

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2017-512350

(P2017-512350A)

(43) 公表日 平成29年5月18日(2017.5.18)

(51) Int.Cl.  
G06F 13/10 (2006.01)F I  
G06F 13/10 340A

テーマコード (参考)

審査請求 未請求 予備審査請求 未請求 (全 54 頁)

(21) 出願番号 特願2016-575306 (P2016-575306)  
 (86) (22) 出願日 平成27年3月6日 (2015.3.6)  
 (85) 翻訳文提出日 平成28年11月8日 (2016.11.8)  
 (86) 国際出願番号 PCT/US2015/019206  
 (87) 国際公開番号 W02015/138245  
 (87) 国際公開日 平成27年9月17日 (2015.9.17)  
 (31) 優先権主張番号 61/950,036  
 (32) 優先日 平成26年3月8日 (2014.3.8)  
 (33) 優先権主張国 米国 (US)  
 (31) 優先権主張番号 62/017,257  
 (32) 優先日 平成26年6月26日 (2014.6.26)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 516270692  
 ディアマンティ インコーポレイテッド  
 アメリカ合衆国 カリフォルニア州 95  
 113 サン ホセ ノース マーケット  
 ストリート 111 スイート 800  
 (74) 代理人 100086771  
 弁理士 西島 孝喜  
 (74) 代理人 100088694  
 弁理士 弟子丸 健  
 (74) 代理人 100094569  
 弁理士 田中 伸一郎  
 (74) 代理人 100067013  
 弁理士 大塚 文昭  
 (74) 代理人 100109070  
 弁理士 須田 洋之

最終頁に続く

(54) 【発明の名称】 集中型ネットワーキング及びストレージのための方法及びシステム

## (57) 【要約】

装置が、物理的ターゲットストレージ媒体コントローラと、物理的ネットワークインターフェイスコントローラと、ストレージ媒体コントローラとネットワークインターフェイスコントローラとの間のゲートウェイを含む集中型入力/出力コントローラを含み、ゲートウェイが、ストレージ媒体コントローラとネットワークインターフェイスコントローラとの間のストレージトラフィック及びネットワークトラフィックのための直接接続を提供する。

【選択図】 図3

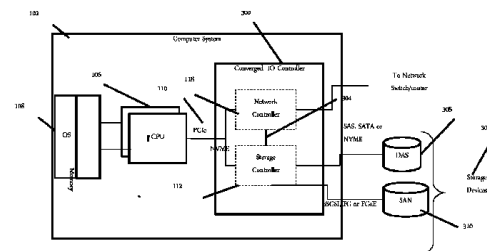


FIG. 3

**【特許請求の範囲】****【請求項 1】**

集中型入力／出力コントローラを備えた装置であって、前記集中型入力／出力コントローラは、

物理的ターゲットストレージ媒体コントローラと、

物理的ネットワークインターフェイスコントローラと、

前記ストレージ媒体コントローラと前記ネットワークコントローラとの間のゲートウェイと、

を備え、前記ゲートウェイは、前記ストレージ媒体コントローラと前記ネットワークインターフェイスコントローラとの間のストレージトラフィック及びネットワークトラフィックのための直接接続を提供する、  
ことを特徴とする装置。

10

**【請求項 2】**

前記ストレージ媒体コントローラによって制御されるストレージ媒体を、該ストレージ媒体の位置に関わらずローカルに接続されたストレージとして提示する仮想ストレージインターフェイスをさらに備える、  
請求項 1 に記載の装置。

**【請求項 3】**

前記ストレージ媒体コントローラによって制御されるストレージ媒体を、該ストレージ媒体の数又はタイプに関わらずローカルに接続されたストレージとして提示する仮想ストレージインターフェイスをさらに備える、  
請求項 1 に記載の装置。

20

**【請求項 4】**

前記ストレージ媒体の動的プロビジョニングを容易にする仮想ストレージインターフェイスをさらに備え、前記物理ストレージは、局所的又は遠隔的に存在することができる、  
請求項 1 に記載の装置。

**【請求項 5】**

前記ストレージ媒体の動的プロビジョニングを容易にする仮想ネットワークインターフェイスをさらに備え、前記物理ストレージは、局所的又は遠隔的に存在することができる、  
請求項 1 に記載の装置。

30

**【請求項 6】**

前記装置は、ホストコンピュータシステム上のコントローラカードとして導入されるように適合される、  
請求項 1 に記載の装置。

**【請求項 7】**

前記ゲートウェイは、前記ホストコンピュータシステムの CPU 上で実行される前記オペレーティングシステムによる介入、ハイパーバイザによる介入、又はその他のソフトウェアによる介入を伴わずに動作する、  
請求項 6 に記載の装置。

40

**【請求項 8】**

前記装置は、前記装置のストレージ機能及びネットワーク機能の少なくとも一方を提供するフィールドプログラマブルゲートアレイ、ASIC 及びネットワークプロセッサのうちの少なくとも 1 つを備える、  
請求項 1 に記載の装置。

**【請求項 9】**

前記装置は、ネットワーク展開されたスイッチとして構成される、  
請求項 1 に記載の装置。

**【請求項 10】**

ストレージ媒体命令を第 1 のプロトコルと少なくとも 1 つの他のプロトコルとの間で変

50

換する、前記装置の機能コンポーネントをさらに含む、請求項 1 に記載の装置。

【請求項 1 1】

ストレージ装置の仮想化方法であって、

第 1 のストレージプロトコルでの命令に応答する物理ストレージ装置にアクセスするステップと、

前記第 1 のストレージプロトコルと第 2 ストレージプロトコルとの間で命令を変換するステップと、

前記第 2 のプロトコルを用いて前記物理ストレージ装置をオペレーティングシステムに提示することにより、前記オペレーティングシステムを使用するホストコンピュータシステムに対して前記物理ストレージ装置が局所的に存在するか、それとも遠隔的に存在するかに関わらず、前記物理ストレージ装置のストレージが動的に展開されるようにするステップと、

を含むことを特徴とする方法。

【請求項 1 2】

前記第 1 のプロトコルは、SATA プロトコル、NVME プロトコル、SAS プロトコル、iSCSI プロトコル、ファイバチャネルプロトコル及びファイバチャネルオーバーサネット(登録商標)プロトコルのうちの少なくとも 1 つである、

請求項 1 1 に記載の方法。

【請求項 1 3】

前記第 2 のプロトコルは、NVMe プロトコルである、

請求項 1 1 に記載の方法。

【請求項 1 4】

オペレーティングシステムと、前記第 1 及び第 2 ストレージプロトコル間における命令の変換を行う装置との間のインターフェイスを提供するステップをさらに含む、

請求項 1 1 に記載の方法。

【請求項 1 5】

前記命令の変換を行う装置と、遠隔地に存在するネットワーク展開されたストレージ装置との間の NVMe オーバーサネット(登録商標)接続を提供するステップをさらに含む、

請求項 1 1 に記載の方法。

【請求項 1 6】

ターゲット物理ストレージ装置に記憶されたアプリケーション、コンテナ及びデータのうちの少なくとも 1 つの移動を容易にする方法であって、

集中型ストレージ及びネットワークングコントローラを提供して、ゲートウェイが、ホストコンピュータの CPU 上で実行されるオペレーティングシステムの介入、ハイパーバイザの介入、又はその他のソフトウェアの介入を必要とせずに前記装置のストレージコンポーネントとネットワークングコンポーネントとの間のネットワーク及びストレージトラフィックのための接続を提供するステップと、

前記少なくとも 1 つのアプリケーション又はコンテナを、前記集中型ストレージ及びネットワークングコントローラによって制御されるターゲット物理ストレージ装置にマッピングすることにより、前記ターゲット物理ストレージ装置に記憶されたアプリケーション、コンテナ又はデータが 1 又は 2 以上の他のコンピュータシステムに移動する際に、前記アプリケーション又はコンテナが、前記ターゲット物理ストレージが接続された前記ホストシステムの CPU 上で実行されるオペレーティングシステムの介入、ハイパーバイザの介入、又はソフトウェアの介入を必要とせずに前記ターゲット物理ストレージにアクセスできるようにするステップと、

を含むことを特徴とする方法。

【請求項 1 7】

前記移動は、Linux(登録商標)コンテナの移動である、

10

20

30

40

50

請求項 16 に記載の方法。

【請求項 18】

前記移動は、ハイパーバイザにおいて実行される仮想マシンの移動である、

請求項 16 に記載の方法。

【請求項 19】

前記移動は、スケールアウトアプリケーションの移動である、

請求項 16 に記載の方法。

【請求項 20】

前記ターゲット物理ストレージは、iSCSI プロトコル、ファイバチャネルプロトコル及びファイバチャネルオーバーサネット(登録商標)プロトコルのうちの少なくとも 1 つを使用するネットワーク展開されたストレージ装置である、

請求項 16 に記載の方法。

【請求項 21】

前記ターゲット物理ストレージは、SAS プロトコル、SATA プロトコル及び NVME プロトコルのうちの少なくとも 1 つを使用する直接接続されたストレージ装置である、

請求項 16 に記載の方法。

【請求項 22】

ネットワークのサービス品質(QoS)を提供する方法であって、

集中型ストレージ及びネットワークングコントローラを提供して、ゲートウェイが、ホストコンピュータの CPU 上で実行されるオペレーティングシステムの介入、ハイパーバイザの介入、又はソフトウェアの介入を伴わずに前記装置のストレージコンポーネントとネットワークングコンポーネントとの間のネットワークトラフィック及びストレージトラフィックのための接続を提供するステップと、

前記集中型ストレージ及びネットワークングコントローラによって処理される前記ストレージトラフィック及び前記ネットワークトラフィックの少なくとも一方に基づいて、ホストコンピュータのオペレーティングシステムの介入を伴わずに、前記ストレージ及びネットワークングコントローラが展開されたデータ経路を有するネットワークに関連する少なくとも 1 つのサービス品質(QoS)パラメータを管理するステップと、  
を含むことを特徴とする方法。

【請求項 23】

前記 QoS パラメータは、帯域幅パラメータ、ネットワークレイテンシパラメータ、I/O 性能パラメータ、スループットパラメータ、ストレージタイプパラメータ及びストレージレイテンシパラメータからなる群から選択される、

請求項 22 に記載の方法。

【請求項 24】

前記 QoS は、前記集中型ストレージ及びネットワークコントローラを介してストレージによってサービスを受けるアプリケーション及びコンテナの少なくとも一方がホストコンピュータから別のコンピュータに移動する際に自動的に維持される、

請求項 22 に記載の方法。

【請求項 25】

前記 QoS は、前記集中型ストレージ及びネットワークコントローラを介してアプリケーション及びコンテナの少なくとも一方にサービスを提供する少なくとも 1 つのターゲットストレージ装置が第 1 の位置から少なくとも 1 つの第 2 の位置に移動する際に自動的に維持される、

請求項 22 に記載の方法。

【請求項 26】

ネットワークトラフィックデータの暗号化、ストレージ内のデータの暗号化、及びネットワークトラフィックデータとストレージ内のデータとの暗号化からなる群からセキュリティ機能が選択される、

請求項 22 に記載の方法。

10

20

30

40

50

## 【請求項 27】

圧縮、保護レベル、RAIDレベル、ストレージ媒体タイプ、グローバル重複排除、並びに目標復旧時点(RPO)及び目標復旧時間(RTO)の少なくとも一方を達成するためのスナップショット間隔からなる群から選択された1又は2以上のストレージ機能が提供される、

請求項22に記載の方法。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

〔関連出願との相互参照〕

10

本出願は、2014年3月8日に提出された「アプリケーション駆動型ストレージアクセスのための方法及びシステム(Method and Apparatus for Application Driven Storage Access)」という名称の米国特許出願第61/950,036号、及び2014年6月26日に提出された「仮想クラスタIOのための装置(Apparatus for Virtualized Cluster IO)」という名称の米国特許出願第62/017,257号に対する優先権を主張するものであり、これらの各仮特許出願はその全体が引用により本明細書に組み入れられる。

## 【0002】

本出願は、ネットワーキング及びデータストレージの分野に関し、具体的には、集中型  
ネットワーキング及びデータストレージ装置の分野に関する。

20

## 【背景技術】

## 【0003】

スケールアウトアプリケーションの蔓延は、このようなアプリケーションを使用する企業に非常に重要な課題をもたらしてきた。通常、企業は、(ハイパーバイザ及びプレミアムハードウェアコンポーネントのようなソフトウェアコンポーネントを伴う)仮想マシンのようなソリューション、及び(一般的にはLinux(登録商標)及びコモディティハードウェアのようなオペレーティングシステムの使用を伴う)いわゆる「ベアメタル」ソリューションを選択する。大まかに言えば、通常、仮想マシンソリューションは、入出力性能が劣り、メモリが不十分であり、性能に一貫性がなく、インフラコストが高い。通常、  
ベアメタルソリューションは、(リソースの変更を困難にしてハードウェアの使用を効率的にする)静的リソース割り当てを行い、容量の計画に課題があり、性能に一貫性がなく、動作が複雑である。いずれの場合にも、性能の非一貫性が既存のソリューションの特徴である。動的リソース割り当てに対処できるとともにコモディティハードウェアを高利用率で利用できるマルチテナント型展開において高性能を発揮するソリューションが必要とされている。

30

## 【0004】

図1に、本明細書に開示するいくつかの実施形態に関連し得る機能及びモジュールを有するサーバなどのコンピュータシステム102の全体的アーキテクチャを示す。従来、(媒体104(例えば、回転媒体又はフラッシュ)などの、サーバ102上のローカルストレージ装置へのアクセスなどの)ストレージ機能と、転送などのネットワーク機能は、ソフトウェアスタック又は(例えば、ネットワーク機能又はストレージ機能のためにそれぞれネットワークインターフェイスコントローラ118又はストレージコントローラ112を伴う)ハードウェア装置のいずれかにおいて別個に行われていた。通常、(オペレーティングシステムと、実施形態によってはコンピュータシステムのストレージ機能及びネットワーク機能に関連する全てのソフトウェアスタックを含むハイパーバイザとを含むことができる)オペレーティングシステムスタック108内では、ソフトウェアストレージスタックが、スモールコンピュータシステムインターフェイス(SCSI)プロトコル、シリアルATA(SATA)プロトコル、又は不揮発性メモリエクスプレス(NVMe)プロトコル(典型的なコンピュータシステム102のPCIエクスプレス(PCIe)バス

40

50

110を通じてソリッドステートドライブ（SSD）のようなディスク接続ストレージ（DAS）にアクセスするためのプロトコル）などの、ストレージ内で使用できる様々なプロトコルの使用を可能にするモジュールを含む。PCIeバス110は、（（単複の）プロセッサ及びメモリを含む）CPU106と、様々なIOカードとの間を相互接続することができる。ストレージスタックは、ボリュームマネージャなどを含むこともできる。ストレージソフトウェアスタック内の動作は、ミラーリング又はRAID、バックアップ、スナップショット、重複排除、圧縮及び暗号化などのデータ保護を含むこともできる。ストレージ機能の一部は、ストレージコントローラ112内にオフロードすることができる。ソフトウェアネットワークスタックは、伝送制御プロトコル/インターネットプロトコル（TCP/IP）、ドメインネームシステムプロトコル（DNS）、アドレス解決プロトコル（ARP）及び転送プロトコルなどの様々なネットワークプロトコルの使用を可能にするためのモジュール及び機能などを含む。ネットワーク機能の一部は、イーサネット（登録商標）接続120などを介してネットワークインターフェイスコントローラ118（又はNIC）又はネットワークファブリックスイッチにオフロードし、（様々なスイッチ及びルータなどを用いて）さらにネットワークにつなげることができる。仮想環境では、PCIエクスプレス標準においてSR-IOVによって指定されるように、NIC118を複数の仮想NICに仮想化することができる。PCIエクスプレス標準によって指定されておらず、それほど一般的でもないが、ストレージコントローラも同様に仮想化することができる。仮想マシンなどの仮想エンティティは、この方法によって固有のプライベートリソースにアクセスすることができる。

10

20

#### 【0005】

図2を参照すると、ハイパーバイザの1つの主要課題は、IO動作の複雑性である。例えば、2つの異なるコンピュータ（図2のコンピュータシステム1とコンピュータシステム2）にわたるデータに関連する動作に対処するには、コンピュータのローカルストレージ装置104、ストレージコントローラ112、CPU106、ネットワークインターフェイスコントローラ118及びハイパーバイザ/オペレーティングシステム108に関連する異なるソフトウェアスタック間をデータが移動する際にデータを幾度となく繰り返しコピーしなければならない、1つのコンピュータから別のコンピュータにデータを移すことや、或いはストレージの構成を変更することなどを伴う行動が行われる度に、IO動作毎に大量の非効率的なデータコピーが生じてしまう。ルート124は、データが2つのコンピュータのソフトウェアスタックを上下に移動して1つのコンピュータから別のコンピュータに至るまでに取ることができる多くの複雑なルート例の1つである。コンピュータシステム2が求めるデータは、当初はコンピュータシステム1のディスクなどのローカルストレージ装置104内に存在することがあり、この場合、ストレージコントローラカード112によってプルされ（IO動作及びコピーを伴い）、PCIeバス110（別のIO動作）を介してCPU108に送られ、ここでコンピュータシステム1のOSスタック108のハイパーバイザ又はその他のソフトウェアコンポーネントによって処理される。次に、このデータを、ネットワークコントローラ118を通じ、ネットワーク122を介して（別の一連のIO動作）コンピュータシステム2に送出（別のIO動作）することができる。このルートはコンピュータシステム2においても継続し、データは、ネットワークコントローラ118を通じてコンピュータシステム2のCPU106に移動した（さらなるIO動作を伴った）後に、PCIeバス110を介してストレージのローカルストレージコントローラ112に送られ、その後ハイパーバイザ/OSスタック108に戻されて実際に使用することができる。これらの動作は、数多くのコンピュータシステムの対にわたって行われ、各交換は、この種の大幅に増加したIO動作を伴う（また他の多くのルートも可能であり、それぞれが多くの動作を伴う）ことがある。企業が次第に採用してきているスケールアウトの状況では、コンピュータシステム間におけるこのような多くの複雑なデータ複製及び転送行為が必要である。例えば、MongoDB（商標）のようなスケールアウトアプリケーションを実装すると、顧客は、リバランス動作中にリアルタイムクエリを繰り返し実行して大規模なデータロードを行わなければならない。このような行為

30

40

50

は、非常に多くの I O 動作を伴い、ハイパーバイザソリューションの性能が低下してしまう。これらのアプリケーションのユーザは、頻繁にリシャードイング（データが配置されたシャードの変更）も行い、やはり 1 つの場所から別の場所へのデータの移動に数多くの I O 動作を伴う多くのコピー動作及び転送動作が関与するので、静的ストレージリソース割り当てを行うベアメタルソリューションにとっての大きな問題となる。これらの問題は、スケールアウトアプリケーションにおいて使用されるデータ量が急速に増加し、（多くの機械を伴うクラウド展開などにおいて）異なるシステム間のつながりが増えるにつれて指数関数的に拡大する。I O 動作の回数及び複雑性を低減し、高価なプレミアムハードウェアを必要とせずにスケールアウトアプリケーションの性能及び拡張性を別様に高めるストレージソリューション及びネットワーキングソリューションが必要とされている。

10

#### 【0006】

引き続き図 2 を参照すると、多くのアプリケーション及び使用事例では、コンピュータシステム 102 間のネットワークを越えてデータ（さらにはストレージ）にアクセスする必要がある。この動作の 3 つの高レベルステップは、一方のコンピュータシステムのストレージ媒体からボックス外にデータを転送し、ネットワーク 122 を越えて移動させ、第 2 のボックス（第 2 のコンピュータシステム 102）内にデータを転送してこの第 2 のコンピュータシステム 102 のストレージ媒体 104 に移動させることを含む。まず、ボックス外転送は、ストレージコントローラ 112 からの介入、OS 108 内のストレージスタックからの介入、OS 108 内のネットワークスタックからの介入、及びネットワークインターフェイスコントローラ 118 からの介入を伴う可能性がある。内部バス（PCIe 110 及びメモリ）を横切る多くのトラバース及びコピー、並びに CPU 106 の処理サイクルが費やされる。これにより、動作性能が低下する（レイテンシ及びスループットの問題が生じる）だけでなく、CPU 上で実行される他のアプリケーションにも悪影響が及ぶ。次に、データがボックス 102 から離れてネットワーク 122 上に移動すると、このデータは他のあらゆるネットワークトラフィックと同様に扱われ、その宛先に転送/ルーティングする必要がある。ポリシーが実行されて決定が行われる。大量のトラフィックが動いている環境では、ネットワーク 122 に輻輳が生じて性能の低下及び利用可能性問題（例えば、パケットの欠損、接続の喪失及び予測不能なレイテンシ）を引き起こす恐れがある。ネットワークは、一時停止機能、逆方向輻輳通知（BCN）、明示的輻輳通知（ECN）などの、輻輳の拡大を避けるための機構及びアルゴリズムを有する。しかしながら、これらは反応的方法、すなわち輻輳形成地点を検出し、発生源において先送りして輻輳を低減するものであり、結果として遅延及び性能への影響を生じる可能性がある。3 番目に、データがその「宛先」コンピュータシステム 102 に到着すると、このデータを処理する必要がある。これにはネットワークインターフェイスコントローラ 118 からの介入、OS 108 内のネットワークスタックからの介入、OS 108 内のストレージスタックからの介入及びストレージコントローラ 112 からの介入が伴う。上述したボックス外動作と同様に、内部バスを横切る多くのトラバース及びコピー、並びに CPU 106 の処理サイクルが費やされる。さらに、データの最終目的地は、さらに異なるボックス内に存在することもある。このことは、より多くのデータ保護（例えば、ミラーリング又はボックス間 RAID）の必要性、又は重複排除の必要性の結果と考えることができる。そうである場合、ネットワークを横切るボックス外及びボックス内へのデータ転送シーケンスを再び繰り返す必要がある。上述したように、この方法を制限するものとしては、本来の性能の低下、予測不能な性能、他のテナント又は動作への影響、利用可能性及び信頼性、並びに非効率的なリソース使用が挙げられる。現行方法の複雑性及び性能への影響を回避するデータ転送システムに対するニーズが存在する。

20

30

40

#### 【発明の概要】

#### 【発明が解決しようとする課題】

#### 【0007】

（管理する各仮想マシンに別個のオペレーティングシステムを提供する）ハイパーバイザの代案として、（単一のオペレーティングシステムによる複数のアプリケーションコン

50

テナの管理を可能にする) L i n u x (登録商標) コンテナなどの技術が開発されてきた。また、アプリケーションをライブラリにパッケージングするプロビジョニングを提供する D o c k e r s などのツールも開発されてきた。本開示全体を通じて説明する他の多くの技術革新の中でも特に、これらの新たな技術の能力を活用してスケールアウトアプリケーションのための改善された方法及びシステムを提供する機会が存在する。

【課題を解決するための手段】

【0008】

本明細書では、イニシエータ、ターゲットストレージ機能及びネットワーク機能を単一のデータ及び制御経路に組み合わせるハードウェアにおける集中型ストレージ及びネットワークコントローラを含み、ホストCPUによる介入を必要とせずにネットワークとストレージとの間の「カットスルー」経路を可能にする方法及びシステムを提供する。参照を容易にするために、本開示では、このソリューションを、集中型ハードウェアソリューション、集中型装置、集中型アダプタ、集中型I/Oコントローラ又は「データワイズ」コントローラなどと様々に呼ぶことができ、このような用語は、文脈において別途示している場合を除き、ターゲットストレージ機能とネットワーク機能とを単一のデータ及び制御経路に組み合わせるハードウェアにおける集中型ストレージ及びネットワークコントローラを含むものとして理解すべきである。

【0009】

この集中型ソリューションは、他の利点の中でも特に、コンピュータリソース及び/又はストレージリソースのクラスタのそのままの性能を高め、クラスタ全体にわたってサービスレベル合意(SLA)を強制するとともに予測可能な性能の保証に役立ち、テナントがその近隣に影響を与えないマルチテナント環境を提供し、ハードウェアの利用度が高いことによってデータセンタの占有面積が減少し、電力が減少し、管理すべきシステムが少なくなる、より高密度なクラスタを提供し、拡張性の高いクラスタを提供し、性能を損なうことなくクラスタ全体にわたってストレージリソースをプールする。

【0010】

本明細書に開示する様々な方法及びシステムは、スケールアウトアプリケーション及び高性能マルチノードプールに必要なリソースを高密度に圧密化する。これらの方法及びシステムは、動的なクラスタ規模のリソース提供、ネットワーク機能及びストレージ機能に対してサービス品質(QoS)、セキュリティ、分離などを保証する能力、並びに生産及び検査/開発のために共有インフラストラクチャを使用する能力を含む複数の顧客利益を提供する。

