

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6055310号
(P6055310)

(45) 発行日 平成28年12月27日(2016.12.27)

(24) 登録日 平成28年12月9日(2016.12.9)

(51) Int.Cl.		F I			
G 0 6 F	9/46	(2006.01)	G 0 6 F	9/46	3 5 0
G 0 6 F	9/50	(2006.01)	G 0 6 F	9/46	4 6 2 Z

請求項の数 18 (全 24 頁)

(21) 出願番号	特願2012-544559 (P2012-544559)	(73) 特許権者	314015767
(86) (22) 出願日	平成22年11月23日 (2010.11.23)		マイクロソフト テクノロジー ライセン
(65) 公表番号	特表2013-514588 (P2013-514588A)		シング, エルエルシー
(43) 公表日	平成25年4月25日 (2013.4.25)		アメリカ合衆国 ワシントン州 9805
(86) 国際出願番号	PCT/US2010/057871		2 レッドモンド ワン マイクロソフト
(87) 国際公開番号	W02011/084257		ウェイ
(87) 国際公開日	平成23年7月14日 (2011.7.14)	(74) 代理人	100140109
審査請求日	平成25年10月9日 (2013.10.9)		弁理士 小野 新次郎
(31) 優先権主張番号	12/640,272	(74) 代理人	100075270
(32) 優先日	平成21年12月17日 (2009.12.17)		弁理士 小林 泰
(33) 優先権主張国	米国 (US)	(74) 代理人	100101373
			弁理士 竹内 茂雄
		(74) 代理人	100118902
			弁理士 山本 修

最終頁に続く

(54) 【発明の名称】 仮想記憶ターゲットオフロード技術

(57) 【特許請求の範囲】

【請求項 1】

一つ又は複数の子パーティションを含むパーティションであって、前記子パーティションは、少なくとも一つの仮想プロセッサとゲストオペレーティングシステムを含む仮想マシンである、パーティションと、

移行可能なストレージサービスを実現する回路であって、前記移行可能なストレージサービスは前記パーティション上で実行可能なものであり、前記移行可能なストレージサービスは前記子パーティションについての仮想ハードディスクの入出力要求を管理するように構成され、前記移行可能なストレージサービスに一意のネットワーク識別子が割り当てられる、移行可能なストレージサービスを実現する回路と、

前記仮想ハードディスクを移行させることなく、前記移行可能なストレージサービスのあるパーティションから少なくとも別のパーティションに移行させる回路と、

前記移行可能なストレージサービスをネットワーク内のストレージターゲットとして設定する回路と

を備えるシステム。

【請求項 2】

前記移行可能なストレージサービスのあるパーティションから少なくとも別のパーティションに移行させる回路は、前記移行可能なストレージサービスをリモートコンピューターシステム上のパーティションに移行させる回路を備える請求項 1 に記載のシステム。

【請求項 3】

10

20

前記子パーティションについての入出力要求に関連付けられるゲスト物理アドレスをシステム物理アドレスへ変換する入出力メモリー管理ユニットを構成する回路をさらに備える請求項 1 又は 2 に記載のシステム。

【請求項 4】

前記子パーティションが第 2 の一意のネットワーク識別子が割り当てられた仮想ネットワークアダプターを含む請求項 1 ~ 3 のいずれか一項に記載のシステム。

【請求項 5】

前記移行可能なストレージサービスが前記一意のネットワーク識別子が割り当てられた仮想ネットワークアダプターを含む請求項 1 ~ 4 のいずれか一項に記載のシステム。

【請求項 6】

前記移行可能なストレージサービスと、前記仮想マシンを管理するように構成される管理サービスを、前記子パーティション以外のパーティションによって実現する回路をさらに備える請求項 1 ~ 5 のいずれか一項に記載のシステム。

【請求項 7】

前記仮想プロセッサにマッピングされる一つ又は複数の論理プロセッサと、前記子パーティションについての入出力要求を受信することに対応して、前記論理プロセッサに前記移行可能なストレージサービスを実行するように通知する回路をさらに備える請求項 1 ~ 6 のいずれか一項に記載のシステム。

【請求項 8】

前記子パーティションについての入出力要求が前記一意のネットワーク識別子と少なくとも 1 つの他の一意のネットワーク識別子との間のものであるときに、前記子パーティションについての入出力要求がセキュリティポリシーに準拠するか否かであることを決定する回路をさらに備える請求項 1 ~ 7 のいずれか一項に記載のシステム。

【請求項 9】

一つ又は複数の子パーティションを含むパーティションであって、前記子パーティションは、少なくとも一つの仮想プロセッサとゲストオペレーティングシステムを含む仮想マシンであるパーティション、を含むシステムで実行される方法であって、

移行可能なストレージサービスであって、前記移行可能なストレージサービスは前記パーティション上で実行可能なものであり、前記移行可能なストレージサービスは前記子パーティションについての仮想ハードディスクの入出力要求を管理するように構成され、前記移行可能なストレージサービスに一意のネットワーク識別子が割り当てられている、移行可能なストレージサービスを実現するステップと、

前記仮想ハードディスクを移行させることなく、前記移行可能なストレージサービスのあるパーティションから少なくとも別のパーティションに移行させるステップと

前記移行可能なストレージサービスをネットワーク内のストレージターゲットとして設定するステップと

を含む方法。

【請求項 10】

前記移行可能なストレージサービスのあるパーティションから少なくとも別のパーティションに移行させるステップは、前記移行可能なストレージサービスをリモートコンピューターシステム上のパーティションに移行させるステップを備える請求項 9 に記載の方法。

【請求項 11】

前記子パーティションについての入出力要求に関連付けられるゲスト物理アドレスをシステム物理アドレスへ変換する入出力メモリー管理ユニットを構成するステップをさらに含む請求項 9 又は 10 に記載の方法。

【請求項 12】

前記子パーティションが第 2 の一意のネットワーク識別子が割り当てられた仮想ネットワークアダプターを含む請求項 9 ~ 11 のいずれか一項に記載の方法。

【請求項 13】

10

20

30

40

50

前記移行可能なストレージサービスが前記一意のネットワーク識別子が割り当てられた仮想ネットワークアダプターを含む請求項 9 ~ 12 のいずれか一項に記載の方法。

【請求項 14】

前記移行可能なストレージサービスと、前記仮想マシンを管理するように構成される管理サービスを、前記子パーティション以外のパーティションによって実現するステップをさらに含む請求項 9 ~ 13 のいずれか一項に記載の方法。

【請求項 15】

前記システムは、前記仮想プロセッサにマッピングされる一つ又は複数の論理プロセッサを更に備え、

前記コンピューター方法は、前記子パーティションについての入出力要求を受信することに対応して、前記論理プロセッサに前記移行可能なストレージサービスを実行するように通知するステップをさらに含む請求項 9 ~ 14 のいずれか一項に記載の方法。

10

【請求項 16】

前記子パーティションについての入出力要求が前記一意のネットワーク識別子と少なくとも 1 つの他の一意のネットワーク識別子との間のものであるときに、前記子パーティションについての入出力要求がセキュリティポリシーに準拠するか否かであることを決定するステップをさらに含む請求項 9 ~ 15 のいずれか一項に記載の方法。

【請求項 17】

請求項 9 ~ 16 のいずれか一項に記載の方法を実行するためのプログラム。

【請求項 18】

20

請求項 9 ~ 16 のいずれか一項に記載の方法を実行するためのプログラムを記録した記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、仮想記憶ターゲットオフロード技術に関する。

【背景技術】

【0002】

[0001]仮想マシン技術は、ワークロードをパッケージ化し、データセンターにそれを移動させるために使用することができる。1つの物理ホストから別の物理ホストへワークロードを移動させるこの機能は、ハードウェアコスト及び管理コストをはるかに低くすることにつながる動的なマシンの統合を可能にするため、ユーザーにとって大きな利点である。仮想マシンは、通常、ハイパーバイザー、管理パーティション又はそれらの組み合わせ内に位置するストレージ仮想化を処理するモジュールを介してストレージにアクセスする。このモデルでは、仮想マシンは、通常、参照により全体としてその内容が本明細書に組み込まれる「パーティションバス」と題した米国特許出願第 11 / 128,647 号に記載される例示的なパーティションバスのようなパーティション間通信バスなどのソフトウェア通信経路を介してモジュールへストレージ入出力要求を送信する。

30

【発明の概要】

【発明が解決しようとする課題】

40

【0003】

仮想マシンとハイパーバイザー（又は管理パーティション）との間の通信は、通信経路の実行及びメッセージを転送するときに発生するコンテキストスイッチによって、CPU サイクルコストを被る。したがって、CPU コストを減少させることによって入出力要求を処理する効率を高めるための技術が望まれる。

