

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5315209号
(P5315209)

(45) 発行日 平成25年10月16日(2013.10.16)

(24) 登録日 平成25年7月12日(2013.7.12)

(51) Int.Cl.

F I

G 0 6 F 13/10 (2006.01)

G 0 6 F 13/10 3 3 0 C

請求項の数 9 (全 28 頁)

(21) 出願番号	特願2009-246376 (P2009-246376)	(73) 特許権者	390009531
(22) 出願日	平成21年10月27日(2009.10.27)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公開番号	特開2010-140471 (P2010-140471A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公開日	平成22年6月24日(2010.6.24)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
審査請求日	平成24年5月11日(2012.5.11)		
(31) 優先権主張番号	12/332957	(74) 代理人	100108501
(32) 優先日	平成20年12月11日(2008.12.11)		弁理士 上野 剛史
(33) 優先権主張国	米国 (US)	(74) 代理人	100112690
			弁理士 太佐 種一
		(74) 代理人	100091568
			弁理士 市位 嘉宏

最終頁に続く

(54) 【発明の名称】 冗長構成を生成するための周辺機器相互接続入出力仮想化デバイスの使用

(57) 【特許請求の範囲】

【請求項 1】

冗長システム構成を生成するためのデータ処理システム実施方法であって、

各仮想機能について仮想機能パス許可テーブルを生成するステップであって、前記仮想機能パス許可テーブルは、当該仮想機能パス許可テーブルに関連付けられた仮想機能が複数のシステムにおける 1 組のアドレス範囲へアクセスする許可を得ているかを定義するエントリを複数有し、前記複数のエントリそれぞれが、無効な機能間アクセスを防止する境界をさらに定義し、前記仮想機能が、シングルルートまたはマルチルート周辺機器相互接続デバイスによって実行される、前記生成するステップと、

データをディスクに書き込むためのディスク書き込み動作のリクエストからのディスク書き込み要求を受信して、要求されたデータを前記仮想機能から提供するステップと、

前記 1 組のアドレス範囲のうちの選択された第 1 のアドレス範囲内に第 1 の受信バッファを生成し、且つ、前記 1 組のアドレス範囲のうちの選択された第 2 のアドレス範囲内に第 2 の受信バッファを生成するステップであって、前記選択された第 2 のアドレスは、前記選択された第 1 のアドレスを有するシステムと異なるシステム上にある、前記生成するステップと、

前記仮想機能のための仮想機能作業待ち行列エントリを生成するステップであって、前記仮想機能作業待ち行列エントリは、前記第 1 の受信バッファ及び前記第 2 の受信バッファの各アドレスを含む、前記生成するステップと、

前記仮想機能パス許可テーブルにおいて、前記各選択されたアドレス範囲の使用を前記

10

20

仮想機能が許可されているかどうかを判定するステップと、

前記仮想機能は許可されていることに応答して、前記要求されたデータを、前記第1の受信バッファ及び前記第2の受信バッファに書き込むステップと、

前記書き込みに応答して、前記リクエストに完了通知を発行するステップと
を含む、前記方法。

【請求項2】

前記仮想機能パス許可テーブル内にエントリが存在していることが、1次パスおよび対応する2次パスによる前記1つまたは複数のシステム内のシステムへの機能のアクセスを許可することである、請求項1に記載の方法。

【請求項3】

前記仮想機能パス許可テーブル内にエントリが存在していないことが、1次パスおよび対応する2次パスによる前記1つまたは複数のシステム内のシステムへの機能のアクセスを許可しないことである、請求項1に記載の方法。

【請求項4】

前記仮想機能パス許可テーブルが、各仮想機能と1組の仮想階層の間の対応、又は、各仮想機能と1組の論理パーティション内の1組のアドレス範囲の間の対応を含む、請求項1～3のいずれか一項に記載の方法。

【請求項5】

前記複数のエントリそれぞれが、前記仮想機能と前記複数のシステムのうちの一つのシステムにおけるアドレス範囲との間の1次通信パスと当該1次通信パスの代替通信パスとを有する、請求項1～4のいずれか一項に記載の方法。

【請求項6】

前記1次通信パスが優先パスであり、当該優先通信パスが利用不可能であることに応答して、前記代替通信パスの1つを使用する、請求項1～5のいずれか一項に記載の方法。

【請求項7】

冗長システム構成を生成するためのデータ処理システムであって、
バスと、

前記バスに接続されたコンピュータ実行可能命令を含むメモリと、

前記コンピュータ実行可能命令を実行する中央プロセッサ・ユニットであって、前記コンピュータ実行可能命令が、請求項1～6のいずれか一項に記載の方法の各ステップを実行するように前記データ処理システムに指令する、前記中央プロセッサ・ユニットと、
を備えている、前記データ処理システム。

【請求項8】

冗長システム構成を生成するためのデータ処理システムであって、
バスと、

前記バスに接続されたコンピュータ実行可能命令を含むメモリと、

前記コンピュータ実行可能命令を実行する中央プロセッサ・ユニットと

を備えており、

前記メモリが、

各仮想機能について仮想機能パス許可テーブルであって、前記仮想機能パス許可テーブルは、当該仮想機能パス許可テーブルに関連付けられた仮想機能が複数のシステムにおける1組のアドレス範囲へアクセスする許可を得ているかを定義するエントリを複数有し、前記複数のエントリそれぞれが、無効な機能間アクセスを防止する境界をさらに定義し、前記仮想機能が、シングルルートまたはマルチルート周辺機器相互接続デバイスによって実行される、前記仮想機能パス許可テーブルと、

前記仮想機能のための仮想機能作業待ち行列エントリであって、前記仮想機能作業待ち行列エントリは、前記1組のアドレス範囲のうちの選択された第1のアドレス範囲内に生成された第1の受信バッファ及び前記1組のアドレス範囲のうちの選択された第2のアドレス範囲内に生成された第2の受信バッファの各アドレスを含み、前記選択された第2のアドレスは、前記選択された第1のアドレスを有するシステムと異なるシステム上にあ

10

20

30

40

50

る、前記仮想機能作業待ち行列エントリと、

を記憶し、

前記コンピュータ実行可能命令が、

データをディスクに書き込むためのディスク書き込み動作のリクエストからのディスク書き込み要求を受信して、要求されたデータを前記仮想機能から提供することと、

前記仮想機能パス許可テーブルにおいて、前記各選択されたアドレス範囲の使用を前記仮想機能が許可されているかどうかを判定することと、

前記仮想機能は許可されていることに応答して、前記要求されたデータを、前記第 1 の受信バッファ及び前記第 2 の受信バッファに書き込むことと、

前記書き込みに応答して、前記リクエストに完了通知を発行すること

を実行するように前記データ処理システムに指令する、

前記データ処理システム。

【請求項 9】

冗長システム構成を生成するためのコンピュータ・プログラムであって、

コンピュータに、請求項 1 ~ 6 のいずれか一項に記載の方法の各ステップを実行させる、前記コンピュータ・プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般に、改良データ処理システムに関し、より具体的には、周辺機器相互接続入出力仮想化構成(peripheral component interconnect input/output virtualization configuration)を使用して冗長構成を生成するための、コンピュータ実施方法、データ処理システム、およびコンピュータ・プログラム製品に関する。

【背景技術】

【0002】

典型的なコンピューティング・デバイスは、元々は 1990 年代にインテル・コーポレーションによって作成され、現在は PCI-SIG によって管理される、周辺機器相互接続(PCI: Peripheral Component Interconnect)規格のあるバージョンまたは実施を利用する、入出力(I/O)アダプタおよびバスを使用する。周辺機器相互接続(PCI)規格は、周辺デバイスをコンピュータ・マザーボードに接続するためのコンピュータ・バスを規定する。PCI エクスプレス(PCI Express)すなわち PCIe は、既存の PCI プログラミング概念を使用する PCI コンピュータ・バスの実施であるが、コンピュータ・バスは、まったく異なるはるかに高速なシリアル物理レイヤ通信プロトコルに基づいている。物理レイヤは、複数のデバイス間で共用できる双方向バスから成るのではなく、正確に 2 つのデバイスに接続される単一の一方向リンクから成る。

【0003】

図 1 を参照すると、周辺機器相互接続エクスプレス仕様に従った周辺機器相互接続エクスプレス(PCIe: peripheral component interconnect express)バスを含むシステムを示す例示的な図が提示されている。図 1 に示された具体的なシステムは、複数のサーバ・ブレード(serverblade) 101 ~ 104 が提供されるブレード・エンクロージャ(blade enclosure)である。サーバ・ブレードは、高密度システム用に設計された内蔵型のコンピュータ・サーバである。サーバ・ブレードは、スペース、電力、および他の考慮点のために多くのコンポーネントが除去されているが、それでも、コンピュータと見なされるためのすべての機能コンポーネントを有している。ブレード・エンクロージャ 100 は、電力、冷却、ネットワーキング、様々な相互接続、およびブレード・エンクロージャ 100 内の様々なサーバ・ブレード 101 ~ 104 の管理などのサービスを提供する。サーバ・ブレード 101 ~ 104 およびブレード・エンクロージャ 100 は、一緒になってブレード・システムを形成する。

【0004】

図 1 に示されるように、周辺機器相互接続エクスプレスが、サーバ・ブレード 101 ~

10

20

30

40

50

104の各々に実装され、周辺機器相互接続エクスプレス・デバイス105～112の1つに接続するために使用される。次にこれらのサーバ・ブレード101～104の各々が、ブレード・エンクロージャ100内のスロットに差し込まれ、次にブレード・エンクロージャが、周辺機器相互接続エクスプレス・イーサネット・デバイス105、107、109、111の出力を、ブレード・エンクロージャ100内のバックプレーン(backplane)を介して、イーサネット・スイッチ113に接続し、次にイーサネット・スイッチが、例えば、ブレード・エンクロージャ100の外部のデバイスへの通信接続などの、外部接続用のイーサネット接続115を生成する。同様に、周辺機器相互接続エクスプレス・ストレージ・デバイス106、108、110、112の各々は、ブレード・エンクロージャ100内のバックプレーンを介して、ストレージ・エリア・ネットワーク・スイッチ114に接続され、次にストレージ・エリア・ネットワーク・スイッチが、外部接続用のストレージ・エリア・ネットワーク接続116を生成する。

10

【0005】

したがって、図1に示されるシステムは、周辺機器相互接続仕様または周辺機器相互接続エクスプレス仕様あるいはその両方が実施された、データ処理システムの例示的な1つのタイプである。周辺機器相互接続仕様または周辺機器相互接続エクスプレス仕様あるいはその両方を使用するデータ処理システムの他の構成も知られている。これらのシステムは、アーキテクチャが多様であり、したがって、その各々を本明細書で詳しく説明することはできない。周辺機器相互接続および周辺機器相互接続エクスプレスに関するより多くの情報については、周辺機器相互接続スペシャル・インタレスト・グループ(PCI-SIG: peripheral component interconnect special interest group)のウェブサイトwww.pcisig.comから入手可能な周辺機器相互接続仕様および周辺機器相互接続エクスプレス仕様を参照されたい。

20

【0006】

周辺機器相互接続仕様および周辺機器相互接続エクスプレス仕様に加えて、周辺機器相互接続スペシャル・インタレスト・グループは、いくつかの論理パーティション(LPAA: logical partition)によって共用できる入出力アダプタ(IOA: input/output adapter)をいかに設計するかについて規定するための、入出力仮想化(IOV: input/output virtualization)規格も規定している。論理パーティションは、コンピュータのプロセッサ、メモリ、およびストレージを、リソースの多数の組に分割して、リソースの各組を、独自のオペレーティング・システム・インスタンスおよびアプリケーションとともに独立に動作させ得るようにする。生成できる論理パーティションの数は、システムのプロセッサ・モデルおよび利用可能なリソースに依存する。一般に、パーティションは、データベース動作、クライアント/サーバ動作、またはテスト環境と生産環境の分離など、様々な目的で使用される。各パーティションは、他のパーティションがあたかも別個のマシンであるかのように、他のパーティションと通信することができる。

30

【0007】

論理パーティションをサポートする現代のシステムでは、いくつかのリソースは、論理パーティション間で共用することができる。上で言及されたように、周辺機器相互接続仕様および周辺機器相互接続エクスプレス仕様では、共用できるそのようなリソースの1つは、入出力仮想化メカニズムを使用する入出力アダプタである。

40

【0008】

さらに、周辺機器相互接続スペシャル・インタレスト・グループは、多数のシステム間で入出力アダプタを共用するための入出力仮想化(IOV)規格も規定している。この機能は、マルチルート(MR: multi-root)入出力仮想化と呼ばれる。

【0009】

図2を参照すると、周辺機器相互接続エクスプレス・マルチルート入出力仮想化を含むシステムを示す例示的な図が提示されている。具体的には、図2は、多数のシステム間で周辺機器相互接続エクスプレス・デバイスを共用するために、図1に示されたアーキテクチャをどのように変更できるかを示している。