【0011】

また、本明細書では、ネットワークを通じたストレージ機能を実行し、ストレージ装置及びネットワーク装置を仮想化して単一テナント環境又はマルチテナント環境における高性能及び決定的性能をもたらす方法及びシステムも提供する。

【0012】

また、本明細書では、N V M e 及び同様のプロトコルを使用するようなストレージ装置を仮想化し、これらの仮想装置を、S A T A を使用するような異なる物理的装置に変換する方法及びシステムも提供する。

【0013】

本明細書に開示する方法及びシステムは、ボックスレベルでの遠隔クレジット管理及び分散スケジューリングアルゴリズムを含むハードウェアを(ネットワークファブリックとは対照的に)ホスト側のみに含むエンドツーエンド輻輳制御のための方法及びシステムも含む。

【0014】

また、本明細書では、ストレージアダプタ、ネットワークアダプタ、コンテナ(例えば、L i n u x (登録商標)コンテナ)又はS o l a r i s ゾーンなどのクラスタを可能にするストレージクラスタ又はその他の要素の仮想化方法及びシステムを含む、集中型ネットワーク/ストレージコントローラによって可能になる様々な方法及びシステムも提供する



。利点の中でも特に、クラスタを仮想化する１つの態様では、物理的クラスタ内でコンテナが位置に依存しないようにすることができる。他の利点の中でも特に、これにより、後述する大いに単純化されたプロセスにおけるマシン間のコンテナの移動が可能になる。

【００１５】

本明細書では、直接接続ストレージ装置（ＤＡＳ）の仮想化方法及びシステムを提供することにより、たとえ物理ストレージ装置が移動して遠隔地に存在する場合でも、オペレーティングシステムスタック１０８がローカルな永続装置を探し求めるようになり、すなわち本明細書では、ＤＡＳの仮想化方法及びシステムを提供する。実施形態では、この方法が、ファブリック全体にわたってＤＡＳを仮想化し、すなわちＤＡＳストレージシステムを選択し、ボックス外に移動させてネットワーク上配置することを含むことができる。実施形態では、この方法が、ＤＡＳを任意の名前空間に分割することを含むことができる。実施形態では、オペレーティングシステムがこの仮想ＤＡＳに、まるで実際のＤＡＳであるかのようにアクセスできるようになり、例えば、ＯＳ１０８がＮＶＭｅを介しＰＣＩｅバスを通じてアクセスすることができる。従って、本明細書では、（ＤＡＳを含む）ストレージ装置を仮想化する能力を提供することにより、たとえ実際にはイーサネット（登録商標）などのネットワークプロトコルを介してストレージ装置にアクセスしている場合でも、ＯＳ１０８がこの仮想ストレージをＤＡＳと見なすようになり、従ってＯＳ１０８は、ローカル物理ストレージ装置の場合に必要なことと異なることを行う必要はない。

10

【００１６】

本明細書では、ＯＳ１０８を何ら修正する必要なくＯＳ１０８に仮想ＤＡＳを露出することを含む、ファブリック全体にわたってＤＡＳを提供する方法及びシステムを提供する。

20

【００１７】

また、本明細書では、（ターゲットストレージシステムを参照する）ストレージアダプタの仮想化方法及びシステムも提供する。

【００１８】

本明細書では、ストレージイニシエーションとストレージターゲットとを単一のハードウェアシステムに組み合わせる方法及びシステムを提供する。実施形態では、これらをＰＣＩｅバス１１０によって接続することができる。単一のルート仮想化機能（ＳＲ－ＩＯＶ）を適用して、いずれかの標準的な装置を選択し、この装置を数百個ものこのような装置であるかのように機能させることができる。本明細書に開示する実施形態は、ＳＲ－ＩＯＶを用いて物理ストレージアダプタの複数の仮想インスタンスを生じることを含む。ＳＲ－ＩＯＶは、Ｉ／Ｏ機能を仮想化するＰＣＩｅ標準であり、ネットワークインターフェイスに使用されてきたが、本明細書に開示する方法及びシステムは、ＳＲ－ＩＯＶをストレージ装置に使用するように拡張する。従って、本明細書では仮想ターゲットストレージシステムが提供される。

30

【００１９】

実施形態は、スイッチフォームファクタ又はネットワークインターフェイスコントローラを含むことができ、本明細書に開示する方法及びシステムは、（ソフトウェア又はハードウェアのいずれかの）ホストエージェントを含むことができる。実施形態は、フロントエンドとバックエンドとの間で仮想化を分割することを含むことができる。

40

【００２０】

実施形態は、集中型ネットワーク及びターゲットストレージコントローラのための様々な展開地点を含むことができる。いくつかの実施形態では、集中型装置をホストコンピュータシステム１０２上に配置し、他の事例では、ディスクを別のボックスに動かす（例えばイーサネット（登録商標）によって、以下の様々なボックスを切り換えるスイッチに接続する）ことができる。仮想化にはレイヤが必要になる場合がある一方で、ストレージリソースとコンピュータリソースとを別個にスケールリングできるようにストレージを分離することができる。また、この時ブレードサーバ（すなわち、ステートレスサーバ）を可能にすることもできる。かつては高価なブレードサーバ及び接続されたストレージエリアネッ

50

トワーク (S A N) を伴っていた設備を代わりにスイッチに接続することができる。実施形態では、この設備が、リソースがラックレベルで分割される「ラックスケール」アーキテクチャを含む。

#### 【0021】

本明細書に開示する方法及びシステムは、様々なタイプの非 D A S ストレージを集中型ネットワーキング / ターゲットストレージ装置内の D A S として仮想化する方法及びシステムを含む。実施形態では、ストレージシステムに対して様々なフロントエンドプロトコルを使用すると同時に、ストレージ装置を O S スタック 1 0 8 に対して D A S として露出して、D A S として望まれるあらゆるストレージを仮想化することができる。

#### 【0022】

本明細書に開示する方法及びシステムは、集中型ネットワーク / ストレージアダプタの仮想化を含む。トラフィックの観点から、システムを 1 つに組み合わせることができる。ストレージアダプタとネットワークアダプタとを組み合わせると仮想化に追加すれば大きな利点を得られる。すなわち、2 つの P C I e バス 1 1 0 を含む単一のホスト 1 0 2 が存在する。P C I e 1 1 0 からルーティングを行うには、R D M A のようなシステムを用いて別のマシン / ホスト 1 0 2 に到達させることができる。これを別個に行った場合、ストレージの R D M A システムとネットワーク R D M A システムとを別個に構成する必要がある。それぞれを結合して 2 つの異なる場所で構成する必要がある。集中型シナリオでは、これが R D M A であり、他のどこかに別のファブリックが存在するとすれば、ストレージとネットワーキングの組み合わせを用いてこれらの 2 つを単一ステップで構成できるので、Q o S を構成するステップ全体がゼロタッチ処理である。すなわち、ストレージが分かれば、Q o S をネットワーク上に別個に構成する必要はない。

#### 【0023】

本明細書に開示する方法及びシステムは、ハードウェア内で、任意に集中型ネットワークアダプタ / ストレージアダプタアプライアンス内で具体化されるネットワーキング機能及びストレージ機能の仮想化及び / 又は間接参照を含む。仮想化は、あるレベルの間接参照であり、プロトコルは、別のレベルの間接参照である。本明細書に開示する方法及びシステムは、多くのオペレーティングシステムがローカルストレージを処理するために使用するのに適した N V M e などのプロトコルを、S A S、S A T A などの別のプロトコルに変換することができる。N V M e などの一貫したインターフェイスを O S 1 0 8 に露出することができ、バックエンドでは、コスト効率の高いストレージ媒体であれば何にでも変換することができる。このことは、ユーザに価格上 / 性能上の利点を与える。コンポーネントが安い / 速い場合には、これらのうちのいずれかのコンポーネントを接続することができる。バックエンドは、N V M e を含む何であってもよい。

#### 【0024】

本明細書では、ネットワーク機能及びストレージ機能のための集中型データ経路を装置に含める方法及びシステムを提供する。別の実施形態は、ネットワーク機能及びストレージ機能のための集中型データ経路をスイッチに提供することができる。

#### 【0025】

実施形態では、本明細書に開示する方法及びシステムが、ストレージ / ネットワークトンネリングを含み、ネットワークを介したストレージシステム間のトンネル経路は、ソースコンピュータ又はターゲットコンピュータのオペレーティングシステムを伴わない。従来のシステムには、別個のストレージ経路とネットワーク経路とが存在し、従ってストレージに遠隔的にアクセスするには、メモリ、I / O バスなどとの間で大々的にコピーを行う必要があった。これらの 2 つの経路を統合するということは、ストレージトラフィックがネットワークの方に真っ直ぐ進むということである。各コンピュータの O S 1 0 8 は、ローカルディスクしか見ない。別の利点は、プログラミングの単純さである。ユーザは、S A N を別個にプログラムする必要がなく、すなわち本明細書に開示する方法は、ワンステッププログラム可能な S A N を含む。ゾーンの発見及び指定などを必要とせずに、暗号化、接続及び分離などを中央でプログラマ的に行うことができる。

10

20

30

40

50

## 【 0 0 2 6 】

本明細書に開示する実施形態は、OS 108がストレージ装置をローカルディスクとして見るようにストレージ装置をOS 108に対して仮想化することを含むことができる。集中型システムは、本明細書に開示する方法及びシステムにおける間接参照のレベルによって、ストレージ媒体の位置だけでなくメディアタイプも隠すことができる。たとえば実際のストレージが遠隔地に存在し、及び/又はSANなどの異なるタイプのものであっても、OSには、ローカルディスクが存在することしか見えない。従って、ストレージが仮想化され、OS 108及びアプリケーションを変更する必要はない。通常は複雑なストレージタイプを裏で構成するために必要な管理、階層化ポリシー、バックアップポリシー及び保護ポリシーなどを全て隠すことができる。

10

## 【 0 0 2 7 】

ストレージの仮想化のどこで間接参照を行うかを選択する方法及びシステムを提供する。いくつかの機能の仮想化は、ハードウェア（例えば、ホスト上のアダプタ、スイッチ、様々なハードウェアフォームファクタ（例えば、FPGA又はASIC））及びソフトウェアで行うことができる。本明細書に開示する方法及びシステムをホストマシン上、トップオブブラックスイッチ上、又はこれらの組み合わせで展開するような異なるトポロジーを利用することができる。選択に通じる要因としては、使いやすさが挙げられる。ステートレスサーバを実行したいと望むユーザは、トップオブブラックを好むと思われる。この方法を気にしないユーザは、ホスト上のコントローラを好むと思われる。

20

## 【 0 0 2 8 】

本明細書に開示する方法及びシステムは、NVMeオーバイーサネット（登録商標）の提供を含む。これらの方法は、装置間で使用されるトンネリングプロトコルの基礎になり得る。NVMeは、従来はローカルPCIeに進むように意図された好適なDASプロトコルである。本明細書に開示する実施形態は、イーサネット（登録商標）を介してNVMeプロトコルトラフィックをトンネリングすることができる。NVMe（不揮発性メモリエクスプレス）は、Linux（登録商標）及びWindowsにおいてPCIeベースのフラッシュストレージへのアクセスを提供するプロトコルである。このプロトコルは、従来のシステムで使用されているソフトウェアスタックを迂回することによって高性能をもたらす。

30

## 【 0 0 2 9 】

本明細書に開示する実施形態は、仮想化されて動的に割り当てられたNVMeデバイスを提供することを含むこともできる。実施形態では、NVMeを重畳させることもできるが、NVMeデバイスを分割し、仮想化して動的に割り当てることもできる。実施形態では、ソフトウェア内にフットプリントは存在しない。オペレーティングシステムは、同じ状態を保つ（集中型ネットワーク/ストレージカードを小型ドライバ）。この結果、仮想ストレージは、直接接続されたディスクのように示されるが、このような装置をネットワークにわたってプールできる点が異なる。

## 【 0 0 3 0 】

本明細書では、ストレージエリアネットワーク（SAN）のような共有の利点と共に、直接接続ストレージ装置（DAS）の単純性を提供する方法及びシステムを提供する。本明細書に開示する様々な実施形態における各集中型アプライアンスはホストとすることができ、あらゆるストレージドライブは、特定のホストに固有のものでありながら（SAN又は他のネットワークアクセス可能ストレージと同様に）他のホストからも見えるようにすることができる。本開示のネットワーク/ストレージコントローラによって可能になる各ボックス内のドライブは、SANと同様に挙動する（すなわち、ネットワーク上で利用可能である）が、管理方法は大幅に単純である。ストレージ管理者がSANを構成する場合、典型的な企業では、「誰が何を見る」など、部門全体がSANのゾーン（例えば、ファイバチャネルスイッチ）を構成することがある。この知識は最初から組み込まれており、ユーザは、SANの管理者にSANを構成する作業を行うように依頼しなければならない。通常のレガシーなSANのアーキテクチャにはプログラマビリティが存在しない。

40

50

本明細書に開示する方法及びシステムは、ネットワーク上に存在するローカルユニットを提供するが、これらのローカルユニットは、ゾーン定義などのような複雑な管理ステップを踏む必要なく、引き続き自機のストレージにアクセスすることができる。これらの装置は、ネットワークとストレージの両方を認識することのみによってS A Nが行うことを実行することができる。従って、これらは、第1のプログラマ的なS A Nを表す。

#### 【0031】

本明細書に開示する方法及びシステムは、集中型ネットワーク及びストレージデータ管理を提供するハードウェアアプライアンスによって可能になる永続的でステートフルな分散的ストレージを含むことができる。

#### 【0032】

本明細書に開示する方法及びシステムは、仮想化のためのコンテナの使用をサポートするように適合された単一アプライアンスにおけるネットワークデータ及びストレージデータ管理の収束を含むこともできる。このような方法及びシステムは、新興的なものではあるがいくつかのさらなる利点を提供するコンテナエコシステムに対応する。

#### 【0033】

本明細書では、N V M eの仮想化を実装する方法及びシステムを開示する。どれほど多くのソースがどれほど多くの宛先に関連しているかに関わらず、ソースからのデータがハブに入る前に最初にシリアル化される限り、ハブは、指定された宛先に順にデータを分配する。そうだとした場合、D M Aエンジンなどのデータトランスポートリソースを1回のみのコピーに低減することができる。この開示は、様々な使用シナリオを含むことができる。1つのシナリオでは、N V M e仮想機能(V F)について、これらが全て同じP C I eバスに接続された場合、どれほど多くのV Fが構成されるかに関わらず、データはこのV Fのプール内に順に到来し、従ってD M Aエンジンは1つしか存在せず、(制御情報のため)1つのストレージブロックしか必要ない。別の使用シナリオでは、分散ディスク/コントローラのプールを有するディスクストレージシステムについて、物理的バス、すなわちP C I eからデータが生じた場合、データはこのディスクのプール内に連続して到来するので、プール内にどれほど多くのディスク/コントローラが存在するかに関わらず、D M Aエンジンなどのトランスポートリソースをコントローラ毎に1つではなく1つのみに低減することができる。

#### 【0034】

様々な例示的かつ非限定的な実施形態によれば、装置が、物理的ターゲットストレージ媒体コントローラと、物理的ネットワークインターフェイスコントローラと、ストレージ媒体コントローラとネットワークコントローラとの間のゲートウェイを含む集中型入力/出力コントローラを含み、ゲートウェイは、ストレージ媒体コントローラとネットワークインターフェイスコントローラとの間のストレージトラフィック及びネットワークトラフィックのための直接接続を提供する。

#### 【0035】

様々な例示的かつ非限定的な実施形態によれば、ストレージ装置の仮想化方法が、第1のストレージプロトコルでの命令に応答する物理ストレージ装置にアクセスするステップと、第1のストレージプロトコルと第2ストレージプロトコルとの間で命令を変換するステップと、第2のプロトコルを用いて物理ストレージ装置をオペレーティングシステムに提示することにより、オペレーティングシステムを使用するホストコンピュータシステムに対して物理ストレージ装置が局所的に存在するか、それとも遠隔的に存在するかに関わらず、物理ストレージ装置のストレージが動的に展開されるようにするステップとを含む。

#### 【0036】

様々な例示的かつ非限定的な実施形態によれば、アプリケーション及びコンテナの少なくとも一方の移動を容易にする方法が、集中型ストレージ及びネットワーキングコントローラを提供して、ゲートウェイが、ホストコンピュータのオペレーティングシステムの介入を伴わずに装置のストレージコンポーネントとネットワーキングコンポーネントとの間

10

20

30

40

50

のネットワーク及びストレージトラフィックのための接続を提供するステップと、少なくとも1つのアプリケーション又はコンテナを、集中型ストレージ及びネットワークコントローラによって制御されるターゲット物理ストレージ装置にマッピングすることにより、アプリケーション又はコンテナが別のコンピュータシステムに移動する際に、アプリケーション又はコンテナが、ターゲット物理ストレージが接続されたホストコンピュータのオペレーティングシステムの介入を伴わずにターゲット物理ストレージにアクセスできるようにするステップとを含む。

#### 【0037】

様々な例示的かつ非限定的な実施形態によれば、ネットワークのサービス品質（QoS）を提供する方法が、集中型ストレージ及びネットワークコントローラを提供して、ゲートウェイが、ホストコンピュータのCPU上で実行されるオペレーティングシステムの介入、ハイパーバイザの介入、又はその他のソフトウェアの介入を伴わずに装置のストレージコンポーネントとネットワークコンポーネントとの間のネットワークトラフィック及びストレージトラフィックのための接続を提供するステップと、集中型ストレージ及びネットワークコントローラによって処理されるストレージトラフィック及びネットワークトラフィックの少なくとも一方に基づいて、ホストコンピュータのCPU上で実行されるオペレーティングシステムの介入、ハイパーバイザの介入、又はその他のソフトウェアの介入を伴わずにストレージ及びネットワークコントローラが展開されたデータ経路を有するネットワークに関連する少なくとも1つのサービス品質（QoS）パラメータを管理するステップとを含む。

#### 【0038】

QoSは、帯域幅パラメータ、ネットワークレイテンシパラメータ、I/O性能パラメータ、スループットパラメータ、ストレージタイプパラメータ及びストレージレイテンシパラメータのうちの1つ又は2つ以上などの様々なパラメータに基づくことができる。QoSは、集中型ストレージ及びネットワークコントローラを介してストレージによってサービスを受けるアプリケーション及びコンテナの少なくとも一方がホストコンピュータから別のコンピュータに移動する際に自動的に維持することができる。同様に、QoSは、集中型ストレージ及びネットワークコントローラを介してアプリケーション及びコンテナの少なくとも一方にサービスを提供する少なくとも1つのターゲットストレージ装置が第1の位置から別の位置又は複数の位置に移動する際に自動的に維持することができる。例えば、要件が増えた際にストレージニーズを満たすように、ストレージ装置をスケーリングし、或いは異なるストレージ媒体タイプを選択することができる。実施形態では、ネットワークトラフィックデータの暗号化、ストレージ装置内のデータの暗号化、又はこれらの両方などのセキュリティ機能を提供することができる。圧縮、保護レベル（例えば、RAIDレベル）、異なるストレージ媒体タイプの使用、グローバル重複排除、並びに目標復旧時点（RPO）及び目標復旧時間（RTO）の少なくとも一方を達成するためのスナップショット間隔などの様々なストレージ機能を提供することもできる。

#### 【0039】

独立した図を通じて同一又は機能的に同様の要素を同じ参照番号によって示す、以下の詳細な説明と共に本明細書に組み込まれて本明細書の一部を成す添付図面は、様々な実施形態をさらに示し、本明細書に開示するシステム及び方法による全ての様々な原理及び利点を説明する役割を果たす。

#### 【図面の簡単な説明】

#### 【0040】

【図1】例示的かつ非限定的な実施形態による全体的アーキテクチャを示す図である。

【図2】例示的かつ非限定的な実施形態によるコンピュータシステムを示す図である。

【図3】例示的かつ非限定的な実施形態による集中型ソリューションを示す図である。

【図4】例示的かつ非限定的な実施形態による、集中型ソリューションによって可能になる2つのコンピュータシステムを示す図である。

【図5】例示的かつ非限定的な実施形態による集中型コントローラを示す図である。

【図 6】例示的かつ非限定的な実施形態による集中型コントローラの配置を示す図である。

【図 7】例示的かつ非限定的な実施形態による複数のシステムを示す図である。

【図 8】例示的かつ非限定的な実施形態によるフィールドプログラマブルゲートアレイ (FPGA) を示すブロック図である。

【図 9】例示的かつ非限定的な実施形態によるコントローラカードのアーキテクチャを示す図である。

【図 10】例示的かつ非限定的な実施形態によるソフトウェアスタックを示す図である。

【図 11】例示的かつ非限定的な実施形態による複数のシステムにわたるアプリケーションコンテナの移動を示す図である。

10

【図 12】例示的かつ非限定的な実施形態による複数のシステムにわたるアプリケーションコンテナの移動を示す図である。

【図 13】例示的かつ非限定的な実施形態による複数のシステムにわたるアプリケーションコンテナの移動を示す図である。

【図 14】例示的かつ非限定的な実施形態による複数のシステムにわたるアプリケーションコンテナの移動を示す図である。

【図 15】例示的かつ非限定的な実施形態による複数のシステムにわたるアプリケーションコンテナの移動を示す図である。

【図 16】例示的かつ非限定的な実施形態によるパケット送信を示す図である。

【図 17】例示的かつ非限定的な実施形態によるストレージアクセススキームを示す図である。

20

【図 18】例示的かつ非限定的な実施形態によるファイルシステムの動作を示す図である。

【図 19】例示的かつ非限定的な実施形態による分散ファイルサーバの動作を示す図である。

【図 20】例示的かつ非限定的な実施形態による高性能分散ファイルサーバ (DFS) を示す図である。

【図 21】例示的かつ非限定的な実施形態によるシステムを示す図である。

【図 22】例示的かつ非限定的な実施形態によるホストを示す図である。

【図 23】例示的かつ非限定的な実施形態によるデータブロックにアクセスするアプリケーションを示す図である。

30

【図 24】例示的かつ非限定的な実施形態によるデータブロックにアクセスするアプリケーションを示す図である。

【図 25】例示的かつ非限定的な実施形態によるシステムを示す図である。

【図 26】例示的かつ非限定的な実施形態による方法を示す図である。

【図 27】例示的かつ非限定的な実施形態による方法を示す図である。

【図 28】例示的かつ非限定的な実施形態による方法を示す図である。

【発明を実施するための形態】

【0041】

当業者であれば、図中の要素は単純性及び明確性を目的として示しており、必ずしも縮尺通りではないと理解するであろう。例えば、図中の要素には、本明細書に開示するシステム及び方法の実施形態をより良く理解する役に立つように他の要素に対して寸法を誇張しているものもある。

40

【0042】

以下、添付図面及び添付書類を参照しながら本開示の様々な例示的かつ非限定的な実施形態を説明することによって本開示を詳細に説明する。しかしながら、本開示は多くの異なる形で具体化することができ、本明細書に示す例示的な実施形態に限定されると解釈すべきではない。むしろ、これらの実施形態は、本開示を徹底的なものとして当業者に本開示の概念を十分に伝えるように提供するものである。本開示の実際の範囲を確認するには、特許請求の範囲を参照すべきである。

50

## 【 0 0 4 3 】

本明細書に開示するシステム及び方法による実施形態を詳細に説明する前に、これらの実施形態は、主に集中型ネットワーキング及びストレージに関連する方法ステップ及び／又はシステムコンポーネントの組み合わせの形で存在することを認められたい。従って、当業者に容易に明らかになる詳細によって本開示が曖昧にならないように、図面には、必要に応じてこれらのシステムコンポーネント及び方法ステップを従来の記号によって表し、本明細書に開示するシステム及び方法の実施形態を理解することに関する特定の詳細のみを示している。

## 【 0 0 4 4 】

図 3 に示すように、集中型ソリューション 3 0 0 は、3 つの重要な態様を含むことができ、ハードウェア、ソフトウェアモジュール及び機能の組み合わせを含むハードウェア装置で実装することができる。第 1 に、ネットワークコントローラ 1 1 8 とストレージコントローラ 1 1 2 の間にカットスルーデータ経路 3 0 4 を提供することにより、ストレージとネットワークとの間のアクセスを、OS スタック 1 0 8、P C I e バス 1 1 0 又は C P U 1 0 6 のいずれの介入も必要とせず直接的に行えるようにすることができる。第 2 に、ストレージとローカルホスト上のエンティティとの間のアクセスなどの、ストレージ装置 3 0 2 などへのカットスルーストレージスタックアクセスを提供して、ストレージにアクセスするために S C S I / S A S / S A T A スタックなどの複雑なレガシーソフトウェアスタックを迂回できるようにすることができる。第 3 に、ネットワークを横切るデータ転送を予約及びスケジューリングする機構などによってネットワークのエンドツーエンド輻輳管理及びフロー制御を提供することにより、ターゲットのデータを遠隔インシエータが利用できることを保証して、中間ネットワークファブリックスイッチを流れる際のトラフィックの輻輳を最小化することができる。第 1 及び第 2 の態様は、データ経路からソフトウェアスタック（従って、C P U 1 0 6 及びメモリ）を削除し、冗長又は不要な動き及び処理を排除する。エンドツーエンド輻輳管理及びフロー制御は、確定的で信頼性の高いデータ転送を提供する。