【課題を解決するための手段】

【0004】

[0002]本開示の 1 つの実施例は方法について記載する。この例において、当該方法は、移行可能なストレージサービスをもたらすステップであって、当該移行可能なストレージサービスは子パーティションについての仮想ハードディスクの入出力要求を管理するよう

50

に構成され、当該移行可能なストレージサービスはネットワークの独自のネットワーク識別子が割り当てられる、ステップと、ネットワーク内のストレージターゲットとして移行ストレージサービスを構成するステップとを含むが、これに限定されない。上記に加えて、他の態様が、特許請求の範囲、図面、及び本開示の一部を形成する記述に記載されている。

【 0 0 0 5 】

[0003]本開示の例示的な実施例は方法について記述する。この例において、方法は、子パーティションについての仮想ハードドライブディスク入出力要求を管理するように構成されたストレージサービスに、ネットワークアダプターの第1の一意のネットワーク識別子をアタッチ (attach) するステップと、ネットワークアダプターによってもたらされる仮想関数を子パーティションにアタッチするステップであって、当該仮想関数が第2の一意のネットワーク識別子を含む、ステップとを含むが、これに限定されない。上記に加えて、他の態様は、特許請求の範囲、図面、及び本開示の一部を形成する記述に記載される。

10

【 0 0 0 6 】

[0004]本開示の例示的な実施例は方法について記述する。この例において、当該方法は、子パーティションにおいてストレージサービスを実行するステップであって、当該ストレージサービスは第2の子パーティションについての仮想ハードドライブディスク入出力要求を管理するように構成され、当該記憶装置にネットワーク内で一意のネットワーク識別子が割り当てられる、ステップを含むが、これに限定されない。上記に加えて、他の態様は、特許請求の範囲、図面、及び本開示の一部を形成する記述に記載される。

20

【 0 0 0 7 】

[0005]当業者であれば、本開示の1つ又は複数の様々な態様が、本開示についての本明細書にて参照される態様を達成する回路及び/又はプログラミングに限定されず、当該回路及び/又はプログラミングが、システム設計者の設計の選択肢に応じて、本明細書にて参照される態様を達成するように構成された、ハードウェア、ソフトウェア、及び/又はファームウェアの実質的に任意の組み合わせであってもよいことを認識することができる。

【 0 0 0 8 】

[0006]上記は概要であり、したがって必然的に、単純化、一般化及び詳細についての省略を含んでいる。当業者であれば、この概要が例示的なものにすぎず、いかなるようにも限定的であるようには意図されていないことを理解するであろう。

30

【図面の簡単な説明】

【 0 0 0 9 】

【図1】[0007]本開示の態様を実施することができる例示的なコンピューターシステムを示す。

【図2】[0008]本開示の態様を実施するための動作環境を示す。

【図3】[0009]本開示の態様を実施するための動作環境を示す。

【図4】[0010]SR-IOVに準拠したネットワークデバイスを含むコンピューターシステムを示す。

40

【図5】[0011]仮想化環境でのメモリー間の関係を示す。

【図6】[0012]本開示の実施例を示す。

【図7】[0013]本開示の態様を説明するための動作環境を示す。

【図8】[0014]本開示の態様を実施するための動作手順を示す。

【図9】[0015]図8の動作手順の代替的な実施例を示す。

【図10】[0016]本開示の態様を実施するための動作手順を示す。

【図11】[0017]図10の動作手順の代替的な実施例を示す。

【図12】[0018]本開示の態様を実施するための動作手順を示す。

【図13】[0019]図12の動作手順の代替的な実施例を示す。

【発明を実施するための形態】

50

【 0 0 1 0 】

[0020]実施例は1つ又は複数のコンピューターシステム上で実行することができる。図1及び以下の説明は、本開示を実施することができる適切なコンピューティング環境の簡潔な一般的な説明を提供することを目的とする。

【 0 0 1 1 】

[0021]本開示にわたって使用される回路という用語は、ハードウェア割り込みコントローラー、ハードドライブ、ネットワークアダプター、グラフィックプロセッサ、ハードウェアベースの動画/音声コーデック、及びこのようなハードウェアを動作させるために使用されるファームウェアなどのハードウェアコンポーネントを含んでもよい。回路という用語はまた、マイクロプロセッサ、特定用途向け集積回路、及び/又は1つもしくは複数の論理プロセッサ、例えば、ファームウェア及び/又はソフトウェアによって構成されたマルチコアの一般的な処理装置の1つ又は複数のコアを含んでもよい。論理プロセッサは、RAM、ROMなどのメモリー、ファームウェア、及び/又はマスメモリからロードされる機能を実行するように動作可能な論理を具体化する命令によって構成することができる。回路がハードウェアとソフトウェアの組み合わせを含む例示的な実施例において、開発者は、論理プロセッサによって実行できるマシン読み取り可能なコードへとその後コンパイルされる論理を具体化するソースコードを書いてもよい。当業者であれば、最先端の技術がハードウェアにより実施される機能とソフトウェアにより実施される機能との間にほとんど差がないところまで進化してきたことを理解することができるので、本明細書に記載される機能を達成するためのハードウェアとソフトウェアの選択は、設計上の選択にすぎない。換言すれば、当業者は、ソフトウェア処理を等価なハードウェア構造へと変換することができ、ハードウェア構造自体を等価なソフトウェア処理へと変換することができることを理解し得るので、ハードウェアによる実施とソフトウェアによる実施との選択は開発者に委ねられる。

【 0 0 1 2 】

[0022]図1を参照すると、例示的なコンピューティングシステム100が示されている。コンピューターシステム100は、論理プロセッサ102、例えば、実行コアのハイパースレッドを含んでもよい。1つの論理プロセッサ102が示されているが、他の実施例では、コンピューターシステム100は、複数の論理プロセッサ、例えば、プロセッサ基板及び/又は各々が複数の実行コアを持つことができる複数のプロセッサ基板あたり、複数の実行コアを有してもよい。図で示すように、様々なコンピューター読み取り可能な記憶媒体110は、様々なシステムコンポーネントを論理プロセッサ102に結合する1つ又は複数のシステムバスによって相互接続することができる。システムバスは、メモリーバスもしくはメモリーコントローラー、周辺バス、及び各種のバスアーキテクチャーうちの任意のものを使用するローカルバスを含む、幾つかの種類のバス構造のうちの任意のものであってもよい。例示的な実施例において、コンピューター読み取り可能な記憶媒体110は、例えば、ランダムアクセスメモリー(RAM)104、例えば電気機械ハードドライブ、固体ハードドライブなどの記憶装置106、ファームウェア108、例えばフラッシュRAM又はROM、例えばCD-ROM、フロッピーディスク、DVD、フラッシュドライブ、外付け記憶装置などの取り外し可能な記憶装置118を含んでもよい。取り外し可能な記憶装置は、磁気カセット、フラッシュメモリーカード、デジタルビデオディスク、ベルヌーイカートリッジなどの、他の種類のコンピューター読み取り可能な記憶媒体を使用することが当業者にとって理解されるべきである。

【 0 0 1 3 】

[0023]コンピューター読み取り可能な記憶媒体110は、プロセッサ実行可能な命令122の不揮発性及び揮発性の記憶装置、データ構造、プログラムモジュール、及びコンピューター100のための他のデータを提供することができる。起動中にコンピューター・システム100内の要素間で情報を転送するのを助ける基本ルーチンを含む基本入出力システム(BIOS)120は、ファームウェア108に格納することができる。多くのプログラムを、ファームウェア108、記憶装置106、RAM104、及び/又は取り

10

20

30

40

50

外し可能な記憶装置 118 に格納することができ、また、オペレーティングシステム及び / 又はアプリケーションプログラムを含む論理プロセッサ 102 によって実行することができる。

【0014】

[0024] コマンド及び情報は、キーボード及びポインティングデバイスを含み得るがこれらに限定されない入力装置 116 を介してコンピューター 100 によって受信されてもよい。他の入力装置は、マイクロフォン、ジョイスティック、ゲームパッド、スキャナーなどを含んでもよい。これら及び他の入力装置は、システムバスに結合され、しばしばユニバーサルシリアルバスポート (USB) などの他のインターフェースによって接続される、シリアルポートインターフェースを介して論理プロセッサ 102 に接続することができる。ディスプレイや他の種類の表示装置はまた、グラフィックプロセッサ 112 の一部であるか又はグラフィックスプロセッサ 112 に接続することができる、ビデオアダプターなどのインターフェースを介してシステムバスに接続することができる。ディスプレイに加えて、コンピューターは、通常、スピーカーやプリンターなどの他の周辺出力装置 (図示せず) を含む。図 1 の例示的なシステムはまた、ホストアダプター、小型コンピューターシステムインターフェース (SCSI) バス、及び SCSI バスに接続された外部記憶装置を含んでもよい。