50

【 0 0 1 0 】

今度はサーバ・ブレード 2 0 1 ~ 2 0 4 は、図 1 のサーバ・ブレード 1 0 1 ~ 1 0 4 で行われたように、サーバ・ブレード 2 0 1 ~ 2 0 4 上に周辺機器相互接続エクスプレス・デバイス自体を含む代わりに、周辺機器相互接続エクスプレス・ルート・ポート (root port) 2 0 5 ~ 2 1 2 を生成し、ブレード・エンクロージャ 2 0 0 のバックプレーンを横断する周辺機器相互接続エクスプレス接続を駆動する。その場合、各サーバ・ブレード 2 0 1 ~ 2 0 4 からの周辺機器相互接続エクスプレス・リンクは、マルチルート周辺機器相互接続エクスプレス・スイッチ 2 1 3 ~ 2 1 4 の一方に接続され、次にマルチルート周辺機器相互接続エクスプレス・スイッチは、周辺機器相互接続エクスプレス・イーサネット / ストレージ・デバイス 2 1 7 ~ 2 2 0 に接続される。周辺機器相互接続エクスプレス・イーサネット / ストレージ・デバイス 2 1 7 ~ 2 2 0 は、外部接続 2 1 5、2 1 6 を介して、外部のイーサネット・デバイスおよびストレージ・デバイスに接続する。そのようにして、周辺機器相互接続エクスプレス・デバイスを、ブレード・エンクロージャ 2 0 0 内で使用することができる。これは、サーバ・ブレード 2 0 1 ~ 2 0 4 間で周辺機器相互接続エクスプレス・デバイス 2 1 7 ~ 2 2 0 が共用されるので、それらの数を最低限に抑えることができる点で、全体的コストを低減する。さらに、これは、周辺機器相互接続エクスプレス・デバイス 2 1 7 ~ 2 2 0 の組み込みを必要としないことによって、サーバ・ブレード 2 0 1 ~ 2 0 4 自体の複雑さおよびコストを低減することができる。

10

【 0 0 1 1 】

周辺機器相互接続スペシャル・インタレスト・グループは、いくつかの論理パーティションによって共用できる入出力アダプタをいかに設計するかについて規定するための規格を提供するが、その仕様は、入出力アダプタをいかにホスト・システムに接続するかについて規定していない。さらに、その規格は、各機能を単一システムにいかに割り当てることができるかについて規定するに過ぎない。

20

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 1 2 】

本発明は、周辺機器相互接続入出力仮想化構成を使用して冗長構成を生成するための、コンピュータ実施方法、データ処理システム、およびコンピュータ・プログラムを提供する。

30

【 課題を解決するための手段 】

【 0 0 1 3 】

本発明の一実施形態によれば、冗長システム構成を生成するためのコンピュータ実施方法が提示される。コンピュータ実施方法は、信用できるエンティティによって 1 組の仮想機能パス許可テーブル (virtual function path authorization table) を生成することであって、エントリが、1 つまたは複数のシステム内の 1 組のアドレス範囲への機能のアクセスを定義し、エントリが、無効な機能間アクセス (crossfunction access) を防止する境界をさらに定義し、仮想機能が、シングルルートまたはマルチルート周辺機器相互接続デバイスによって実行される、生成することと、リクエストから要求を受信して、要求されたデータを仮想機能から提供することと、1 組のアドレス範囲において選択されたアドレス範囲内に受信バッファ (receivebuffer) を生成することと、選択されたアドレス範囲内の受信バッファのアドレスを含む、仮想機能のための仮想機能作業待ち行列エントリ (virtualfunction work queue entry) を生成することと、を含む。さらに、コンピュータ実施方法は、1 組の仮想機能パス許可テーブルにおいて、選択されたアドレス範囲の使用を仮想機能が許可されているかどうかを判定することと、仮想機能は許可されているという判定に回答して、要求されたデータを、1 つまたは複数のシステム内の各アドレス範囲の受信バッファに書き込むことと、要求されたデータの書き込みに応答して、リクエストに完了通知を発行することと、を含む。

40

【 0 0 1 4 】

別の実施形態によれば、冗長システム構成を生成するためのデータ処理システムが提示

50

される。データ処理システムは、バスと、バスに接続されたコンピュータ実行可能命令を含むメモリと、中央プロセッサ・ユニットとを含む。中央プロセッサ・ユニットは、コンピュータ実行可能命令を実行し、コンピュータ実行可能命令は、信用できるエンティティによって1組の仮想機能バス許可テーブルを生成することであって、エントリが、1つまたは複数のシステム内の1組のアドレス範囲への機能のアクセスを定義し、エントリが、無効な機能間アクセスを防止する境界をさらに定義し、仮想機能が、シングルルートまたはマルチルート周辺機器相互接続デバイスによって実行される、生成することと、リクエストから要求を受信して、要求されたデータを仮想機能から提供することと、1組のアドレス範囲において選択されたアドレス範囲内に受信バッファを生成することと、選択されたアドレス範囲内の受信バッファのアドレスを含む、仮想機能のための仮想機能作業待ち行列エントリを生成することと、1組の仮想機能バス許可テーブルにおいて、選択されたアドレス範囲の使用を仮想機能が許可されているかどうかを判定することと、仮想機能は許可されているという判定に回答して、要求されたデータを、1つまたは複数のシステム内の選択されたアドレス範囲の受信バッファに書き込むことと、要求されたデータの書き込みに応答して、リクエストに完了通知を発行することと、を実行するようにデータ処理システムに指令する。

10

【0015】

別の実施形態によれば、冗長システム構成を生成するためのコンピュータ・プログラム製品が提示される。コンピュータ・プログラム製品は、コンピュータ実行可能命令が保存されたコンピュータ可読媒体を含む。コンピュータ実行可能命令は、信用できるエンティティによって1組の仮想機能バス許可テーブルを生成するためのコンピュータ実行可能命令であって、エントリが、1つまたは複数のシステム内の1組のアドレス範囲への仮想機能のアクセスを定義し、エントリが、無効な機能間アクセスを防止する境界をさらに定義し、仮想機能が、シングルルートまたはマルチルート周辺機器相互接続デバイスによって実行される、コンピュータ実行可能命令と、リクエストから要求を受信して、要求されたデータを仮想機能から提供するためのコンピュータ実行可能命令と、1組のアドレス範囲において選択されたアドレス範囲内に受信バッファを生成するためのコンピュータ実行可能命令と、を含む。コンピュータ実行可能命令は、選択されたアドレス範囲内の受信バッファのアドレスを含む、仮想機能のための仮想機能作業待ち行列エントリを生成するためのコンピュータ実行可能命令と、1組の仮想機能バス許可テーブルにおいて、選択されたアドレス範囲の使用を仮想機能が許可されているかどうかを判定するためのコンピュータ実行可能命令と、仮想機能は許可されているという判定に回答して、要求されたデータを、1つまたは複数のシステム内の選択されたアドレス範囲の受信バッファに書き込むためのコンピュータ実行可能命令と、要求されたデータの書き込みに応答して、リクエストに完了通知を発行するためのコンピュータ実行可能命令と、をさらに含む。

20

30

【図面の簡単な説明】

【0016】

【図1】周辺機器相互接続エクスプレス規格を実施するシステム・アーキテクチャのブロック図である。

【図2】周辺機器相互接続マルチルート入出力仮想化を含む図1のシステムのブロック図である。

40

【図3】周辺機器相互接続マルチルート入出力ファブリックを利用する分散コンピューティング・システムのブロック図である。

【図4】本発明の例示的な実施形態が実施できる、多数の論理パーティションを使用するシステム・リソースの仮想化のブロック図である。

【図5】例示的な一実施形態による、周辺機器相互接続エクスプレス・マルチルート入出力仮想化対応エンドポイントのブロック図である。

【図6】周辺機器相互接続エクスプレス・マルチルート対応周辺機器相互接続エクスプレス・スイッチのブロック図である。

【図7】例示的な一実施形態による、仮想機能作業待ち行列エントリのブロック図である

50

。

【図 8】例示的な一実施形態による、マルチルート・デバイス内の任意の与えられた仮想階層にアクセスするための仮想機能の権限を確認するためのテーブルのブロック図である。

。

【図 9】例示的な一実施形態による、マルチルート・デバイスの冗長パス実施のための代替ルート仮想階層を指定するためのテーブルのブロック図である。

【図 10】例示的な一実施形態による、許可アドレスと仮想機能の関係を指定するためのテーブルのブロック図である。

【図 11】例示的な一実施形態による、図 10 のアドレスを使用する仮想機能作業待ち行列エントリのブロック図である。

10

【図 12】例示的な一実施形態による、マルチルート・デバイスおよびマルチルート・スイッチを使用する冗長システムの構成のブロック図である。

【図 13】例示的な一実施形態による、シングルルート・デバイスを使用する冗長論理パーティションの構成のブロック図である。

【図 14】例示的な一実施形態による、マルチルート・マルチシステム構成のマルチルート・ファブリック構成のプロセスのフローチャートである。

【図 15】例示的な一実施形態による、パートナ・システムに伝達するために必要とされる仮想階層番号をシステムが決定するためのプロセスのフローチャートである。

【図 16】例示的な一実施形態による、仮想機能作業待ち行列エントリを設定するためのプロセスのフローチャートである。

20

【図 17】例示的な一実施形態による、入出力ファブリック・パス動作ステータスおよび代替パス使用を動的に決定するためのプロセスのフローチャートである。

【図 18】例示的な一実施形態による、データ・レプリケーションのためのデュアル・パスを使用して書き込み動作を実行するプロセスのフローチャートである。

【発明を実施するための形態】

【0017】

当業者によって理解されるように、本発明は、システム、方法、またはコンピュータ・プログラム製品として実施することができる。したがって、本発明は、全面的にハードウェアの実施形態、(ファームウェア、常駐ソフトウェア、マイクロコードなどを含む)全面的にソフトウェアの実施形態、またはソフトウェア態様とハードウェア態様を組み合わせた実施形態の形態を取ることができ、本明細書では、それらのすべてを一般に、「回路」、「モジュール」、または「システム」と呼ぶことができる。さらに、本発明は、媒体内に具体化されたコンピュータ使用可能プログラム・コードを有する表現である、任意の有形な媒体内に具体化されたコンピュータ・プログラム製品の形態を取ることができる。

30

【0018】

1つまたは複数のコンピュータ使用可能またはコンピュータ可読媒体の任意の組合せを利用することができる。コンピュータ使用可能またはコンピュータ可読媒体は、例えば、電子、磁気、光、電磁気、赤外線、または半導体のシステム、装置、デバイス、または伝播媒体とすることができるが、それらに限定されない。コンピュータ可読媒体のより具体的な例(非網羅的なリスト)は、1つもしくは複数の電線を有する電気接続、ポータブル・コンピュータ・ディスク、ハード・ディスク、ランダム・アクセス・メモリ(RAM)、リード・オンリ・メモリ(ROM)、消去可能プログラマブル・リード・オンリ・メモリ(EPROMもしくはフラッシュ・メモリ)、光ファイバ、ポータブル・コンパクト・ディスク・リード・オンリ・メモリ(CDROM)、光記憶デバイス、インターネットもしくはイントラネットをサポートするような伝送媒体、または磁気記憶デバイスを含む。コンピュータ使用可能またはコンピュータ可読媒体は、例えば、紙または他の媒体の光学的スキャニングによって、プログラムが電子的に獲得でき、次に必要であれば、コンパイル、インタプリテーション、または他の適切な方法で処理でき、次にコンピュータ・メモリ内に保存できるように、プログラムが印刷された紙または別の適切な媒体とすることさえできることに留意されたい。本文書との関連では、コンピュータ使用可能またはコ

40

50

ンピュータ可読媒体は、命令実行システム、装置、もしくはデバイスによって使用するため、またはそれらに関連して使用するために、プログラムを格納、保存、伝達、伝播、または転送できる任意の媒体とすることができる。コンピュータ使用可能媒体は、ベースバンド内または搬送波の一部として、コンピュータ可読プログラム・コードが具体化された伝播データ信号を含むことができる。コンピュータ使用可能プログラム・コードは、無線、有線、光ファイバ・ケーブル、RFなどを含むが、それらに限定されない、任意の適切な媒体を使用して、伝送することができる。

【0019】

本発明の動作を実施するためのコンピュータ・プログラム・コードは、Java、Smalltalk、またはC++などのオブジェクト指向プログラミング言語、および「C」プログラミング言語または類似のプログラミング言語などの従来の手続き型プログラミング言語を含む、1つまたは複数のプログラミング言語の任意の組合せで書くことができる。プログラム・コードは、スタンドアロン・ソフトウェア・パッケージとして、全面的にユーザのコンピュータ上で、部分的にユーザのコンピュータ上で、一部はユーザのコンピュータ上、一部はリモート・コンピュータ上で、または全面的にリモート・コンピュータもしくはサーバ上で実行することができる。後者のシナリオでは、リモート・コンピュータは、ローカル・エリア・ネットワーク(LAN)もしくはワイド・エリア・ネットワーク(WAN)を含む任意のタイプのネットワークを介して、ユーザのコンピュータに接続することができ、または接続は、(例えば、インターネット・サービス・プロバイダを使用してインターネットを介して)外部コンピュータに対して行うことができる。

