーンスロットに挿入)された場合に、プラグインのイベントが発生したことを感知するか、又はU.2 GPU250が現在のピン(ミッドプレーンスロット)から除去された場合に、プラグアウトのイベントが発生したことを感知する。一実施形態において、BMC234はNVMe管理インターフェース(NVMe-MI: NVMe Management Interface)プロトコルを用いてU.2 GPU250と通信するためにPCIe又はSMBusインターフェースを利用する。

10

【0033】

520ステップで、BMC234がイベントをパワーオンのイベントと判別した場合、BMC234はBMC234に連結されたU.2 GPU250の能力を判別する。例えば、図6に示すように、524ステップで、BMC234はBMC234に連結されたU.2 GPU250に対する必須製品のデータ(VPD: Vital Product Data)をリード(read)する。526ステップで、BMC234は、U.2 GPU250の能力にに関する情報を集めるために、追加的に又は代替的にNVMe“識別(Identify)”コマンドを利用する。522ステップで、BMC234がU.2 GPU250の能力を判別するために、スキャンされないBMC234に連結された一つ以上のU.2 GPU250が有ると判別した場合、BMC234は、残っているU.2 GPU又はU.2 GPUに対して524ステップ及び526ステップを反復する。BMC234がBMC234に連結された各U.2 GPU250がその能力を判別するためにスキャンされたと判別した場合、手順は530ステップに進行する。

20

【0034】

530ステップで、BMC234がU.2 GPUの判別された能力に基づいてドライブ242で能力テーブルをアップデートする。いくつかの実施形態において、判別された能力はU.2 GPUの現在の使用効率、U.2 GPUの特徴、U.2 GPU世代レーション(generation)、U.2 GPUの処理能力、スレッド(thread)プロセッサの個数、U.2 GPUのDRAMサイズ、帯域幅、遅延時間、精密度、入出力サイズ、及びMHz動作速度の中の一つ以上を含む。例えば、BMC234がスロット-3でU.2 GPU250がドライブ242に対する変換-1又は変換-2を遂行すると判別した場合、ドライブ242に対する能力テーブルは、変換-1及び変換-2に対応するものとしてスロット-3のアドレスを含むようにアップデートされる。528ステップで、BMC234がBMC234に連結される一つ以上のドライブ242が能力テーブルをアップデートしなかったと判別した場合、BMC234は、530ステップに戻り、次のドライブ242をアップデートする。BMC234がBMC234に連結された各ドライブ242が能力テーブルをアップデートしたと判別した場合、手順は510ステップに戻り、BMC234は、他のパワーオン、プラグイン、プラグアウトのイベントの感知の待機を再開する。

30

40

【0035】

BMC234が510ステップで感知されたイベントが、パワーオンのイベントではないと判別した場合、手順は532ステップに進行する。532ステップで、BMC234がイベントをプラグインのイベントであると判別した場合、BMC234はプラグインのイベントを誘発したBMC234に連結されたU.2 GPU250の能力を判別する。このプラグインのイベントは、U.2 GPUが現場交換可能な装置として具現され、現場交換可能なU.2 GPUが既に一つ以上のドライブ242を含むシャーシに連結され

50

る場合に対応する。U . 2 GPU 250の能力を判別する例として、図6に示すように、534ステップで、BMC 234はU . 2 GPU 250に対する必須製品のデータ (GPU . VPD) をリード (read) する。526ステップで、BMC 234は追加的に又は代替的にU . 2 GPU 250の能力に関連する情報を集めるためにNVMe “識別” コマンドを利用する。プラグインのイベントを誘発したU . 2 GPU 250の能力が判別された場合、手順は540及び538ステップに進行する。538及び540ステップで、530及び528ステップに対して上述したように、BMC 234は、新たなU . 2 GPU 250の能力に基づいてBMC 234に連結されたドライブ242の能力テーブルをアップデートする。