10

【0015】

[0025] コンピューターシステム 100 は、リモートコンピューターへの論理接続を使用してネットワーク環境で動作することができる。リモートコンピューターは、別のコンピューター、サーバー、ルーター、ネットワーク PC、ピアデバイス又は他の共通ネットワークノードであってもよく、通常、コンピューターシステム 100 に関連して上述した要素のうちの多く又はすべてを含んでもよい。

20

【0016】

[0026] LAN 又は WAN ネットワーキング環境で使用される場合、コンピューターシステム 100 は、ネットワークインターフェースカード 114 を介して LAN 又は WAN に接続することができる。内蔵であっても又は外付けであってもよい NIC 114 は、論理プロセッサに接続することができる。ネットワーク環境において、コンピューターシステム 100 に関連して示したプログラムモジュール又はその一部は、リモートメモリー記憶装置に格納することができる。ここで説明するネットワーク接続が例示的なものであって、コンピューター間の通信リンクを確立する他の手段を使用することができることが理解されるであろう。さらに、本発明の多数の実施例は特にコンピューター化されたシステムに適していることを想定しているが、本明細書のいかなる記載も本発明をそのような実施例に限定することを意図されてはいない。

30

【0017】

[0027] 図 2 及び 3 を参照すると、仮想マシンをもたらすように構成されたコンピューターシステム 200 及び 300 の高レベルのブロック図を示す。本開示の例示的な実施例において、コンピューターシステム 200 及び 300 は、図 1 に記載した要素を含んでもよく、仮想マシンを達成するように動作可能なコンポーネントを含んでもよい。図 2 に移ると、1つのそのようなコンポーネントは、当技術分野で仮想マシンモニターと呼ばれることもあるハイパーバイザー 202 である。図示された実施例のハイパーバイザー 202 は、コンピューターシステム 100 のハードウェアへのアクセスを制御して解決するように構成されてもよい。概して、ハイパーバイザー 202 は、パーティションと呼ばれる実行環境、例えば仮想マシン、を生成することができる。実施例において、子パーティションは、ハイパーバイザー 202 によってサポートされる分離の基本単位と見なすことができる。すなわち、各々の子パーティション (246 及び 248) は、ハイパーバイザー 202 及び / 又は親パーティションの制御下にあるハードウェア資源の組、例えばメモリー、デバイス、論理的なプロセッササイクルなど、にマッピングすることができ、ハイパーバイザー 202 は、別のパーティションのリソースへのアクセスから 1つのパーティション内の処理を分離することができ、例えば、1つのパーティション内のゲストオペレーテ

40

50

イングシステムは、別のパーティションのメモリーから分離することができる。実施例において、ハイパーバイザー 202 は、スタンドアロンのソフトウェア製品、オペレーティングシステムの一部であってもよく、マザーボードのファームウェアに埋め込まれてもよく、特殊な集積回路であってもよく、又はこれらの組み合わせであってもよい。

【0018】

[0028]図示された例において、コンピューターシステム 100 は、オープンソースコミュニティにおけるドメイン0に類似したものとも考えることができる、親パーティション 204 を含む。親パーティション 204 は、通常はオープンソースコミュニティのバックエンドドライバーと呼ばれる仮想化サービスプロバイダー 228 (VSP に) を使用して、子パーティションにおいて実行されるゲストオペレーティング・システムにリソースを提供するように構成することができる。この例のアーキテクチャーにおいて、親パーティション 204 は、基盤となるハードウェアへのアクセスをゲート制御することができる。概して、VSP 228 は、仮想サービスクライアント (VSC) (通常、オープンソースコミュニティのフロントエンドドライバーと呼ばれる) によって、ハードウェア資源へのインターフェースを多重化するのに使用することができる。各々の子パーティションは、ゲストオペレーティングシステム 220 から 222 がその上で実行するスレッドを管理しスケジュールすることができる、仮想プロセッサ 230 から 232 などの 1 つ又は複数の仮想プロセッサを含んでもよい。一般に、仮想プロセッサ 230 から 232 は、特定のアーキテクチャーを備えた物理プロセッサの表現を提供する実行可能な命令及び関連する状態情報である。例えば、1 つの子パーティションが Intel x86 プロセッサの特徴を有する仮想プロセッサを有し得る一方、別の仮想プロセッサは PowerPC プロセッサの特徴を有してもよい。この例における仮想プロセッサは、命令についての仮想プロセッサの実行が論理プロセッサによってバックアップされるように、コンピューターシステムの論理プロセッサにマッピングすることができる。したがって、これらの例示的な実施例では、たとえば、別の論理プロセッサがハイパーバイザーの命令を実行している間に、複数の仮想プロセッサが同時に実行することができる。仮想プロセッサ、様々な VSC、及びパーティション内のメモリーの組み合わせを、仮想マシンと考えることができる。

【0019】

[0029]ゲストオペレーティングシステム 220 から 222 は、例えば、マイクロソフト (登録商標)、アップル (登録商標)、オープンソースコミュニティからオペレーティングシステムなどの任意のオペレーティングシステムを含んでもよい。ゲストオペレーティングシステムは、動作のユーザー / カーネルモードを使用することができ、スケジューラー、メモリーマネージャーなどを含み得るカーネルを有してもよい。各ゲストオペレーティングシステム 220 から 222 は、ターミナルサーバー、電子商取引サーバー、電子メールサーバーなど、及びゲストオペレーティングシステム自体などの、格納されたアプリケーションを有し得る関連付けられたファイルシステムを有してもよい。ゲストオペレーティングシステム 220 から 222 は仮想プロセッサ 230 から 232 上で実行するスレッドをスケジュールすることができ、そのようなアプリケーションのインスタンスを達成することができる。

【0020】

[0030]ここで図 3 を参照すると、図 2 における上記のものに対する代替的なアーキテクチャーを示す。図 3 は図 2 と同様のコンポーネントを示す。しかし、この例示的な実施例において、ハイパーバイザー 202 は、仮想化サービスプロバイダー 228 及びデバイスドライバー 224 を含んでもよく、親パーティション 204 は設定ユーティリティ 236 を含んでもよい。このアーキテクチャーにおいて、ハイパーバイザー 202 は、図 2 のハイパーバイザー 202 と同一又は類似の機能を実行することができる。図 3 のハイパーバイザー 202 は、スタンドアロンのソフトウェア製品、オペレーティングシステムの一部であってもよく、マザーボードのファームウェア内に埋め込まれてもよく、また、ハイパーバイザー 202 の一部は、特殊な集積回路によって達成することができる。この例で

は、親パーティション 204 は、ハイパーバイザー 202 を設定するために使用できる命令を有してもよいが、ハードウェアへのアクセス要求は親パーティション 204 に渡される代わりにハイパーバイザー 202 によって扱われてもよい。

【0021】

[0031] 本開示の実施例において、参照によりその全体が明確に本明細書に組み込まれる「単ルート入出力仮想化仕様」リビジョン 1.0 に準拠したネットワークアダプターは、図面に記載されるものなどのコンピューター・システムにインストールすることができる。例示的なアダプターは、インテル（登録商標）の「ギガビット E T デュアルポートサーバーアダプター」であってもよい。SR-IOV 対応のネットワークデバイスは、物理的な機能に対するインターフェースを仮想化することによって、例えば仮想マシン間の入出力アダプター又は任意の他の処理を共有することができるハードウェアデバイスである。仮想関数（VF）としても知られる各々の仮想インターフェースは、おおまかに言って、コンピューターシステムの PCI-Express バス上の別個のネットワークインターフェースカードのように見える。たとえば、各仮想関数は、エミュレートされた PCI 設定スペースと、独自のネットワーク識別子、例えば、メディアアクセス制御アドレス（MAC アドレス）、ワールドワイドな名前（world wide name）などを有してもよい。したがって、各仮想関数は、物理関数へアクセスするために、一意にアドレス指定され、強力に分離された別個のパスをサポートすることができる。