【0020】

本発明は、本発明の実施形態による方法、装置(システム)、およびコンピュータ・プログラム製品のフローチャート図またはブロック図あるいはその両方を参照しながら以下で説明される。フローチャート図またはブロック図あるいはその両方の各ブロック、およびフローチャート図またはブロック図あるいはその両方におけるブロックの組合せは、コンピュータ・プログラム命令によって実施することができることが理解されよう。

【0021】

これらのコンピュータ・プログラム命令は、汎用コンピュータ、専用コンピュータ、またはマシンを生成するための他のプログラマブル・データ処理装置のプロセッサに提供することができ、コンピュータまたは他のプログラマブル・データ処理装置のプロセッサによって実行される命令は、フローチャートまたはブロック図あるいはその両方の1つまたは複数のブロックで指定される機能/動作を実施するための手段を生成する。特定の方法で機能するようにコンピュータまたは他のプログラマブル・データ処理装置に指令することができる、これらのコンピュータ・プログラム命令はまた、コンピュータ可読媒体内に保存することができ、コンピュータ可読媒体内に保存された命令は、フローチャートまたはブロック図あるいはその両方の1つまたは複数のブロックで指定される機能/動作を実施する命令手段を含む製造品を生成する。

【0022】

コンピュータ・プログラム命令はまた、一連の動作ステップがコンピュータまたは他のプログラマブル装置上で実行されて、コンピュータ実施プロセスを生成するように、コンピュータまたは他のプログラマブル・データ処理装置にロードすることができ、コンピュータまたは他のプログラマブル装置上で実行される命令は、フローチャートまたはブロック図あるいはその両方の1つまたは複数のブロックで指定される機能/動作を実施するためのプロセスを提供する。

【0023】

例示的な実施形態は、入出力仮想化機能から別個のシステムへの多数のパスを可能にする、マルチルート入出力仮想化(MR-IOV)アダプタおよび入出力ファブリック(input/output fabric)の構成のためのメカニズムを提供する。例示的な実施形態は、周辺機器相互接続エクスプレス(PCIe)アダプタまたはエンドポイントに関して説明されるが、本発明は、そのようなものに限定されない。むしろ、例示的な実施形態のメカニズ

ムは、入出力アダプタ内で入出力仮想化をサポートする任意の入出力ファブリックにおいて実施することができる。

【 0 0 2 4 】

さらに、例示的な実施形態は、ハイパーバイザ (hypervisor) が利用される実施に関して説明されるが、本発明は、そのようなものに限定されない。反対に、ソフトウェア、ハードウェア、またはソフトウェアとハードウェアの任意の組合せのどれで実施されようと、現在知られている、または今後開発される、ハイパーバイザ以外の他のタイプの仮想化プラットフォームが、本発明の主旨および範囲から逸脱することなく使用できる。

【 0 0 2 5 】

ここで図を、特に図 3 を参照すると、周辺機器相互接続マルチルート入出力ファブリック (peripheral component interconnect multi-root input output fabric) を利用する分散コンピューティング・システムのブロック図が、本発明の例示的な一実施形態に従って示されている。図 3 は、システム・ノードを共用入出力アダプタと接続するための周辺機器相互接続ファブリックを追加して、図 1 および図 2 の構成を増強する。図 3 に示されるように、分散コンピューティング・システム 3 0 0 は、周辺機器相互接続マルチルート入出力ファブリック 3 4 4 に結合された複数のルート・ノード (root node) 3 6 0 ~ 3 6 3 を含み、次に周辺機器相互接続マルチルート入出力ファブリック 3 4 4 は、マルチルート入出力ファブリック構成マネージャ (multi-root input output fabric configuration manager) 3 6 4 と、周辺機器相互接続入出力アダプタまたはエンドポイント 3 4 5 ~ 3 4 7 に結合される。各ルート・ノード 3 6 0 ~ 3 6 3 は、1 つまたは複数の対応するルート・コンプレックス (root complex) 3 0 8、3 1 8、3 2 8、3 3 8、3 3 9 を含み、それらは、入出力リンク 3 1 0、3 2 0、3 3 0、3 4 2、3 4 3 をそれぞれ介して、周辺機器相互接続マルチルート入出力ファブリック 3 4 4 に接続され、さらにルート・ノード (RN) 3 6 0 ~ 3 6 3 のメモリ・コントローラ 3 0 4、3 1 4、3 2 4、3 3 4 に接続される。入出力ファブリック 3 4 4 は、リンク 3 5 1、3 5 2、3 5 3 を介して、入出力アダプタ 3 4 5、3 4 6、3 4 7 に接続される。入出力アダプタ 3 4 5、3 4 6、3 4 7 は、周辺機器相互接続エクスプレス入出力アダプタ 3 4 5 におけるような非入出力仮想化対応アダプタ、周辺機器相互接続エクスプレス入出力アダプタ 3 4 6 におけるようなシングルルート (SR) 入出力仮想化アダプタ、または周辺機器相互接続エクスプレス入出力アダプタ 3 4 7 におけるようなマルチルート入出力仮想化アダプタとすることができる。

【 0 0 2 6 】

示されるように、ルート・コンプレックス 3 0 8、3 1 8、3 2 8、3 3 8、3 3 9 は、ルート・ノード 3 6 0、3 6 1、3 6 2、3 6 3 の一部である。ルート・ノード 3 6 3 に示されるように、ルート・ノード当たり 2 つ以上のルート・コンプレックスが存在してもよい。ルート・コンプレックスは、中央プロセッサ / メモリを入出力アダプタに接続する入出力階層 (input/output hierarchy) のルートである。ルート・コンプレックスは、ホスト・ブリッジ (host bridge) と、ゼロまたは 1 つ以上のルート・コンプレックス組み込みエンドポイントと、ゼロまたは 1 つ以上のルート・コンプレックス・イベント・コレクタ (root complex event collector) と、1 つまたは複数のルート・ポートを含む。各ルート・ポートは、別個の入出力階層をサポートする。入出力階層は、ルート・コンプレックス、例えば、ルート・コンプレックス 3 0 8 と、ゼロまたは 1 つ以上の相互接続スイッチまたはブリッジ (スイッチもしくは周辺機器相互接続マルチルート入出力ファブリック 3 4 4 などの周辺機器相互接続エクスプレス・ファブリックを含む) あるいはその両方と、周辺機器相互接続エクスプレス入出力アダプタまたはエンドポイント 3 4 5 ~ 3 4 7 などの 1 つまたは複数のエンドポイントから成ることができる。

【 0 0 2 7 】

ルート・コンプレックスに加えて、各ルート・ノードは、1 つまたは複数の中央処理装置 3 0 1、3 0 2、3 1 1、3 1 2、3 2 1、3 2 2、3 3 1、3 3 2 と、メモリ 3 0 3、3 1 3、3 2 3、3 3 3 と、メモリ・コントローラ 3 0 4、3 1 4、3 2 4、3 3 4 か

10

20

30

40

50

ら成る。メモリ・コントローラ304、314、324、334は、中央処理装置301、302、311、312、321、322、331、332を、バス305、306、307、315、316、317、325、326、327、335、336、337によって、メモリ303、313、323、333と接続し、バス309、319、329、340、341によって、入出力ルート・コンプレックス308、318、328、338、339と接続する。メモリ・コントローラは一般に、メモリのためのコヒーレンシ・トラフィック (coherency traffic) の処理などの機能を実行する。ルート・ノード360、361は、1つのコヒーレンシ・ドメイン (coherency domain) を形成するために、それぞれのメモリ・コントローラ304、314を介して、接続359で互いに接続することができる。したがって、ルート・ノード360~361は、単一の対称マルチプロセッシング (SMR: symmetric multi-processing) システムとして動作することができ、またはルート・ノード362、363におけるように、別個のコヒーレンシ・ドメインを有する独立ノードとすることもできる。

10

【0028】

マルチルート入出力ファブリック構成マネージャ364は、ルート・ノードの他の動作から分離されることができ、したがって、入出力ファブリック344に別個に接続されるように示されている。しかし、これはシステムに費用を追加し、したがって、本明細書で開示される実施形態は、1つまたは複数のルート・ノード360、361、362、363の一部として、この機能を含むことができる。構成マネージャ364は、マルチルート入出力ファブリック344のリソースを構成し、リソースをルート・ノード360、361、362、363に割り当てる。

20

【0029】

図3に示されたハードウェアが様々であり得ることを当業者であれば理解されよう。例えば、光ディスク・ドライブなどの他の周辺デバイスを、示されたハードウェアに追加して、またはそれと交換して使用することができる。示された例は、本発明に関するアーキテクチャの限定を暗示することを意図していない。

【0030】

図3の分散コンピューティング・システム300の例を使用して、例示的な実施形態は、入出力仮想化デバイスの単一機能に、多数のシステムにアクセスする能力を提供する。その能力は、冗長バスを有する入出力サブシステムを構成マネージャ364によって構成することを可能にし、入出力仮想化デバイスの単一機能が、多数のシステム間の多数の通信バスを介して多数のシステムにアクセスすることを可能にする。

30

【0031】

例示的な実施形態は、入出力 (I/O) ファブリック344が、ルート・ノード360、361、362、363のシステムもしくは論理パーティション (LPAR) など、2つ以上のシステムによって共用される状況、各システムまたは論理パーティションが潜在的に、周辺機器相互接続エクスプレス入出力アダプタもしくはエンドポイント345~347などの入出力アダプタ (IOA) を、他の論理パーティションと共用できる状況、または多数のシステムが、マルチルート入出力仮想化ファブリックの使用によって、入出力アダプタを共用できる状況に対処する。例示的な実施形態は、ルート・ノードの多数のシステムまたは論理パーティションにアクセスすることを周辺機器相互接続エクスプレス入出力アダプタ347などの入出力仮想化アダプタの単一機能に許可する一方で、またアクセスを許可すべきではないシステムへのアクセスを防止するための、メカニズムを規定する。したがって、エンドポイント345~347とマルチルート・ノードのメモリ303、313、323、333の間に冗長バスを確立する目的で、単一の入出力仮想化機能は、多数の仮想階層 (VH: virtual hierarchy) またはバス、およびマルチルート入出力ファブリック344の仮想化機能へのアクセスを許可される。

40

【0032】

図4を参照すると、本発明の例示的な実施形態が実施できる多数の論理パーティションを使用するシステム・リソースの仮想化のブロック図が提示されている。論理パーティシ

50

ョン化プラットフォーム (logical partitioned platform) 400 内のハードウェアは、例えば、図3のルート・ノード360、361、362、363内で実施することができ、ルート・ノードに割り当てられた、マルチルート入出力ファブリック344および入出力アダプタ345～347の部分をさらに含むことができる。

【0033】

論理パーティション化プラットフォーム400は、パーティション化ハードウェア430と、オペレーティング・システム402、404、406、408と、プラットフォーム・ファームウェア410のパーティション管理を含む。オペレーティング・システム402、404、406、408は、論理パーティション化プラットフォーム400上で同時に動作する、単一のオペレーティング・システムの多数のコピー、また多数の異種のオペレーティング・システムとすることができる。

10

【0034】

オペレーティング・システム402、404、406、408は、パーティション403、405、407、409内に配置される。ハイパーバイザ・ソフトウェアまたはファームウェアは、プラットフォーム・ファームウェア410のパーティション管理を実施するために使用できるソフトウェアの一例である。ファームウェアは、電力がなくてもその内容を保持するメモリ・チップ内に、例えば、リード・オンリ・メモリ (ROM)、プログラマブル・リード・オンリ・メモリ (PROM)、消去可能プログラマブル・リード・オンリ・メモリ (EPROM)、電氣的消去可能プログラマブル・リード・オンリ・メモリ (EEPROM)、および不揮発性ランダム・アクセス・メモリ (NVRAM) 内に保存された「ソフトウェア」である。

20

【0035】

加えて、パーティション403、405、407、409は、パーティション・ファームウェア411、413、415、417も含む。パーティション・ファームウェア411、413、415、417は、初期ブート・ストラップ・コード、例えば、電気電子学会 (IEEE) 1275規格のオープン・ファームウェア (Open Firmware)、およびランタイム・アブストラクション・ソフトウェア (RTAS: runtime abstraction software) を使用して、実施することができる。パーティション403、405、407、409がインスタンス化される場合、ブート・ストラップ・コードのコピーが、プラットフォーム・ファームウェア410によって、パーティション403、405、407、409にロードされる。その後、制御は、ブート・ストラップ・コードに移され、ブート・ストラップ・コードを用いて、オープン・ファームウェアおよびランタイム・アブストラクション・ソフトウェアをロードする。その後、パーティション・ファームウェア411、413、415、417を実行するために、パーティション403、405、407、409に関連する、または割り当てられたプロセッサが、パーティションのメモリにディスパッチされる。