【0036】

BMC 234がイベントをプラグインのイベントではないと判別した場合、例えばBMC 234がイベントをプラグアウトのイベントであると判別した場合、手順は540ステップに進行する。このプラグアウトのイベントは、U . 2 GPUが現場交換可能な装置として具現され、U . 2 GPUがシャシー231から除去された場合に対応する。540及び538ステップで、BMC 234はU . 2 GPU 250に対応するプラグアウトのイベントに基づいてドライブ242の能力テーブルをアップデートする。例えば、BMC 234は能力テーブルからU . 2 GPU 250の以前のアドレスを除去する。

【0037】

いくつかの実施形態において、パワーオン、プラグイン、及びプラグアウト以外に、BMC 234は、イベントがストレージ管理イベントであると判別する。これはストレージ管理者 (例えば、ホストサーバ210のアプリケーション212) が、どのドライブ242がどのU . 2 GPU 250を利用するかを変更したことを示す。イベントがストレージ管理のイベントの場合、BMC 234は、これによって能力テーブルをアップデートする。例えば、ストレージ管理のイベントは特定のU . 2 GPU 250が特定のドライブ242を支援できないことを示し、BMC 234は特定のドライブ242の能力テーブルから特定のU . 2 GPU 250を除去する。

【0038】

いくつかの実施形態において、BMC 234がドライブ242の能力テーブルをアップデートした場合、BMC 234は、どのドライブに対して及びどの変換のために利用可能なものとしてリストするためのU . 2 GPUを判別するに当たって、各U . 2 GPUの能力以外の考慮事項を考慮する。例えば、一実施形態において、BMC 234は、特定のドライブ242を特定のU . 2 GPUに割り当てることで、U . 2 GPU 250に対する負荷の均衡をとる。このような実施形態において、BMC 234は、割り当てられたU . 2 GPU 250のアドレスを含むように、与えられたドライブ242の能力テーブルのみをアップデートする。他の実施形態において、BMC 234は特定のU . 2 GPU 250の能力に基づいて特定の変換を処理するために特定のU . 2 GPU 250を割り当てる。このような実施形態において、BMC 234は割り当てられたU . 2 GPU又は特定の変換のためのU . 2 GPUのアドレスを含むように能力テーブルをアップデートするが、U . 2 GPU能力が他の変換を処理するのに十分であるとしても、他の変換のためのアドレスを含まないこともある。例えば、BMC 234は変換 - 3に対してスロット - 2に位置する最も強力なU . 2 GPU 250を予約する。スロット - 2でU . 2 GPUは、変換 - 1 ~ 変換 - 9を処理するために適切な能力を有する。BMC 234は変換 - 3に対応するものとしてスロット - 2を含むように能力テーブルをアップデートするが、他の変換に対してスロット - 2を含まないことも有り、従って変換 - 3以外の他の変換に対してドライブ242がスロット - 2でU . 2 GPUを利用できないようにする。いくつかの実施形態において、BMC 234は、利用可能なU . 2 GPUの特徴、現在利用可能なU . 2 GPUの特徴 (例えば、専有されないか又は完全に割り当てられない)、及び/又は同時に実行する並列アルゴリズムの個数と比較して、変換に関するアルゴリズムの類型に基づいてロードバランシング (load balancing) を遂行する。

10

20

30

40

50

【 0 0 3 9 】

本明細書で使用した用語は、特定の実施形態を説明するための目的であり、本発明の思想を制限しようとするものではない。本明細書で使用したように、文脈が明確に違うことを意味しない限り、単数形は複数形を含むものとして意図する。“含む”、“含み”の用語は、本明細書で使用した場合、明示された特徴、整数、ステップ、動作、要素及び／又は部品の存在を明示するが、一つ以上の他の特徴、整数、ステップ、動作、要素、部品及び／又はグループの存在又は追加を排除しない。本明細書で使用したように、“及び／又は”の用語は一つ以上の関連する列挙項目の全ての組合せ及び特定の組合せを含む。“少なくとも一つ”のような表現は、要素のリストの前にある場合、全体要素のリストを修正し、リストの個別要素を修正しない。

10

【 0 0 4 0 】

本明細書で使用したように、“できる”の使用は、本発明の実施形態を説明する場合、“本発明の一つ以上の実施形態”を指し示す。本明細書で使用したように、“使用”、“使用する”、“使用された”の用語は“利用”、“利用する”、“利用された”の用語に各々同義語と見なされ得る。なお、“例示的な”の用語は例示又は図面を指し示すものとして意図する。