10

【0022】

[0032] 図 4 に目を向けると、SR-IOV に準拠したアダプター 402（「アダプター」）を含むコンピューターシステム 400 を示す。上記のものと同様に、コンピューターシステム 400 は、図 1-3 に関して上記と同様のコンポーネントを含んでもよい。アダプター 402 は、ネットワーク及び内部ルーター 412 に接続することができる、ポートに対応できる物理関数 410 を含んでもよい。内部ルーター 412 は、例えば、各々が仮想ポートを備えた仮想アダプターなど、仮想関数 404 又は 406 に割り当てられるものなどの、アダプター 402 のネットワーク識別子 420 から 424 へ、又はそれらから、データをルーティングするように構成することができる。

20

【0023】

[0033] 例示的な実施例において、ネットワークアダプター 402 はイーサネット（登録商標）アダプターであってもよく、仮想関数は仮想イーサネットアダプターであってもよい。この例では、仮想関数の一意の識別子はイーサネット MAC アドレスである。ファイバーチャネルの例では、アダプター 402 はファイバーチャネルホストバスアダプターであってもよく、仮想関数は、ワールドワイドノード名及びワールドワイドポート名を含むワールドワイド名を有する仮想ファイバーチャネルホストバスアダプターであってもよい。インフィニバンド（商標）の例では、仮想関数は、グローバル識別子を有する仮想インフィニバンドエンドポイントであってもよい。

30

【0024】

[0034] ネットワーク識別子 424 は、ファイバーチャネルホストバスアダプター又はイーサネットアダプターなどの特定のネットワークアダプターが複数の一意の識別子が単一の物理ポートを共有することを可能にすることを示し、破線で示される。ファイバーチャネルでは、この能力は、NポートID仮想化又はNPIVと呼ばれ、イーサネットにおいて、アダプターはプロミスキヤス（promiscuous）モードと呼ばれるものにおいて動作し、埋め込まれた仮想スイッチを含み、又は、メモリーバッファを分離するために、特定のMACアドレスに対してアドレス指定されるデータをフィルタリングしルーティングしてもよい。

40

【0025】

[0035] 各々のネットワーク識別子は、ネットワークを介して送信できるように情報をフォーマットするように構成されたソフトウェアプロトコルスタック（414-418）に関連付けることができる。特定のTCP/IPの例において、プロセスは、アプリケーション層のポートを介してTCP/IPスタックのアプリケーション層のインスタンスに結

50

合することができる。最終的に、プロトコルスタックの異なる機能によって処理される情報は、ファブリックを介して送信できるデータのフレームの組立を担当しているメディアアクセス制御層として知られるものに存在する、一群の関数によって処理することができる。プロトコルスタックのこの層は、ネットワーク上で送信されるフレームに、仮想関数のためのメディアアクセス制御アドレスを追加する。その後、プロトコルスタックは、フレーム内の情報を電気信号に変換してネットワークへ当該フレームを送信するように構成される物理層に、組み立てられたフレームを渡す。

【 0 0 2 6 】

[0036]入出力メモリー管理ユニット 4 2 6 (I / O - M M U) は、 P C I - E x p r e s s インターコネクトなどの直接メモリーアクセス動作を R A M に対して実行できる入出力相互接続を結合するために使用することができる。本開示の実施例において、 I / O - M M U 4 2 6 は、パーティションからのゲスト物理アドレスをシステム物理アドレスに変換するハイパーバイザー 2 0 2 からのページテーブルを含んでもよい。 I / O - M M U 4 2 6 は、コンピューターシステム 4 0 0 内の複数の位置に存在することができることを示し、破線で示される。たとえば、 I / O - M M U は、マザーボード上のチップ又は論理プロセッサのコンポーネントであってもよい。

【 0 0 2 7 】

[0037]図 5 は、本開示の実施例における、ゲスト物理アドレスとシステム物理アドレスとの間の関係を示す。ゲストメモリーはハイパーバイザー 2 0 2 によって制御されるメモリーの表示である。ゲストメモリーはゲストオペレーティングシステムに割り当てられ、それらのメモリーマネージャーによって制御することができる。ゲスト物理アドレスは、システムの物理アドレス (S P A)、例えばハイパーバイザー 2 0 2 によって管理される物理的なコンピューターシステムのメモリー、によってバックアップすることができる。図面によって示すように、実施例において、 G P A 及び S P A は、メモリー・ブロック、例えば、メモリーの 1 つ又は複数のページ、へと配置することができる。 G P A と S P A との関係は、その内容が参照により全体として本明細書に組み込まれる、「強化されたシャドウページテーブルアルゴリズム」なるタイトルの、同一出願人による米国特許出願第 1 1 / 1 2 8 , 6 6 5 号に記載されるものなどのシャドウページテーブルによって維持することができる。動作において、ゲスト・オペレーティングシステムが G P A ブロック 1 内にデータを格納するとき、データは実際にはシステム上のブロック 6 などの異なる S P A に格納することができる。本開示の実施例において、 I / O - M M U 4 2 6 は、 1 つの G P A 空間から別の G P A 空間に直接的にストレージデータを移動するために、入出力動作中に変換を行うことができる。この実施例では、論理プロセッササイクルは、これらの変換を達成するための命令をハイパーバイザーにおいて実行する必要なくして、保存することができる。

【 0 0 2 8 】

[0038]図 6 は仮想ストレージターゲットオフロード技術を記述する高レベルの動作環境を示す。図 6 は、 S R - I O V ネットワークアダプター 4 0 2 とその仮想関数 4 0 6 を介してストレージの仮想化クライアント 6 0 4 と通信する、仮想マシンストレージサービス 6 0 2 を示す。図に示すように、本開示のこの実施例において、 S R - I O V ネットワークアダプター 4 0 2 は、ソフトウェア通信経路をバイパスすることにより、仮想マシンと仮想マシンのストレージサービスとの間の入出力を伝送するために使用することができる。これは、仮想マシンの入出力を実行するために使用される C P U サイクルの量を減らし、ストレージ・サービス 6 0 2 に移行する能力を高め、潜在的に、親パーティションにおいて実行されるホストオペレーティングシステムの負担及び / 又はハイパーバイザー 2 0 2 の負担を低減する。

【 0 0 2 9 】

[0039]仮想マシンストレージサービス 6 0 2 は、 S A N によって提供される論理ユニット番号 (L U N)、例えば、子パーティションの代わりに他のストレージ仮想化技術によって既に仮想化されてもよいダーク (d i r k)、などの物理的な記憶装置と通信するように

10

20

30

40

50

構成することができる。1つの例では、これは、仮想マシンから入出力要求を受信しLUNヘルディングするように、仮想マシンストレージサービス602を構成することを含んでもよい。別の例では、LUNがサブ割り当てされる場合、仮想マシンストレージサービス602は、仮想ハードドライブを生成し、仮想マシンにそれらを公開(エクスポート)し、LUN又は物理ドライブ上の仮想ハードドライブ(VHD)ファイルとしてそれらを格納するように構成することができる。VHDファイルは単一のファイル内にカプセル化することができる仮想マシンハードディスクを表す。仮想マシンストレージサービス602は、ファイルを解析し、物理的なストレージであるかのようにゲストオペレーティングシステム220にさらすことができるディスクをもたらしすることができる。仮想マシンストレージサービス602によって生成された仮想ハードディスクは、ローカルであるように見えるような方法でゲストオペレーティングシステムにアクセス可能であるバスに対して表すことができる。

10

【0030】

[0040]本開示の実施例において、仮想マシンストレージサービス602は、一意のネットワーク識別子を仮想マシンストレージサービス602に付加し、例えば、データセンター内のストレージターゲットとして仮想マシンストレージサービス602をアドバタイズするために使用されるストレージターゲットパラメーターを設定することによって、ネットワーク内のファイバーチャネルターゲット又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ース(iSCSI)ターゲットなどの、ストレージターゲットであるように構成することができる。iSCSIの例示的な環境において、仮想マシンストレージサービス602は、インターネットプロトコルを介して子パーティションにアクセス可能なLUNをもたらしことによって、iSCSIターゲットを実施することができる。仮想マシンストレージクライアント604又はゲストオペレーティングシステムは、仮想マシンストレージサービス602のアドレスを取得することができ、SCSIハードディスクへの接続をエミュレートする接続を設定することができる。SCSI又はハードドライブ及び仮想マシンストレージサービス602が子パーティションに仮想ハードドライブを提供するのと同じ方法で、仮想マシンストレージクライアント604は仮想マシンストレージサービス602を扱うことができる。この例では、仮想マシンストレージクライアント604は、ネットワークファイルシステム環境で行われるようにリモートディレクトリーをマウントする必要なく、仮想マシンストレージサービス602によって提供される仮想ディスク上のファイルシステムを直接的に作成及び管理することができる。ゲストOS220の観点から、それはハードドライブと同様の方法で機能する1つ又は複数の論理ユニットに結合されるネットワークに結合されたネットワークアダプターを有する。