30

【0036】

パーティション化ハードウェア430は、複数のプロセッサ432、434、436、438と、複数のシステム・メモリ・ユニット440、442、444、446と、複数の入出力アダプタ448、450、452、454、456、458、460、462と、ストレージ・ユニット470と、不揮発性ランダム・アクセス・メモリ・ストレージ498を含む。プロセッサ432、434、436、438、メモリ・ユニット440、442、444、446、不揮発性ランダム・アクセス・メモリ・ストレージ498、および入出力アダプタ448、450、452、454、456、458、460、462、またはそれらの部分の各々は、論理パーティション化プラットフォーム400内の多数のパーティションの1つに割り当てることができ、その各々は、オペレーティング・システム402、404、406、408の1つに対応する。

40

【0037】

プラットフォーム・ファームウェア410は、論理パーティション化プラットフォーム400のパーティショニングを生成および実施するために、パーティション403、40

50

5、407、409のための多数の機能およびサービスを実行する。プラットフォーム・ファームウェア410は、基礎をなすハードウェアと同じファームウェア実施仮想マシンを含むことができる、パーティション管理ファームウェアを含むことができる。したがって、プラットフォーム・ファームウェア410内のパーティション管理ファームウェアは、論理パーティション化プラットフォーム400のハードウェア・リソースを仮想化することによって、独立のオペレーティング・システム・イメージ402、404、406、408の同時実行を可能にする。

【0038】

サービス・プロセッサ490は、パーティション403、405、407、409におけるプラットフォーム・エラーの処理など、様々なサービスを提供するために使用することができる。これらのサービスは、エラーをベンダに報告するためのサービス・エージェントとしても動作することができる。パーティション403、405、407、409の動作は、ハードウェア管理コンソール480などのハードウェア管理コンソールを介して制御することができる。ハードウェア管理コンソール480は、異なるパーティションへのリソースの再アロケーションを含む様々な機能をシステム管理者がそこから実行できる、別個の分散コンピューティング・システムである。制御され得る動作は、パーティションが動作中であってなくても、パーティションに割り当てられるコンポーネントに関連したパーティションの構成などを含む。

【0039】

論理パーティション(LPAR)環境では、1つのパーティション内のリソースまたはプログラムが別のパーティションにおける動作に影響を及ぼすことは許可されない。さらに、有用なこととして、リソースの割り当ては、細かい粒度(fine-grained)とする必要がある。例えば、特定の周辺機器相互接続ホスト・ブリッジ(PHB: peripheral component interconnect host bridge)下のすべての入出力アダプタを同じパーティションに割り当てることは、パーティション間でリソースを動的に移動させる能力を含むシステムの構成可能性を制限するので、しばしば許可されない。

【0040】

したがって、個々の入出力アダプタまたは入出力アダプタの部分を別個のパーティションに割り当てることができ、同時に、他のパーティションのリソースにアクセスすることなどによって、割り当てられたリソースが他のパーティションに影響を及ぼすことを防止することができるように、いくつかの機能が、入出力アダプタを入出力バスに接続するブリッジ内で必要とされる。

【0041】

図5を参照すると、周辺機器相互接続エクスプレス・マルチルート入出力仮想化対応エンドポイントのブロック図が提示されている。図5に示されるように、図3のマルチルート周辺機器相互接続エクスプレス入出力アダプタ347などの周辺機器相互接続エクスプレス・マルチルート入出力仮想化エンドポイント500は、周辺機器相互接続エクスプレス・ファブリックの周辺機器相互接続エクスプレス・スイッチなどとの通信がそれを介して実行され得る周辺機器相互接続エクスプレス・ポート501を含む。内部ルーティング502は、構成管理機能503および構成管理機能509、ならびに複数の仮想機能(VF: virtual function)504~506への通信経路を提供する。構成管理機能503は、仮想機能504~506とは反対に、物理機能(PF: physical function)とすることができ、構成管理機能509は、基本機能(BF: base function)とすることができる。周辺機器相互接続仕様において使用される用語としての物理「機能」は、単一の構成空間によって表される1組のロジックである。言い換えると、物理「機能」は、例えば、分離不能リソース507において提供され得るような、メモリ内の機能関連構成空間に保存されたデータに基づいて構成可能な回路ロジックである。基本「機能」509についても同様のことが言える。

【0042】

構成管理機能503は、仮想機能504~506を構成するために使用することができ

10

20

30

40

50

る。仮想機能は、周辺機器相互接続エクスプレス・マルチルート入出力仮想化エンドポイント 500 の共用可能リソース・プール 508 において提供され得る、例えば、リンクなどの 1 つまたは複数の物理エンドポイント・リソースを、例えば、別の機能と共用する、入出力仮想化対応エンドポイント内の機能である。仮想機能は、ハイパーバイザによるランタイム介入なしに、直接的に、システム・イメージからの入出力およびメモリ動作のためのシンクとすることができ、システム・イメージへのダイレクト・メモリ・アクセス (DMA : direct memory access)、完了、および割り込み動作のソースとすることができる。

【0043】

マルチルート入出力仮想化エンドポイント 500 は、例えば、図 3 のルート・ノード 360 ~ 363 など、多数のルート・ノード間で共用することもできる。基本機能としての構成管理機能 509 は、例えば、どのルート・ノードが各物理機能にアクセスするかなど、物理機能の特性を構成するために使用することができる。

【0044】

周辺機器相互接続エクスプレス・エンドポイントは、周辺機器相互接続エクスプレス・エンドポイントによってサポートされる「機能」に関して、多くの異なるタイプの構成を有することができる。例えば、エンドポイントは、単一の物理機能、多数の独立物理機能をサポートすることができ、または多数の従属物理機能さえもサポートすることができる。ネイティブ入出力仮想化をサポートするエンドポイントでは、エンドポイントによってサポートされる各物理機能は、1 つまたは複数の仮想機能に関連付けることができ、仮想機能自体は、他の物理機能に関連付けられた仮想機能に従属することができる。ルート・ノードに割り当てられる入出力仮想化エンドポイントの単位は、物理機能であり、マルチルート入出力仮想化対応エンドポイントは、多数の物理機能を含む。

【0045】

一実施形態では、仮想機能 (VF) - 仮想階層 (VH) 許可テーブル (virtual function (VF) to virtual hierarchy (VH) authorization table) 510 は、図 3 の構成マネージャ 364 が、多数の仮想階層へのアクセスを各機能に与えることを可能にする。この態様は後で説明される。やはりさらに説明される仮想機能作業待ち行列 511 が、仮想機能のために、デバイス・ドライバ・ソフトウェアによって設定され、仮想機能によって実行される動作を指定する。テーブル内の仮想機能作業待ち行列エントリは、要求された特定の動作のために使用する 1 つまたは複数の仮想階層番号も含む。

【0046】

図 6 を参照すると、周辺機器相互接続エクスプレス・マルチルート対応周辺機器相互接続エクスプレス・スイッチのブロック図が提示されている。周辺機器相互接続エクスプレス・スイッチ 520 は、例えば、図 3 の周辺機器相互接続マルチルート入出力ファブリック 344 において、周辺機器相互接続マルチルート入出力仮想化仕様によって規定されるように、使用することができる。スイッチ 520 は論理的に、ルート・ノードに接続されるポート当たり 1 つの、多数の仮想プレーン (virtual plane) から成る。例えば、ルート・ノード 521 は、周辺機器相互接続エクスプレス・リンク 524 によって、論理周辺機器相互接続 - 周辺機器相互接続 (P2P : peripheral component interconnect to peripheral component interconnect) ブリッジ 527 に接続され、それは論理的に、周辺機器相互接続 - 周辺機器相互接続ブリッジ 536 ~ 538 に対するスイッチに内部的に接続される。同様に、ルート・ノード 522 は、周辺機器相互接続エクスプレス・リンク 525 によって、論理周辺機器相互接続 - 周辺機器相互接続ブリッジ 528 に接続され、それは論理的に、周辺機器相互接続 - 周辺機器相互接続ブリッジ 530 ~ 532 に対するスイッチに内部的に接続され、またルート・ノード 523 は、周辺機器相互接続エクスプレス・リンク 526 によって、論理周辺機器相互接続 - 周辺機器相互接続ブリッジ 529 に接続され、それは論理的に、周辺機器相互接続 - 周辺機器相互接続ブリッジ 533 ~ 535 に対するスイッチに内部的に接続される。

【0047】

10

20

30

40

50

その場合、周辺機器相互接続 - 周辺機器相互接続ブリッジ 530、533、536は、マルチルート周辺機器相互接続エクスプレス・デバイス 542のリソースを共用できるように、周辺機器相互接続エクスプレス・マルチルート・リンク 539を共用する。同様に、その場合、周辺機器相互接続 - 周辺機器相互接続ブリッジ 531、534、537は、マルチルート周辺機器相互接続エクスプレス・デバイス 543のリソースを共用できるように、周辺機器相互接続エクスプレス・マルチルート・リンク 540を共用し、その場合、周辺機器相互接続 - 周辺機器相互接続ブリッジ 532、535、538は、マルチルート周辺機器相互接続エクスプレス・デバイス 544のリソースを共用できるように、周辺機器相互接続エクスプレス・マルチルート・リンク 541を共用する。

【0048】

スイッチ 520を設定するための制御点が、基本機能 (BF) 545である。この出力仮想化構成メカニズム、例えば、基本機能 545は、マルチルート周辺機器相互接続マネージャ (MR-PCIM: multi-root peripheral component interconnect manager) ・プログラムが、スイッチ 520内の論理構造を決定することを可能にする。例えば、図 6は、各ルート・ノード 521~523が、各周辺機器相互接続エクスプレス・マルチルート・デバイス 542~544の一部にアクセスする、かなり対称な構成を示している。通常のシステムでは、システムを使用するユーザの必要を満たすために、システム管理者は、あまり対称的ではない方法で、入出力を設定することを望むことがある。

【0049】

基本機能 545、509は、マルチルート周辺機器相互接続マネージャ・プログラムによってアクセスされる。このプログラムが存在する場所は、周辺機器相互接続スペシャル・インタレスト・グループ入出力仮想化仕様によって規定されていない。プログラムは、例えば、ルート・ノード 521~523の1つによって示されるような、マルチルート周辺機器相互接続マネージャにもっぱら専用され、ルート・ポート・ノードの1つに接続されるノード内に存在することができ、または別個のプロセッサ、例えば、図 4の 490のようなサービス・プロセッサが接続される、ベンダ独自のポートを介して提供することができる。マルチルート周辺機器相互接続マネージャがどこで実行されるかに関わらず、主要な要件は、このプログラムが堅牢であり、システム内の他のアプリケーションの動作または動作失敗によって影響され得ないことである。

【0050】

例示的な実施形態は、図 5に示された周辺機器相互接続エクスプレス・マルチルート入出力仮想化エンドポイント 500など、2つ以上のシステムにアクセスするための入出力仮想化アダプタの構成のためのメカニズムを提供する。例示的な実施形態のメカニズムは、図 6の周辺機器相互接続エクスプレス・スイッチ 520など、1つまたは複数の周辺機器相互接続エクスプレス・スイッチを含むことができる入出力ファブリックが、2つ以上のシステム、例えば、図 3のルート・ノード 362および 363によって共用される状況に対処する。

【0051】

図 7を参照すると、例示的な一実施形態による、仮想機能 (VF) 作業待ち行列エントリのブロック図が提示されている。提供された例は、図 5における仮想機能作業待ち行列エントリ 511を表している。仮想機能作業待ち行列エントリ 601のフィールド 605、607は、周辺機器相互接続エクスプレス・ファブリック仮想階層番号を含む。フィールド 605、607は、動作のためのダイレクト・メモリ・アクセス・データを、どのシステムに送信すべきか、またはどのシステムから受信すべきかを、仮想機能に指示する。このフィールドは、デバイス・ドライバ・ソフトウェアが、同じデータを多数のシステムに送信することを可能にする。例えば、ネットワーク・アダプタ上で仮想機能ネットワーク接続から受信したパケットは、多数のシステムのシステム・メモリに書き込むことができる。例えば、データがディスクなどの不揮発性ストレージ媒体にコミットされる前にシステムが故障した場合にデータのコピーを維持するために、図 3のシステム 362、363内のシステム・メモリ 323、333に書き込むことができる。冗長システムではデ

10

20

30

40

50

ータは安全であるので、冗長性を有することによって、システムは、例えば、データが実際にディスクに書き込まれる前に、データがディスクにコミットされたことを通知することができる。データが物理的にディスクに書き込まれる前に書き込み動作の完了を通知することは、より速い応答、より短い待ち時間を、ディスク書き込み動作のリクエストに提供する。