## 【 0 0 4 5 】

上述したように、集中型ソリューション 3 0 0 の 1 つの利点は、オペレーティングシステムスタック 1 0 8 が従来の P C I e 1 1 0 又は同様のバスを介して集中型ソリューション 3 0 0 に接続することにより、たとえ物理ストレージが遠隔地に存在する場合でも、OS スタック 1 0 8 が、集中型ソリューション 3 0 0 と、ストレージ装置 3 0 2 へのカットスルーを通じて制御するあらゆるストレージとを 1 又は 2 以上のローカルな永続装置として見るようになる点である。とりわけ、この利点は、D A S 3 0 8 を仮想化する能力を含み、この能力は、ファブリックを介して D A S 3 0 8 を仮想化し、すなわち D A S 3 0 8 ストレージシステムを取り出し、コンピュータシステム 1 0 2 の外部に動かしてネットワーク上に置くことを含むことができる。集中型ソリューション 3 0 0 のストレージコントローラ 1 1 2 は、S A S、S A T A 又は N V M e などの様々な既知のプロトコルを介してネットワーク 1 2 2 上の D A S 3 0 8 に接続して制御することができる。実施形態では、仮想化が、D A S 3 0 8 を任意の名前空間に分割することを含むことができる。実施形態では、オペレーティングシステムが、仮想 D A S 3 0 8 に実際のローカル物理 D A S であるかのようにアクセスすることができ、例えば OS 1 0 8 が、N V M e などの標準プロトコルを介し、P C I e バス 1 1 0 を介して集中型ソリューション 3 0 0 のストレージコントローラ 1 1 2 にアクセスすることができる。この場合も、OS 1 0 8 は、ソリューション 3 0 0 全体を D A S などのローカル物理装置として見る。従って、本明細書では、たとえば実際にはストレージがネットワーク 1 2 2 を介してアクセスされる場合でも、OS 1 0 8 があらゆるストレージタイプを D A S と見なし、OS 1 0 8 がローカル物理ストレージの場合に必要とされる以外のことを行う必要がないようにする（D A S 及び S A N 3 1 0 などの他のストレージタイプを含む）ストレージ仮想化能力を提供する。ストレージ装置 3 0 2 が S A N 3 1 0 ストレージである場合、集中型ソリューションのストレージコントローラ 1 1 2 は、インターネットスモールコンピュータシステムインターフェイス（i S

10

20

30

40

50

C S I )、ファイバチャネル ( F C )、又はファイバチャネルオーバーサネット (登録商標) ( F C o E ) などのストレージエリアネットワークに使用するのに適したプロトコルを通じて S A N 3 1 0 を制御することができる。従って、集中型ソリューション 3 0 0 は、O S スタック 1 0 8 を、とりわけイーサネット (登録商標)、S A S、S A T A、N V M e、i S C S I、F C 又は F C o E などのストレージにおいて使用されている他のプロトコルのいずれかから、異なるストレージタイプ及びプロトコルが P C I e 1 1 0 を介してアクセスできるローカルストレージのように見えるようにする N V M e のような単純なプロトコルに変換する。さらに、この変換は、(あらゆる種類のターゲットストレージシステムを参照する) ストレージアダプタの仮想化も可能にする。従って、本明細書に開示する方法及びシステムは、様々なタイプの非 D A S ストレージを集中型ネットワーキング / ターゲットストレージ装置 3 0 0 内の D A S として仮想化する方法及びシステムを含む。実施形態では、ストレージシステムに対して様々なプロトコルを使用すると同時に、ストレージ装置を O S スタック 1 0 8 に対して D A S として露出して、D A S として望まれるあらゆるストレージを仮想化することができる。従って、本明細書では、N V M e 及び同様のプロトコルを使用するようなストレージ装置を仮想化し、これらの仮想装置を、S A T A を使用するような異なる物理装置に変換する方法及びシステムを提供する。

10

20

30

40

50

#### 【 0 0 4 6 】

ネットワーク 1 2 2 を介したストレージシステム間のトンネル経路がソース又はターゲットコンピュータのオペレーティングシステムに関与しないストレージ / ネットワークトンネリング 3 0 4 は、いくつかの利点を可能にする。従来のシステムには、別個のストレージ経路とネットワーク経路とが存在し、従ってストレージに遠隔的にアクセスするには、メモリ、I / O バスなどとの間で大々的にコピーを行う必要があった。これらの2つの経路を統合するということは、ストレージトラフィックがネットワークの方に真っ直ぐ進むということである。プログラミングの単純さが利点である。ユーザは、S A N 3 1 0 を別個にプログラムする必要がなく、すなわち本明細書に開示する方法は、ワンステッププログラム可能な S A N 3 1 0 を可能にする。ゾーンの発見及び指定などを必要とせずに、構成、暗号化、接続及び分離などを中央でプログラマ的行うことができる。一例として、典型的な S A N は、「イニシエータ」、「ターゲット」、及びイニシエータとターゲットとを接続するスイッチファブリックで構成される。通常、どのイニシエータがどのターゲットを見るかは、「ゾーン」と呼ばれるファブリックスイッチによって規定 / 制御される。従って、イニシエータ又はターゲットが移動した場合には、ゾーンを更新する必要がある。通常、S A N の第 2 の制御部分は「ターゲット」と共に存在する。これらは、どのイニシエータポートがどの論理ユニット番号 ( L U N ) (ターゲットによって露出されたストレージユニット) を見ることができるかを制御することができる。通常、この制御は、L U N マスキング及び L U N マッピングと呼ばれる。この場合も、イニシエータが移動した場合には「ターゲット」を再プログラムする必要がある。このような環境では、アプリケーションが (フェイルオーバ、ロードリバランシングなどに起因して) 1 つのホストから別のホストに移動した場合にゾーニング及び L U N マスキング / マッピングを更新する必要がある。或いは、全てのイニシエータが全てのターゲットを見るように S A N を予めプログラムしておくこともできる。しかしながら、これを行うと、S A N をスケールングできずセキュアでなくなる。本開示を通じて説明する別のソリューションでは、アプリケーション、コンテナ又はストレージ装置のこのような移動に S A N の再プログラミングが必要なく、ゼロタッチソリューションとなる。集中型ソリューション 3 0 0 によって維持され実行されるマッピングは、アプリケーション又はコンテナ、ターゲットストレージ媒体、又はこれらの両方の、ホスト C P U 上で実行される O S、ハイパーバイザ、又はその他のソフトウェアによる介入を伴わない (複数の位置を含めて) 個別の移動及びスケールングを可能にすることができる。

#### 【 0 0 4 7 】

O S 1 0 8 がストレージをローカルディスクとして見ることにより、ストレージの単純な仮想化が可能になる。集中型システム 3 0 0 は、本明細書に開示する方法及びシステム



における間接参照のレベルによって、ストレージ媒体の位置だけでなくメディアタイプも隠すことができる。たとえば実際のストレージが遠隔地に存在し、及び/又はS A N 3 1 0などの異なるタイプのものであっても、O S 1 0 8には、ローカルディスクが存在することしか見えない。従って、集中型ソリューション3 0 0を通じてストレージが仮想化され、O S 1 0 8及びアプリケーションを変更する必要はない。通常は複雑なストレージタイプを裏で構成するために必要な管理、階層化ポリシー、バックアップポリシー及び保護ポリシーなどを全て隠すことができる。

#### 【0 0 4 8】

集中型ソリューション3 0 0は、ストレージエリアネットワーク(S A N)の利点と共に、直接接続ストレージ(D A S)の単純性を可能にする。本明細書に開示する様々な実施形態における各集中型アプライアンス3 0 0はホストとして機能することができ、あらゆるストレージ装置3 0 2は、特定のホストに固有のものでありながら(S A N 3 1 0又は他のネットワークアクセス可能ストレージと同様に)他のホストからも見えるようにすることができる。本開示のネットワーク/ストレージコントローラによって可能になる各ボックス内のドライブは、S A N 3 1 0と同様に挙動する(例えば、ネットワーク上で利用可能である)が、管理方法は大幅に単純である。ストレージ管理者が通常通りにS A Nを構成する場合、典型的な企業では、「誰が何を見る」など、部門全体がS A N 3 1 0のゾーン(例えば、ファイバチャネルスイッチ)を構成することがある。この知識は最初から組み込まれていなければならない。ユーザは、S A Nの管理者にS A Nを構成する作業を行うように依頼しなければならない。通常のレガシーなS A N 3 1 0のアーキテクチャにはプログラマビリティが存在しない。本明細書に開示する方法及びシステムは、ネットワーク上に存在するローカルユニットを提供するが、これらのローカルユニットは、ゾーン定義などのような複雑な管理ステップを踏む必要なく、引き続き自機のストレージにアクセスすることができる。これらの装置は、ネットワークとストレージの両方を認識することのみによってS A Nが行うことを実行することができる。従って、これらは、第1のプログラマ的なS A Nを表す。

#### 【0 0 4 9】

ソリューション3 0 0は、ストレージ媒体3 0 2及びネットワーク1 2 2の両方を制御する「集中型I Oコントローラ」として説明することができる。この集中型コントローラ3 0 0は、ストレージコントローラ1 1 2とネットワークコントローラ(N I C)1 1 8とを単純に統合したものではない。実際のストレージ機能及びネットワーク機能は、ネットワークインターフェイスとの間をデータがトラバースする際にストレージ機能が実行されるように融合される。これら機能は、後述するF P G A(1又は2以上)又はA S I C(1又は2以上)などのハードウェアソリューションにおいて提供することができる。

#### 【0 0 5 0】

図4を参照すると、集中型ソリューション3 0 0によって可能になる2又は3以上のコンピュータシステム1 0 2は、それぞれのストレージターゲットのホストとしての役割を果たすことができ、ストレージとネットワークを融合して両インターフェイスを制御することにより、集中型ソリューション3 0 0によって可能になる別のコンピュータシステム1 0 2へのポイントツーポイント経路4 0 0又はイーサネット(登録商標)スイッチ4 0 2などによって内部バス又はC P U /ソフトウェアの機能をトラバースすることなく、ネットワーク1 2 2を介してストレージ装置3 0 2に遠隔的に直接アクセスすることができる。最高の性能(高I O P及び低レイテンシ)を達成することができる。さらに、この時点でクラスタ全体にわたってストレージリソース3 0 2をプールすることができる。図4には、このことを点線の楕円形4 0 0によって概念的に示している。

#### 【0 0 5 1】

実施形態では、集中型ソリューション3 0 0を、図1に示すような従来のコンピュータシステムの様々なコンポーネント、及び図3に関して説明した集中型I Oコントローラ3 0 0と共にホストコンピュータシステム1 0 2に含めることができる。図5を参照すると、別の実施形態では、集中型コントローラ3 0 0をトップオブブラックスイッチなどのスイ

ッチ内に配置し、従ってストレージ対応スイッチ500を可能にすることができる。このスイッチは、ネットワーク122上に存在することができ、従来のコンピュータシステム102のネットワークコントローラなどのネットワークコントローラ118によってアクセスすることができる。

#### 【0052】

図6を参照すると、1又は2以上のホストコンピュータシステム102、並びに集中型ソリューション300によって可能になるシステム102と非可能システム102とに接続できるストレージ対応スイッチ500の両方に集中型コントローラ300が配置されたシステムを展開することができる。上述したように、ホストコンピュータシステム102及びストレージ対応スイッチ500における(単複の)集中型コントローラ300のためのターゲットストレージ302は、ネットワークを越えて互いに見えており、例えば仮想化ソリューションなどへの統一リソースとして扱われる。要するに、本開示の様々な別の実施形態におけるホストシステム、スイッチ、又はこれらの両方に、同じ装置上における集中型ネットワーク及びストレージトラフィックの処理を含む知性を配置することができる。

10

#### 【0053】

従って、本明細書に開示する実施形態は、スイッチフォームファクタ又はネットワークインターフェイスコントローラを含むことができ、或いはこれらの両方は、(ソフトウェア又はハードウェアのいずれかの)ホストエージェントを含むことができる。これらの様々な展開は、ホスト上及び/又はスイッチ上、及び/又はフロントエンドとバックエンドとの間などにおける仮想化能力の分散を可能にする。いくつかの機能を仮想化するためにレイヤが必要になることもある一方で、ストレージリソースとコンピュータリソースとを別個にスケールできるようにストレージを分離することもできる。また、ブレードサーバ(すなわち、ステートレスサーバ)を可能にすることもできる。かつては高価なブレードサーバ及び接続されたストレージエリアネットワーク(SAN)を伴っていた設備を、代わりにストレージ対応スイッチ500に接続することができる。実施形態では、この設備が、リソースがラックレベルで分割される「ラックスケール」アーキテクチャを含む。

20

#### 【0054】

ストレージの仮想化のどこで間接参照を行うかを選択する方法及びシステムを提供する。いくつかの機能の仮想化は、ハードウェア(例えば、ホスト102上の集中型アダプタ300、ストレージ対応スイッチ500、様々なハードウェアフォームファクタ(例えば、FPGA又はASIC))及びソフトウェアで行うことができる。本明細書に開示する方法及びシステムをホストマシン102上、トップオブラックスイッチ500上、又はこれらの組み合わせで展開するような異なるトポロジーを利用することができる。仮想化を行うべき場所の選択に通じる要因としては、使いやすさが挙げられる。ステートレスサーバを実行したいと望むユーザは、トップオブラックストレージ対応スイッチ500を好むと思われる。この方法を気にしないユーザは、ホスト102上の集中型コントローラ300を好むと思われる。

30

#### 【0055】

図7は、2つのコンピュータシステム102(コンピュータシステム1及びコンピュータシステム2)を含む集中型コントローラ300と、ストレージ対応スイッチ500とを用いて可能になる1組のシステムのさらに詳細な図である。DAS308及びSAN310などのストレージ装置302は、集中型コントローラ300又はストレージ対応スイッチ500によって制御することができる。DAS308は、SASプロトコル、SATAプロトコル又はNVMeプロトコルのいずれかを使用する場合に制御することができる。SAN310は、iSCSI、FC又はFCoEのいずれかを使用する場合に制御することができる。ストレージコントローラ300を有するホスト102間の接続は、ポイントツーポイント経路400を介したものの、イーサネット(登録商標)スイッチ402を介したものの、又は従来のコンピュータシステムへの接続も提供できるストレージ対応スイッチ5

40

50

00を介したものとすることができる。上述したように、知的集中型コントローラ300を含む複数のシステムの各々は、ホストとしての役割と、他のホストが見るストレージターゲットロケーションとしての役割とを果たすことにより、コンピュータシステム102のオペレーティングシステム108のために単一のストレージクラスタとして取り扱われるオプションを提供することができる。

#### 【0056】

本明細書に開示する方法及びシステムは、ハードウェア集中型コントローラ300内で、任意に集中型ネットワークアダプタ/ストレージアダプタアライアンス300内で具体化されるネットワーキング機能及びストレージ機能の仮想化及び/又は間接参照を含む。仮想化は、あるレベルの間接参照であり、プロトコルは、別のレベルの間接参照である。本明細書に開示する方法及びシステムは、多くのオペレーティングシステムがローカルストレージを処理するために使用するのに適したNVMeなどのプロトコルを、SAS、SATAなどの別のプロトコルに変換することができる。NVMeなどの一貫したインターフェイスをOS108に露出することができ、集中型コントローラ300の他方の側では、コスト効率の高いストレージ媒体302であれば何にでも変換することができる。このことは、ユーザに価格上/性能上の利点を与える。コンポーネントが安い/速い場合には、これらのうちのいずれかのコンポーネントを接続することができる。集中型コントローラ300の側は、NVMeを含むあらゆる種類のストレージに面することができる。さらに、ストレージ媒体タイプは、以下に限定されるわけではないが、HDD、(SLC、MLC又はTLCフラッシュベースの)SSD、RAM、その他、又はこれらの組み合わせのいずれかとするすることができる。

10

20

#### 【0057】

実施形態では、集中型コントローラを、NVMe仮想機能を仮想化し、イーサネット(登録商標)スイッチ402を介してNVMeを通じてストレージ対応スイッチ500に接続されたような遠隔ストレージ装置302へのアクセスを提供するように適合させることができる。従って、集中型ソリューション300は、NVMeオーバイーサネット(登録商標)700、すなわちNVMe over Eの使用を可能にする。従って、本明細書に開示する方法及びシステムは、NVMeオーバイーサネット(登録商標)の提供を含む。これらの方法は、集中型コントローラ300及び/又はストレージ対応スイッチ500によって可能になるホストコンピュータシステム102などの装置間で使用されるトンネリングプロトコルの基礎になり得る。NVMeは、従来ローカルPCIe110に進むように意図されている好適なDASプロトコルである。本明細書に開示する実施形態は、イーサネット(登録商標)を介してNVMeプロトコルトラフィックをトンネリングすることができる。NVMe(不揮発性メモリエクスプレス)は、Linux(登録商標)及びWindowsにおいてPCIeベースのフラッシュへのアクセスを提供するプロトコルである。このプロトコルは、従来のシステムで使用されているソフトウェアスタックを迂回することによって高性能をもたらすと同時に、イーサネット(登録商標)を介して他の装置にトンネリングされる(OSスタック108によって使用されるような)NVMe及びトラフィックを変換する必要性を回避する。

30

#### 【0058】

図8は、IOコントローラカード上に存在して集中型ソリューション300の実施形態を可能にすることができるFPGA800のブロック図である。なお、単一のFPGA800を示しているが、様々な機能ブロックは、複数のFPGA、1又は2以上の顧客特定用途向け集積回路(ASIC)などにまとめることもできる。例えば、別個の(ただし相互接続されている)FPGA又はASICにおいて、様々なネットワーキングブロック及び様々なストレージブロックを取り扱うことができる。本開示を通じて、FPGA800についての言及は、文脈によって別途示されている場合を除き、図8に反映する機能及び同様の機能を可能にすることができる他の形態ハードウェアも含むと理解されたい。また、ネットワーキング機能及び/又はストレージ機能などのいくつかの機能グループは、マーチャントシリコン内に具体化することもできる。

40

50

## 【 0 0 5 9 】

図 8 の F P G A 8 0 0 の実施形態は、4つの主なインターフェイスを有する。第 1 に、ホストコンピュータ 1 0 2 の P C I e バス 1 1 0 などへの P C I e インターフェイスが存在する。従って、このカードは P C I e エンドポイントである。第 2 に、D R A M / N V R A M インターフェイスが存在する。例えば、外部 D R A M 又は N V R A M に、組み込み C P U、メタデータ及びデータ構造、並びにパケット / データバッファリングによって使用される D D R インターフェイスを提供することができる。第 3 に、D A S 3 0 8 及び S A N 3 1 0 などの、媒体へのストレージインターフェイスが存在する。ストレージインターフェイスは、S A S、S A T A、N V M e、i S C S I、F C 及び / 又は F C o E のためのインターフェイスを含むことができ、実施形態では、S A N 3 1 0 のようなネットワーク対応ストレージに固有の、又はネットワーク対応ストレージへのカットスルーを介した、回転媒体、フラッシュ又はその他の永続的形態のストレージへのいずれかのインターフェイスとすることができる。第 4 に、ネットワークファブリックへのイーサネット (登録商標) などのネットワークインターフェイスが提供される。N V M e オーバイーサネット (登録商標) を可能にするには、ストレージインターフェイス及びネットワークインターフェイスを部分的に使用することができる。

10

## 【 0 0 6 0 】

F P G A 8 0 0 の内部機能は、集中型ソリューション 3 0 0 のいくつかの有効化機能、及び全体を通じて説明する本開示の他の態様を含むことができる。ホストには、一連の仮想エンドポイント (v N V M e) 8 0 2 を提供することができる。これにより、ネットワークインターフェイスに使用される S R - I O V プロトコルと同様に、ホストに仮想ストレージターゲットが提供される。この F P G A 8 0 0 の実施形態では、N V M e が低ソフトウェアオーバーヘッドの利点を有し、これによってさらに高性能がもたらされる。仮想 N V M e デバイス 8 0 2 には、割り当て / 割り当て解除 / 移動及びサイズ変更を動的に行うことができる。S R - I O V と同様に、P C I e ドライバ 1 1 0 (以下を参照) にインターフェイス接続する 1 つの物理的機能 (P F) 8 0 6 と、各々が N V M e デバイスのように見える複数の仮想機能 8 0 7 (V F) とが存在する。

20

## 【 0 0 6 1 】

F P G A 8 0 2 の機能には、本明細書では D M A エンジン 8 0 4 と呼ぶこともある 1 又は 2 以上の読み取り及び書き込み直接メモリアクセス (D M A) キュー 8 0 4 も提供される。これらは、割り込みキュー、ドアベル及びその他の標準機能を含んでホストコンピュータシステム 1 0 2 との間で D M A を実行することができる。

30

## 【 0 0 6 2 】

F P G A 8 0 0 の装置マッピング機構 8 0 8 は、仮想 N V M e デバイス 8 0 2 の位置を特定することができる。この位置は、任意に局所的に存在する (すなわち、図示のストレージ媒体インターフェイス 8 2 4 の 1 つに接続される) ことも、或いはストレージコントローラ 3 0 0 の別のホスト 1 0 2 上に遠隔的に存在することもある。遠隔 v N V M e デバイスにアクセスするには、トンネル 8 2 8 を通じてネットワーク 1 2 2 に進む必要がある。

40

## 【 0 0 6 3 】

N V M e 仮想化機構 8 1 0 は、N V M e プロトコル命令及び動作を、D A S 3 0 8 を使用する場合には S A S 又は S A T A などの (バックエンドストレージ媒体 3 0 2 の N V M e を使用する場合には変換は不要)、或いはバックエンドで S A N 3 1 0 ストレージを使用する場合には i S C S I、F C 又は F C o E などの、バックエンドストレージ媒体 3 0 2 の対応するプロトコル及び動作に変換することができる。本明細書におけるバックエンドについての言及は、集中型コントローラ 3 0 0 のホスト 1 0 2 とは逆の側を意味する。

## 【 0 0 6 4 】

データ変換機能 8 1 2 は、ストレージ媒体 3 0 2 にデータが記憶された時に、このデータをフォーマットすることができる。これらの動作は、書き換え、変換、圧縮、(R A I D などの) 保護、暗号化、及び適用可能なタイプのターゲットストレージ媒体 3 0 8 によ

50

って処理されるために必要なあらゆる方法でのデータフォーマットの変更を含むその他の機能を含むことができる。いくつかの実施形態では、ストレージ媒体 308 が遠隔地に存在することができる。

#### 【0065】

実施形態では、ストレージ読み取りキュー及び書き込みキュー 814 が、転送中の中継データのデータ構造又はバッファリングを含むことができる。実施形態では、(FPGA 800 から離れて存在することができる) NVRAM の DRAM などの一時ストレージ装置を使用してデータを一時的に記憶することができる。

#### 【0066】

ローカルストレージスケジューラ及びシェイパー 818 は、ストレージ媒体 302 へのアクセスを優先させて制御することができる。スケジューラ及びシェイパー 818 では、ストリクトプライオリティー、重み付きラウンドロビンスケジューリング、IOPシェイパー、並びにキュー毎、イニシエータ毎、ターゲット毎又はcグループ毎などに適用できるポリサーを含むことができる、ローカルストレージのためのあらゆる適用可能なSLAポリシーを実行することができる。

10

#### 【0067】

データ配置機構 820 は、ストレージ媒体 302 上でどのようにデータをレイアウトするかを決定するアルゴリズムを実装することができる。このアルゴリズムは、媒体全体にわたるストライピング、単一の装置 302 へのローカライジング、装置 302 のサブセットの使用、又は装置 302 の特定のブロックへのローカライジングなどの、当業者に周知の様々な配置スキームを伴うことができる。

20

#### 【0068】

ストレージメタデータ管理機構 822 は、データ配置、ブロックiノード及びオブジェクトiノード、圧縮、重複排除及び保護のためのデータ構造を含むことができる。メタデータは、FPGA 800 から離れたNVRAM / DRAM、又はストレージ媒体 302 のいずれかに記憶することができる。

#### 【0069】

複数の制御ブロック 824 は、ストレージ媒体へのインターフェイスを提供することができる。これらの制御ブロックは、他の考えられる制御ブロックの中でも特に、適当なタイプのターゲットストレージ媒体 302 に必要な場合毎に、SAS、SATA、NVMe、PCIe、iSCSI、FC 及び / 又はFCoEを含むことができる。