20

30

【0031】

[0041]図7は、本開示の態様を実施するための例示的な動作環境を示す。図6と同様に、1つ又は複数のSR-IOVネットワークアダプターが、仮想マシンと仮想マシンストレージサービスとの間の入出力を転送するために使用することができ、これによって、ソフトウェア通信パスを使用して入出力を送信する必要性をなくすることができる。これにより、仮想マシンの入出力を実行するために使用されるCPUサイクルの量が低減され、ストレージサービス602に移行する能力が高められ、ホストオペレーティングシステム上の負担及び/又はハイパーバイザー202上の負担が潜在的に低減される。

40

【0032】

[0042]この例の環境においては、2つのコンピューターシステム700及び702を含むデータセンターがスイッチ704に接続されて図示されている(2つのコンピューターシステムが示されるが、当業者であれば、データセンターがより多くのコンピューターシステムを有し得ることを理解することができる)。コンピューターシステム700及び702は図1-4に記載されるものと同様のコンポーネントを有してもよく、スイッチ704は相互接続されたスイッチ及びルーターの全体的なインフラストラクチャーであってもよい。さらに、コンピューターシステム700及び702は、本明細書に開示される技術をより明確に説明するために、特定の機能を含むように示されており、本発明は、示され

50

たトポロジーにおいて実施されるものには限定されない。

【 0 0 3 3 】

[0043] コンピュータシステム 7 0 0 は、本明細書に記載された技術に従ってストレージサービス 6 0 2 に移行するように構成されるマネージャー 2 5 0 を含んでもよく、したがって、仮想マシンストレージサービス 6 0 2 は、同一の又は異なるコンピュータシステムにおいて 1 つのパーティションから別のパーティションに移行することができることを示すように点線で示される。仮想関数 7 0 6 及び 7 0 8 は、特定の実施例において仮想マシンストレージサービス 6 0 2 が仮想関数を介してアクセスする必要なしに S R - I O アダプター 4 0 2 とインターフェースすることができることを示すように、点線で示される。この例示的な実施例において、親パーティション 2 0 4 及び 7 1 2 は物理ハードウェアの制御を有することができ、仮想関数は必要とされない。

10

【 0 0 3 4 】

[0044] 続けて図面について概括すると、仮想マシンストレージサービス 6 0 2 は、本開示の実施例において、割り当てられた一意の識別子を抽出し、任意の必要な状態情報と共に当該識別子を異なるパーティションに移動することによって、移行させることができる。1 つの例では、このプロセスは、マネージャー 2 5 0 を実行する論理プロセッサが一意の識別子を抽出すること、マネージャー 2 5 0 を実行する論理プロセッサが、アダプター (4 0 2 又は 7 1 8) に、異なるパーティションにおける仮想関数に一意の識別子を加えるよう指示すること、及び、マネージャー 2 5 0 を実行する論理プロセッサが、仮想マシンストレージサービス 6 0 2 のインスタンスに、それ自身を仮想関数に加えることを含んでもよい。別の例において、このプロセスは、マネージャー 2 5 0 を実行する論理プロセッサが一意の識別子を抽出すること、マネージャー 2 5 0 を実行する論理プロセッサが、アダプター (4 0 2 又は 7 1 8) に、当該アダプター (4 0 2 又は 7 1 8) に一意の識別子を付加するよう指示すること、及び、マネージャー 2 5 0 を実行する論理プロセッサが、ファブリック上で通信するために一意の識別子を使用するよう、異なるパーティションにおいてインスタンスを作成された仮想マシンストレージサービス 6 0 2 のインスタンスに指示することを含んでもよい。

20

【 0 0 3 5 】

[0045] 以下は動作手順を示す一連のフローチャートである。理解を容易にするために、フローチャートは、最初のフローチャートが全体の「全体像」の視点を介して実施を示し、後続のフローチャートがさらなる追加及び / 又は詳細を提供するように構成される。さらに、当業者であれば、破線で描かれた動作が任意のものであると考えられることを理解することができるであろう。

30

【 0 0 3 6 】

[0046] 図 8 を参照すると、本開示の態様を実施するための動作手順が示される。図に示すように、動作 8 0 0 が動作手順を開始し、動作 8 0 2 は移行ストレージサービスを達成する (effectuate) ことを示し、当該移行可能なストレージサービスは子パーティションについての仮想ハードディスク入出力要求を管理するように構成され、移行可能なストレージサービスにはネットワークについての一意のネットワーク識別子が割り当てられる。例えば、図 6 に移ると、仮想マシンストレージサービス 6 0 2 などの移行可能なストレージサービスは、コンピュータシステムによって達成することができる。つまり、仮想マシンストレージサービス 6 0 2 を示す命令は論理プロセッサによって実行することができる。仮想マシンストレージサービス 6 0 2 は、一意のネットワーク識別子に接続され、1 つのパーティションから別のものへと、それ自体によって、すなわち他の管理モジュールを移動させることなく、移動させることができるので、仮想マシンストレージサービス 6 0 2 は移行可能であると考えられる。

40

【 0 0 3 7 】

[0047] 例示的な実施例において、仮想マシンストレージサービス 6 0 2 はネットワーク上で一意の識別子を排他的に使用することができる、例えば、それは、データセンター内の一意のネットワークアドレスを使用して通信する唯一のプロセスであってもよい。この

50

例において、仮想マシンストレージサービス602は、状態情報が異なるパーティションに送信されて仮想マシンストレージサービス602の別のインスタンスを構成するのに使用できるように、自身の状態をシリアル化する(serialize)ように構成することができる。別の例示的な実施例において、仮想マシンストレージサービス602は、仮想関数に接続される仮想マシンにおいて実行することができる。この例では、仮想マシンストレージサービス602はまた、一意の識別子を使用してネットワークにおいて排他的に通信することができる。仮想マシンストレージサービス602を移行させることは、仮想マシンストレージサービス602を含む仮想マシンの状態をシリアル化すること、及びそれを別のパーティションに送信することを含んでもよい。

【0038】

[0048]具体的な例では、図7に移ると、仮想マシンストレージサービス602は、親パーティション204から子パーティション246へ移行させることができる。この特定の例では、論理プロセッサはマネージャー250を実行することができる、すなわち、論理プロセッサは、マネージャー250を示す命令を実行することができ、データセンター内で通信するために、仮想マシンストレージサービス602によって使用される一意の識別子を抽出することができる。その後、当該一意の識別子は子パーティション246に送信することができ、仮想マシンストレージサービス602のインスタンスを開始することができる。アダプター402中のルーティングテーブルを更新することができ、入出力要求は、親パーティション204の代わりに、子パーティション246へと、アダプター402によってルーティングすることができる。この例では、子パーティション246は、既に使用されている任意の他の一意の識別子に加えて一意の識別子を使用するように構成することができる。

【0039】

[0049]図8の説明を続けると、動作804は、移行ストレージサービスをネットワーク内のストレージターゲットとして設定することを示す。例えば、本開示の実施例において、仮想マシンストレージサービス602は、データセンター内のストレージターゲットであるように構成することができる。上記と同様に、仮想マシンストレージサービス602は、ネットワーク内で一意のネットワーク識別子に接続され、ゲストOS220によってストレージターゲットとして検出することができる。通信セッションをゲストOS220と仮想マシンストレージサービス602との間に開くことができ、ゲストOS220は、仮想マシンストレージサービス602によって公開される仮想ハードドライブを検出し、ローカルハードドライブであるかのように仮想ディスクを使用することができる。特定の例において、仮想マシンストレージサービス602は、上述のようにiSCSIターゲットをエミュレートすることができる。この例では、仮想マシンストレージサービス602は、物理ディスクの代わりに仮想ディスクを公開することができ、LUN又は物理ディスクに対して読み取り又は書き込みをすることによって仮想マシンからの入出力を処理することができる。

【0040】

[0050]図9に移ると、図8の動作手順の代替的な実施例が示されている。動作906は、移行可能なストレージサービスをリモートコンピューターシステムへ移行させることを示す。例えば、図6に目を向けると、実施例において、移行可能なストレージサービス、たとえば仮想マシンストレージサービス602は、データセンター内のリモートコンピューターシステムに移行させることができる。例えば、1つの実施例において、リモートコンピューターシステムは、仮想マシンストレージサービス602を現在ホスティングしているコンピューターシステムよりも多くの入出力帯域幅を利用可能にしてもよく、ストレージサービス602を移動する決定をすることができる。この例では、論理プロセッサは、マネージャー250を実行し、ストレージサービス602に割り当てられる一意の識別子を抽出し、それをリモートコンピューターに送信することができる。その後、リモートコンピューターのマネージャー250は、ストレージサービス602のインスタンスに一意の識別子を付加することができる。