【 0 0 5 2 】

仮想機能作業待ち行列エントリ 6 0 1 の他のフィールドは、動作タイプ 6 0 2 と、転送長 6 0 3 と、動作アドレス 6 0 4、6 0 6 を含む。動作タイプ 6 0 2 は、どの動作を実行すべきかを仮想機能に指示する。例えば、ネットワーク・アダプタの場合、動作は、受信バッファの設定とすることができる。この場合、受信バッファは、各システムにつき 1 つの対の、動作アドレスと周辺機器相互接続エクスプレス・ファブリック仮想階層番号のフィールド対を 2 つ以上使用して、2 つ以上のシステムにおいて設定することができる。例えば、6 0 4 と 6 0 5、6 0 6 と 6 0 7 など、これらのフィールドの対は、受信データを送信すべき各システムにつき 1 つ存在する。この場合、転送長 6 0 3 は、バッファ長となるように設定される。

【 0 0 5 3 】

動作タイプおよびフィールド・タイプは、アダプタによって提供される機能によって様々になり得ることを当業者であれば理解されよう。正しいシステムにデータを送るために、周辺機器相互接続エクスプレス・ファブリック仮想階層番号が、各アドレスに提供される。

【 0 0 5 4 】

図 8 を参照すると、例示的な一実施形態による、マルチルート・デバイス内の任意の与えられた仮想階層にアクセスするための仮想機能の権限を確認するためのテーブルのブロック図が提示されている。仮想機能 - 仮想階層許可テーブル 6 1 0 は、図 5 の仮想機能 - 仮想階層許可テーブル 5 1 0 の一例である。マルチルート・デバイスでは、アダプタは、異なるシステムによってアクセスされ得る機能間にファイアウォールの等価物を提供する。図 7 の周辺機器相互接続エクスプレス・ファブリック仮想階層番号フィールド 6 0 5、6 0 7 は、異なるシステムによって制御される異なる機能に通常は割り当てられる仮想階層番号を使用して、ファイアウォールをトンネリングするためのメカニズムを提供する。図 7 の周辺機器相互接続エクスプレス・ファブリック仮想階層番号フィールド 6 0 5、6 0 7 は、1 つのシステム内のデバイス・ドライバ・ソフトウェアによって設定されるので、アダプタ上の許可されたファイアウォールだけをトンネリングするように、システムが関連する仮想機能を設定できるように、使用される仮想階層番号が確認されることは重要である。必要な機能は、仮想機能 - 仮想階層許可テーブル 6 1 0 を介して提供される。アダプタ内の各仮想機能について仮想機能 - 仮想階層番号許可テーブル (virtual function to virtual hierarchy number authorization table) 6 1 1、6 1 5 が存在する。この例では、テーブルは、多数のエントリ 6 1 2 ~ 6 1 4、6 1 6 ~ 6 1 8 を含むことができ、テーブルが適用される仮想機能がアクセスを許可された各仮想階層につき 1 つのエントリを含むことができる。仮想機能が仮想機能作業待ち行列エントリ 6 0 1 を処理できるようにする前に、仮想機能が仮想階層番号にアクセスする許可を得ていることを確認するために、周辺機器相互接続エクスプレス・ファブリック仮想階層番号フィールド 6 0 5、6 0 7 は、適切な仮想機能 - 仮想階層許可テーブルと照合される。許可されていない場合、仮想機能作業待ち行列エントリ 6 0 1 の処理は許可されず、エラーが、デバイス・ドライバ・ソフトウェアに通知される。仮想機能 - 仮想階層許可テーブル 6 1 0 は、信用できるソフトウェアによって設定される。例えば、信用できるソフトウェアは、図 3 のマルチルート入出力ファブリック構成マネージャまたはマルチルート周辺機器相互接続マネージャ 3 6 4 とすることができる。テーブルは、システム内のデバイス・ドライバ・ソフトウェアによって変更することができず、したがって、トンネリング・プロセスの制御を安全にする。これらのテーブルの使用のさらなる説明は後で説明される。

【 0 0 5 5 】

図 9 を参照すると、例示的な一実施形態による、マルチルート・デバイスの冗長パス実施のための代替ルート仮想階層を指定するためのテーブルのブロック図が提示されている。仮想機能に許可された仮想階層番号のテーブル (authorized virtual hierarchy number for virtual function table) 6 2 0 は、図 5 の内部ルーティング 5 0 2 内に維持される 1 次および 2 次パス・エントリの対応である。仮想機能 - 仮想階層許可テーブル 6 1 0 において仮想階層番号によって指定されたパスの 1 つが利用不可能になった場合、冗長かつ堅牢な構成は、動作のために所望されるシステムへの代替パスを使用するための能力を提供する。この場合、仮想機能 - 仮想階層許可テーブル 6 1 0 の代わりに、仮想機能に許可された仮想階層番号の拡張テーブル 6 2 0 が使用できる。仮想機能に許可された仮想階層番号のテーブル 6 2 0 における相違は、各エントリ 6 2 1、6 2 3、6 2 5 に対して、動作不能の仮想階層番号の代わりに使用する代替仮想階層番号を指定する代替エントリ 6 2 2、6 2 4、6 2 6 が存在することである。例えば、エントリ 6 2 1 が仮想階層番号「1」を指定し、エントリ 6 2 2 が仮想階層番号「3」を指定した場合、仮想機能作業待ち行列エントリ 6 0 1 が仮想階層番号「1」を指定し、仮想階層番号「1」が動作不能として検出されたときは、仮想階層番号「1」を用いて利用可能だった同じシステム内の同じシステム・メモリにアクセスするために、仮想階層番号「3」が使用できる。したがって、例えば、同じ動作のための多数のシステムに同じデータを送信することによって、全システム故障を回避する方法が存在するばかりでなく、冗長性を通して入出力ファブリック故障を回避する方法も存在する。仮想階層番号は、パス識別子として使用することができる。

【 0 0 5 6 】

図 1 0 を参照すると、例示的な一実施形態による、許可アドレスと仮想機能の関係を指定するためのテーブルのブロック図が示されている。仮想機能 - アドレス範囲許可テーブル (virtual function to address range authorization table) 6 2 8 は、許可を必要とする各仮想機能のためのテーブルを含む。各機能について、関連する仮想機能がアクセスを許可されたアドレスの範囲を表すテーブル 6 3 0、6 4 0 内のエントリを用いて、1 組の許可アドレスが提供される。この例では、第 1 の仮想機能 6 3 0 のためのテーブルは、1 組の関連するアドレスを有する。第 1 の仮想機能が使用を許可されたアドレス 6 3 0 は、許可アドレス 6 3 2 ~ 6 3 8 として列挙される。同様に、最後の仮想機能「VFn」も、テーブル・エントリ 6 4 0 によって示される 1 組のエントリを有する。仮想機能 - アドレス範囲許可テーブル 6 2 8 の機能は、例えば、同じルート・ノードの異なる論理パーティション内のアドレス範囲などのリソースへの仮想機能によるアクセスを許可する点において、図 8 の仮想機能 - 仮想階層許可テーブル 6 1 0 の機能と同様である。仮想機能 - アドレス範囲許可テーブル 6 2 8 は、仮想階層を使用する代わりにアドレスを用いて、仮想機能 - 仮想階層許可テーブル 6 1 0 と同様の能力を提供する。

【 0 0 5 7 】

図 1 1 を参照すると、例示的な一実施形態による、図 1 0 のアドレスを使用する仮想機能 (VF) 作業待ち行列エントリのブロック図が提示されている。提供された例は、図 5 の仮想機能作業待ち行列 5 1 1 のさらなる例としての、図 7 の仮想機能作業待ち行列エントリ 6 0 1 を表している。この例では、仮想機能作業待ち行列エントリ 6 4 2 は、以前のように、動作タイプ 6 4 4 および転送長 6 4 6 を含む多数のフィールドを含む。図 7 の先の仮想機能作業待ち行列エントリとの相違は、仮想階層番号が存在しないことである。仮想階層番号の代わりに、動作アドレス 6 4 8 から動作アドレス 6 5 0 が見出される。動作アドレスは、データに関連付けられたロケーション、例えば、同じルート・ノードの異なる論理パーティション内のアドレスを指定する。

【 0 0 5 8 】

図 1 2 を参照すると、例示的な一実施形態による、マルチルート・デバイスおよびマルチルート・スイッチを使用する冗長システムの構成のブロック図が提示されている。図 3 の分散コンピューティング・システム 3 0 0 の例を使用して、マルチルート・デバイスと、コンピュータ電子コンプレックス (CEC: computer electronic complex) を使用し

てマルチルート・デバイス 1 727 およびマルチルート・デバイス 2 728 などのコンピュータ電子コンプレックス通信デバイスに接続されるマルチルート・スイッチとを使用するシステム 700 の構成が定義される。

【0059】

コンピュータ電子コンプレックス 1 701 およびコンピュータ電子コンプレックス 2 702 を含む 2 つのコンピュータ・システムが示されているが、3 ウェイ以上の冗長システムも構成できることは当業者であれば理解されよう。コンピュータ電子コンプレックスは、図 3 のルート・ノードに対応し、周辺機器相互接続ホスト・ブリッジ (PHB) は、図 3 のルート・コンプレックスに対応する。

【0060】

2 つのコンピュータ電子コンプレックスは、論理パーティションの組を形成するために、図 4 におけるように、パーティショニングすることもできる。2 つのコンピュータ電子コンプレックスは、システム・メモリ 703、704 と、各々 3 つの周辺機器相互接続ホスト・ブリッジ 705 ~ 707、708 ~ 710 から成る。マルチルート周辺機器相互接続マネージャ 711 は、図 3 の構成マネージャ 364 に対応する。これは冗長性の高いシステムであるので、1 次マルチルート周辺機器相互接続マネージャ 711 の故障、コンピュータ電子コンプレックス 1 の故障、周辺機器相互接続ホスト・ブリッジ 1 (PHB 1) 705 の故障、またはマルチルート周辺機器相互接続マネージャ 711 がマルチルート入出力ファブリック動作を制御することを妨げる他の任意の故障の場合には、1 次マルチルート周辺機器相互接続マネージャ 711 を引き継ぐことができる、バックアップ・マルチ

【0061】

マルチルート周辺機器相互接続マネージャ 711、712 は、周辺機器相互接続ホスト・ブリッジ 1 (PHB 1) 705 およびマルチルート・スイッチ 1 719 への周辺機器相互接続エクスプレス・リンク 713 を介して、また周辺機器相互接続ホスト・ブリッジ 6 (PHB 6) 710 およびマルチルート・スイッチ 2 720 への周辺機器相互接続エクスプレス・リンク 716 を介して、周辺機器相互接続エクスプレス・マルチルート入出力仮想化仕様によって管理仮想階層として定義されるマルチルート・ファブリックの仮想階層 0 に接続される。他の周辺機器相互接続ホスト・ブリッジは、マルチルート・ファブリックへの 1 次仮想階層接続および 2 次仮想階層接続を形成する。具体的には、コンピュータ電子コンプレックス 1 の 1 次仮想階層は、仮想階層 1 であり、コンピュータ電子コンプレックス 1 は、周辺機器相互接続ホスト・ブリッジ 2 (PHB 2) 706、さらにマルチルート・スイッチ 1 719 への周辺機器相互接続エクスプレス・リンク 714 を介して、仮想階層 1 に接続する。コンピュータ電子コンプレックス 1 の 2 次仮想階層接続は、周辺機器相互接続ホスト・ブリッジ 3 (PHB 3) 707、さらにマルチルート・スイッチ 2 720 への周辺機器相互接続エクスプレス・リンク 718 を介して仮想階層 3 に接続する。同様に、コンピュータ電子コンプレックス 2 の 1 次仮想階層は、周辺機器相互接続ホスト・ブリッジ 5 (PHB 5) 709、さらにマルチルート・スイッチ 2 720 への周辺機器相互接続エクスプレス・リンク 715 を介して仮想階層 4 に接続する。コンピ

【0062】

「2 次」リンクは、必ずしもバックアップ目的だけとは限らず、デバイスがその下に配置されるスイッチに応じて、デバイスとの通信用にも使用される。一般に、デバイスからコンピュータ電子コンプレックスへの最短パスが使用され、それは、動作待ち時間を短縮するために最小数のスイッチを通過するパスである。その場合、多数のスイッチを通過するパスは、バックアップ目的で残して置かれる。周辺機器相互接続エクスプレス・リンク 721、722 は、代替パスを提供するためにスイッチ間接続を提供する。

【 0 0 6 3 】

各マルチルート・スイッチの下には、2つのマルチルート・デバイスが示されている。マルチルート・デバイス1 727は、接続738によってネットワークに接続し、周辺機器相互接続エクスプレス・リンク723を介してマルチルート・スイッチ1 719に接続するネットワーク・デバイスとして示されている。同様に、マルチルート・デバイス2 728は、接続739によってネットワークに接続し、周辺機器相互接続エクスプレス・リンク726を介してマルチルート・スイッチ2に接続するネットワーク・デバイスとして示されている。