30

#### 【0070】

FPGA 800 のストレージネットワークトンネル 828 は、本開示を通じて集中型ソリューション 300 に関連して説明するトンネリング/カットスルー能力を提供することができる。とりわけ、トンネル 828 は、ストレージトラフィックとネットワークトラフィックとの間のゲートウェイを提供する。これには、カプセル化/脱カプセル化又はストレージトラフィック、データの書き換え及びフォーマット、及びデータ転送のエンドツーエンド調整が含まれる。この調整は、ホストコンピュータシステム 102 又は複数のコンピュータシステム 102 内の、図 4 に関連して説明したポイントツーポイント経路 404 などのためのノードを横切るFPGA 800 間の調整とすることができる。シーケンス番号、パケット喪失、タイムアウト及び再送信などの様々な機能を実行することができる。トンネリングは、FCoE 又はNVMe over Ethernet によるものを含め、イーサネット(登録商標)を介して行うことができる。

40

#### 【0071】

仮想ネットワークインターフェイスカード機構 830 は、仮想ネットワークインターフェイスカードとして提供されるホスト 102 への複数のSR-IOVエンドポイントを含むことができる。1つの物理的機能(PF) 836 は、PCIeドライバ110(以下のソフトウェアの説明を参照)と、各々がネットワークインターフェイスカード(NIC) 118のように見える複数の仮想機能(VF) 837 とに接続することができる。

#### 【0072】

50

一連の受信 / 送信 D M A キュー 8 3 2 は、割り込みキュー、ドアベル及びその他の標準機能を含んでホスト 1 0 2 との間で D M A を実行することができる。

【 0 0 7 3 】

分類子及びフロー管理機構 8 3 4 は、典型的には I E E E 標準 8 0 2 . 1 Q サービスクラス ( C O S ) マッピング、又はその他の優先順位レベルへの標準的なネットワークトラフィック分類を実行することができる。

【 0 0 7 4 】

アクセス制御及び書き換え機構 8 3 8 は、典型的にはイーサネット (登録商標) タブル ( M A C S A / D A 、 I P S A / D A 、 T C P ポートなど) 上でパケットの再分類又は書き換えを行うように動作するアクセス制御リストを含むアクセス制御リスト ( A C L ) 及び書き換えポリシーに対処することができる。

10

【 0 0 7 5 】

転送機能 8 4 0 は、レイヤ 2 ( L 2 ) 又はレイヤ 3 ( L 3 ) 機構などを通じてパケットの宛先を決定することができる。

【 0 0 7 6 】

一連のネットワーク受信及び送信キュー 8 4 2 は、データ構造又はネットワークインターフェイスへのバッファリングに対処することができる。パケットデータには、 F P G A 8 0 0 から離れた D R A M を使用することができる。

【 0 0 7 7 】

ネットワーク / リモートストレージスケジューラ及びポリサー 8 4 4 は、優先順位を提供してネットワークインターフェイスへのアクセスを制御することができる。ここでは、ストリクトプライオリティー、重み付きラウンドロビン、I O P シェイパー及び帯域幅シェイパー、並びにキュー毎、イニシエータ毎、ターゲット毎、C グループ毎、又はネットワークフローベース毎などのポリサーを含むことができる、リモートストレージ及びネットワークトラフィックのための S L A ポリシーを実行することができる。

20

【 0 0 7 8 】

ローカルネットワークスイッチ 8 4 8 は、宛先が F P G A 8 0 0 又はホスト 1 0 2 に固有のものである場合、トラフィックが F P G A 8 0 0 から出てネットワークファブリック 1 2 2 に進まなくても済むように、F P G A 内のキュー間でパケットを転送することができる。

30

【 0 0 7 9 】

エンドツーエンド輻輳制御 / クレジット機構 8 5 0 は、ネットワーク輻輳を防ぐことができる。この動作は、2つのアルゴリズムを用いて行われる。第 1 に、リモート F P G A 8 0 0 を含むエンドツーエンド予約 / クレジット機構が存在することができる。この機構は、直ぐにデータを受け入れることができる場合にリモート F P G A 8 0 0 がストレージ転送を可能にする S C S I 転送準備機能に類似することができる。同様に、F P G A 8 0 0 が転送を要求した場合、ローカル F P G A 8 0 0 は、リモート F P G A 8 0 0 にクレジットを割り当てる。ここでも、リモートストレージのための S L A ポリシーを実行することができる。第 2 に、Nick McKown 著、「入力キュースイッチのための i S L I P スケジューリングアルゴリズム ( The i S L I P Scheduling Algorithm for Input - Queues Switches )」、I E E E / A C M T R A N S A C T I O N S O N N E T W O R K I N G 第 7 巻、第 2 号、1 9 9 9 年 4 月、において提案されている入力キューのための i S L I P アルゴリズムなどの反復ラウンドロビンアルゴリズムなどの分散スケジューリングアルゴリズムが存在することができる。このアルゴリズムは、中間ネットワークファブリックをクロスバーとして用いてクラスタ全体にわたって実行することができる。

40

【 0 0 8 0 】

書き換え、タグ及び C R C 機構 8 5 2 は、適当なタグ及び C R C 保護を用いてパケットをカプセル化 / 脱カプセル化することができる。

【 0 0 8 1 】

50

MAC インターフェイスなどの一連のインターフェイス 854 は、イーサネット(登録商標)へのインターフェイスを提供することができる。

【0082】

一連の組み込み CPU 及びキャッシュコンプレックス 858 は、ローカルホスト及びネットワークリモート FPGA 800 との間でプロセス制御プラン、例外処理、及びその他の通信を実行することができる。

【0083】

DDR コントローラなどのメモリコントローラ 860 は、外部 DRAM / NVRAM のコントローラとして機能することができる。

【0084】

本明細書では、FPGA 800 による 1 つの例において具体化される集中型ソリューション 300 によって提供される機能を統合した結果、ストレージの開始とストレージのターゲットとを単一のハードウェアシステムに組み合わせる方法及びシステムを提供する。実施形態では、これらを PCI e バス 110 によって接続することができる。単一のルート仮想化機能 (SR - IOV) などを適用し、いずれかの標準的な装置 (例えば、いずれかのストレージ媒体 302 装置) を選択して数百個のこのような装置であるかのように機能させることができる。本明細書に開示する実施形態は、SR - IOV のようなプロトコルを用いて物理ストレージアダプタの複数の仮想インスタンスを生じることを含む。SR - IOV は、I / O 機能を仮想化する PCI e 標準であり、ネットワークインターフェイスに使用されてきたが、本明細書に開示する方法及びシステムは、SR - IOV をストレージ装置に使用するように拡張する。従って、本明細書では仮想ターゲットストレージシステムが提供される。実施形態では、仮想ターゲットストレージシステムが、異なる媒体を DAS 310 などの 1 又は複数のディスクであるかのように取り扱うことができる。

【0085】

本明細書に開示する方法及びシステムの実施形態は、実施形態によって可能になる FPGA 800 のように、仮想化されて動的に割り当てられた NVMe デバイスを提供することを含むこともできる。実施形態では、通常の NVMe プロトコルを重畳させることもできるが、NVMe デバイスを分割し、仮想化して動的に割り当てることもできる。実施形態では、ソフトウェア内にフットプリントは存在しない。オペレーティングシステム 108 は、同じ又はほぼ同じ状態を保つ (場合によっては、集中型ネットワーク / ストレージカード 300 を見る小型ドライバを有する)。この結果、仮想ストレージは、直接接続されたディスクのように見えるが、このようなストレージ装置 302 をネットワーク 122 にわたってプールの点で異なる。

【0086】

本明細書では、NVMe の仮想化を実装する方法及びシステムを開示する。どれほど多くのソースがどれほど多くの宛先に関連しているかに関わらず、ソースからのデータがハブに入る前に最初にシリアル化される限り、ハブは、指定された宛先に順にデータを分配する。そうだとした場合、DMA キュー 804、832 などのデータトランスポートリソースを 1 回のみのコピーに低減することができる。この開示は、様々な使用シナリオを含むことができる。1 つのシナリオでは、NVMe 仮想機能 (VF) について、これらが全て同じ PCI e バス 110 に接続された場合、どれほど多くの VF 807 が構成されるかに関わらず、データはこの VF のプール 807 内に順に到来し、従って DMA エンジン 804 は 1 つしか存在せず、(制御情報のため) 1 つのストレージブロックしか必要ない。

【0087】

別の使用シナリオでは、分散ディスク / コントローラのプールを有するディスクストレージシステムについて、物理的バス、すなわち PCI e 110 からデータが生じた場合、データはこのディスクのプール内に連続して到来するので、プール内にどれほど多くのディスク / コントローラが存在するかに関わらず、DMA エンジン 804 などのトランスポートリソースをコントローラ毎に 1 つではなく 1 つのみに低減することができる。

10

20

30

40

50

## 【 0 0 8 8 】

本明細書に開示する方法及びシステムは、集中型ネットワーク/ストレージアダプタ 300 の仮想化を含むこともできる。トラフィックの観点から、システムを 1 つに組み合わせることができる。ストレージアダプタとネットワークアダプタとを組み合わせることで仮想化に追加すれば大きな利点を得られる。すなわち、2 つの P C I e バス 1 1 0 を含む単一のホスト 1 0 2 が存在する。P C I e 1 1 0 からルーティングを行うには、リモートダイレクトメモリアクセス ( R D M A ) のようなシステムを用いて別のマシン/ホスト 1 0 2 に到達させることができる。これを別個に行った場合、ストレージの R D M A システムとネットワーク R D M A システムとを別個に構成する必要がある。それぞれを結合して 2 つの異なる場所で構成する必要がある。集中型ソリューション 3 0 0 では、これが R D M A であり、他のどこかに別のファブリックが存在するとすれば、ストレージとネットワーキングの組み合わせを用いてこれらの 2 つを単一ステップで構成できるので、Q o S を構成するステップ全体がゼロタッチ処理である。すなわち、ストレージが分かれば、Q o S をネットワーク上に別個に構成する必要はない。従って、集中型ソリューション 3 0 0 により、R D M A ソリューションのためのネットワーク及びストレージのシングルステップ構成が可能になる。

10

## 【 0 0 8 9 】

再び図 4 を参照すると、図 8 に関して説明したような F P G A 8 0 0 又は同様のハードウェアによってリモートアクセスが可能になる。図 4 では、仮想化境界を点線 4 0 8 によって示している。この線よりも左側では、オペレーティングシステム 1 0 8 に仮想ストレージ装置 (例えば N V M e 8 0 2 ) 及び仮想ネットワークインターフェイス 8 3 0 が示される。オペレーティングシステムは、これらが仮想装置であることが分からない。仮想化境界 4 0 8 の右側には、(例えば、上述した S A T A 又はその他のプロトコルを使用する) 物理ストレージ装置 3 0 2 及び物理的ネットワークインターフェイスが存在する。ストレージ仮想化機能は、図 8 の v N V M e 8 0 2 及び N V M e 仮想化機構 8 1 0 によって実行される。ネットワーク仮想化機能は、v N I C 機構 8 3 0 によって実行される。物理ストレージ媒体の位置も、オペレーティングシステム 1 0 8 から隠されている。サーバ全体にわたる物理的ディスク 3 0 2 をプールして遠隔的にアクセスすることが効率的である。オペレーティングシステム 1 0 8 は、ストレージ媒体 3 0 2 (このストレージ媒体 3 0 2 は仮想装置であるが、オペレーションシステム 1 0 8 は物理的装置として見る) に読み取り又は書き込みトランザクションを発行する。物理ストレージ媒体 3 0 2 がたまたま遠隔地に存在する場合、読み取り/書き込みトランザクションは、正しい物理的位置にマッピングされ、カプセル化され、イーサネット (登録商標) を通じてトンネリングされる。このプロセスは、図 8 の装置マッピング機構 8 0 8 、N V M e 仮想化機構 8 1 0 、データ変換機構 8 1 2 及びストレージネットワークトンネル 8 2 8 によって実行することができる。ターゲットサーバ (第 2 のコンピュータシステム) は、ストレージ読み取り/書き込みをアントンネリングして、そのローカルストレージ媒体 3 0 2 に直接アクセスする。トランザクションが書き込みである場合、媒体 3 0 2 にデータが書き込まれる。トランザクションが読み取りである場合、データが準備され、オリジンサーバにマッピングされ、カプセル化され、イーサネット (登録商標) を通じてトンネリングされる。オリジンオペレーティングシステム 1 0 2 にトランザクション完了が到着する。従来のシステムでは、これらのステップが、ストレージ要求を処理するためのソフトウェアの介入、データフォーマット及びネットワークアクセスを必要とする。図示のように、これらの複雑なソフトウェアステップは全て回避される。

20

30

40

## 【 0 0 9 0 】

図 9 は、本開示を通じて説明する集中型ソリューション 3 0 0 の 1 つの実施形態としてのコントローラカード 9 0 2 のアーキテクチャの簡略ブロック図である。コントローラカード 9 0 2 は、例えば、G e n 3 x 1 6 カードなどの標準的なフルハイト、ハーフレングスの P C I e カードとすることができる。しかしながら、非標準的なカードサイズも容認可能であり、様々なタイプのターゲットシャーシに収まるようなサイズであることが好ま

50



しい。P C I e のフォームファクタは、P C B 上で使用するスタックアップ及びレイヤを制限する。

【 0 0 9 1 】

コントローラカード 9 0 2 は、2 R U、4 ノードシャーシなどの汎用シャーシのアドオンカードとして使用することができる。シャーシの（スレッドと呼ばれる）各ノードの典型的な幅は、1 R U 及び 6 . 7 6 " である。通常、このマザーボードは、背面付近に P C I e G e n 3 x 1 6 コネクタを提供することができる。コントローラカード 9 0 2 は、ライザーカードを用いてマザーボードの上部に取り付けることができ、従ってカードとマザーボードの間の間隙は、概ねスロット幅に制限することができる。

【 0 0 9 2 】

実施形態では、P C I コネクタによって供給される最大電力が 7 5 W である。コントローラカード 9 0 2 は、約 6 0 W 以下を消費することができる。

【 0 0 9 3 】

このシャーシは良好な空気流を提供することができるが、この例では、デュアル X e o n プロセッサ及び 1 6 D I M M によって空気が温められるので、カードは、周囲温度が 1 0 上昇することを予想すべきである。ほとんどのサーバの最大周囲温度は 3 5 であり、従ってコントローラカード 9 0 2 における大気温度は、状況によっては 4 5 以上になる可能性がある。熱対策の一部として、カスタムヒートシンク及びバッフルを検討することができる。

【 0 0 9 4 】

図 9 に示すコントローラカード 9 0 2 の実施形態には、データバス F P G A 又はデータバスチップ 9 0 4、及びネットワーキング F P G A 又はネットワーキングチップ 9 0 8 という 2 つの F P G A を示している。

【 0 0 9 5 】

データバスチップ 9 0 4 は、P C I e コネクタ 1 1 0 を介したホストコンピュータ 1 0 2 への接続性を提供する。ホストプロセッサから見れば、コントローラカード 9 0 2 は複数の N V M e デバイスのように見える。データバスチップ 9 0 4 は、N V M e を標準的な S A T A / S A S プロトコルにブリッジし、この実施形態では、S A T A / S A S リンクを介して最大 6 個の外部ディスクドライブを制御する。なお、S A T A は、最大 6 . 0 G b p s をサポートし、S A S は、最大 1 2 . 0 G b p s をサポートする。

【 0 0 9 6 】

ネットワーキングチップ 9 0 8 は、N I C デバイス 1 1 8 の 2 つの 1 0 G イーサネット（登録商標）ポート及び e C P U 1 0 1 8 を 2 つの外部 1 0 G イーサネット（登録商標）ポートに切り換える。ネットワーキングチップ 9 0 8 は、仮想化に使用する多くのデータ構造も含む。

【 0 0 9 7 】

通常、ホスト 1 0 2 のマザーボードは、I n t e l チップセット内の 2 つの別個の P C I e G e n 3 x 8 バスに分割できる P C I e G e n 3 x 1 6 インターフェイスを提供する。P C I e G e n 3 x 8 バス 1 1 0 の一方は、I n t e l N I C デバイス 1 1 8 に接続される。第 2 の P C I e G e n 3 x 8 バス 1 1 0 は、P L X P C I e スイッチチップ 1 0 1 0 に接続される。スイッチチップ 1 0 1 0 のダウンストリームポートは、2 つの P C I e G e n 3 x 8 バス 1 1 0 として構成される。バス 1 1 0 の一方は、e C P U に接続され、第 2 のバス 1 1 0 は、データバスチップ 9 0 4 に接続される。

【 0 0 9 8 】

データバスチップ 9 0 4 は、データストレージとして外部メモリを使用する。単一の x 7 2 D D R 3 チャンネル 1 0 1 2 は、ほとんどの状況にとって十分な帯域幅を提供すべきである。ネットワーキングチップ 9 0 8 も、データストレージとして外部メモリを使用し、単一の x 7 2 D D R 3 チャンネルは、ほとんどの状況にとって十分と思われる。また、このデータ構造は、不揮発性 D I M M（典型的には、データ保持用のエネルギー貯蔵要素としての内蔵電力切り替え回路及びスーパーコンデンサを有する N V D I M M）などの、高性能及

10

20

30

40

50

び十分な密度を提供するような不揮発性メモリの使用を必要とする。

#### 【0099】

eCPU1018は、2組のインターフェイスを用いてネットワーキング908と通信する。eCPU108は、NVMeと同様の通信のためのPCIe Gen2x4インターフェイスを有する。eCPU1018は、L2スイッチなどを介してネットワーキングチップ908に接続する2つの10Gイーサネット(登録商標)インターフェイスも有する。

#### 【0100】

2つのチップ904、908の内部設計全体を通じて、AXIバス1020(ARMチップセットのバス仕様)が使用される。AXIバス1020は、データバスチップ904とネットワーキングチップ908との間のシームレスな通信を可能にするためにチップ間接続に使用される。シリアルインターフェイスであるXilinx Aurora(商標)プロトコルを物理的レイヤとして使用することができる。

10

#### 【0101】

FPGA構成の重要要件は、(1)PCIe構成の開始前にデータバスチップ904を準備する(QSPIフラッシュメモリ(4倍SPIバスインターフェイスを含むシリアルフラッシュメモリ)が十分に高速である)こと、及び(2)好ましくはチップがフィールドアップグレード可能であること、である。構成されるフラッシュメモリは、構成ビットストリームの少なくとも3つのコピーを記憶できるほど十分に大きいことが好ましい。このビットストリームは、Xilinx(商標)FPGAによって使用される構成メモリパターンを意味する。通常、ビットストリームは不揮発性メモリに記憶され、初期起動中にFPGAを構成するために使用される。eCPU1018は、構成フラッシュメモリの読み取り及び書き込みを行う機構を有することができる。ホスト102のプロセッサには、新たなビットストリームが存在することができる。eCPU1018は、フラッシュメモリをアップグレードしようと試みる前にセキュリティ及び認証に対処することができる。

20

#### 【0102】

ネットワーキングサブシステムでは、コントローラカード902が、ホストプロセッサと外部との間の全てのネットワークトラフィックを処理することができる。ネットワーキングチップ908は、NIC118及び外部からの全てのネットワークトラフィックを傍受することができる。

30

#### 【0103】

この実施形態におけるIntel NIC118は、ネットワーキングチップ908に2つの10GigE XFIインターフェイス1022を接続する。組み込みプロセッサも同じことを行う。ネットワーキングチップ908は、L2スイッチング機能を実行して、2つの外部10GigEポートにイーサネット(登録商標)トラフィックを送出する。同様に、着信10GigEトラフィックは、NIC118、eCPU1018、又はネットワーキングチップ908の内部ロジックに直接進む。

#### 【0104】

コントローラカード902は、2つの外部10Gイーサネット(登録商標)ポートにSFP+光学コネクタを使用することができる。他の実施形態では、カードが、外部PHY及びRJ45コネクタを用いて10GBASE-Tをサポートすることができるが、SFP+とRJ45の切り換えを可能にするために別個のカード又はカスタムパドルカード構成が必要になり得る。

40

#### 【0105】

ネットワーキングチップ908は、LEDの動作を含む外部ポート及び光学素子の全ての管理を制御することができる。従って、PRST、I2C/MADIOなどの信号は、NIC108の代わりにネットワーキングチップ908に接続することができる。

#### 【0106】

ストレージサブシステムでは、データバスチップ904が、ミニSAS HDコネクタを直接駆動することができる。図10に示すような実施形態では、最新のSAD規格をサ

50

ポートするように、信号を 12 Gbps で動作するように設計することができる。

【0107】

ボードスペースを有効に使用するために、2つのx4ミニSAS HDコネクタを使用することができる。たとえ同時に6組の信号しか使用できない場合でも、8組の信号全てをデータバスチップ904に接続することができる。

【0108】

シャーシ上では、高速銅ケーブルを用いてミニSAS HDコネクタをマザーボードに接続することができる。ミニSAS HDコネクタの配置では、様々なシャーシの物理的スペース及びケーブルのルーティングを考慮することができる。

【0109】

コントローラカード902への電力は、PCIex16コネクタによって供給される。外部電源接続を使用する必要はない。PCIx16コネクタは、PCIe仕様により、電源投入後に最大25Wの電力しか供給することができない。コントローラカード902は、PCIeの構成後まで25W未満を引き出すように設計することができる。従って、初期電源投入後には、いくつかのインターフェイス及びコンポーネントをリセット状態に保持する必要がある。コネクタは、構成後に最大75Wの電力を供給することができ、75Wが3.3Vのレールと12Vのレールに分割されるように構成することができる。

【0110】

図10に、集中型ソリューション300にインターフェイス接続するドライバ1002を含む、FPGA800によって可能になるようなソフトウェアスタック1000を示す。NVMeコントローラ1004は、NVMeコントローラの機能を提供してホストに仮想装置1012を割り当てるハードウェア(例えば、FPGA800)の機能セットである。図10では、dev1、dev2、dev3が仮想装置1012の例であり、それぞれコンテナ1018LXC1、LXC2及びLXC3に動的に割り当てられる。NVMe-SATAブリッジ1008は、仮想装置1012(dev1、dev2、dev3)を変換してストレージ装置302(例えば、図のSSD)にマッピングするハードウェアサブシステム(例えば、FPGA800)の一部である。接続1010は、(上述した他の考えられる接続オプションから)SATA接続を提供するハードウェアシステムの一部である。イーサネット(登録商標)リンク120は、仮想装置1012(すなわち、dev1、dev2、dev3)を、ストレージトンネリングプロトコルを用いてイーサネット(登録商標)リンク120を介して接続された他の(単複の)ホスト102に露出することができる。PCI-E(NVMeドライバ)1002は、ストレージ側のハードウェアサブシステムをプログラムして駆動することができる。このドライバ1002は、ホスト上でオペレーティングシステム(例えば、この例ではLinux(登録商標)OS)の一部として実行することができる。ブロックレイヤ1014は、集中型ソリューションのPCIeドライバ1002にインターフェイス接続して仮想ストレージ装置1012を露出することができる、Linux(登録商標)オペレーティングシステムの従来のSCSIサブシステムとすることができる。コンテナ1018(LXC1、LXC2、LXC3)は、仮想ストレージ装置1012(それぞれdev1、dev2及びdev3)を要求してこれらに動的に割り当てることができる。

【0111】

図11~図15に、複数のシステム102にわたるアプリケーションコンテナ1018(例えば、Linux(登録商標)コンテナ)の移動例を、最初に集中型ソリューション300が存在しない場合について、次に集中型ソリューション300が存在する場合について示す。図11には、従来のストレージコントローラ112と、OS/ハイパーバイザスタック108内の仮想ソフトウェアをホストするネットワークコントローラ118とを有する2つの従来のコンピュータシステム102の例を示す。コンピュータシステム1(C1)は、CPU、メモリ、従来のストレージコントローラ112及びネットワークコントローラ118を含む、図1に示すものと同様の構成を有する。このシステムは、Linux(登録商標)、Microsoft Windows(商標)などのオペレーティング

10

20

30

40

50

システム 108、及び/又は Xen、VMware などのハイパーバイザソフトウェアを実行して、ネイティブに又は仮想マシン又はコンテナなどの仮想環境を通じて複数のアプリケーションをサポートする。このコンピュータシステム 102 では、仮想マシン VM 1104 の内部でアプリケーション App 1102 が実行される。仮想化コンテナ LXC 1110 及び LXC 2114 内では、アプリケーション App 2108 及び App 3112 がそれぞれ実行される。これらのアプリケーションに加え、オペレーティングシステム 108 を通じてアプリケーション App 4118 がネイティブに実行される。典型的なことではあるが、実際のシナリオでは、(3 つ全てではなく) 仮想マシン又はコンテナ又はネイティブアプリケーションのいずれかしか存在しないこともあり、ここでは仮想環境の全ての事例に対応するように意図的にこれらを組み合わせた形で示している。コンピュータシステム 2 (C 2) 102 は、コンテナ内で及びネイティブに App 5 及び App 6 をそれぞれサポートする同様の構成を有する。これらのアプリケーションの各々は、互いに無関係にそれぞれのストレージ装置 302 にアクセスし、すなわち APP 1 が S 1 を使用し、App 2 が S 2 を使用し、以下同様である。これらの (S 1 ~ S 6 で示す) ストレージ装置 302 は、独立した物理エンティティに限定されるものではない。これらは、必要と認められる際には、1 又は 2 以上の物理的ストレージ要素から論理的に分割することもできる。図示のように、(各ストレージ装置 302 からアプリケーションへの矢印によって示す)、ストレージ装置 302 とアプリケーション 1102、1108、1112、1118 との間のデータフローは、アプリケーションに到達する前にストレージコントローラ 112 及びオペレーティングシステム/ハイパーバイザスタック 108 を通過し、図 1 に関連して説明した課題を伴う。