【 0 0 4 1 】

本明細書で説明した本発明の実施形態による電子又は電気装置及び／又は他の関連する装置又は部品は、任意の適切なハードウェア、ファームウェア（例えば、応用注文型集積回路（ASIC: application-specific integrated circuit））、ソフトウェア、又はソフトウェア、ファームウェア、及びハードウェアの組合せを利用して具現される。例えば、このような装置の多様な部品は、一つの集積された回路（IC）チップ又は別途のICチップに形成される。なお、このような装置の多様な部品は、フレキシブル印刷回路フィルム、テープキャリアパッケージ（TCP: Tape Carrier Package）、印刷回路基板（PCB: Printed Circuit Board）上に具現されるか又は一つの基板に形成される。なお、このような装置の多様な部品は、一つ以上のコンピューティング装置で、本明細書で説明した多様な機能を遂行するためコンピュータプログラムの命令を実行し、他のシステム部品と相互作用する一つ以上のプロセッサで実行されるプロセス又はスレッドである。コンピュータプログラムの命令は、例えばランダムアクセスメモリー（RAM）のような標準メモリー装置を用い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、コンピュータプログラムの命令は、例えばCD-ROM、フラッシュドライブなどのような他の非一時的なコンピュータ読み取り可能なメディアに格納される。また、当業者（通常の技術者）は多様なコンピューティング装置の機能性が一つのコンピューティング装置に結合されるか又は統合されることを認識しなければならない。また、本発明の例示的な実施形態の思想及び範囲を逸脱せずに特定のコンピューティング装置の機能性が一つ以上の他のコンピューティング装置に分散されることを認識しなければならない。

20

30

【 0 0 4 2 】

別途に定義しない場合、本明細書で使用した全ての用語（技術的及び科学的用語を含む）は、本発明が属する技術における当業者により、一般的に理解されるものと同一の意味を有する。用語は、一般的に利用される辞典で定義されるように、関連する技術及び／又は本明細書の文脈上の意味と一致する意味を有するものとして解釈されるべきであり、本明細書で明確に定義しない限り、理想的であるか形式的な感覚により解釈してはならない。

40

【 0 0 4 3 】

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。

【産業上の利用可能性】

【 0 0 4 4 】

50

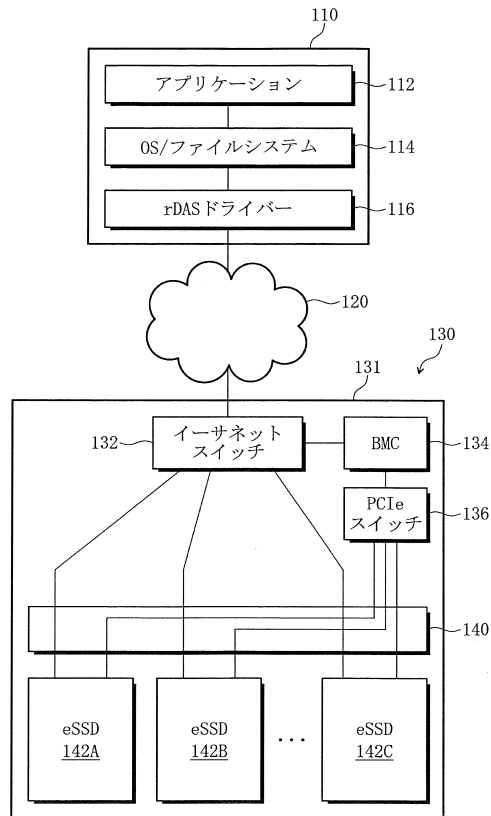
本発明は、ストレージシステムにおいて、複数のグラフィック処理装置を利用してマシンラーニングアルゴリズムを実行するデータストレージ及び処理システムとその動作方法に有用である。