【 0 0 6 4 】

冗長性の理由で、2つのネットワーク・アダプタ接続738、739は、外部ネットワーク・スイッチに接続する可能性が最も高く、両方のデバイスは、同じネットワークにアクセスする。そのようにして、一方のネットワーク・アダプタが故障した場合でも、両方の中央電子コンプレックスは依然として、残りのアダプタを介してネットワークにアクセスする。ネットワーク・アダプタ接続727、728に加えて、2つのディスク・アダプタが、マルチルート・デバイス3 729およびマルチルート・デバイス4 730として存在し、これらのデバイスの両方は、同じ1組のディスク・ドライブ731にアクセスする。マルチルート・デバイス3 729は、周辺機器相互接続エクスプレス・リンク724を介してマルチルート・スイッチ1 719に接続し、マルチルート・デバイス4 730は、周辺機器相互接続エクスプレス・リンク725を介してマルチルート・スイッチ2 720に接続する。

【 0 0 6 5 】

この例では、マルチルート・デバイス1 727は、4つの仮想階層、すなわち、仮想階層1 732、仮想階層2 733、仮想階層3 734、および仮想階層4 735にアクセスする。これらの仮想階層の各々は、通常は別個の周辺機器相互接続エクスプレス機能に関連付けられる。例えば、1つの仮想機能が別の仮想機能の仮想階層にアクセスできないように機能の各々がファイヤウォール737によって分離されることになる仮想機能。ファイヤウォール・トンネル736を仮想階層1 732と仮想階層2 733の間（例えば、マルチルート・デバイス1 727の仮想機能1と仮想機能2の間）に生成して、マルチルート・デバイス1 727が、仮想階層の異なる組に接続された両方のコンピュータ電子コンプレックス内のメモリ703およびメモリ704にダイレクト・メモリ・アクセスによってデータを送ることを可能にすることができる。

【 0 0 6 6 】

マルチルート・デバイス1 727は、図5に示された周辺機器相互接続エクスプレス・マルチルート入出力仮想階層エンドポイント500と論理的に同様である。そのようなものとして、マルチルート・デバイス1 727は、図5および図8の仮想機能 - 仮想階層許可テーブル510、610と、図7の仮想機能作業待ち行列エントリ601を有する図5の仮想機能作業待ち行列511を含む。マルチルート周辺機器相互接続マネージャ711におけるような信用できるソフトウェアが、仮想機能が仮想階層1 732および仮想階層2 733の両方にアクセスすることを、基本的にファイヤウォール・トンネル736を抜けるトンネルを形成して可能にするように、仮想機能 - 仮想階層許可テーブルを設定する。

【 0 0 6 7 】

ファイヤウォールを抜けるトンネルの他の実施形態も使用することができる。例えば、1つの仮想機能が、何らかの手段によって別の仮想機能への通信パスを生成し、情報を、データ上で実行する動作とともに、他の仮想機能へ渡す能力が提供できる。他の手段は、ファイヤウォールを抜けるトンネルが信用できるコードによって制御され得るように、説明されたメカニズムのような手段を設定する安全な方法も必要とする。

【 0 0 6 8 】

以下では、ディスクに書き込まれる予定のデータをネットワーク・リンク738から受信する動作について説明する。仮想機能の処理を担当するコンピュータ電子コンプレック

10

20

30

40

50

ス1 701内のデバイス・ドライバが、システム・メモリ703内に受信バッファを生成する。加えて、コンピュータ電子コンプレックス1 701は、例えば、2つのコンピュータ電子コンプレックス間のネットワーク接続を使用することによって、コンピュータ電子コンプレックス2 702内の対応するドライバと通信する。対応するコンピュータ電子コンプレックス2 702のドライバは、システム・メモリ704内に対応する受信バッファをアロケートし、その後、受信バッファのアドレスをコンピュータ電子コンプレックス1 701内のドライバに伝達する。その後、コンピュータ電子コンプレックス1 701内のドライバは、仮想階層1 732を介してコンピュータ電子コンプレックス1 701の受信バッファを、仮想階層2 733を介してコンピュータ電子コンプレックス2 702の受信バッファを指し示す仮想機能作業待ち行列エントリを、マルチルート・デバイス1 727の仮想機能内に設定する。仮想機能によって通信パケットを受信すると、仮想機能は、バッファがどこに配置されているかを識別するために仮想機能作業待ち行列エントリ内に与えられた情報を使用し、その後、仮想機能に対応する図8の仮想機能 - 仮想階層許可テーブル610を使用して、ファイヤウォールをトンネリングするための仮想機能の権限を検証する。権限が合格すると、マルチルート・デバイス1 727は、次に仮想階層1 732を介してコンピュータ電子コンプレックス1 701のシステム・メモリにダイレクト・メモリ・アクセスによってデータを送り、次に仮想階層2 733を介してコンピュータ電子コンプレックス2 702のシステム・メモリ704にダイレクト・メモリ・アクセスによってデータを送る。これらのダイレクト・メモリ・アクセスの正常完了時に、デバイス・ドライバは、マルチルート・デバイス1 727から割り込みによって通知され、両方のコンピュータ電子コンプレックスに対して正常に完了した動作を検出する。この時点で、データは、コンピュータ電子コンプレックスの一方における故障から安全であり、動作は、正常に完了したと見なすことができる。この時点で、例えば、システム・クラッシュからデータは安全であるので、データがまだディスク上になくても、ディスク書き込み動作のオリジネータは、動作が完了したと通知される。その後、コンピュータ電子コンプレックス1 701内のデバイス・ドライバは、マルチルート・デバイス3 729を介するディスク書き込み動作を待ち行列に入れる。

【0069】

図12をさらに参照すると、2つのマルチルート・スイッチから成るマルチルート・ファブリックを通る多数のパスが提示されている。例えば、周辺機器相互接続エクスプレス・リンク717の故障が存在する場合、マルチルート・デバイスは、システム・メモリ704に対する書き込みを、上で説明されたように実行することができない。マルチルート・デバイスが、図9に示された冗長テーブルを実装している場合、そのリンクを介するマルチルート・デバイス1 727からコンピュータ電子コンプレックス2 702へのパスが動作可能でないとき、仮想階層2 733の代わりに仮想階層4 735による代替パスが存在することを決定するために、図9に示されたテーブルを使用することができ、データは、リンク723、次にマルチルート・スイッチ1 719、次に周辺機器相互接続エクスプレス・リンク721、722、次にマルチルート・スイッチ2 720、次に周辺機器相互接続エクスプレス・リンク715、次に周辺機器相互接続ホスト・ブリッジ5 (PHB5) 709を介して、システム・メモリ704まで流れる。

【0070】

図13を参照すると、例示的な一実施形態による、シングルルート・デバイスのみを使用する冗長論理パーティションの構成のブロック図が提示されている。図4の論理パーティション化プラットフォーム400のさらなる一例として、多数の仮想階層という概念が存在しない構成800が提示されている。別個の仮想階層を有する代わりに、仮想機能にダイレクト・メモリ・アクセス・アドレス範囲を割り当てるという概念が存在する。単一のシステム801は、多数の論理パーティション802~803から成り、その各々は、1つまたは複数の中央処理装置804~807を有し、各中央処理装置は、メモリ808~809を有する。論理パーティションは、1つまたは複数の周辺機器相互接続ホスト・ブリッジ (PHB) 810~811を共用し、シングルルート・デバイス814~815

10

20

30

40

50

は、周辺機器相互接続エクスプレス・リンク 8 1 2 ~ 8 1 3 を介して、周辺機器相互接続ホスト・ブリッジに接続される。シングルルート・デバイスは、ネットワーク・アダプタとして示されており、デバイスは、リンク 8 1 6 ~ 8 1 7 を介してネットワークに接続される。図 1 2 におけるように、他方のリンクが故障した場合でも、両方のリンクが同じネットワークにアクセスできるように、2 つのネットワーク・リンクが、ネットワーク・スイッチ（図示されず）に接続される。また図 1 2 におけるように、仮想機能 8 1 8 ~ 8 2 1 は、ファイアウォール 8 2 3 によって分離され、仮想機能が別の仮想機能の論理パーティション・メモリにアクセスすることを許可するために、ファイアウォール・トンネル 8 2 2 が生成される。アクセスは、各仮想機能がただ 1 つだけの論理パーティションのメモリにアクセスすることを要求する、標準的な周辺機器相互接続エクスプレス入出力仮想化仕様とは異なる。

10

【 0 0 7 1 】

シングルルート・トンネリングを可能にするデータ構造は、図 8 に示されたマルチルートの場合に必要なものに類似している。仮想階層を含むテーブルの代わりに、許可された各仮想階層番号は、許可された周辺機器相互接続エクスプレス・ダイレクト・メモリ・アクセス・アドレス範囲によって置き換えられる。マルチルートの場合のマルチルート周辺機器相互接続マネージャに類似するシングルルート周辺機器相互接続マネージャ（図示されず）は、周辺機器相互接続エクスプレス・ダイレクト・メモリ・アクセス・アドレス範囲をアロケートし、図 1 0 の仮想機能 - アドレス範囲許可テーブル 6 2 8 を設定する。マルチルートの場合の仮想階層のためのように、明示的に設定されない限り、1 つの論理パーティションが別の論理パーティションのメモリにアクセスできないように、論理パーティション内のソフトウェアは、テーブルにアクセスさせてもらえない。図 1 2 のマルチルートの場合におけるように、2 つの論理パーティションは、受信バッファを設定するために、マルチルートの場合のコンピュータ電子コンプレックスにおいてソフトウェアが行ったのと同じ方法で通信する。図 1 1 の仮想機能作業待ち行列エントリ 6 4 2 は、この場合、仮想階層番号を含むが、図 1 1 の 6 4 8 のような動作アドレスをパス識別子として使用することができる。

20

【 0 0 7 2 】

図 1 3 には、2 つの論理パーティションが示されているが、論理パーティションの 3 ウェイ以上の冗長組も構成できることは当業者であれば理解されよう。示されるように、仮想機能は、マルチルート周辺機器相互接続デバイス仮想機能およびシングルルート周辺機器相互接続デバイス仮想機能の一方とすることができる。

30

【 0 0 7 3 】

図 1 4 を参照すると、例示的な一実施形態による、マルチルート・マルチシステム構成のマルチルート・ファブリック構成のプロセスのフローチャートが提示されている。構成プロセス 9 0 0 は、図 1 2 に示されたような構成を提供する、図 3 の構成マネージャ 3 6 4 の構成プロセスの一例である。構成プロセス 9 0 0 が開始し（ステップ 9 0 2 ）、マルチルート周辺機器相互接続マネージャが、マルチルート・ファブリックを構成する（ステップ 9 0 4 ）。マルチルート・ファブリックの構成は、冗長性のための任意の所望の代替ルートを含む、デバイスからルート・コンプレックスまでの正しいルートを生成する。マルチルート周辺機器相互接続マネージャは、仮想階層番号 - 周辺機器相互接続ホスト・ブリッジ（PHB）相関をルート・コンプレックスから利用可能にする（ステップ 9 0 6 ）。マルチルート周辺機器相互接続マネージャは、任意の代替相関を含む仮想機能 - 仮想階層番号許可テーブルを設定するために、デバイス物理機能のためのデバイス・ドライバを起動し（ステップ 9 0 8 ）、その後、構成プロセス 9 0 0 は終了する。

40

【 0 0 7 4 】

図 1 5 を参照すると、例示的な一実施形態による、パートナ・システムに伝達するために必要とされる仮想階層番号をシステムが決定することを可能にするプロセスのフローチャートが提示されている。プロセス 1 0 0 0 は、ルート・ノード 3 6 0 およびルート・ノード 3 6 1 によって図 1 2 の構成を使用する、または図 1 3 の構成ならびに図 4 の論理パ

50

ーティション 4 0 3 および論理パーティション 4 0 5 を使用するプロセスの一例である。

【 0 0 7 5 】

プロセス 1 0 0 0 が開始し (ステップ 1 0 0 2)、それぞれのパートナーおよびパートナーに関連付けられた仮想階層番号を発見するために、コンピュータ電子コンプレックスが互いに通信し、または論理パーティションが使用される場合、論理パーティションが互いに通信する (ステップ 1 0 0 4)。コンピュータ電子コンプレックスまたは論理パーティションの各々は、それぞれのコンプレックスまたはパーティションに関連付けられたデバイスを発見し、それぞれの発見デバイスのためのデバイス・ドライバをロードし、それぞれの仮想機能のための仮想機能 - 仮想階層番号許可テーブルを読む (ステップ 1 0 0 6)。デバイス・ドライバは今では、図 7 の適切な仮想機能作業待ち行列エン트리 6 0 1 を設定するために必要な仮想階層番号を有している。その後、プロセス 1 0 0 0 は終了する (ステップ 1 0 0 8)。