10

20

#### 【0112】

図 12 に示すように、アプリケーション又はコンテナが C 1 から C 2 に移動すると、その対応するストレージ装置も移動する必要がある。この移動が必要になり得る理由は、既存のアプリケーション (App 1 ~ App 4) 内の挙動変化などに起因して、これらのアプリケーションをサポートするための C 1 の (CPU、メモリなどの) リソースが一定期間にわたって不足している可能性があるからである。

#### 【0113】

通常は、アプリケーションの状態及びストレージのサイズが極端でない限り、十分な時間内に移動を行う方が容易である。一般に、ストレージ強化アプリケーションは大量 (例えば、複数テラバイト) のストレージを使用することがあり、この場合には、許容時間内にストレージ 302 を移動させることが実際的でないこともある。この場合、図 13 に示すように、ストレージは相変わらず元の場所に留まり、ソフトウェアレベルの入れ替え/トンネリングを行ってストレージに遠隔的にアクセスすることができる。

30

#### 【0114】

図 13 に示すように、App 2108 は、コンピュータシステム C 2 に移動した後も、両システム C 1 及び C 2 のオペレーティングシステム又はハイパーバイザ 108 をトラバースすることにより、コンピュータシステム C 1 に位置する元々のストレージ S 2 にアクセスし続ける。この理由は、ネットワークコントローラ 118 を介したストレージコントローラ 112 及びこれに接続されたストレージ装置 302 へのストレージアクセスのマッピングが、メイン CPU 内で実行されるオペレーティングシステム又はハイパーバイザソフトウェアスタック 108 によって行われるからである。

40

#### 【0115】

図 13 に示すように、App 2108 は、C 2 に移動した後、両システム C 1 及び C 2 のオペレーティングシステム又はハイパーバイザ 108 をトラバースすることにより、C 1 に位置する元々のストレージ S 2 にアクセスし続ける。この理由は、ネットワークコントローラ 118 を介した C 2 から C 1 への、さらには C 1 のストレージコントローラ 112 へのストレージアクセスのマッピングが、各コンピュータシステムのメイン CPU 内で実行されるオペレーティングシステム又はハイパーバイザソフトウェア 108 によって行われるからである。

50

## 【 0 1 1 6 】

図 1 4 に示すように集中型コントローラ 3 0 0 を適用した場合の同様のシナリオについて検討する。図示のように、このシナリオは、別個のストレージコントローラ 1 1 2 及びネットワークコントローラ 1 1 8 が集中型 I O コントローラ 3 0 0 に置き換わっている点を除いて図 1 1 とほぼ同一である。この場合、( 図 1 5 に示すように ) A p p 2 1 1 0 8 がそのコンテナ L X C 1 と共に C 2 に移動してもストレージ S 2 は移動せず、コンピュータシステム C 1 内に存在するメイン C P U において実行されるいずれかのソフトウェア ( オペレーティングシステム、ハイパーバイザ 1 0 8、又は他のいずれか ) がトラバースされるのを避けることによってアクセスが最適化される。

## 【 0 1 1 7 】

従って、本明細書では、ストレージ装置が存在するメイン C P U を迂回する新規の方法を提供し、この方法はさらに、( a ) 複数のコンピュータシステムを横切ってストレージにアクセスする際のレイテンシ及び帯域幅を大幅に低減することができ、( b ) ストレージが存在するマシンからアプリケーションを移動させる必要がある状況を大幅に単純化して改善する。

## 【 0 1 1 8 】

イーサネット (登録商標) ネットワークは、最善努力の原則で挙動し、従って本質的に損失が多くバースト性である。いずれのパケットも永久に失われる可能性があり、或いはバッファリングされて他のパケット共にバースト式で遅延して送出される可能性がある。一方で、典型的なストレージ中心のアプリケーションは損失及びバーストに敏感であり、ストレージトラフィックがいつイーサネット (登録商標) ネットワークを介して送信されるかが重要である。

## 【 0 1 1 9 】

バス / ネットワークを介した従来のストレージアクセスは、信頼性の高い予測可能な方法で行われる。例えば、ファイバチャネルネットワークは、エンドシステムによって行われるアクセス数を制限するようにクレジットベースのフロー制御を採用する。そして、エンドシステムに与えられるクレジット数は、必要なレイテンシ及び帯域幅ニーズを満たす予測可能時間内にストレージ要求を受け取って満足させるのに十分なコマンドバッファをストレージ装置が有しているかどうかに基づく。以下の図に、S A T A、ファイバチャネル ( F C )、S C S I、S A S などの異なるタイプのバスが採用するいくつかのクレジットスキームを示す。

## 【 0 1 2 0 】

図 1 6 に示すように、イーサネット (登録商標) ネットワークは、最善努力の原則で挙動し、従って本質的に損失が多くバースト性になりがちである。いずれのパケットも永久に失われる可能性があり、或いはバッファリングされて他の多くのパケットと共に輻輳誘導バーストの形で遅延して送出される可能性がある。典型的なストレージ集中のアプリケーションは損失及びバーストに敏感であり、ストレージトラフィックがいつバス及びイーサネット (登録商標) を介して送信されるかが重要であるため、これらは、整合性を維持するために信頼性の高い予測可能な方法で行われる。例えば、従来、ファイバチャネルネットワークは、エンドシステムによって同時に行われるアクセス数を制限するようにクレジットベースのフロー制御を採用する。エンドシステムに与えられるクレジット数は、必要なレイテンシ及び帯域幅要件を満たす予測可能時間内にストレージ要求を受け取って満足させるのに十分なコマンドバッファをストレージ装置 3 0 2 が有しているかどうかに基づくことができる。図 1 6 に、他のタイプのこのようなスキームの中でも特に、S A T A バス 1 6 0 2、ファイバチャネル ( F C ) 1 6 0 4 及び S C S I / S A S 接続 1 6 0 8 などの異なるタイプのバスが採用するクレジットスキームの一部を示す。

## 【 0 1 2 1 】

図示のように、例えば F C コントローラ 1 6 1 0 は、F C ベースのストレージ装置 1 6 1 2 に送信を行う前に最大「N」個のストレージコマンドまでの固有のバッファリングを有することができるが、F C 装置 1 6 1 2 は、例えばこの例では「N」より大きくするこ

10

20

30

40

50

とも、「N」と等しくすることも、又は「N」より小さくすることもできる「M」個という異なるバッファ制限を有することができる。典型的なクレジットベースのスキームは、ターゲットレベル（この例では、FC装置1602などのストレージ装置302のうちの1つがターゲットである）を用いてクレジットを集約し、これらのクレジットに関する情報が、ターゲット302にアクセスしようと試みている様々なソース（この例では、FCコントローラ1610などのコントローラがソースである）に伝播される。例えば、「N」のキュー深度を有するターゲットに2つのソースがアクセスしている場合、ソースに与えられるクレジットの合計は、いかなる時点においてもターゲットが「N」個よりも多くのコマンドを受け取らないように「N」を超えることはない。ソース間のクレジットの分配は任意とすることも、或いは様々なタイプのポリシー（例えば、コスト/価格設定又はSLAなどに基づく優先順位）に基づくこともできる。コマンド要求を満たすことによってキューが処理されると、必要に応じてソースにおいてクレジットを補充することができる。この種のクレジットベースのストレージアクセスに従うことにより、ターゲットのキューがどうしようもなくなくなることによって生じる損失を避けることができる。

#### 【0122】

FCOE及びiSCSIなどのイーサネット（登録商標）を介した典型的なストレージアクセスは、ターゲット志向、クレジットベースのコマンドフルフィルメントを、イーサネット（登録商標）リンクを介した転送に拡張することができる。このような場合、これらは、ソース志向よりもむしろターゲット装置志向とすることができる。本明細書では、代わりにいずれの又はどの種類のソースがどれほど多くのクレジットを獲得するかに基づくことができる新たなクレジットベースのスキームを提供する。例えば、図17に示すように、ネットワークをストレージに直接インターフェイス接続させる上述の集中型ソリューション300は、マルチプレクサを用いて、ソース志向のクレジットベースのスケジューリングスキームをターゲット装置志向のクレジットベースのスキームにマッピングすることができる。

#### 【0123】

図17に示すように、イーサネット（登録商標）上に4つのソースが存在し、2つのターゲットストレージ装置302が存在する。典型的なターゲット志向のクレジットベースのスキームは、（1ターゲット当たり1つの）2つのキュー、又は1ソース当たり2つの接続を各ターゲットに露出する。代わりに、図17に示すように、キュー（Q1、Q2、Q3、Q4）1702は1ソース当たりのものであり、マルチプレクサ（S）1708を横切って2つのターゲット志向キュー（Q5、Q6）1704にマッピング/多重化される。この種のソース志向のクレジットベースのスキームを採用することにより、ターゲットストレージ装置302の数に関わらずにアクセス帯域幅及び予測可能なアクセスレイテンシを保証することができる。一例として、1つのタイプのマルチプレクサは、Q1がそのソースによって圧倒されないように、Q1のキューサイズ「P」がQ5及びQ6の「L+M」を超えないことを確実にするものである。

#### 【0124】

実施形態では、ストレージ装置302からデータブロックにアクセスできるようにする方法及びシステムについて説明する。具体的には、アプリケーションがそのデータにアクセスして特定の一連のアクセス要件を満たせるようにする新規方法について説明する。

#### 【0125】

本明細書で使用する「アプリケーション駆動型データストレージ」（ADS）は、アプリケーションに対するアプリケーションデータの記憶方法、アクセス方法、転送方法、キャッシュ方法及び分配方法に関してあらゆるアプリケーションに透明性をもたらすストレージを含む。ADSは、アプリケーションがこれらの個々の相を制御して特定のアプリケーションの特定のニーズに対応できるようにする。一例として、あるアプリケーションは、アプリケーション自体の複数のインスタンスと、ネットワークを越えて複数のLinux（登録商標）ノードにわたって拡散した複数のプロセスとを含むことができる。これらのプロセスは、これらの間で共有的又は排他的に複数のファイルにアクセスすることができ

る。これらのプロセスは、アプリケーションがこれらのファイルをどのように処理したいと望むかに基づいて、異なるファイル部分に頻繁にアクセスしたいと望み、クイックアクセス又は使い捨てを必要とすることができる。これらの基準に基づいて、必要時にキャッシュ及び/又はストレージの異なる階層の特定のファイル部分を迅速なアクセスのためにセッション毎又はファイル毎にプリフェッチ及び/又は保持したいと望むことができる。これらのアプリケーション固有の要件は、ファイルシステム全体のディスクストライピング、先読み連続ブロックのプリフェッチ、サーバ又はLRU又はFIFOベースのファイルコンテンツキャッシングにおける物理メモリの確保などの一般的な方法で満たすことができない。

#### 【0126】

10

アプリケーション駆動型データストレージI/Oは、単純にストレージエンティティのみに適応可能なものではない。アプリケーション駆動型データストレージI/Oは、ストレージスタック全体に複数の方法で影響を与える。まず、アプリケーション駆動型データストレージI/Oは、Linux(登録商標)ページングシステム、バッファリング、基本ファイルシステムクライアント、TCP/IPスタック、分類、QoS処理及びネットワークングハードウェア及びソフトウェアによって提供されるパケットキューイングを含むアプリケーションが実行されているコンピュータノードのストレージI/Oスタックに影響を与える。次に、アプリケーションノードとそのストレージとを相互接続し、イーサネット(登録商標)セグメント、最適経路選択、スイッチのバッファリング、レイテンシに敏感なストレージトラフィックの分類及びQoS処理、並びにストレージI/Oに関連するインプロージョン問題を含むネットワークングインフラストラクチャに影響を与える。また、基本ファイルレイアウト、冗長性、アクセス時間、様々なタイプのストレージ及びリモートリポジトリの階層化を含むファイルの観点からデータを記憶して維持するストレージインフラストラクチャにも影響を与える。

20

#### 【0127】

本明細書に開示する方法及びシステムは、アプリケーションノード内の典型的なアプリケーションに影響を与える要素、及び集中型ソリューション300がいくつかの重要なアプリケーション要件に対処するためにステータクオーを変更する方法に関するものを含む。

#### 【0128】

30

従来のLinux(登録商標)スタックは、汎用メモリ割り当て、プロセススケジューリング、ファイルアクセス、メモリマッピング、ページキャッシングなどの単純なビルディングブロックを含むことができる。これらは、Linux(登録商標)上で実行されるアプリケーションにとって必須であるが、NoSQLなどの入力/出力(I/O)集中的な特定のカテゴリのアプリケーションにとっては最適でない。NoSQLアプリケーションは非常にI/O集中的であり、そのデータアクセスを一般的な方法で予測することは困難である。ユーティリティコンピュータ環境内でアプリケーションを展開する必要がある場合、Linux(登録商標)がこれらのビルディングブロックの汎用最小実装を提供することは理想的でない。これらのビルディングブロックは、非常にフレキシブルであるとともに、(単複の)アプリケーションから制御可能なアプリケーション関連の特徴を有することが好ましい。

40

#### 【0129】

全てのアプリケーションは独自の特定の要件を有するが、例示的な実施形態では、アプリケーションのNoSQLクラスが、Linux(登録商標)スタックによって処理された時にNoSQLアプリケーション及びその他のI/O集中的アプリケーションの性能を大幅に向上させる以下の要件を有する。まず、これらの要件は、ファイルレベル優先順位を使用する点である。Linux(登録商標)ファイルシステムは、最低限でも異なるファイル間でのアクセスレベル優先順位を提供すべきである。例えば、(他方よりも優先順位の高いデータベース/テーブル/インデックスなどの)一方が他方よりも高い優先順位を与えられた2つの異なるファイルにアクセスする(複数のスレッドからなる)アプリケーショ

50

ンプロセス。このプロセスでは、アクセスされるデータに基づいて貴重なストレージ I / O リソースを優先的に利用することができる。これには、1つのスレッド / プロセスを高優先順位又は低優先順位で実行することによって間接的に対処できるとの主張があると思われるが、これらのプロセスレベルの優先順位は、ファイルシステム又はストレージコンポーネントに通信されない。プロセス又はスレッドレベルの優先順位は、CPU リソースの利用のみを対象とする。さらに、これらの2つのファイルに同じスレッドがアクセスし、従ってどのデータ (ファイル) がアクセスされているに基づいて2つの異なるレベルでストレージソースを利用することも可能である。次に、アクセスレベルのプリファレンス要件が存在することもできる。Linux (登録商標) ファイルシステムは、ファイル (オープンファイル) のセッション中に、ファイルセッション間の優先順位、ブロックのバッファリング量、様々なブロックの保持 / 有効期限プリファレンス、リソース閾値及び競合のためのアラート、並びに性能統計などの様々なプリファレンス (主に S L A ) を提供することができる。一例として、MongoDB 又は Cassandra などの NoSQL アプリケーションが、書き込み及び読み取りのための2又は3以上のスレッドを有している時に、読み取りよりも書き込みのプリファレンスの方が高い場合、同じファイルの読み取りファイルセッションよりも書き込みファイルセッションの方に高プリファレンスを与える必要があり得る。この能力により、同じファイルの2つのセッションが、2つの異なる優先順位を有することができる。

10

20

30

40

50

#### 【0130】

NoSQL アプリケーションの多くは、異なるタイプのデータを同じファイルに記憶し、例えば MongoDB は、ユーザコレクション及び (b ツリー) インデックスコレクションを同じデータベースファイルの組に記憶する。MongoDB は、インデックスページ (b ツリー及びコレクション) をユーザコレクションページよりも優先してメモリに保持しておきたいと望むことができる。MongoDB は、これらのファイルがオープンになると、Linux (登録商標)、ファイル及びストレージシステムに影響を与えて、アプリケーションの要件を知らない汎用 FIFO 又は LRU に基づいてページを処理するのは対照的に、MongoDB ポリシーに従ってこれらのページを処理したいと望むことができる。

#### 【0131】

リソースアラート及び性能統計は、NoSQL データベースが基本ファイル及びストレージシステムの挙動を理解できるようにし、これに従ってデータベースクエリにサービスを提供し、或いはデータベースシャーディング、又は同じホストにおいて実行される (バックアップ、シャーディング、ナンバー読み取り / 書き込みクエリサービスなどの) 他のジョブのファイル I / O プリファレンスの減分 / 増分などの実行すべき動作をトリガすることができる。例えば、IOP 及びレイテンシの最小数、最大数及び平均数、並びに一定時間にホストメモリとの間でスラッシュイン及びスラッシュアウトされた上位10個の候補ページに関する性能統計は、アプリケーションが上述したパラメータを動的に調整してアプリケーション自体を微調整できるようにする。

#### 【0132】

キャッシング及び階層化プリファレンスのための要件も存在することができる。Linux (登録商標) ファイルシステムは、アプリケーションがそのファイルにアクセスしている間に動的に構成可能なキャッシングポリシーを有する必要があり得る。通常、現在の Linux (登録商標) ファイルシステムは、アプリケーションがストリームのような連続的方法でファイルを読み取ることを希望して連続ファイルブロックをプリフェッチする。しかしながら、実際には、ウェブサーバ及びビデオストリーマのような多くのレガシーアプリケーションの場合、台頭しつつある NoSQL アプリケーションは厳密に連続読み取りに従わない。これらのアプリケーションは、ブロックをランダムに読み取る。一例として、MongoDB は、キーの検索時にキーを発見するまでランダムにブロックにアクセスする、ファイルの一部にフラットにレイアウトされた b ツリーの形でインデックステーブルにドキュメントキーを記憶する。さらに、これらのファイルは、このような b ツリーベ



ースのインデックステーブル専用のもではない。これらのファイルは、ユーザドキュメント及びシステムインデックスファイルなどの様々なタイプのテーブル（コレクション）間で共有される。このため、Linux（登録商標）ファイルシステムは、効率的なメモリ使用などのために、ファイルのどの部分をキャッシュし、先読みし、スワップアウトする必要があるかを予測することができない。

#### 【0133】

本明細書に開示する方法及びシステムの実施形態では、ストレージ要件において様々なアプリケーションにわたる共通スレッドが存在する。具体的には、特定の必要な時点及び場所における特定のタイプのデータのレイテンシ及びI/Oは、これらのアプリケーションの性能に対して非常に強い影響力を有する。

10

#### 【0134】

例えば、本明細書では、上述したホストレベルの要件に対処するために、アプリケーションがホスト内及び他の場所におけるプリファレンスに従ってデータの記憶、検索、保持及び階層化に完全に影響を与えてこれらを制御できるようにする十分に微調整されたファイルシステムクライアントのための方法及びシステムを開示する。

#### 【0135】

図18に示すように、ファイルシステム（FS）クライアント1802は、別個のファイルセッション（fd1及びfd2）のために別個のバッファプールを保持する。ファイルシステム（FS）クライアント1802は、アプリケーション又は一連のプロセス毎に集約メモリプールの事前割り当て及び維持も行う。SLAブローカ1804は、ファイルI/Oが実行されたプロセス/スレッド内で内部的に、又は別の一連のプロセスから外部的にアプリケーションが実行して、FSクライアント1802に影響を与えて適切なストレージI/O SLAを動的に提供することができる。SLAを外部プロセスから制御すると、アプリケーション自体を修正することなく、これらの新たなストレージ制御機能についての知識を有していないレガシーアプリケーションが直ちに可能になる。

20

#### 【0136】

本明細書に開示する方法及びシステムは、ネットワーク及びホストを横切るデータ検索のための広範にわたる階層化サービスを提供することができる。以下の図19に示すように、高性能分散ファイルサーバ（DFS）1902は、キャッシュ形態ストレージ形態のどの媒体（DRAM、NVRAM、SSD又はHDD）にファイルのどの部分が存在すべきかを動的に決定して実行するようにプラットフォーム1904内でアプリケーションがコンテナ化された形で実行されることを可能にする。これらのアプリケーションコンテナ1908は、ファイルのストライピング、ミラーリング、レイディング及び障害回復（DR）を行う必要があるかどうかなどの他のストレージポリシーを決定することができる。

30

#### 【0137】

本明細書に開示する方法及びシステムは、高性能DFS1902内のアプリケーションコンテナが、特定のファイルページをローカルストレージ及び/又は遠隔位置から積極的に検索し、これらのページを必要時に後で高速検索できるように特定の場所にプッシュできるような幅広いキャッシングサービスも提供する。例えば、これらの方法及びシステムは、アプリケーションを実行するホストのメモリ及びSSDの使用を局所化し、アプリケーションの関心ページをこれらのホストのローカルメモリ/SSDのいずれかに積極的にプッシュすることができる。これらの方法及びシステムは、その後の必要時にアプリケーションによって超低レイテンシ検索を行えるように、この目的でDFSプラットフォームに提供されたこれらのメモリ、SSD及びHDDの階層を使用することができる。

40

#### 【0138】

アプリケーションのホストにわたるキャッシュの拡大使用は限らない。例えば、MongoDBでは、ワーキングセットが一時的にそのローカルホストのメモリを越えて成長するとスラッシングが発生し、これによってクエリ処理の性能が大幅に低下する。この理由は、新たなページがクエリを満たすようにするために必要なファイルデータページが廃棄され、その後に元々のページに戻す必要がある場合に、システムがディスクサブシステム

50

から再び新たにページを読み込まなければならない、これによってクエリを完了する上で大きなレイテンシを伴うからである。アプリケーション駆動型ストレージアクセスは、廃棄されたページのキャッシュを、MongoDBが再びそのページを必要とするまでネットワークの他の場所に（別のアプリケーションホストのメモリ/SSD、又は高性能DFSシステム1902のストレージのローカル階層に）一時的に保持することによってこれらの種類のシナリオを避けることにより、クエリを完了する上でレイテンシを大幅に低減する。

#### 【0139】

図20に示すように、高性能DFS1902は、必要時に、及びアプリケーションによって影響され制御された際にアプリケーションデータのキャッシング及びサービスを行うために、単一の統合されたRAM及びSSDベースの階層/キャッシュ2002内のアプリケーションホスト全体にわたってDRAM及びSSDリソースを利用する。

#### 【0140】

本明細書では、図21に示すような、一連のホスト(H1~HN)、ファイル又はブロックサーバ2102及びストレージサブシステム2104を含むシステムを開示する。通常、ホストH1~HNは、ストレージに恒久的又は一時的に記憶されたデータへのアクセスを必要とするアプリケーションを実行するコンピュータである。ファイル又はボリュームサーバ2102は、典型的には中央処理装置(CPU)、メモリ及び専用ハードウェアを含むハードウェアを実行してネットワーク装置及びストレージ装置などの外部装置に接続するデータオーガナイザ及びデータサーバとすることができる。ファイル又はボリュームサーバ2102は、ユーザデータをブロックと呼ばれる複数の固定数又は可変数のバイトに関して構造化する。ファイル又はボリュームサーバ2102は、これらのデータブロックを内部又は外部ストレージに記憶する。ランダムではあるが論理的に関連する一連のブロックをファイル又はボリュームに構造化する。1又は2以上のホストH1~HNは、アプリケーションプログラミングインターフェイス(API)又は他のいずれかのプロトコルを通じてこれらのファイル又はボリュームにアクセスすることができる。ファイル又はボリュームサーバは、1又は2以上のファイル及びボリュームを1又は2以上のホストに提供することができる。なお、ホスト、及びファイル又はボリュームサーバは、直接又はネットワークを介して接続された2つの異なる物理エンティティの形を取ることも、或いは単一の物理コンピュータの形で論理的に共存することもできる。

#### 【0141】

ストレージ2104は、データ片を一時的又は恒久的に保持できる一群のエンティティとすることができる。通常、ストレージ2104は、静的又は動的ランダムアクセスメモリ(RAM)、ソリッドステートストレージ(SSD)、ハードディスクドライブ(HDD)又はこれらの全ての組み合わせを含む。ストレージは、リンク又はネットワークを介してファイル又はボリュームサーバ2102に接続された独立した物理エンティティとすることができる。ストレージは、ファイル又はボリュームサーバ2102と単一の物理エンティティの形で統合することもできる。従って、ホストH1~HN、ファイル又はボリュームサーバ2102及びストレージ2104は、単一のハードウェアエンティティの形で物理的にまとめて配置することができる。