【 0 0 4 1 】

[0051]特定の例において、図 7 に目を向けると、仮想マシンストレージサービス 6 0 2 は、子パーティション 2 4 6 から親パーティション 7 1 2 へと移行させることができる。この特定の例において、コンピューターシステム 7 0 0 のマネージャー 2 5 0 は、仮想マシンストレージサービス 6 0 2 に付加された一意の識別子を抽出し、それをコンピューターシステム 7 0 2 に送信することができる。コンピューターシステム 7 0 2 のマネージャー 2 5 0 は、論理プロセッサ上で実行することができ、親パーティション 7 1 2 において実行される仮想マシンストレージサービス 6 0 2 のインスタンスに一意の識別子を付加することができる。この例において、仮想ストレージサービス 6 0 2 は、仮想関数 7 0 8 を使用して又は使用せずに子パーティション 2 4 6 において仮想マシンストレージサービス 6 0 2 によってサービスを提供されたクライアントからの入出力を送信 / 受信するときに、一意の識別子を使用することができる。

10

【 0 0 4 2 】

[0052]この特定の例において、仮想マシンストレージサービス 6 0 2 の状態情報及びプロトコルスタックは、入出力サービスを中断することができないように、コンピューターシステム 7 0 2 に送信することができる。たとえば、コンピューターシステム 7 0 0 のプロトコルスタックの少なくとも機能的に等価な状態を反映するよう、コンピューターシステム 7 0 2 のマネージャー 2 5 0 がプロトコルスタックを構成することを可能にするのに十分な情報を、コンピューターシステム 7 0 2 に送信することができる。状態情報は、送信されることになっている次のパケットの数 (番号) 、使用されるソケット番号、最大バッファサイズ、サーバーのポート番号、クライアントのポート番号などを含んでもよい。状態情報はまた、より高レベルのプロトコル情報などの情報を含んでもよい。他の例は、使用される暗号化プロトコルに関連する情報であってもよい。

20

【 0 0 4 3 】

[0053]この例示的な実施例においては、クライアントの観点から接続が断たれる代わりに休止されたため、クライアントに対するサービスは中断されずに動作する。例えば、仮想マシンストレージサービス 6 0 2 が移行される場合、プロトコルスタックは、例えば現在の動作を完了させるか又はキャンセルすることによって、それが実行している現在の動作を終えることができ、必要に応じて、プロトコルが短い期間の間情報を送信することを控えることを要求する仮想マシンストレージクライアント 6 0 4 にバインドされるプロトコルに対して、バックオフメッセージを送信することができる。コンピューターシステム 7 0 2 上のプロトコルスタックがインスタンス化される場合、それは、コンピューターシステム 7 0 0 のプロトコルスタックと同等の状態を有することができる、以前にコンピューターシステム 7 0 0 に関連付けられていた一意の識別子を用いてネットワーク上で通信することができる。コンピューターシステム 7 0 2 上で新たに設定されたプロトコルスタックは、必要に応じて再開 (resume) メッセージを送信するよう構成することができ、仮想マシンストレージクライアント 6 0 4 にサービス提供するプロトコルは入出力の送信を再開することができる。スイッチ 7 0 4 は、コンピューターシステム 7 0 2 上の仮想マシンストレージサービス 6 0 2 にプロトコルメッセージが送信されるように、ルーティングを解決することができる。

30

40

【 0 0 4 4 】

[0054]図 9 の説明を続けると、動作 9 0 8 は、子パーティションについての入出力要求に関連付けられるゲスト物理アドレスをシステムの物理アドレスへ変換するように、入出力メモリー管理ユニットを構成することを示す。例えば、図 7 を参照すると、本発明の実施例において、コンピューターシステム 7 0 0 の入出力メモリー管理ユニット 4 2 6 は、ゲスト物理アドレスをシステム物理アドレスに変換するために使用することができる。たとえば、ゲストオペレーティングシステム 2 2 0 が入出力動作、例えば読み取り又は書き込み、を開始する場合、ゲストオペレーティングシステム 2 2 0 は、システム物理アドレスに変換する必要があるゲスト物理アドレスを含むコマンドを生成する。例示的な実施例において、これらの変換は MMU の代わりに I / O - MMU 4 2 6 において生じてモ

50

よい。I/O - MMU 426 にメモリー変換をオフロードすることによって、ハイパーバイザー 202 及び/又は親パーティション 204 の負担が低減される。たとえば、ゲスト OS 220 は、ディスクオフセットをゲストメモリーアドレスへ読み込む要求を含む、読み取り動作を発行することができる。この例では、入出力メモリー管理ユニット 426 は、子パーティション 248 のゲストメモリーアドレスをシステムアドレスにマッピングするテーブルを使用することができ、ゲストが読み取りを始めることを望むゲストメモリーアドレスを物理的に支持するシステムアドレスへとゲストメモリーアドレスに変換することができる。仮想マシンストレージサービス 602 は、要求を受信し、クライアントが要求する情報を取得し、以前に要求されたデータを含む応答メッセージを提供することができる。当該応答はゲストメモリーアドレスとして指定されたバッファ内に提供することができ、その場合、アダプター 402 及び I/O - MMU 426 は提供されたゲストメモリーアドレスをシステム物理アドレスに変換することができ、次いで、アダプター 402 は、クライアントの要求を満たすために、応答バッファからの応答データを要求バッファにコピーすることができる。

【0045】

[0055]この技術は、クライアントが仮想マシンストレージサービス 602 と同じ物理コンピュータ上にある場合に周辺装置によって実行されるメモリー間直接メモリーアクセス (DMA) 動作に類似する。この例示的な実施例において、ネットワークアダプター 402 は、システム物理アドレスの 1 つのブロックから情報を取り出し、それを、仮想マシンストレージクライアント 604 又は仮想マシンストレージサービス 602 に代わってシステム物理アドレスの別のブロックに移動するので、入出力動作はメモリー間の DMA 動作に類似したものとなり得る。特定の例は、仮想マシンストレージクライアント 604 によって発行される読み取り動作を含んでもよい。この例では、仮想マシンストレージクライアント 604 は、それが制御するメモリーページに読み込むことを望むストレージデータのページを指定する、読み取り動作を発行することができる。この例では、データのページは、要求を満たすために、仮想マシンストレージサービス 602 によって使用されるページへとコピーされ、次いで、仮想マシンストレージクライアント 604 によって指定されたメモリーページに当該データをコピーする。

【0046】

[0056]図 9 の説明を続けると、動作 910 は、子パーティションから入出力ジョブ要求を受信することを示し、当該子パーティションは、ネットワークの第 2 の一意のネットワーク識別子を含む仮想関数に接続される。たとえば、図 6 に示すように、実施例において、子パーティション 248 は仮想関数 406 を含んでもよい。この例では、子パーティション 248 は仮想関数 406 を介して SR - IOV アダプター 402 に排他的にインターフェースすることができ、入出力要求を送信することができる。アダプター 402 は、コマンドが仮想マシンストレージサービス 602 に関連付けられた一意の識別子にアドレス指定されていると判断して、それに対してコマンドを送信することができる。この場合、子パーティション 248 からの入出力コマンドは、ハイパーバイザー 202 又はパーティション間の通信インターフェースを介して要求を送信させることなく、仮想マシンストレージサービス 602 に送信することができる。さらに、アダプター 402 は、どのメモリーページをバッファとして使用するか、従って、どのアドレス空間の間でデータをコピーするか、を決定する際に、クライアント 604 及び仮想マシンストレージサービス 602 の一意の識別子を使用することができる。

【0047】

[0057]具体的な例において、入出力要求は、(ゲスト物理アドレスにおける)データの位置及びデータが書き込まれるべき仮想ハードドライブ上の位置を指定する、書き込み動作であってもよい。この例において、ストレージ仮想化クライアント 604 は、仮想マシンストレージサービス 602 の一意の識別子に宛てた情報の 1 つ又は複数のパケットに、当該要求を配置することができる。この例では、アダプター 402 は、要求を受信して仮想マシンストレージサービス 602 に送信することができる。アダプター 402 は、さら