【符号の説明】

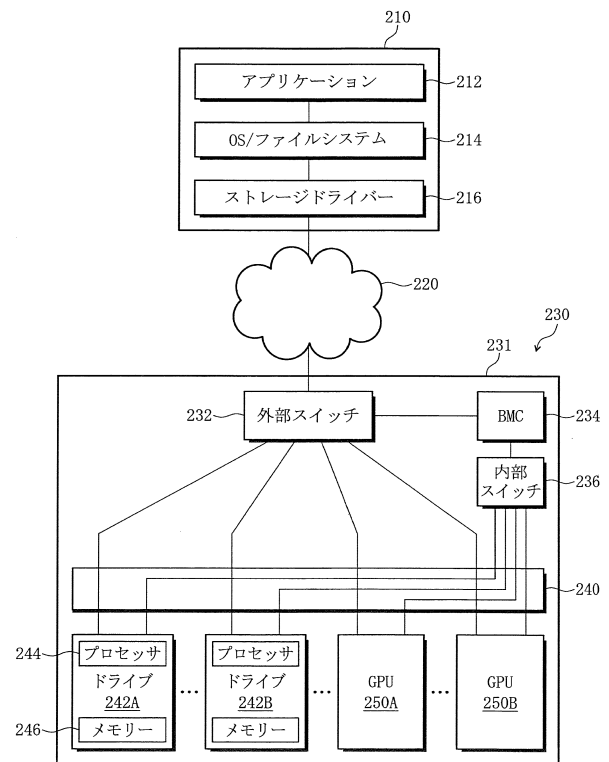
【 0 0 4 5 】

1 1 0、2 1 0	ホストサーバ	
1 1 2、2 1 2	アプリケーション	
1 1 4、2 1 4	OS / ファイルシステム	
1 1 6	遠隔直接アクセスストレージ (r D A S : r e m o t e D i r e c t A	
c c e s s S t o r a g e)	ドライバ	10
1 2 0、2 2 0	ネットワーク	
1 3 0、2 3 0	ストレージ部	
1 3 1、2 3 1	シャシー (c h a s s i s)	
1 3 2	イーサネット (e t h e r n e t) スイッチ	
1 3 4、2 3 4	ベースボード管理制御器 (B M C : b a s e b o a r d m a n a	
g e m e n t c o n t r o l l e r)		
1 3 6	P C I e スイッチ	
1 4 0、2 4 0	シャシーミッドプレーン (c h a s s i s m i d p l a n e)	
1 4 2 A ~ 1 4 2 C	e S S D (e t h e r n e t S o l i d S t a t e D r	
i v e s)		20
2 1 6	ストレージドライバ (遠隔直接アクセスストレージドライバ)	
2 3 2	外部スイッチ (イーサネットスイッチ)	
2 3 6	内部スイッチ (P C I e スイッチ)	
2 4 2、2 4 2 A ~ 2 4 2 B	ドライブ	
2 4 4、3 6 0	プロセッサ	
2 4 6	メモリー	
2 5 0、2 5 0 A ~ 2 5 0 B、3 5 0	U . 2 グラフィック処理装置 (G P U)	
3 5 2	コネクタ	
3 5 4	イーサネットインターフェース	
3 5 6	P C I e インターフェース	30
3 5 8	システム管理バス (S M B u s) インターフェース	
3 6 2	D R A M	
3 6 4	不揮発性メモリー (N V M)	
3 6 6	電源供給器 / 電源調節器	
5 0 0	能力テーブル	
6 2 0	アルゴリズム	