#### 【0142】

通常、ホストは、図22に示すような複数の論理エンティティを含む。通常、アプリケーション2202はホスト内で実行され、ホストのローカルオペレーティングシステム2204又は代替りの他のいずれかのエンティティによって提供されるAPIを介してデータ要素にアクセスする。通常、オペレーティングシステム2204は、そのファイルシステムクライアント2206を介してファイルシステムにインターフェイス接続するための標準的なAPIインターフェイスを有する。ファイルシステムクライアント2206は、遠隔的に又は局所的に位置するファイル又はボリュームサーバ2210にインターフェイス接続するようにホスト内で実行されるソフトウェアエンティティである。ファイルシステムクライアント2206は、単一の又は複数のファイル又はボリューム内に存在するア

アプリケーション 2202 が必要とするデータ要素をファイル又はボリュームサーバ 2210 から検索し、アプリケーションがデータ要素の処理を完了するまでこれらのデータ要素をホストのメモリ 2208 に保持することによってデータ要素を提供する。典型的なアプリケーションのシナリオでは、必要に応じて特定のデータ片が複数回読み取られ及び/又は修正される。また、複数のデータ要素からなるファイル又はボリューム全体が、特定のタイプのアプリケーションのローカルメモリ 2208 のサイズより潜在的にはるかに大きいことも一般的である。これにより、将来的にアプリケーション 2202 がこれらのデータブロックにアクセスする可能性があるか、それ伴いかの予測に基づいてどのデータブロックをメモリ 2208 に保持し、又はメモリ 2208 から排出すべきかを決定するための、オペレーティングシステム 2204 及びファイルシステムクライアント 2206 の実装が複雑になる。これまで、既存の実装では、ファイル又はボリュームサーバ 2210 からの新たなデータブロックを取り入れるために、先入れ先出し (FIFO) 又は最長未使用時間 (LRU) などのいくつかの汎用なアプリケーション非依存方法を実行して、メモリにおけるデータブロックの保持又は排出を行っていた。さらに、別のデータブロックを記憶するために、あるデータブロックによって占められているメモリを再要求する場合、元々のデータは、将来的な使用について考慮せずに単純に消去される。通常、ディスクサブシステムは非常に低速であり、データブロックを読み出してファイル又はボリュームサーバ 2210 によってファイルシステムクライアント 2206 及びメモリ 2208 に転送する際に高レイテンシを伴う。従って、元々のデータブロックが消去されると、アプリケーションは、近い将来に元々のデータにアクセスしようと試みた場合に長く待つ必要性が生じ得る。この種の実装に関する主な問題点は、データアクセス経路内のモジュール、すなわちオペレーティングシステム 2204、ファイルシステムクライアント 2206、メモリ 2208、ブロックサーバ 2210 及びストレージのいずれもが、どのデータブロックがいつ、どのような頻度でアプリケーション 2202 によってアクセスされ予定であるかを全く知らない点である。

#### 【0143】

図 23 に、アプリケーション 2202 がストレージ 2212 からのデータブロックにアクセスする様子を描いた例示的なシナリオを示す。番号付きの円は、データブロックへのアクセス過程に関与するステップを示す。以下、これらのステップについて説明する。第 1 に、アプリケーション 2202 が、ファイル又はオペレーティングシステム 2204 の API を用いてデータブロックにアクセスする。オペレーティングシステム 2204 は、ファイルシステムクライアント 2206 のための同等の API を呼び出して同じデータブロックにアクセスする。第 2 に、ファイルシステムクライアント 2206 が、この目的のためのローカルメモリバッファ内にデータが存在するかどうかを発見しようと試みる。発見された場合、以下のステップ (3) ~ (7) をスキップする。第 3 に、ブロックサーバ 2210 からデータ検索コマンドを送信する。第 4 に、ブロックサーバ 2210 が、ストレージからデータブロックを読み取るための読み取りコマンドをストレージ 2212 に送信する。第 5 に、ストレージ 2212 が、ストレージからの読み取り後にデータブロックをブロックサーバ 2210 に戻す。第 6 に、ブロックサーバ 2210 が、データブロックをファイルシステムクライアント 2206 に戻す。第 7 に、ファイルシステムクライアント 2206 が、将来的にアクセスできるようにメモリ 2208 のメモリバッファにデータを保存する。第 8 に、ファイルシステムクライアント 2206 が、要求されたデータをアプリケーション 2202 に戻す。

#### 【0144】

本明細書に開示する方法及びシステムでは、NoSQL 及びビッグデータの領域における最新のクラスのアプリケーションによるデータアクセスに関連する性能要件に対処するために、オペレーティングシステム 2204、ファイルシステムクライアント 2206、メモリ 2208、ブロックサーバ 2210 及びストレージ 2212 を含むデータブロックアクセスのコンポーネントがあらゆるアプリケーション 2202 によって制御されることを提案する。すなわち、本発明者らは以下を提案する。第 1 に、オペレーティングシステ

ム 2 2 0 4 が、アプリケーションによるファイルシステムクライアント 2 2 0 6 の制御を可能にするさらなる A P I を提供できるようにする。2 番目に、( a ) ファイル又はボリュームが、他のファイル又はボリュームのために共有又は削除されない固有のデータを保持する専用のメモリバッファプールを有するという意味において、アプリケーション 2 2 0 2 が特定のファイル又はボリュームのための専用メモリプールをメモリ 2 2 0 8 内に作成できるようにすることと、( b ) アプリケーション 2 2 0 2 が、ファイル又はボリュームとの 2 つの独立セッションがこれらのデータを保持するための独立したメモリバッファを有するように、ファイル又はボリュームとの特定のセッションのための専用メモリプールをメモリ 2 2 0 8 内に作成できるようにすることをサポートするようにファイルシステムクライアント 2 2 0 6 を拡張する。一例として、非常に重要なファイルセッションは、迅速かつ頻繁なアクセスのためにより多くの存在するデータを利用できるように、メモリ 2 2 0 8 内に多くのメモリバッファを有することができる一方で、同じファイルとの第 2 のセッションには非常に少ないバッファを割り当て、従って様々なファイル部分にアクセスするために遅延及びそのバッファの再使用を伴う必要があるようにし、( c ) 迅速なアクセスのために他のホスト又はブロックサーバ 2 2 1 0 全体にわたってメモリ 2 2 0 8 を越えて拡張されたバッファプールをアプリケーション 2 2 0 2 が作成できるようにする。これにより、他のホストのメモリ 2 2 0 8、及びファイル又はブロックサーバ 2 2 1 0 内に存在するあらゆるメモリ 2 4 0 2 にデータブロックを保持することができ、( d ) アプリケーション 2 2 0 2 が、いずれかのデータブロックを、他のファイル、ボリューム又はセッションのためのデータブロックに対してメモリ 2 2 0 8 内に長く存在させることができるようにする。これにより、アプリケーションが、あるデータブロックを迅速にアクセスできるように常に利用可能となるように選別し、オペレーティングシステム 2 2 0 4 又はファイルシステムクライアント 2 2 0 6 が固有の排出ポリシーに基づいてそのデータブロックを排出しないようにすることができ、( e ) アプリケーション 2 2 0 2 が、いずれかのデータブロックを、他のファイル、ボリューム又はセッションのためのデータブロックに対してメモリ 2 2 0 8 内に長く存在させないことができるようにする。これにより、アプリケーションは、オペレーティングシステム 2 2 0 4 及びファイルシステムクライアント 2 2 0 6 に、そのデータブロックのバッファを排出し、選択時には再使用することを知らせることができるようになる。このことは、他の正常なデータブロックを長期にわたって維持する役に立つ。第 3 に、( a ) アプリケーションコンテナ 2 4 0 0 がアプリケーション 2 2 0 2 の関心データブロックを事前にフェッチして、後で素早くアクセスできるようにローカルメモリ 2 4 0 2 に記憶するとともに、ストレージ 2 2 1 2 に関連するレイテンシペナルティを回避できるようにする能力と、( b ) ホストのメモリ 2 2 0 8 から排出されたページを後でアプリケーション 2 2 0 2 がアクセスできるようにローカルメモリ 2 4 0 2 に記憶できるようにする能力とを有する図 2 4 に示すアプリケーションコンテナ 2 4 0 0 の観点から、ブロックサーバ 2 2 1 0 がアプリケーション固有のモジュールをホストできるようにする。

#### 【 0 1 4 5 】

上記の ( 2 ) ( c ) のアプリケーション駆動型の特徴については、さらなる解説が必要である。2 つのシナリオが存在する。1 つのシナリオは、データブロックをブロックサーバ 2 2 1 0 のメモリから検索するものである。他方のシナリオは、同じデータブロックを別のホストから検索するものである。本明細書に開示する方法及びシステムは、2 つのホスト ( H 1 ) 及び ( H 2 ) がストレージ 2 2 1 2 から全く同じブロックデータを読み取るものと仮定して、図 2 5 に示すようなシステムを提供する。データブロックが別のホスト ( H 2 ) に存在する旨が通知されると、ブロックサーバ 2 2 1 0 にこのデータブロックをストレージ 2 2 1 2 から検索するように求める代わりに、このデータブロックはファイルシステムクライアント 2 2 0 6 によってホスト ( H 2 ) から直接検索され、これによって速度が落ちるとともに高レイテンシを伴う。

#### 【 0 1 4 6 】

実施形態では、ファイルシステムクライアント 2 2 0 6 が、同じ場所にさらに重要なデ

ータブロックを記憶するという理由で(D 1)からデータブロックを排出すると決定した場合、このファイルシステムクライアント2206は、排出されたデータブロックをファイルシステムクライアント2206'に送信して、このファイルシステムクライアント2206'のためにメモリ2208'に記憶することができる。

【0147】

なお、ホストに障害が起きた場合に上述した技術を適用して高速フェイルオーバを実現することもできる。さらに、上述したキャッシング技術を、特にRAMに関連して使用して、ウォームキャッシュを用いたフェイルオーバを実現することもできる。図25は、ウォームキャッシュを用いた高速フェイルオーバシステムの例を示すものである。最終的に、ノードの障害中に、新たなノード上のエンドアプリケーションが、(RAM内の)キャッシュが温まる前の時間を経験せず、これによって低アプリケーション性能の期間を伴わなくなる。

10

【0148】

本明細書では、プロセッサと、アプリケーションのニーズに従ってストレージアクセスを制御するアプリケーション固有のモジュールを含むファイルサーバとを有するシステム及び方法を提供する。

【0149】

本明細書では、アプリケーション固有のモジュールがアプリケーションのニーズに従ってストレージアクセスを制御できるようにするプロセッサと(固定サイズのバイトブロック、可変数のバイトを有する同様の又は異なるオブジェクトを構成する)データストレージとを有するシステム及び方法を提供する。

20

【0150】

また、本明細書では、アプリケーションによる使用目的でそれまで維持されていた古いファイル又はストレージデータブロックをホストのメモリ及び/又はその一時的又は永続的ストレージ要素から検索し、アプリケーションが後でできるようにこれらを別のホストのメモリ又は及び/又はその一時的又は永続的ストレージ要素に記憶するシステム及び方法も提供する。

【0151】

また、本明細書では、アプリケーションによる使用目的でそれまで維持されていたいずれかのファイル又はストレージデータブロックをホストのメモリ及び/又はその一時的又は永続的ストレージ要素から検索し、アプリケーションが後でできるようにこれらを別のホストのメモリ又は及び/又はその一時的又は永続的ストレージ要素に記憶するシステム及び方法も提供する。

30

【0152】

また、本明細書では、データアクセスのレイテンシを低減する目的で、ホストのメモリ及び/又はその一時的又は永続的ストレージ要素を利用して、別のホストにおいて実行されるアプリケーションが後でアクセスすると思われるいずれかのファイル又はストレージデータブロックを記憶するシステム及び方法も提供する。

【0153】

アプリケーションによる使用目的のためにそれまで維持されていた、ホストのメモリ及び/又はその一時的又は永続的ストレージ要素からのファイル又はストレージデータブロックを、後でアプリケーションが使用できるように別のホストのメモリ又は及び/又はその一時的又は永続的ストレージ要素に記憶することができる。

40

【0154】

アプリケーションによる使用目的のためにそれまで維持されていたファイル又はストレージデータブロックを、ネットワークを介してホストのメモリ及び/又はその一時的又は永続的ストレージ要素から別のホストに転送する機構。

【0155】

様々な例示的かつ非限定的な実施形態によれば、物理的ターゲットストレージ媒体コントローラと、物理的ネットワークインターフェイスコントローラと、ストレージ媒体コン

50

コントローラとネットワークインターフェイスコントローラとの間のゲートウェイとを含む集中型入力/出力コントローラを含む装置が開示され、ゲートウェイは、ストレージ媒体コントローラとネットワークインターフェイスコントローラとの間のストレージトラフィック及びネットワークトラフィックのための直接接続を提供する。

【0156】

いくつかの実施形態によれば、装置は、ストレージ媒体コントローラによって制御されるストレージ媒体を、ストレージ媒体の位置に関わらずローカルに接続されたストレージとして提示する仮想ストレージインターフェイスをさらに含むことができる。さらに他の実施形態によれば、装置は、ストレージ媒体コントローラによって制御されるストレージ媒体を、ストレージ媒体のタイプに関わらずローカルに接続されたストレージとして提示する仮想ストレージインターフェイスをさらに含むことができる。さらに他の実施形態によれば、装置は、ストレージ媒体の動的プロビジョニングを容易にする仮想ストレージインターフェイスをさらに含むことができ、物理ストレージ装置は、局所的又は遠隔的に存在することができる。

【0157】

さらに他の実施形態によれば、装置は、ストレージ媒体の動的プロビジョニングを容易にする仮想ネットワークインターフェイスをさらに含むことができ、物理ストレージは、局所的又は遠隔的に存在することができる。さらに他の実施形態によれば、装置を、ホストコンピュータシステム上のコントローラカードとして導入されるように適合することができ、ゲートウェイは、ホストコンピュータシステムのオペレーティングシステムによる介入を伴わずに動作する。

【0158】

さらに他の実施形態によれば、装置は、装置のストレージ機能及びネットワーク機能の少なくとも一方を提供する少なくとも1つのフィールドプログラマブルゲートアレイを含むことができる。さらに他の実施形態によれば、装置を、ネットワーク展開されたスイッチとして構成することができる。さらに他の実施形態によれば、装置は、ストレージ媒体命令を第1のプロトコルと少なくとも1つの他のプロトコルとの間で変換する、装置の機能コンポーネントをさらに含むことができる。

【0159】

図26に、例示的かつ非限定的なストレージ装置の仮想化方法を示す。まず、ステップ2600において、第1のストレージプロトコルでの命令に応答する物理ストレージ装置にアクセスする。次にステップ2602において、第1のストレージプロトコルと第2ストレージプロトコルとの間で命令を変換する。最後に、ステップ2604において、第2のプロトコルを用いて物理ストレージ装置をオペレーティングシステムに提示することにより、オペレーティングシステムを使用するホストコンピュータシステムに対して物理ストレージ装置が局所的に存在するか、それとも遠隔的に存在するかに関わらず、物理ストレージ装置のストレージが動的に展開されるようにする。

【0160】

様々な実施形態によれば、第1のプロトコルは、SATAプロトコル、NVMeプロトコル、SASプロトコル、iSCSIプロトコル、ファイバチャネルプロトコル及びファイバチャネルオーバーサネット(登録商標)プロトコルのうちの少なくとも1つである。他の実施形態では、第2のプロトコルがNVMeプロトコルである。

【0161】

いくつかの実施形態では、方法は、オペレーティングシステムと、第1及び第2ストレージプロトコル間における命令の変換を行う装置との間のインターフェイスを提供するステップ、及び/又は命令の変換を行う装置と、遠隔地に存在するネットワーク展開されたストレージ装置との間のNVMeオーバーサネット(登録商標)接続を提供するステップをさらに含むことができる。

【0162】

図27に、アプリケーション及びコンテナの少なくとも一方の移動を容易にする例示的

10

20

30

40

50

かつ非限定的な方法を示す。まずステップ 2700 において、集中型ストレージ及びネットワークコントローラを提供し、ゲートウェイが、ホストコンピュータのオペレーティングシステムの介入を伴わずに装置のストレージコンポーネントとネットワークコンポーネントとの間のネットワーク及びストレージトラフィックのための接続を提供する。次にステップ 2702 において、少なくとも 1 つのアプリケーション又はコンテナを、集中型ストレージ及びネットワークコントローラによって制御されるターゲット物理ストレージ装置にマッピングすることにより、アプリケーション又はコンテナが別のコンピュータシステムに移動する際に、アプリケーション又はコンテナが、ターゲット物理ストレージが接続されたホストシステムのオペレーティングシステムの介入を伴わずにターゲット物理ストレージにアクセスできるようにする。

10

#### 【0163】

様々な実施形態によれば、この移動は、Linux (登録商標) コンテナ又はスケールアウトアプリケーションの移動である。

#### 【0164】

さらに他の実施形態によれば、ターゲット物理ストレージは、iSCSI プロトコル、ファイバチャネルプロトコル及びファイバチャネルオーバーサネット (登録商標) プロトコルのうちの少なくとも 1 つを使用するネットワーク展開されたストレージ装置である。さらに他の実施形態では、ターゲット物理ストレージが、SAS プロトコル、SATA プロトコル及び NVMe プロトコルのうちの少なくとも 1 つを使用するディスク接続されたストレージ装置である。

20

#### 【0165】

図 28 に、ネットワークのサービス品質 (QoS) を提供する例示的かつ非限定的な方法を示す。まずステップ 2800 において、集中型ストレージ及びネットワークコントローラを提供し、ゲートウェイが、ホストコンピュータのオペレーティングシステムの介入を伴わずに装置のストレージコンポーネントとネットワークコンポーネントとの間のネットワーク及びストレージトラフィックのための接続を提供する。次にステップ 2802 において、集中型ストレージ及びネットワークコントローラによって処理されるストレージトラフィック及びネットワークトラフィックの少なくとも一方に基づいて、ホストコンピュータのオペレーティングシステムの介入を伴わずに、ストレージ及びネットワークコントローラが展開されたデータ経路を有するネットワークに関連する少なくとも 1 つのサービス品質 (QoS) パラメータを管理する。

30

#### 【0166】

本開示のわずかな実施形態の図示及び説明しかしていないが、当業者には、以下の特許請求の範囲に示す本開示の思想及び範囲から逸脱することなく多くの変更及び修正を行い得ることが明らかであろう。本明細書で引用した全ての外国及び国内の特許出願及び特許、並びに他の全ての出版物は、法律の許す最大限までその全体が本明細書に組み入れられる。

#### 【0167】

本明細書で説明した方法及びシステムは、プロセッサ上でコンピュータソフトウェア、プログラムコード及び / 又は命令を実行する機械を介して部分的に又は全体的に展開される。本開示は、機械における方法として、機械の一部としての又は機械に関するシステム又は装置として、或いは 1 又は 2 以上の機械において実行されるコンピュータ可読媒体内に具体化されたコンピュータプログラム製品として実装することができる。実施形態では、プロセッサを、サーバ、クラウドサーバ、クライアント、ネットワークインフラストラクチャ、モバイルコンピュータプラットフォーム、定置式コンピュータプラットフォーム、又はその他のコンピュータプラットフォームの一部とすることができる。プロセッサは、プログラム命令、コード及びバイナリ命令などを実行できるあらゆる種類のコンピュータ装置又は処理装置とすることができる。プロセッサは、単一プロセッサ、デジタルプロセッサ、組み込みプロセッサ、マイクロプロセッサ、或いは記憶されているプログラムコード又はプログラム命令の実行を直接的又は間接的に容易にすることができるコプロセッ

40

50

サ（マスコプロセッサ、グラフィックコプロセッサ及び通信コプロセッサなど）などのいずれかの変種とすることができ、又はこれらを含むことができる。また、プロセッサは、マルチプログラム、スレッド及びコードの実行を可能にすることもできる。これらのスレッドは、プロセッサの性能を高めてアプリケーションの同時動作を容易にするように同時に実行することができる。一実装として、本明細書で説明した方法、プログラムコード、プログラム命令などを1又は2以上のスレッドの形で実装することもできる。スレッドは、これらに関連する優先順位を割り当てられることができる他のスレッドを引き起こすことができ、プロセッサは、プログラムコードとして提供される命令に基づく優先順位又は他のいずれかの順序に基づいてこれらのスレッドを実行することができる。プロセッサ、又はプロセッサを利用するあらゆる機械は、本明細書及び他の場所で説明した方法、コード、命令及びプログラムを記憶する非一時的メモリを含むことができる。プロセッサは、本明細書及び他の場所で説明した方法、コード及び命令を記憶できる非一時的ストレージ媒体にインターフェイスを介してアクセスすることができる。プロセッサに関連する、コンピュータ装置又は処理装置によって実行できる方法、プログラム、コード、プログラム命令又は他のタイプの命令を記憶するストレージ媒体としては、以下に限定されるわけではないが、CD-ROM、DVD、メモリ、ハードディスク、フラッシュドライブ、RAM、ROM及びキャッシュなどのうちの1つ又は2つ以上を挙げることができる。

10

#### 【0168】

プロセッサは、マルチプロセッサの速度及び性能を高めることができる1又は2以上のコアを含むことができる。実施形態では、プロセスを、デュアルコアプロセッサ、クアッドコアプロセッサ、及び2又は3以上の（チップと呼ばれる）独立コアを組み合わせたその他のチップレベルマルチプロセッサなどとすることができ。

20

#### 【0169】

本明細書で説明した方法及びシステムは、サーバ、クライアント、ファイヤウォール、ゲートウェイ、ハブ、ルータ、又はその他のこのようなコンピュータハードウェア及び/又はネットワーキングハードウェア上でコンピュータソフトウェアを実行する機械を通じて部分的に又は全体的に展開することができる。ソフトウェアプログラムは、ファイルサーバ、プリントサーバ、ドメインサーバ、インターネットサーバ、イントラネットサーバ、クラウドサーバ、並びに二次サーバ、ホストサーバ及び分散サーバなどの他の変種を含むことができるサーバに関連することができる。サーバは、メモリ、プロセッサ、コンピュータ可読媒体、ストレージ媒体、（物理及び仮想）ポート、通信装置、並びに有線又は無線媒体を介して他のサーバ、クライアント、機械及び装置などにアクセスできるインターフェイスのうちの1つ又は2つ以上を含むことができる。本明細書及び他の場所で説明した方法、プログラム又はコードは、サーバによって実行することもできる。また、本出願で説明した方法の実行に必要な他の装置は、サーバに関連するインフラストラクチャの一部と見なすことができる。

30

#### 【0170】

サーバは、以下に限定されるわけではないが、クライアント、他のサーバ、プリンタ、データベースサーバ、プリントサーバ、ファイルサーバ、通信サーバ、分散サーバ及びソーシャルネットワークなどを含む他の装置へのインターフェイスを提供することができる。また、この結合及び/又は接続は、ネットワークを越えた遠隔的なプログラムの実行を容易にすることもできる。これらの装置の一部又は全部のネットワーキングは、本開示の範囲から逸脱することなく、1又は2以上の位置におけるプログラム又は方法の並行処理を容易にすることができる。また、インターフェイスを介してサーバに接続された装置は、いずれも方法、プログラム、コード及び/又は命令を記憶できる少なくとも1つのストレージ媒体を含むことができる。異なる装置上で実行されるプログラム命令は、中央リポジトリによって提供することができる。この実装では、遠隔リポジトリが、プログラムコード、命令及びプログラムのためのストレージ媒体として機能することができる。

40

#### 【0171】

ソフトウェアプログラムは、ファイルクライアント、プリントクライアント、ドメイン

50



クライアント、インターネットクライアント、イントラネットクライアント、並びに二次クライアント、ホストクライアント及び分散クライアントなどの他の変種を含むことができるクライアントに関連することができる。クライアントは、メモリ、プロセッサ、コンピュータ可読媒体、ストレージ媒体、（物理及び仮想）ポート、通信装置、並びに有線又は無線媒体を介して他のクライアント、サーバ、機械及び装置などにアクセスできるインターフェイスのうちの1つ又は2つ以上を含むことができる。本明細書及び他の場所で説明した方法、プログラム又はコードは、クライアントによって実行することもできる。また、本出願で説明した方法の実行に必要な他の装置は、クライアント関連するインフラストラクチャの一部と見なすことができる。

#### 【0172】

クライアントは、以下に限定されるわけではないが、サーバ、他のクライアント、プリンタ、データベースサーバ、プリントサーバ、ファイルサーバ、通信サーバ及び分散サーバなどを含む他の装置へのインターフェイスを提供することができる。また、この結合及び/又は接続は、ネットワークを越えた遠隔的なプログラムの実行を容易にすることもできる。これらの装置の一部又は全部のネットワークングは、本開示の範囲から逸脱することなく、1又は2以上の位置におけるプログラム又は方法の並行処理を容易にすることができる。また、インターフェイスを介してクライアントに接続された装置は、いずれも方法、プログラム、アプリケーション、コード及び/又は命令を記憶できる少なくとも1つのストレージ媒体を含むことができる。異なる装置上で実行されるプログラム命令は、中央リポジトリによって提供することができる。この実装では、遠隔リポジトリが、プログラムコード、命令及びプログラムのためのストレージ媒体として機能することができる。