に、データを子パーティションのゲスト物理アドレスから仮想マシンストレージサービス 602 に割り当てられるシステム物理アドレスに移動することができる。つまり、アダプター 402 及び I/O - MMU 426 は、送信バッファ及び受信バッファの両方をゲスト物理アドレスからシステム物理アドレスへ変換するようにこせいすることができ、アダプター 402 は、その後、システム物理アドレスの観点から内部で、データを内部の送信バッファから受信バッファにコピーすることができる。その後、仮想マシンストレージサービス 602 は、その仮想ハードドライブの実施と一貫した適切な位置にデータを格納することができる。当業者であれば理解することができるように、これは仮想ハードドライブファイルを使用することを含んでもよく、LUN 上にデータを格納することを含んでもよく、又は、おそらくは冗長的に、データを格納するための他の技術及び位置を含んでもよい。

10

【0048】

[0058]図9の説明を続けると、動作912は、第1のパーティションにおいて移行可能なストレージサービスを実行すること、第2のパーティションにおいて仮想マシンを管理するように構成される管理サービスを実行することを示し、子パーティションは第3のパーティションである。例えば、1つの実施例において、仮想マシンストレージサービス602は子パーティション246などの第1のパーティションにおいて実行することができ、親パーティション204は管理サービスを実行することができ、仮想マシンストレージクライアント604はパーティション248において実行することができる。この例示的な実施例では、仮想マシンストレージサービス602は、管理プロセスとは別個のパーティションにある。この構成において、子パーティション246は、SANターゲットのように動作する専用のストレージパーティションのように効果的に動作してもよい。この構成により、ハイパーバイザー202及び親パーティション上の負担を低減することができる。たとえば、親パーティションからストレージサービスを分離することによって、オペレーティングシステム内のロッキング(locking)を低減することができる。さらに、このようにしてコンピューターシステムを構成することにより、ハイパーバイザースケジューラー上の負担は、パーティション間で送信する必要があるメッセージの数を減らすことによって低減される。

20

【0049】

[0059]図9の説明を続けると、動作914は、移行可能なストレージサービスを一意のネットワーク識別子を含むネットワークアダプターの仮想関数と関連付けること、及びネットワークアダプターの第2の仮想関数に子パーティションを接続することを示す。例えば、図7に目を向けると、実施例において、仮想マシンストレージサービス602は、仮想関数404などの仮想関数に関連付けることができる。仮想マシンストレージサービス602が子パーティション246において実行される例において、仮想関数404は、制御された方法で、すなわち、子パーティション246における任意のプロセスがそのパーティションの外側にあるデータにアクセスしないことを確実にする方法で、アダプター402にアクセスするように使用することができる。さらに、仮想マシンのスナップショット動作を、仮想マシンストレージサービス602を移行するために使用することができる。

30

40

【0050】

[0060]図9の説明を続けると、動作916は、子パーティションから入出力ジョブ要求を受信することに応答して論理プロセッサに通知を送信すること、及び論理プロセッサが移行可能なストレージサービスを実行していることを決定することを示す。例えば、1つの実施例において、入出力ジョブがソフトウェア処理を必要とする場合、ハイパーバイザー202は割り込みを受信して実行することができる。ハイパーバイザー202は、仮想マシンストレージサービス602を実行しているか又は実行することになっている論理プロセッサを識別することができ、割り込み又は軽量の(lightweight)通知を送信することによって、その論理プロセッサに通知することができる。仮想マシンストレージサービス602が子パーティションに位置する場合、割り込みは、メッセージを処理す

50

るために管理パーティションを開始する必要なしに、論理プロセッサに送信することができる。仮想マシンストレージサービス 602 が現在コンテキストを実行している場合、仮想マシンストレージサービス 602 への切り替えは、生じる必要はなく、代わりに軽量の通知を使用することができるように、中断することもない。

【0051】

[0061]図9の説明を続けると、動作918は、入出力トラフィックがネットワークアダプターを介して独自のネットワーク識別子と少なくとも1つの他の一意のネットワーク識別子との間で運ばれるときに、入出力トラフィックがセキュリティポリシーに準拠していることを決定することを示す。例えば、実施例において、アダプター402はネットワークトラフィックのセキュリティポリシーを含んでもよい。この例示的な実施例において、アダプター402は、仮想マシンストレージサービス602と例えば仮想マシンに付加される別の一意の識別子との間で送信される入出力トラフィックがセキュリティポリシーに従うことを決定するように構成することができる。特定の例において、セキュリティポリシーは、すべての入出力トラフィックが暗号化されることを必要とするかもしれない。この例では、アダプター402は、仮想ハードドライブへの書き込みが平文であるか暗号化されているかを決定するように構成することができる。別の例では、セキュリティポリシーは、仮想ローカルエリアネットワークが完全に別々に保持され、異なる仮想ローカルエリアネットワーク内のエンドポイント間でデータトラフィックが許可されないことを要求してもよい。

【0052】

[0062]次に図10に移ると、動作1000、1002、及び1004を含む本開示の態様を実施するための動作手順を示す。動作1000は動作手順を開始し、動作1002は、ネットワークアダプターの第1の一意のネットワーク識別子を、子パーティションについての仮想ハードドライブディスク入出力要求を管理するように構成されたストレージサービスへ付加することを示す。例えば、図6に目を向けると、本開示の実施例において、SR-IOVアダプター402は、複数のネットワーク識別子をもたらすことができ、それらのうちの1つを仮想マシンストレージサービス602に割り当てることができる。ファイバーチャネルの例では、ファイバーチャネルホストバスアダプターは、複数の一意の識別子が同じポートで使えるようにするために、NポートID仮想化又は(NPIV)を使用することができる。このファイバーチャネルの例では、仮想マシンストレージサービス602は、ファブリック上で通信するために、割り当てられたNPIVアドレスを排他的に使用することができる。

【0053】

[0063]図10の説明を続けると、動作1004は、ネットワークアダプターによってもたらされる仮想関数を子パーティションに付加することを示し、仮想関数は第2の一意のネットワーク識別子を含む。例えば、再び図6を参照すると、SR-IOVアダプター402は、一意のネットワーク識別子を含む仮想関数406をインスタンス化して、仮想マシンに付加することができる。この例示的な実施例において、アダプター402は、ハイパーバイザー202をバイパスしてアダプターを介して入出力要求をストレージサービス602へルーティングするスイッチ、又は別個のパーティション間通信機構として機能するように構成される。今度はこれによって、パーティションに通知しパーティションを切り替えるために論理プロセッサ上で命令を実行するのに使用される時間を減らすことができる。

【0054】

[0064]次に図11に移ると、追加の動作1106、1108、1110、1112、1114を含む図10の動作手順の代替的な実施例を示す。動作1106は、第1の一意のネットワーク識別子を含むように第2の仮想関数を構成する要求を、第2のネットワークアダプターを含むリモートコンピューターシステムへ送信することを示す。例えば、1つの実施例において、論理プロセッサは、マネージャー250において命令を実行することができ、仮想マシンストレージ・サービス602に付加された一意のネットワーク識別

子を含む別のアダプターを有するリモートコンピューターシステムにおいて仮想関数を構成する要求を生成することができる。図 7 に目を向けると、特定の例において、コンピューターシステム 700 のマネージャー 250 は、アダプター 718 を有するコンピューターシステム 702 へ、生成された要求を送信することができる。この例における要求は、仮想関数 710 をインスタンス化し、仮想マシンストレージサービス 602 のインスタンスに関連付けられた一意の識別子を含むように命じるために、コンピューターシステム 702 のマネージャー 250 によって使用することができる。

【0055】

[0065] 図 11 の説明を続けると、動作 1108 は、子パーティションにストレージサービスを移行すること、及び子パーティションに割り当てられる第 2 の仮想関数を、第 1 の一意のネットワーク識別子を使用するように構成することを示す。例えば、図 7 に目を向けると、論理プロセッサは、マネージャー 250 を実行し、例えば親パーティション 204 から子パーティション 246 へと仮想マシンストレージサービス 602 を移行させることができる。この例では、論理プロセッサは、マネージャー 250 を実行し、仮想マシンストレージサービス 602 に関連付けられる一意の識別子を抽出し、それをアダプター 402 に送信することができる。アダプター 402 は、仮想関数 404 をインスタンス化し、それに一意の識別子を付加することができる。その後、マネージャー 250 は、仮想マシンストレージサービス 602 のインスタンスに一意の識別子を付加することができる。この例示的な実施例において、仮想マシンストレージサービス 602 は、管理プロセスとは別個のパーティションにあり、事実上、iSCSI ターゲットのように動作する専用のストレージパーティションとなっている。