10

【 0 0 7 6 】

図 1 6 を参照すると、例示的な一実施形態による、仮想機能作業待ち行列エントリを設定するためのプロセスのフローチャートが提示されている。プロセス 1 1 0 0 は、図 1 2 の C E C 1 7 0 1 または図 1 3 の L P A R 1 8 0 2 などの中央電子コンプレックスによって、図 5 の仮想機能作業待ち行列エントリ 5 1 1 を確立するためのプロセスの一例である。プロセス 1 1 0 0 が開始し (ステップ 9 0 2)、マスタ・コンピュータ電子コンプレックスまたは論理パーティションが、システム仮想機能内の仮想機能作業待ち行列エントリを設定する (ステップ 1 1 0 4)。生成されたエントリは、動作が適用可能なすべてのコンピュータ電子コンプレックスまたは論理パーティションに対して仮想階層番号を指定する。マスタ・コンピュータ電子コンプレックスまたは論理パーティションは、特定の動作のためのデバイス・ドライバが存在する場所である。すべてのコンピュータ電子コンプレックスまたは論理パーティションは、同時に実行するマスタ動作を有することができる。すなわち、様々なコンピュータ電子コンプレックスまたは論理パーティション間に作業負荷を分散させるために、1 つのコンピュータ電子コンプレックスまたは論理パーティションは、作業負荷の一部を担い、その部分を制御することができ、別のコンピュータ電子コンプレックスまたは論理パーティションは、作業負荷の別の一部を担うことができる。

20

【 0 0 7 7 】

デバイスは、要求された動作を実行し、動作のための仮想機能作業待ち行列エントリ内の仮想階層番号およびアドレスを使用して、データをすべての適切なコンピュータ電子コンプレックスまたは論理パーティションに送信する (ステップ 1 1 0 6)。その後、プロセス 1 1 0 0 は終了する (ステップ 1 1 0 8)。

30

【 0 0 7 8 】

図 1 7 を参照すると、例示的な一実施形態による、入出力ファブリック・パス動作ステータスおよび代替パス使用を動的に決定するためのプロセスのフローチャートが提示されている。構成プロセス 1 2 0 0 は、パス利用可能性を決定するための、図 1 2 の M R デバイス 1 7 2 7 などのデバイスのプロセスの一例である。構成プロセス 1 2 0 0 が開始し (ステップ 1 2 0 2)、デバイスが、システム・メモリへの仮想階層パスの動作ステータスを定期的に決定し、仮想階層パスが利用可能でない場合にはフラグを設定する (ステップ 1 2 0 4)。例えば、デバイスは、ダイレクト・メモリ・アクセスを介してシステム・メモリ内のロケーションを読み、読み取り時に動作タイムアウトなどのエラーを受信した場合、パスを利用不能としてマーク付ける。一般に、デバイスは、1 次パスが利用可能である場合、そのパス上で動作を開始し、それ以外の場合、デバイスは、代替パスを使用する (ステップ 1 2 0 6)。その後、プロセス 1 2 0 0 は終了する (ステップ 1 2 0 8)。

40

【 0 0 7 9 】

図 1 8 を参照すると、例示的な一実施形態による、データ・レプリケーションのためのデュアル・パスを使用して書き込み動作を実行するプロセスのフローチャートが提示されている。構成プロセス 1 3 0 0 は、図 1 2 の C E C 1 7 0 1 または図 1 3 の L P A R 1

50

802などの中央電子コンプレックスのデバイス・ドライバ部分の一例である。プロセス1300は、デュアル・パス・ディスク書き込み動作の一例であるが、より多くまたは多数のパスも同様に使用することができる。例えば、書き込み動作は、1組のパスを含むことができ、その組は3つ以上のパスを含む。ストレージ・サーバ・ディスク書き込み動作のプロセス1300が開始し(ステップ1302)、デバイス・ドライバが、各中央電子コンプレックスに1つの受信バッファが配置される、2つの通信ネットワーク受信バッファを指し示すための仮想機能作業待ち行列エントリを設定する(ステップ1304)。ディスクに書き込まれるデータが、通信ネットワークを介してホストから受信され、両方の中央電子コンプレックス内のバッファに書き込まれる(ステップ1306)。各受信バッファへの書き込み動作は、同時またはほぼ同時に発生する。ディスク書き込み動作は、データがまだディスクに書き込まれていなくても、完了したものとしてホストに通知される(ステップ1308)。ディスク書き込みの正常完了の後、データは物理的にディスクに書き込まれ、両方の中央電子コンプレックスのメモリから廃棄される(ステップ1310)。その後、プロセス1300は終了する(ステップ1312)。

【0080】

したがって、例示的な実施形態は、入出力仮想化デバイスの単一の機能が、多数のシステム間の多数のパスを介して多数のシステムにアクセスするための能力を提供する。特に、単一の機能は、冗長通信パスを確立するために、入出力ファブリックの多数の仮想階層へのアクセスを許可されることができる。冗長システムの確立は、データ・インテグリティの理由で多数のシステムの2つ以上のシステムにデータが送信されること、または同じ入出力仮想化アダプタの同じ機能から2つ以上のシステム内のデータにアクセスすることを可能にする。例示的な一実施形態では、許可は、仮想機能 - 仮想階層許可対応テーブルまたは仮想機能 - アドレス範囲許可テーブルの使用を通して確立される。対応は、機能が、リソースに関連付けられた別の仮想機能または仮想階層のリソースを使用するために、仮想機能または仮想階層を分離するファイアウォールをトンネリングすることを特に許可する。

【0081】

図のフローチャートおよびブロック図は、本発明の様々な実施形態による、システム、方法、およびコンピュータ・プログラム製品の可能な実施の、アーキテクチャ、機能、および動作を示す。この点で、フローチャートまたはブロック図の各ブロックは、規定された論理機能を実施するための1つまたは複数の実行可能命令を含む、モジュール、セグメント、またはコードの部分を表すことができる。いくつかの代替実施では、ブロック内で述べられた機能は、図内で述べられた順序とは異なる順序で発生し得ることに留意されたい。例えば、連続して示される2つのブロックは、実際には実質的に同時に実行されることがあり、またはブロックは、含まれる機能に応じて、時には逆の順序で実行されることもある。ブロック図またはフローチャート図あるいはその両方の各ブロック、およびフローチャート図またはブロック図あるいはその両方におけるブロックの組合せは、指定された機能もしくは動作を実行する専用ハードウェア・ベースのシステムによって、または専用ハードウェアおよびコンピュータ・プログラム命令の組合せによって実施できることに留意されたい。

【0082】

添付の特許請求の範囲におけるすべての手段またはステップおよび機能要素に対応する構造、材料、動作、および均等物は、具体的に特許請求されるように他の特許請求要素と組み合わせて機能を実行するための任意の構造、材料、または動作を含むことが意図されている。本発明の説明は、例示および説明の目的で提示されたものであり、網羅的であること、または開示された形態の本発明に限定することは意図していない。本発明の範囲および主旨から逸脱しない多くの修正および変形が、当業者には明らかであろう。本発明の原理および実際の応用を最もよく説明し、企図された特定の使用に適した様々な修正を有する様々な実施形態のために当業者が本発明を理解できるように、実施形態が選択され、説明された。

【 0 0 8 3 】

本発明は、全面的にハードウェアの実施形態、全面的にソフトウェアの実施形態、またはハードウェア要素とソフトウェア要素の両方を含む実施形態の形態を取ることができる。好ましい一実施形態では、本発明は、ファームウェア、常駐ソフトウェア、マイクロコードなどを含むが、それらに限定されない、ソフトウェアで実施される。

【 0 0 8 4 】

さらに、本発明は、コンピュータもしくは任意の命令実行システムによって使用される、またはそれらに関連して使用されるプログラム・コードを提供するコンピュータ使用可能またはコンピュータ可読媒体から取得可能なコンピュータ・プログラム製品の形態を取ることができる。この説明のため、コンピュータ使用可能またはコンピュータ可読媒体は、命令実行システム、装置、もしくはデバイスによって使用される、またはそれらに関連して使用されるプログラムを含み、保存し、伝達し、伝播し、または転送することができる任意の有形な装置とすることができる。

10

【 0 0 8 5 】

媒体は、電子、磁気、光、電磁気、赤外線、または半導体システム（または装置もしくはデバイス）、あるいは伝播媒体とすることができる。コンピュータ可読媒体の例は、半導体または固体メモリ、磁気テープ、着脱可能コンピュータ・ディスク、ランダム・アクセス・メモリ（RAM）、リード・オンリ・メモリ（ROM）、固定磁気ディスク、および光ディスクを含む。光ディスクの現在の例は、コンパクト・ディスク・リード・オンリ・メモリ（CD-ROM）、コンパクト・ディスク・リード/ライト（CD-R/W）、およびDVDを含む。

20

【 0 0 8 6 】

プログラム・コードの保存または実行あるいはその両方を行うのに適したデータ処理システムは、システム・バスを介して直接的または間接的にメモリ要素に結合された少なくとも1つのプロセッサを含む。メモリ要素は、プログラム・コードの実際の実行中に利用される局所メモリ、大容量記憶装置、および実行中に大容量記憶装置からコードを取り出さなければならない回数を減らすために、少なくとも一部のプログラム・コードの一時記憶を提供するキャッシュ・メモリを含むことができる。

【 0 0 8 7 】

（キーボード、ディスプレイ、ポインティング・デバイスなどを含むが、それらに限定されない）入出力すなわちI/Oデバイスは、直接的に、または介在するI/Oコントローラを介してシステムに結合することができる。

30

【 0 0 8 8 】

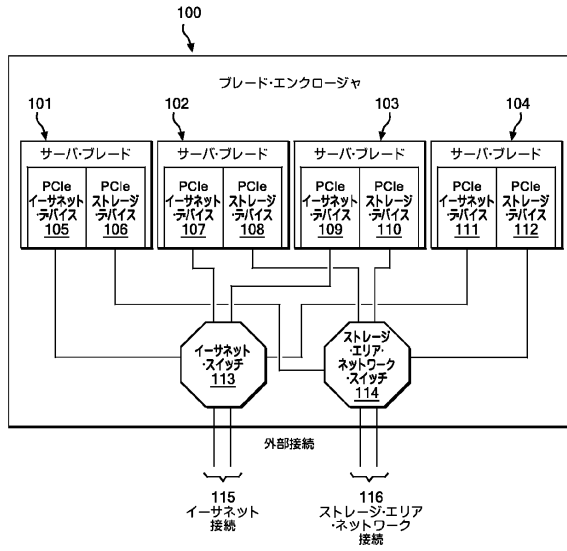
データ処理システムが、介在する私設ネットワークまたは公共ネットワークを介して、他のデータ処理システムまたはリモート・プリンタもしくはストレージ・デバイスに結合できるように、ネットワーク・アダプタも、システムに結合することができる。モデム、ケーブル・モデム、およびイーサネット・カードは、現在利用可能なタイプのネットワーク・アダプタのうちのいくつかに過ぎない。

【 0 0 8 9 】

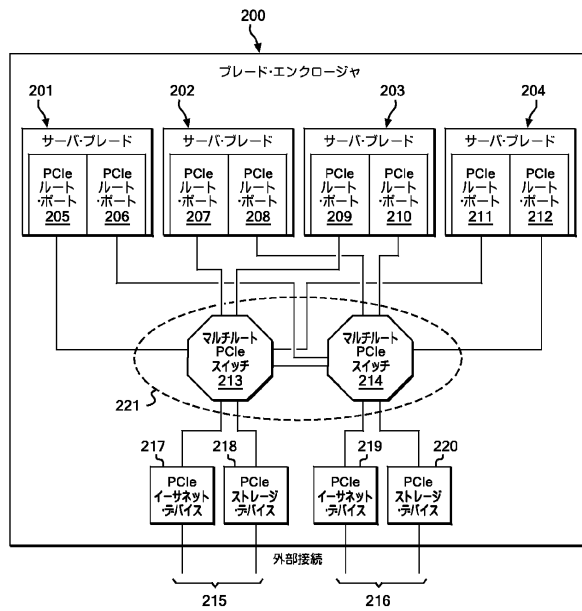
本発明の説明は、例示および説明の目的で提示されたものであり、網羅的であること、または開示された形態の本発明に限定することは意図していない。多くの修正および変形が、当業者には明らかであろう。本発明の原理および実際の応用を最もよく説明し、企図された特定の使用に適した様々な修正を有する様々な実施形態のために当業者が本発明を理解できるように、実施形態が選択され、説明された。