#### 【0173】

本明細書で説明した方法及びシステムは、ネットワークインフラストラクチャを通じて部分的に又は全体的に展開することができる。ネットワークインフラストラクチャは、当業で周知のコンピュータ装置、サーバ、ルータ、ハブ、ファイアウォール、クライアント、パーソナルコンピュータ、通信装置、ルーティング装置、及びその他の能動装置及び受動装置、モジュール及び/又はコンポーネントなどの要素を含むことができる。ネットワークインフラストラクチャに関連するコンピュータ装置及び/又は非コンピュータ装置は、他のコンポーネントとは別に、フラッシュメモリ、バッファ、スタック、RAM及びROMなどのストレージ媒体を含むことができる。本明細書及び他の場所で説明した処理、方法、プログラムコード、命令は、これらのネットワークインフラストラクチャ要素のうちの1つ又は2つ以上によって実行することができる。本明細書で説明した方法及びシステムは、サービス型ソフトウェア（SaaS）、サービス型プラットフォーム（PaaS）、及び/又はサービス型インフラストラクチャ（IaaS）の特徴を伴うものを含むあらゆる種類のプライベートネットワーク、コミュニティネットワーク、又は混成クラウドコンピューティングネットワーク又はクラウドコンピューティング環境と共に使用するように適合することができる。

#### 【0174】

本明細書及び他の場所で説明した方法、プログラムコード及び命令は、セルラーネットワークh a a送信者制御コンタクト媒体コンテンツアイテムマルチセルにおいて実装することができる。セルラーネットワークは、周波数分割多重アクセス（FDMA）ネットワーク又は符号分割多重アクセス（CDMA）ネットワークのいずれかとすることができる。セルラーネットワークは、モバイル装置、セルサイト、基地局、リピータ、アンテナ及びタワーなどを含むことができる。セルラーネットワークは、GSM（登録商標）、GPRS、3G、EVDO、メッシュ又はその他のネットワークタイプを含むことができる。

#### 【0175】

本明細書及び他の場所で説明した方法、プログラムコード及び命令は、モバイル装置上で又はモバイル装置を介して実装することができる。モバイル装置は、ナビゲーション装置、セルラー電話機、携帯電話機、携帯情報端末、ラップトップ、パームトップ、ネットブック、ページャ、電子ブックリーダー及び音楽プレーヤなどを含むことができる。これら

の装置は、他のコンポーネントとは別に、フラッシュメモリ、バッファ、RAM、ROM及び1又は2以上のコンピュータ装置などのストレージ媒体を含むことができる。モバイル装置に関連するコンピュータ装置は、記憶されているプログラムコード、方法及び命令を実行することができる。或いは、モバイル装置は、他の装置と協調して命令を実行するように構成することもできる。モバイル装置は、サーバにインターフェイス接続されてプログラムコードを実行するように構成された基地局と通信することができる。モバイル装置は、ピア・ツー・ピアネットワーク、メッシュネットワーク又はその他の通信ネットワーク上で通信することができる。プログラムコードは、サーバに関連するストレージ媒体に記憶し、サーバに組み込まれたコンピュータ装置によって実行することができる。基地局は、コンピュータ装置及びストレージ媒体を含むことができる。ストレージ装置は、基地局に関連するコンピュータ装置によって実行されるプログラムコード及び命令を記憶することができる。

10

#### 【0176】

コンピュータソフトウェア、プログラムコード及び/又は命令は、何らかの時間間隔にわたって計算に使用されるデジタルデータを保持するコンピュータコンポーネント、装置及び記録媒体、ランダムアクセスメモリ(RAM)として知られている半導体ストレージ、光学ディスク、ハードディスク、テープ、ドラム、カード及びその他のタイプ様の磁気ストレージの形態などの、通常はさらに永続的な記憶のための大容量ストレージ、プロセッサレジスタ、キャッシュメモリ、揮発性メモリ、不揮発性メモリ、CD、DVDなどの光学ストレージ、フラッシュメモリ(例えば、USBスティック又はキー)、フロッピーディスク、磁気テープ、紙テープ、パンチカード、スタンドアロン型RAMディスク、Zipドライブ、取り外し可能大容量ストレージ及びオフラインなどの取り外し可能媒体、動的メモリ、静的メモリ、読み取り/書き込みストレージ、可変ストレージ、リードオンリ、ランダムアクセス、順次アクセス、位置アドレス指定可能、ファイルアドレス指定可能、コンテンツアドレス指定可能、ネットワーク接続ストレージ、ストレージエリアネットワーク、バーコード、磁気インクなどの他のコンピュータメモリを含むことができる機械可読媒体に記憶し及び/又は機械可読媒体においてアクセスすることができる。

20

#### 【0177】

本明細書で説明した方法及びシステムは、物理項目及び/又は無形項目を1つの状態から別の状態に変換することができる。本明細書で説明した方法及びシステムは、物理項目及び/又は無形項目を表すデータを1つの状態から別の状態に変換することもできる。

30

#### 【0178】

図を通じたフローチャート及びブロック図に示したものを含む、本明細書で説明し図示した要素は、要素間の論理的境界を示唆するものである。しかしながら、ソフトウェア又はハードウェア工学の慣例に従い、図示の要素及びその機能は、コンピュータ実行可能媒体h a a送信者制御コンタクト媒体コンテンツアイテムを介して機械上の実装することができる。プロセッサは、記憶されているプログラム命令を、モノリシックソフトウェア構造、独立ソフトウェアモジュール、或いは外部ルーチン、コード及びサービスなど、又はこれらのいずれかの組み合わせを採用するモジュールとして実行することができる。全てのこのような実装は、本開示の範囲に含まれるものとする。このような機械の例としては、以下に限定されるわけではないが、携帯情報端末、ラップトップ、パーソナルコンピュータ、携帯電話機、その他のハンドヘルドコンピュータ装置、医療機器、有線又は無線通信装置、トランスデューサ、チップ、計算機、衛星、タブレットPC、電子ブック、ガジェット、電子装置、装置h a a送信者制御コンタクト媒体コンテンツアイテム人工知能、コンピュータ装置、ネットワーキング装置、サーバ及びルータなどを挙げることができる。また、フローチャート及びブロック図に示した要素又は他のいずれかの論理コンポーネントは、プログラム命令を実行できる機械上の実装することもできる。従って、上述した図面及び説明には、開示するシステムの機能面を示しているが、明確に示しているか、或いは文脈から別様に明らかでない限り、これらの説明からこれらの機能面を実装するソフトウェアの特定の構成を推測すべきではない。同様に、上記で識別し説明した様

40

50

々なステップは変更することができ、ステップの順序は、本明細書に開示した技術の特定の応用に適合することができると理解されるであろう。全てのこのような変形及び修正は、本開示の範囲内に含まれるものとする。従って、様々なステップの順番の図示及び／又は説明については、特定の用途に必要な限り、又は明確に示していない限り、或いは文脈から別様に明らかでない限り、これらのステップの特定の実行順を必要とするものであると理解すべきではない。

【0179】

上述した方法及び／又はプロセス、及びこれらに関連するステップは、ハードウェア、ソフトウェア、又は特定の応用に適したハードウェアとソフトウェアのあらゆる組み合わせで実現することができる。ハードウェアは、汎用コンピュータ及び／又は専用コンピュータ装置、或いは特定のコンピュータ装置又は特定のコンピュータ装置の特定の態様又はコンポーネントを含むことができる。これらのプロセスは、1又は2以上のマイクロプロセッサ、マイクロコントローラ、組み込み式マイクロコントローラ、プログラマブルデジタルシグナルプロセッサ又は他のプログラマブル装置において、内部及び／又は外部メモリと共に実現することができる。これに加えて、又はこれとは別に、これらのプロセスは、特定用途向け集積回路、プログラマブルゲートアレイ、プログラマブルアレイロジック、又は電子信号を処理するように構成できる他のいずれかの装置又は装置の組み合わせで具体化することもできる。さらに、これらのプロセスの1つ又は2つ以上を、機械可読媒体上で実行できるコンピュータ実行可能コードとして実現することもできると理解されるであろう。

10

20

【0180】

コンピュータ実行可能コードは、Cなどの構造化プログラミング言語、C++などのオブジェクト志向プログラミング言語、或いは上記の装置のうちの1つ、並びにプロセッサの異種の組み合わせ、プロセッサアーキテクチャ、又は異なるハードウェアとソフトウェアとの組み合わせ、又はプログラム命令を実行できる他のいずれかの機械において実行されるように記憶し、コンパイルし、又は解釈することができる（アセンブリ言語、ハードウェア記述言語、及びデータベースプログラミング言語及び技術を含む）他のいずれかの高水準又は低水準プログラミング言語を用いて作成することができる。

【0181】

従って、1つの態様では、上述した方法及び方法の組み合わせを、1又は2以上のコンピュータ装置において実行された時に方法のステップを実行するコンピュータ実行可能コードで具体化することができる。別の態様では、これらの方法を、そのステップを実行するシステム内で具体化し、複数の形で装置全体に分散させることも、又は全ての機能を専用のスタンドアロン型装置又はその他のハードウェアに統合することもできる。別の態様では、上述したプロセスに関連するステップを実行する手段が、上述したハードウェア及び／又はソフトウェアのいずれかを含むことができる。全てのこのような置換及び組み合わせは、本開示の範囲に含まれるものとする。

30

【0182】

詳細に図示し説明した好ましい実施形態に関連して開示を行ったが、当業者には、本開示に対する様々な修正及び改善が容易に明らかになるであろう。従って、本開示の思想及び範囲は、上記の例によって限定されるものではなく、法律によって認められる限りの広い意味で理解すべきである。

40

【0183】

本開示を説明する文脈における（特に以下の特許請求の範囲の文脈における）「1つの（英文不定冠詞）」及び「その（英文定冠詞）」という用語の使用、並びに同様の言及は、本明細書で別途示していない限り、或いは文脈によって明らかに矛盾しない限り、単数形及び複数形の両方を含むと解釈されたい。「含む（comprising）」、「has送信者制御コンタクト媒体コンテンツアイテム」、及び「含む（including、containing）」という用語は、別途示していない限り、非制限的な用語（すなわち、「含むけれども限定されるわけではない（including、but not

50

limited to)」として解釈されたい。本明細書における値の範囲の記述は、本明細書において別途示していない限り、その範囲内に収まる各別個の値を個別に示す簡略表現方法の役割を果たすものにすぎず、各別個の値は、本明細書に個別に示されているかのように明細書に組み入れられる。本明細書で説明した全ての方法は、本明細書において別途示していない限り、又は文脈によって明らかに矛盾しない限り、あらゆる好適な順序で実行することができる。本明細書に示すありとあらゆる例の使用、又は例示的言語（例えば、「～などの（such as）」）は、本開示の理解を容易にするためのものにすぎず、別途主張していない限り、本開示の範囲に限定をもたらすものではない。本明細書におけるあらゆる表現は、本開示を実施するために必須のいずれかの非請求要素を示すものとして解釈すべきではない。

10

**【0184】**

上述した明細書は、現在のところ本発明の最良の形態であると考えられるものを当業者が実施して使用できるようにするものであり、当業者であれば、本明細書における特定の実施形態、方法及び実施例の変種、組み合わせ及び同等物の存在を理解して評価するであろう。従って、本開示の限定は、上述した実施形態、方法及び実施例によってではなく、本開示の範囲及び思想に含まれる全ての実施形態及び方法によって行うべきである。

**【0185】**

本明細書で参照した全ての文献は、引用により本明細書に組み入れられる。

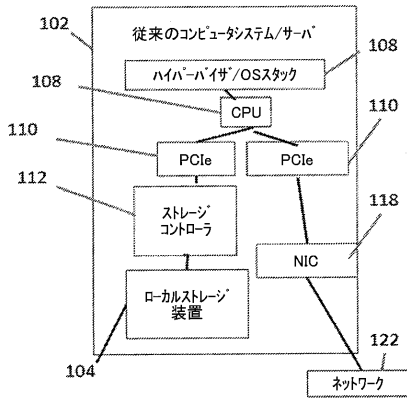
**【符号の説明】****【0186】**

102 コンピュータシステム  
106 CPU  
108 メモリ  
110 PCIバス  
112 ストレージコントローラ  
118 ネットワークコントローラ  
300 集中型IOコントローラ  
302 ストレージ装置  
304 カットスルーデータ経路  
308 DAS  
310 SAN

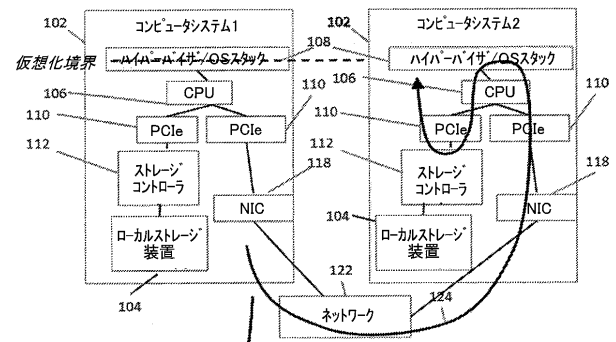
20

30

【図 1】

Fig. 1  
先行技術

【図 2】

Fig. 2  
先行技術

【図 3】

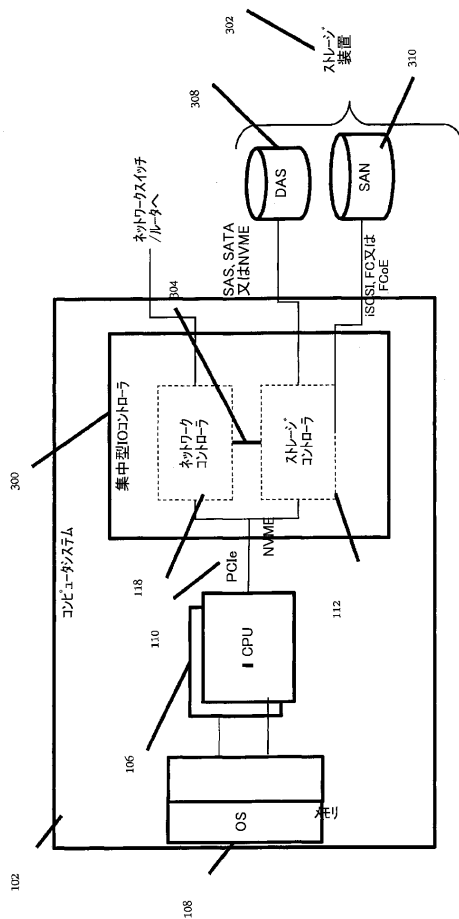


Fig. 3

【図 4】

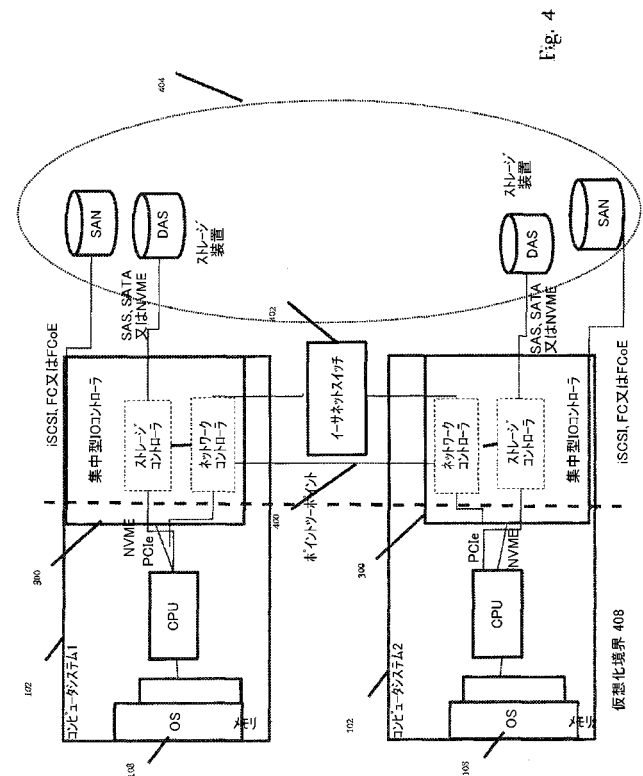
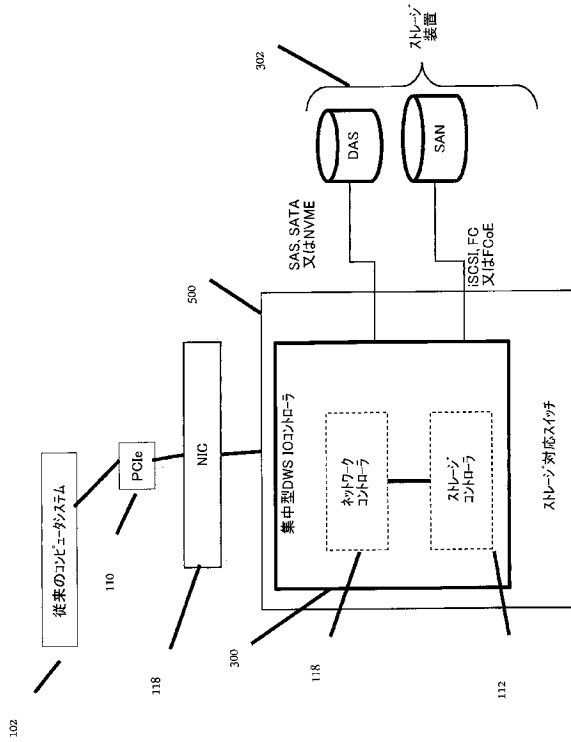