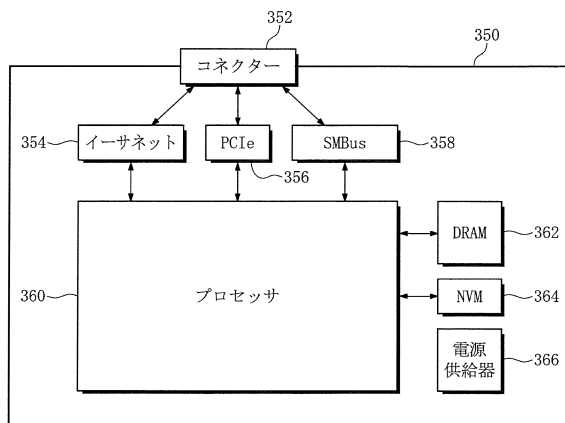
【図 1】



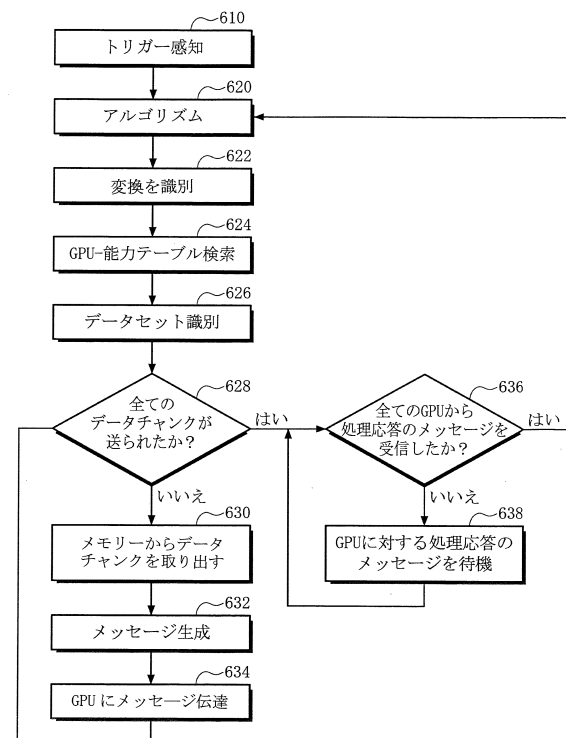
【図 2】



【図 3】



【図 4】

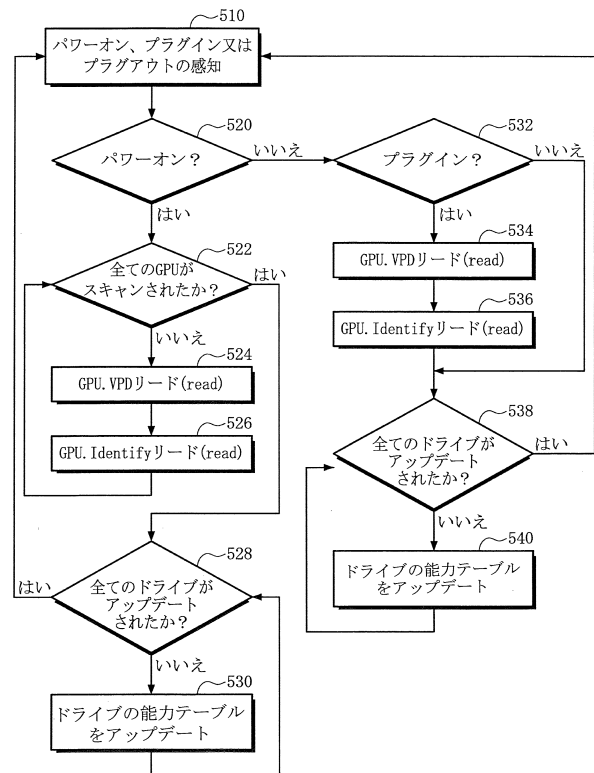


【図 5】

500

変換	GPU 住所	GPU 住所	GPU 住所	GPU 住所	GPU 住所
変換 1	スロット-3	スロット-4			
変換 2	スロット-3	スロット-5	スロット-6	スロット-13	
変換 n	スロット-7				

【図 6】



フロントページの続き

- (72)発明者 オラリゲ, ソンボン ポール
アメリカ合衆国, 94566, カリフォルニア州, プレザントン, パセオ グラナダ 3
050
- (72)発明者 シュワドラー, デイビッド
アメリカ合衆国, 95070, カリフォルニア州, サラトガ, パセオ プレサード 13
165

審査官 松平 英

- (56)参考文献 米国特許出願公開第2014/0129753(US, A1)
米国特許出願公開第2011/0292058(US, A1)
国際公開第2016/185542(WO, A1)
米国特許出願公開第2011/0295967(US, A1)
米国特許出願公開第2014/0176583(US, A1)

- (58)調査した分野(Int.Cl., DB名)
- | | |
|------|-------|
| G06F | 3/06 |
| G06F | 3/08 |
| G06F | 13/10 |
| G06N | 20/00 |