40

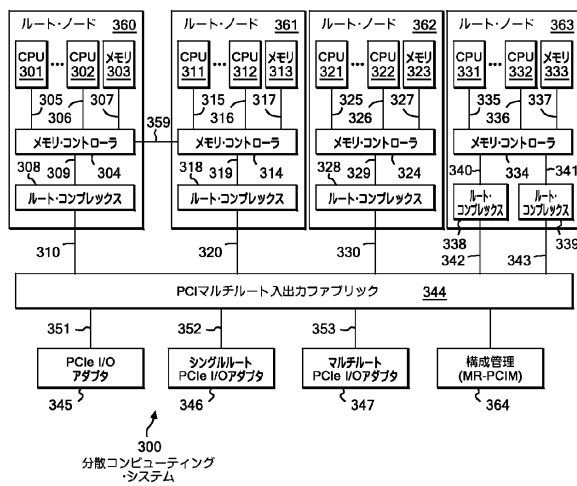
【図 1】



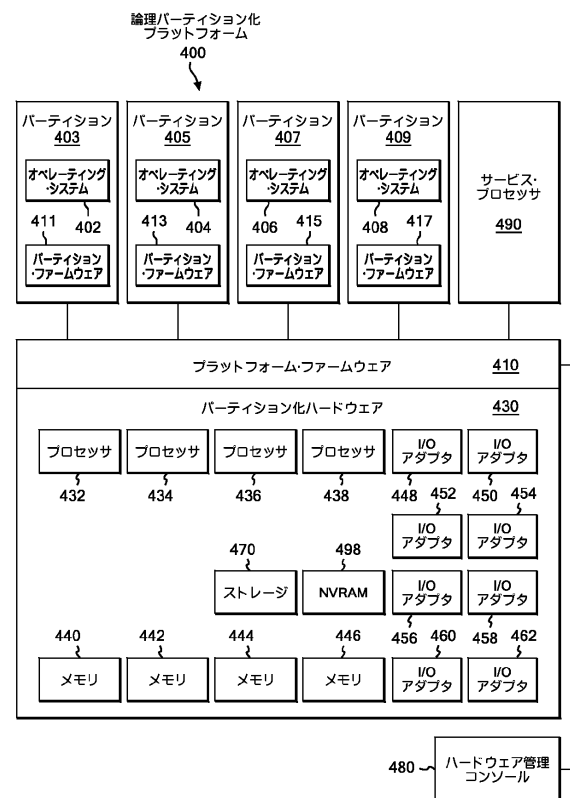
【図 2】



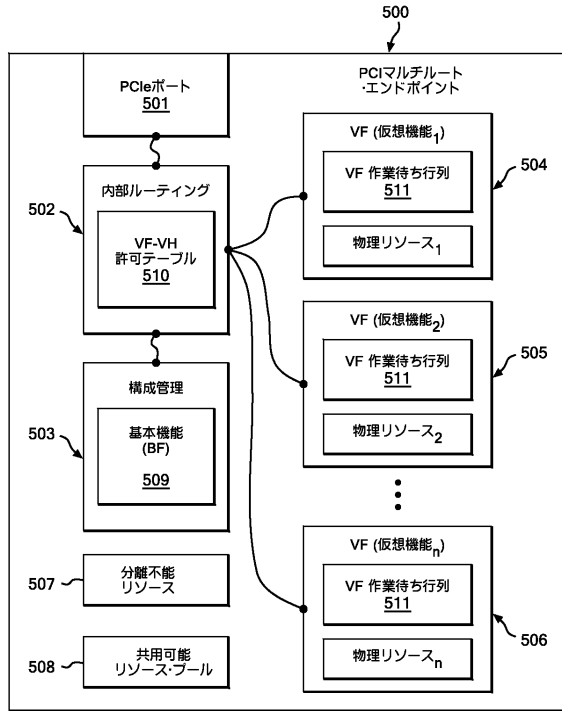
【図 3】



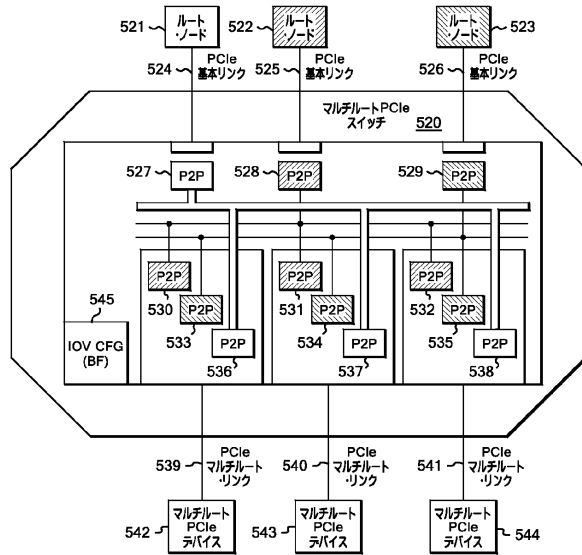
【図 4】



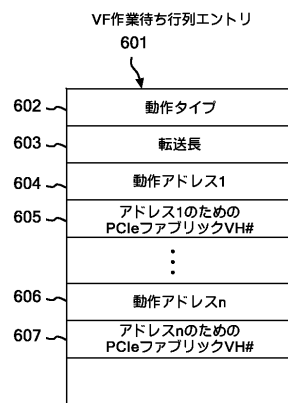
【図 5】



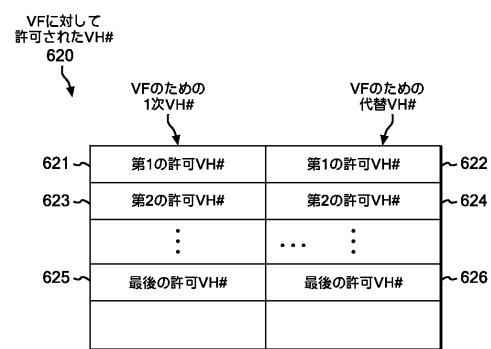
【図 6】



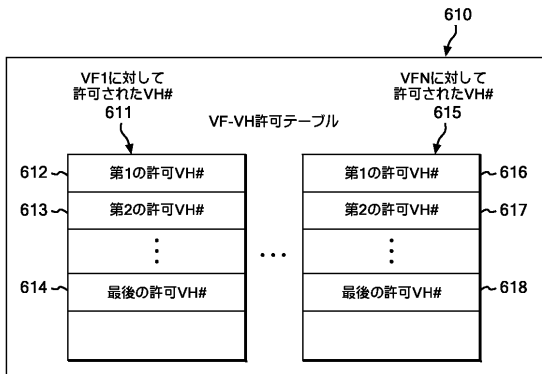
【図 7】



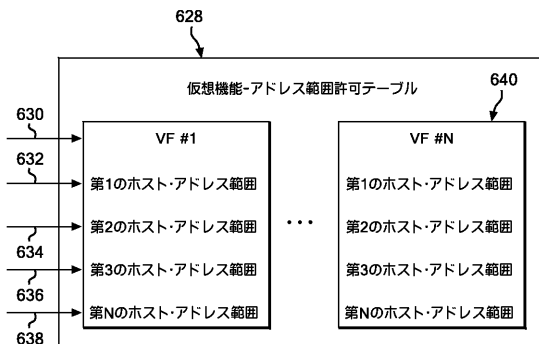
【図 9】



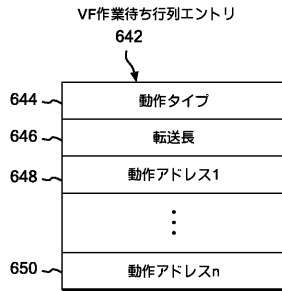
【図 8】



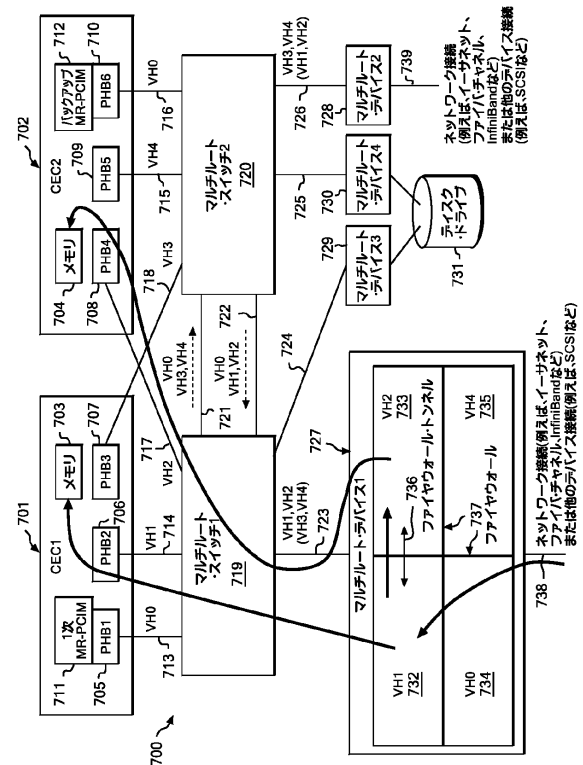
【図 10】



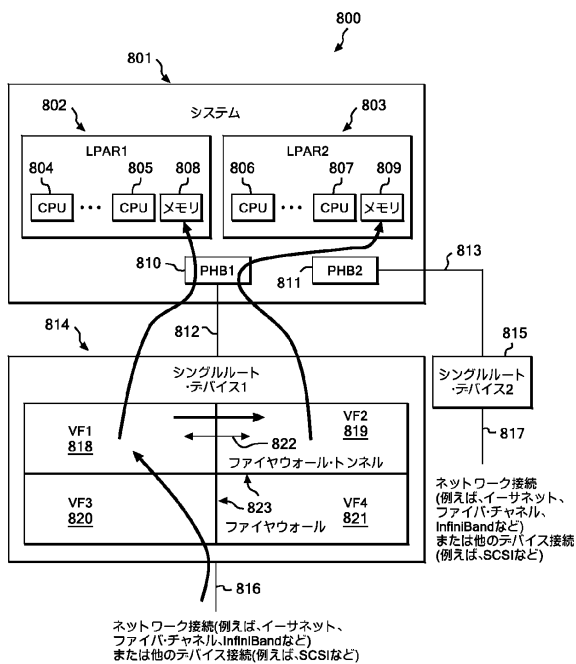
【 ㄨ 1 1 】



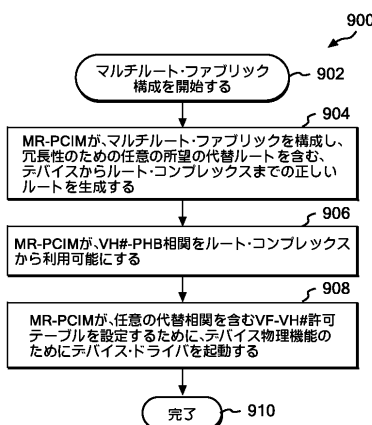
【 図 1 2 】



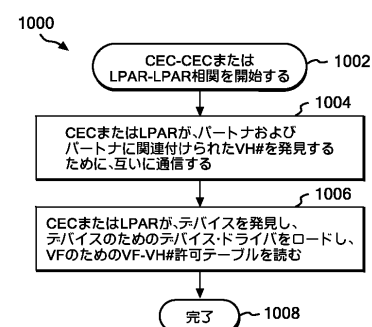
【 図 1 3 】



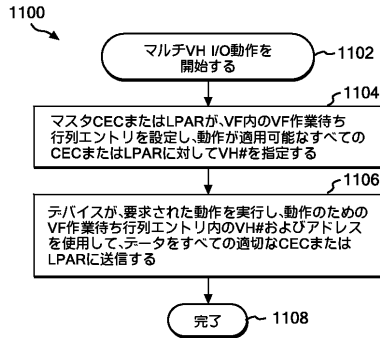
【 図 1 4 】



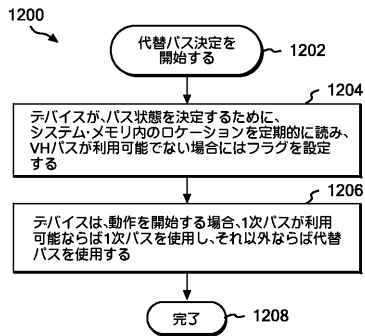
【 図 1 5 】



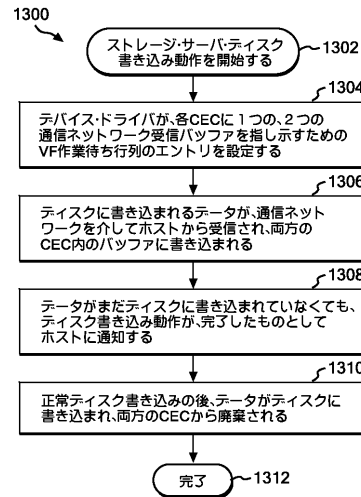
【図 16】



【図 17】



【図 18】



フロントページの続き

(72)発明者 スティーブ・サーバー

アメリカ合衆国 78758 テキサス州 オースティン バーネットロード 11501

(72)発明者 ダグラス・エム・フレイマス

アメリカ合衆国 10532 ニューヨーク州 ホーゾン スカイラインドライブ 19

審査官 坂東 博司

(56)参考文献 特開2008-171413(JP,A)

特開2008-065561(JP,A)

特開2004-013899(JP,A)

特開2008-009980(JP,A)

特開2009-093316(JP,A)

米国特許出願公開第2005/0281191(US,A1)

米国特許出願公開第2003/0088698(US,A1)

米国特許出願公開第2009/0094403(US,A1)

米国特許第07398337(US,B1)

米国特許第07464218(US,B1)

(58)調査した分野(Int.Cl., DB名)

G06F 13/10