Fig. 4

【図 5】



【図 6】

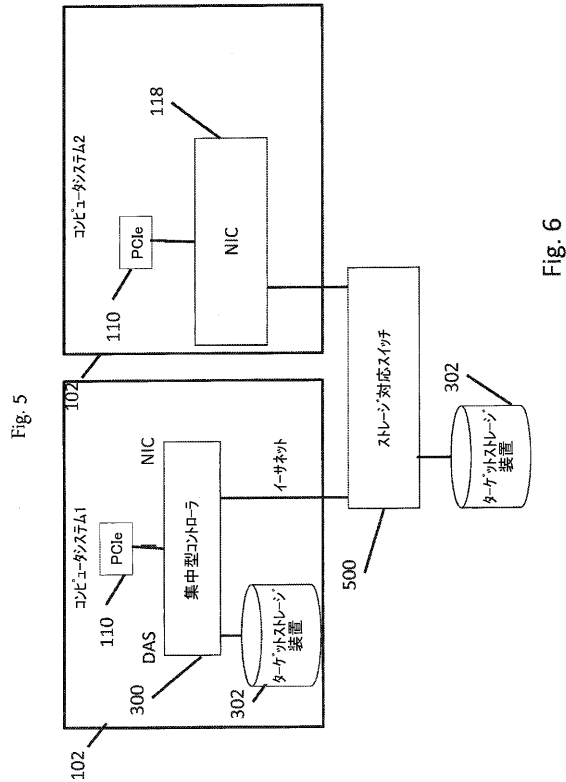


Fig. 6

【図 7】

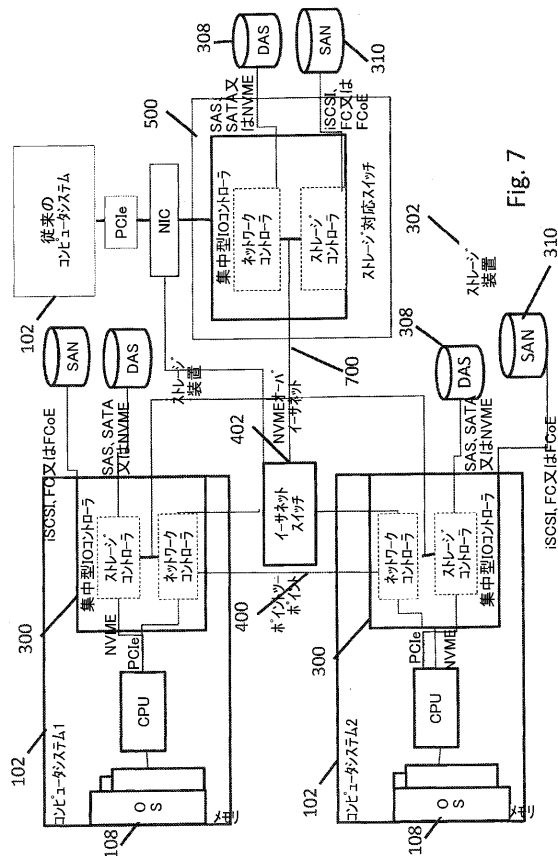


Fig. 7

【図 8】

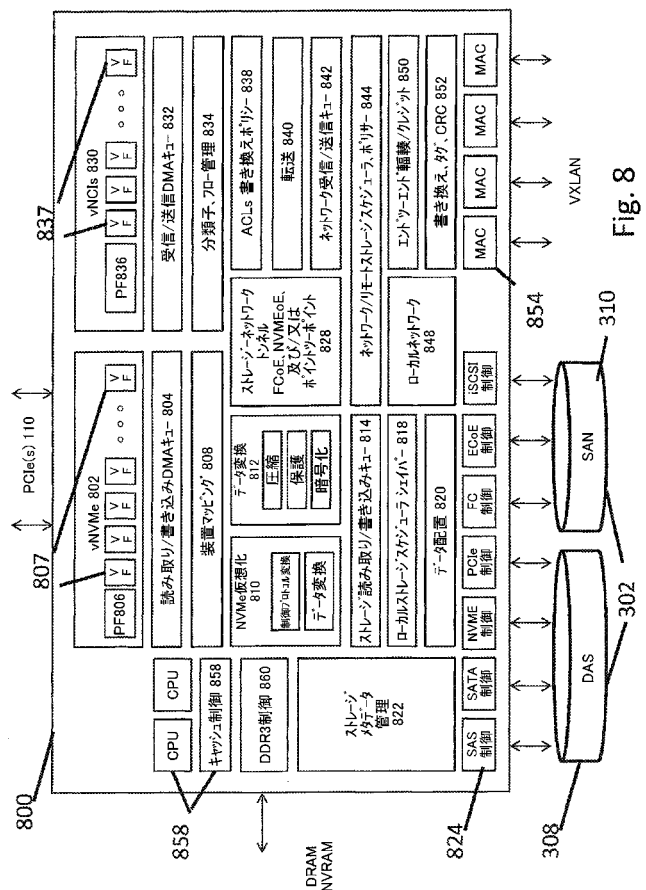


Fig. 8

【図 9】

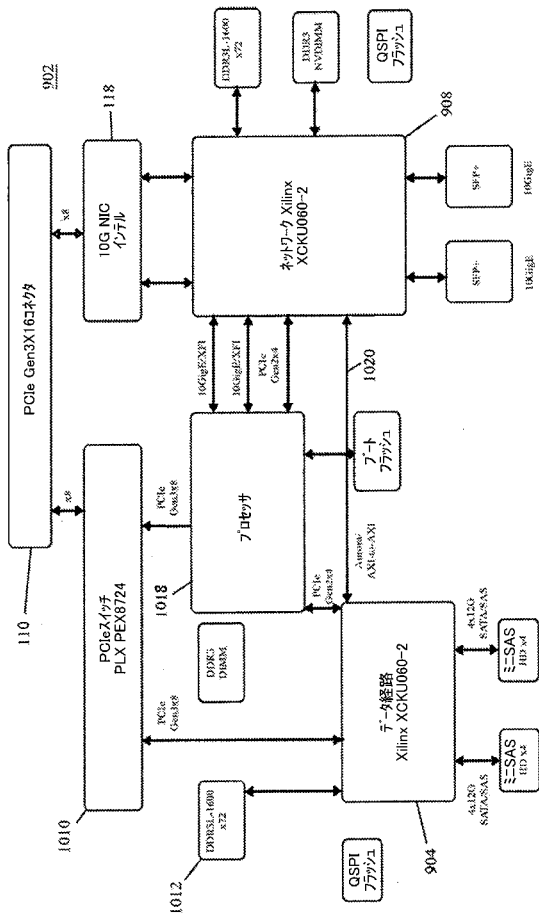


Fig. 9

【図 10】

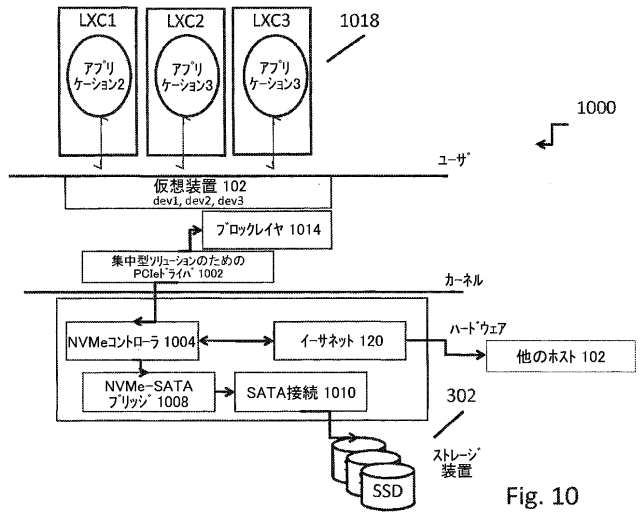


Fig. 10

【図 11】

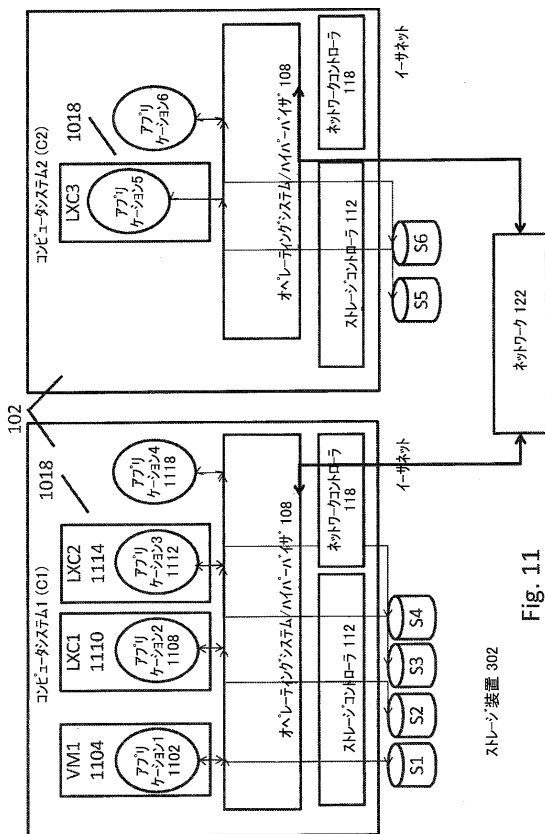


Fig. 11

【図 12】

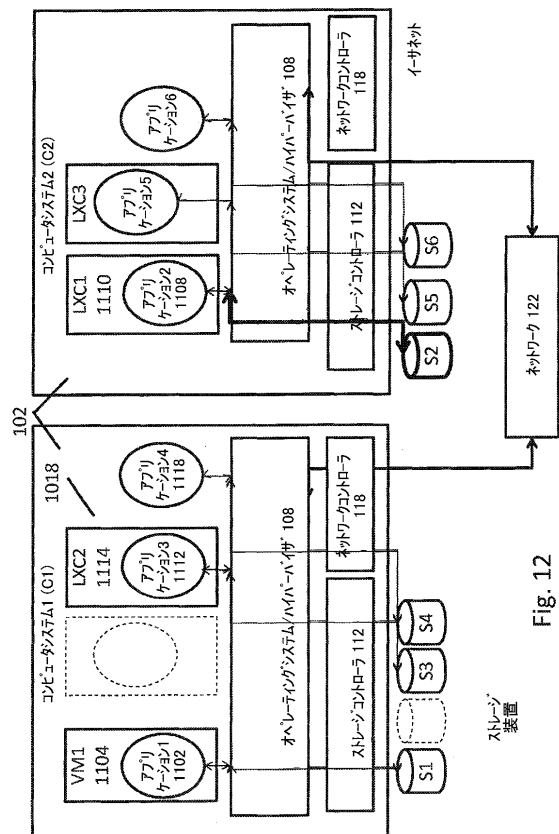


Fig. 12

【図 13】

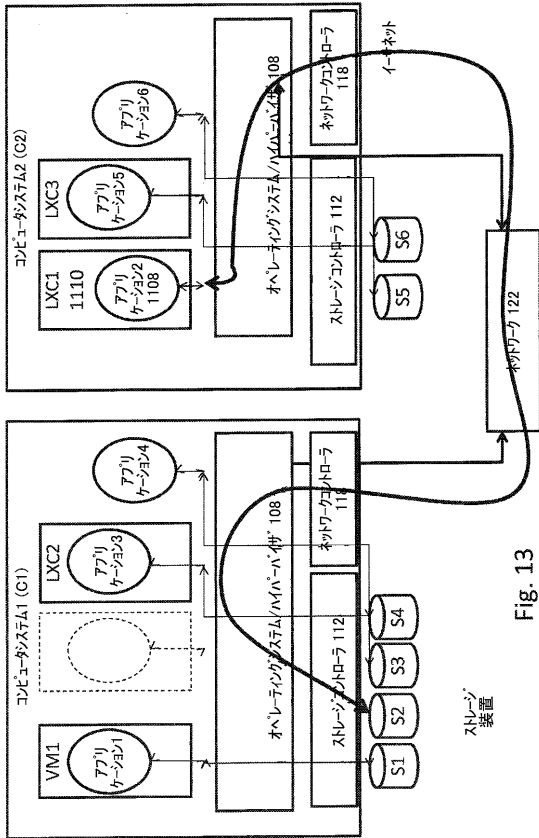


Fig. 13

【図 14】

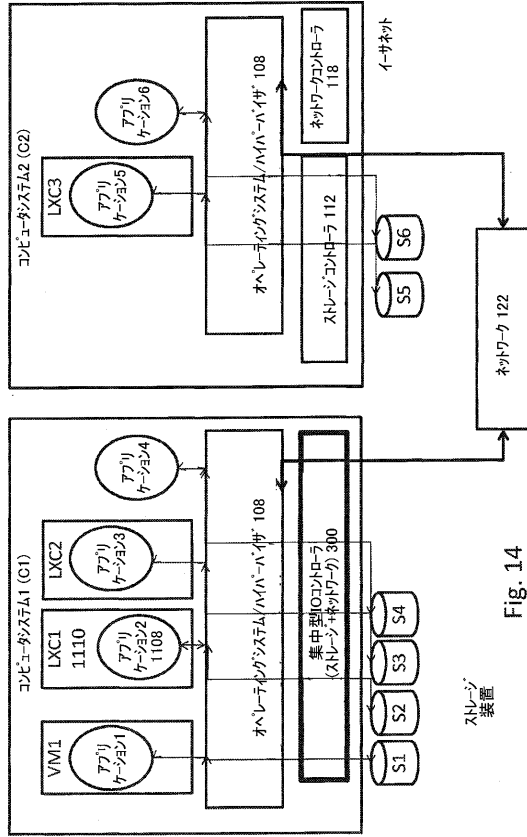


Fig. 14

【図 15】

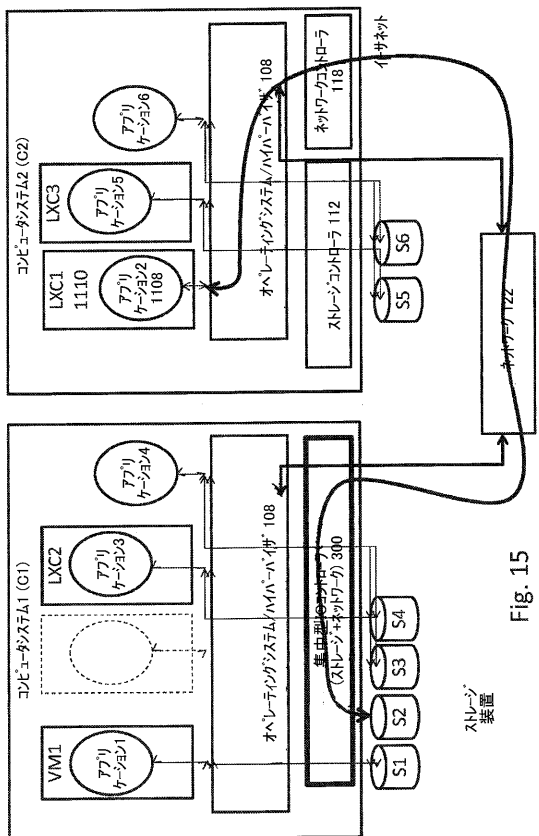


Fig. 15

【図 16】

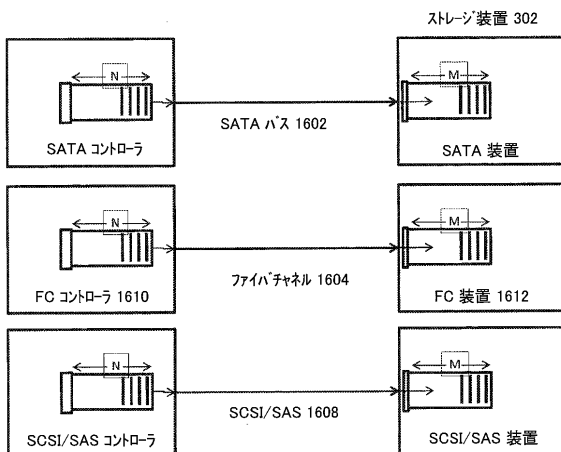


Fig. 16



【図 17】

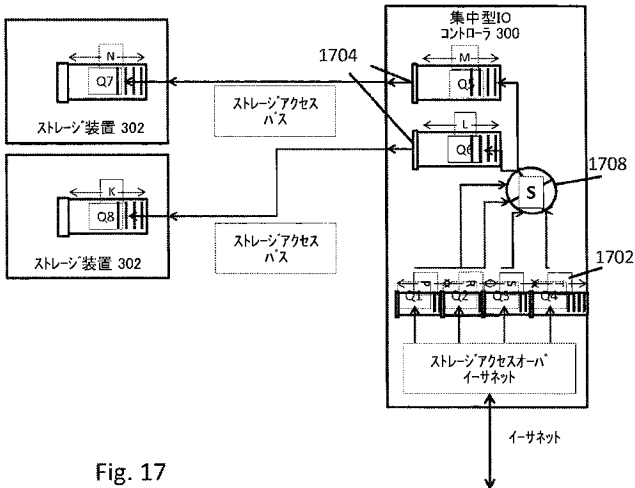


Fig. 17

【図 18】

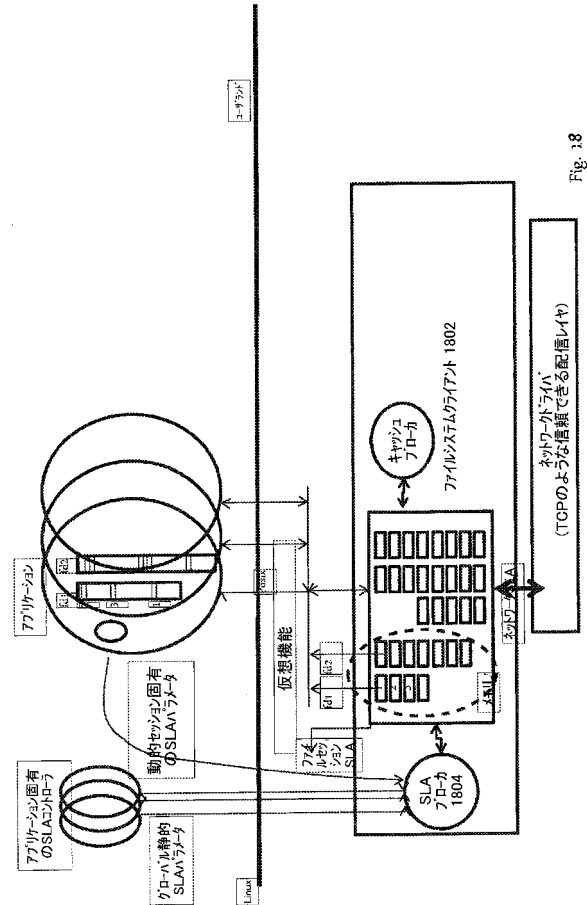


Fig. 18

【図 19】

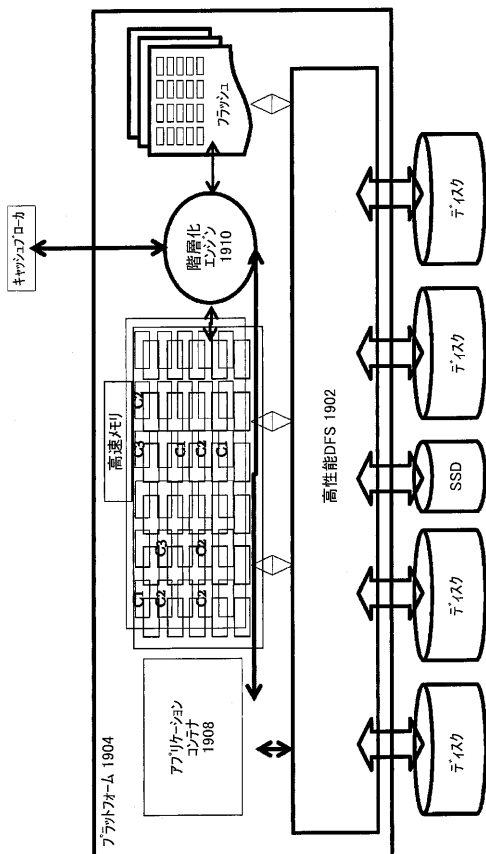


Fig. 19

【図 20】

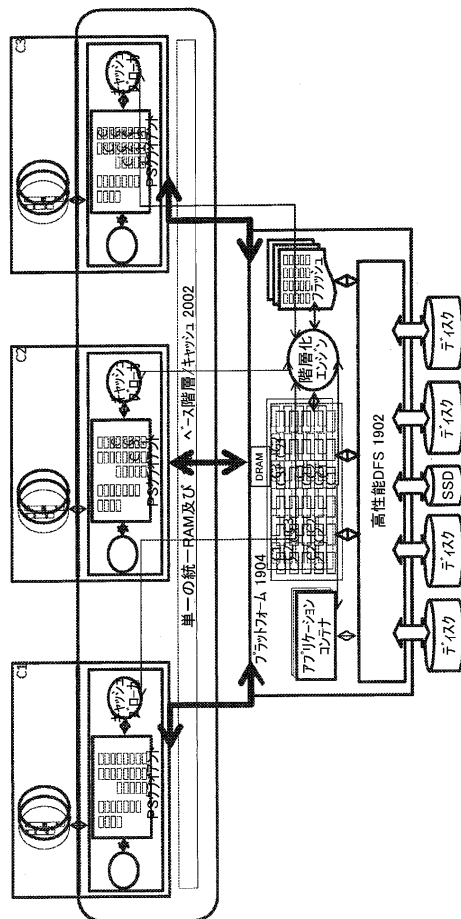
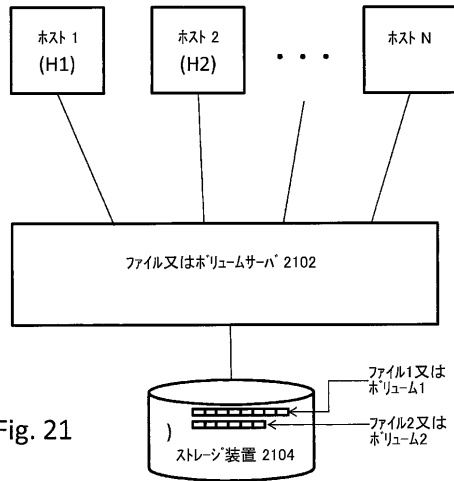
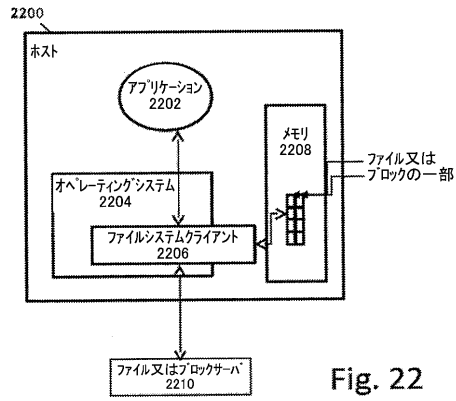


Fig. 20

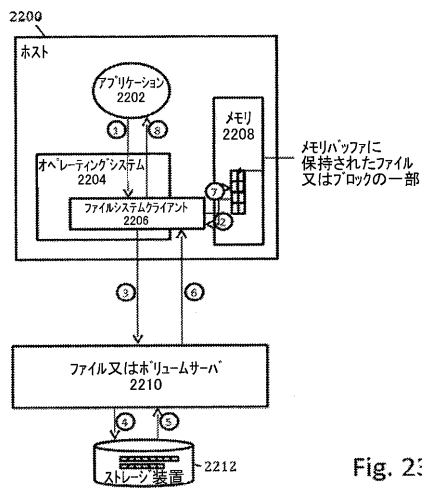
【図 2 1】



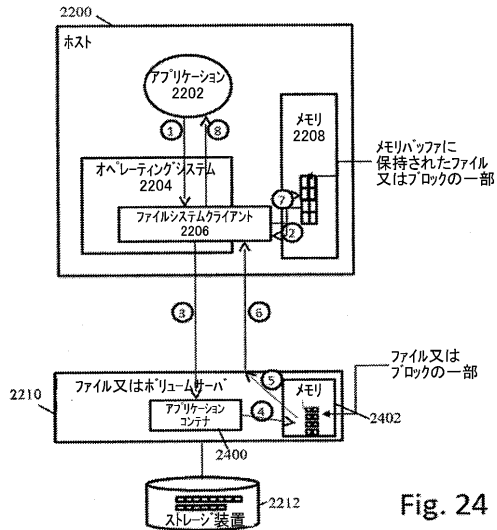
【図 2 2】



【図 2 3】



【図 2 4】



【図 25】

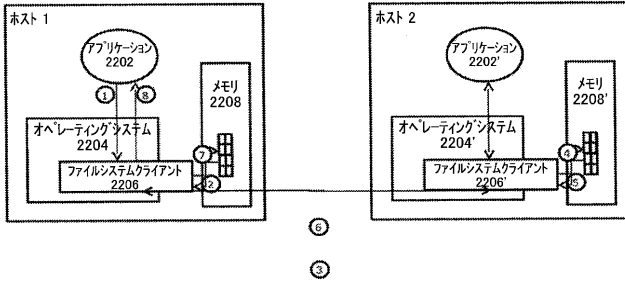


Fig. 25

【図 26】

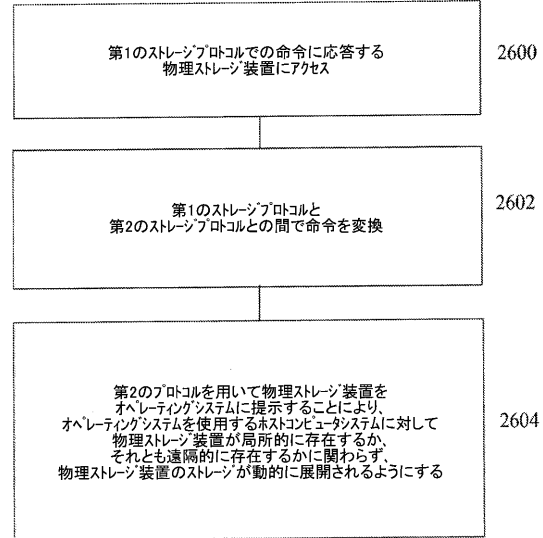


Fig. 26

【図 27】

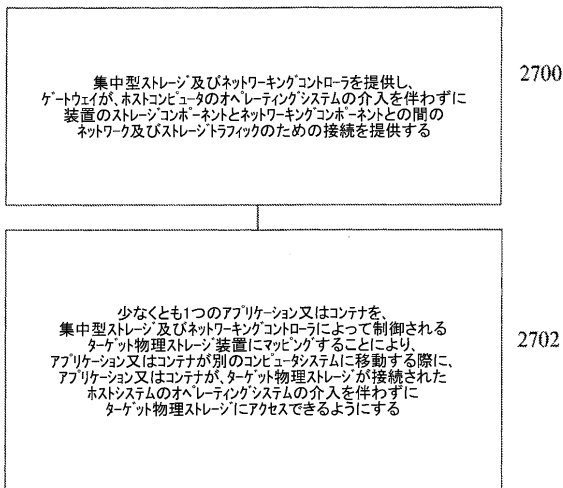


Fig. 27

【図 28】

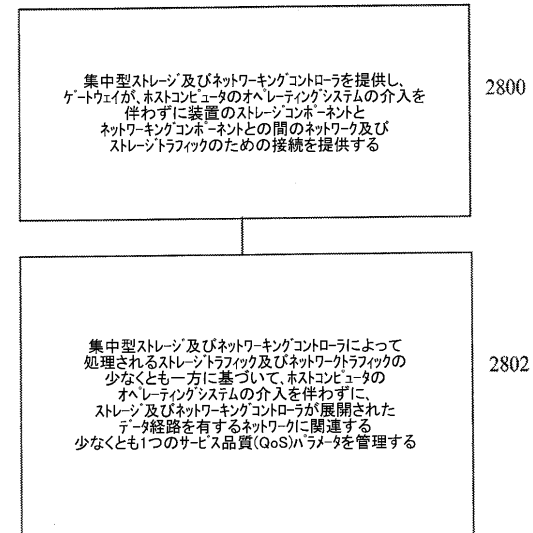




Fig. 28

## 【 国際調査報告 】

<b>INTERNATIONAL SEARCH REPORT</b>		International application No. <b>PCT/US2015/019206</b>
<b>A. CLASSIFICATION OF SUBJECT MATTER</b> <b>H04L 29/06(2006.01)i, G06F 3/06(2006.01)i</b>		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b>		
Minimum documentation searched (classification system followed by classification symbols) H04L 29/06; G06F 12/06; G06F 15/167; H04L 12/28; G06F 15/173; H04L 12/24; G06F 3/06		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Korean utility models and applications for utility models Japanese utility models and applications for utility models		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) cKOMPASS(KIPO internal) & Keywords: converged, IO, network, storage, controller, gateway, QoS, SATA, NVMe, protocol		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 2010-0005234 A1 (ILANGO S. GANGA et al.) 07 January 2010 See paragraphs [0012]-[0017], [0021]-[0028], [0036], [0039]; and figure 1.	1-10
Y		11-27
Y	US 2013-0198312 A1 (ELIEZER TAMIR et al.) 01 August 2013 See paragraphs [0021], [0023], [0037], [0062], [0084], [0087], [0100]-[0101]; and figures 5-6, 9.	11-27
A	US 2009-0003361 A1 (RANGA BAKTHAVATHSALAM) 01 January 2009 See paragraphs [0040]-[0045]; and figures 3a-4c.	1-27
A	US 2013-0232267 A1 (KEVIN D. SHATZKAMER et al.) 05 September 2013 See paragraphs [0012]-[0013], [0033]-[0034]; and figures 1-2.	1-27
A	KR 10-2008-0052846 A (ELECTRONICS AND TELECOMMUNICATIONS RESEARCH INSTITUTE) 12 June 2008 See paragraphs [0026]-[0029]; and figure 2.	1-27
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 28 May 2015 (28.05.2015)		Date of mailing of the international search report 28 May 2015 (28.05.2015)
Name and mailing address of the ISA/KR  International Application Division Korean Intellectual Property Office 189 Cheongsa-ro, Seo-gu, Daejeon Metropolitan City, 302-701, Republic of Korea Facsimile No. +82-42-472-7140		Authorized officer KIM, Seong Woo Telephone No. +82-42-481-3348 

**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No.

**PCT/US2015/019206**

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 2010-0005234 A1	07/01/2010	CN 101621410 A DE 102009031126 A1 JP 2010-016819 A JP 5064447 B2 US 8359408 B2	06/01/2010 25/02/2010 21/01/2010 31/10/2012 22/01/2013
US 2013-0198312 A1	01/08/2013	CN 104246742 A DE 112013000601 T5 US 2013-198311 A1 US 2014-325013 A1 WO 2013-109640 A1	24/12/2014 18/12/2014 01/08/2013 30/10/2014 25/07/2013
US 2009-0003361 A1	01/01/2009	US 7917682 B2	29/03/2011
US 2013-0232267 A1	05/09/2013	None	
KR 10-2008-0052846 A	12/06/2008	None	

## フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US

(74)代理人 100109335

弁理士 上杉 浩

(74)代理人 100120525

弁理士 近藤 直樹

(74)代理人 100151987

弁理士 谷口 信行

(72)発明者 チョウ ジェフリー

アメリカ合衆国 カリフォルニア州 9 4 3 0 6 パロ アルト チマルス ドライヴ 6 6 8

(72)発明者 シャルマ ゴパール

アメリカ合衆国 カリフォルニア州 9 5 1 4 8 サン ホセ プレイス デ ルイス 3 6 3 1

(72)発明者 グハ アミタヴァ

アメリカ合衆国 カリフォルニア州 9 5 1 2 9 サン ホセ ロイヤル アン ドライヴ 6 1 2 7

(72)発明者 フォン ケヴィン

アメリカ合衆国 ネバダ州 8 9 1 3 4 ラス ベガス アイアン ヒル レーン 1 4 1 7

(72)発明者 ガネーシュ ジャヤセナン スンダラ

アメリカ合衆国 カリフォルニア州 9 5 0 1 4 クパチーノ ホリーヘッド レーン 1 1 0 1