【0056】

[0066] 動作 1110 に移ると、入出力メモリー管理ユニットによって、子パーティションからの入出力要求に関連付けられるゲスト物理アドレスをシステム物理アドレスに変換することを示す。例えば、図 7 を参照すると、本開示の実施例において、コンピューターシステム 700 の入出力メモリー管理ユニット 426 は、ゲスト物理アドレスをシステム物理アドレスに変換するために使用することができる。たとえば、ゲストオペレーティングシステム 220 が入出力動作、例えば読み取り又は書き込みを開始するとき、ゲストオペレーティングシステム 220 はゲスト物理アドレスを含むコマンドを生成する。この例では、入出力メモリー管理ユニット 426 は、子パーティション 246 のゲストメモリーアドレスを親パーティション 204 によって使用されるシステムアドレスにマッピングするテーブルを使用することができる。アダプター 402 及び I/O - MMU 426 は、送信バッファー及び受信バッファーの両方をゲスト物理アドレスからシステム物理アドレスに変換するように構成することができ、次いで、アダプター 402 は、内部の送信バッファーから内部の受信バッファーへデータをコピーすること及びその逆を行うことができる。

【0057】

[0067] 動作 1112 に移ると、入出力トラフィックが一意のネットワーク識別子と第 2 の一意のネットワーク識別子との間で運ばれるときに入出力トラフィックのセキュリティポリシーへの準拠を監視するようにネットワークアダプターを構成することを示す。例えば、実施例では、アダプター 402 はネットワークトラフィックのセキュリティポリシーを含んでもよい。この例示的な実施例において、アダプター 402 は、仮想マシンストレージサービス 602 と例えば仮想マシンに付加されるものなどの別の一意の識別子との間で送信される入出力トラフィックがセキュリティポリシーに準拠していることを決定するように構成することができる。特定の例は、特定の仮想マシンがネットワークにおいて特定の一意の識別子を使用して入出力を送信することを要求するセキュリティポリシーを含んでもよい。この例におけるアダプター 402 は、仮想マシンからの情報のパケットを監視し、それらがセキュリティポリシーに準拠しているかどうかを決定することができる。

【0058】

[0068] 動作 1114 に移ると、所定のしきい値を超える量の入出力要求がリモートコン

10

20

30

40

50

ピューターシステムから受信されたという決定に応答して、ストレージサービスをインスタンス化してストレージサービスに第1の一意のネットワーク識別子を割り当てる要求を、リモートコンピューターシステムに送信することを示す。例えば、本開示の1つの実施例において、論理プロセッサは、マネージャー250を示す命令を実行して、仮想マシンストレージサービス602のインスタンスを作成してそれを一意の識別子に付加するようコンピューターシステム702などのリモートコンピューターシステムに指示する要求を送信することができる。論理プロセッサは、入ってくる入出力要求に関連付けられた一意の識別子を監視し、しきい値を超える数の要求がコンピューターシステム702から受信されたことを決定した後に、この要求を生成することができる。特定の例において、マネージャー250は、過去30分にわたって入出力要求の60%がコンピューターシステム702に現在関連付けられる一意の識別子から受信されたことと決定したかもしれない。この例では、マネージャー250は、仮想マシンストレージサービス602がコンピューターシステム702上でローカルに実行して、それを移行させる場合に、データセンターの性能を増加することができることと決定することができる。

【0059】

[0069]ここで図12に移ると、動作1200及び1202を含む動作手順が示される。動作1200は動作手順を開始し、動作1202は子パーティションにおいてストレージサービスを実行することを示し、記憶装置は第2の子パーティションについての仮想ハードドライブディスク入出力要求を管理するように構成され、記憶装置にはネットワーク内の独自のネットワーク識別子が割り当てられる。例えば、1つの実施例において、仮想マシンストレージサービス602は、子パーティション、例えば子パーティション246においてもたすことができ、ネットワーク内の一意の識別子、例えばワールドワイドな名前を割り当てることができる。この例示的な実施例における子パーティション246は、ハイパーバイザー202及び/又は親パーティション204によって制御することができる。この構成において、子パーティション246は、事実上、iSCSIターゲットのように動作する専用のストレージパーティションとなり得る。

【0060】

[0070]ここで図13に移ると、動作1304、1306、1308、1310、及び1312を含む図12の動作手順の代替的な実施例が示される。動作1304を見ると、所定のしきい値を超える量の入出力要求がリモートコンピューターシステムから受信されたという決定に応答して、ストレージサービスをインスタンス化してストレージサービスに第1の一意のネットワーク識別子を割り当てる要求をリモートコンピューターシステムに送信することを示す。例えば、本開示の1つの実施例において、論理プロセッサは、マネージャー250を示す命令を実行し、仮想マシンストレージサービス602のインスタンスを作成してそれを一意の識別子に付加することをコンピューターシステム702などのリモートコンピューターシステムに指示する要求を送信することができる。論理プロセッサは、入ってくる入出力要求に関連付けられた一意の識別子を監視し、しきい値を超える数の要求がコンピューターシステム702から受信されたことと決定した後に、この要求を生成することができる。特定の例において、マネージャー250は、過去30分にわたる入出力要求の60%がコンピューターシステム702に現在関連付けられる一意の識別子から受信されたことと決定するかもしれない。この例では、マネージャー250は、仮想マシンストレージサービス602がコンピューターシステム702上でローカルに実行して、それを移行させる場合に、データセンターの性能を増加することができることと決定することができる。

【0061】

[0071]図13の説明を続けると、動作1306はストレージサービスをハイパーバイザーへ移行することを示す。例えば、図7に目を向けると、実施例において、仮想マシンストレージサービス602はハイパーバイザー202に移行することができる。この例示的な実施例において、コンピューターシステム702は図3に示すものと同様のアーキテクチャーを有してもよく、ストレージ・サービス602を子パーティション246からハイ

パーバイザー 202 に移動させる決定がなされてもよい。この例では、論理プロセッサは、マネージャー 250 を実行して仮想マシンストレージサービス 602 に関連付けられる一意の識別子を抽出することができ、ハイパーバイザー 202 は、仮想マシンストレージサービス 602 のインスタンスにそれを付加することができる。例示的な実施例において、ハイパーバイザー 202 はハードウェアを制御するので、それは、アダプター 402 の物理的な機能にアクセスするように構成することができる。ファイバーチャネルの例においては、ファイバーチャネルホストバスコントローラーは、アダプター 402 を介して入出力コマンドを送受信するべく、一意の識別子を使用するために NPIV を使用することができる。

【0062】

10

[0072] 図 13 の説明を続けると、動作 1308 は、ストレージサービスを親パーティションに移行することを示す。例えば、図 7 に目を向けると、実施例において、仮想マシンストレージサービス 602 は子パーティション 246 から親パーティション 204 又は 712 に移行することができる。この例では、論理プロセッサは、マネージャー 250 を実行し、仮想マシンストレージサービス 602 に関連付けられる一意の識別子を抽出し、また、リモートコンピューター又はローカルコンピューターシステム上の親パーティション 204 にそれを送信することができる。その後、一意の識別子はストレージサービス 602 のインスタンスに付加することができる。

【0063】

[0073] 図 13 の説明を続けると、動作 1310 は、ハイパーバイザーにストレージサービスを割り当てることを示す。この例示的な実施例において、仮想関数 406 は、子パーティション 248 に接続することができ、ネットワーク上の第 2 の一意のネットワーク識別子を有することができる。図に示すように、この例示的な実施例において、子パーティション 246 及び 248 の両方は同じ SR-IOV アダプター 402 に接続することができる。したがって、この例示的な実施例において、入出力要求は、ハイパーバイザー 202 を介する代わりに SR-IOV アダプター 402 を介して、又はスイッチ 704 を介して入出力を送信する必要なしにパーティション間通信機構を介して、渡すことができる。

【0064】

[0074] 図 13 の説明を続けると、動作 1312 は、子パーティションについての入出力要求に関連付けられるゲスト物理アドレスをシステム物理アドレスへ変換するように入出力メモリー管理ユニットを構成することを示す。例えば、図 7 を参照すると、本開示の実施例において、コンピューターシステム 700 の入出力メモリー管理ユニット 426 は、ゲスト物理アドレスをシステム物理アドレスに変換するために使用することができる。たとえば、ゲストオペレーティングシステム 220 が入出力動作、例えば読み取り又は書き込み、を開始する場合、ゲストオペレーティングシステム 220 はゲスト物理アドレスを含むコマンドを生成する。この例では、入出力メモリー管理ユニット 426 は、子パーティション 248 のゲストメモリーアドレスを親パーティションによって使用されるシステムアドレスへマッピングするテーブルを使用することができる。アダプター 402 及び I/O-MMU 426 は、送信バッファー及び受信バッファーの両方をゲスト物理アドレスからシステム物理アドレスに変換するように構成することができ、次いで、アダプター 402 は、内部の送信バッファーから内部の受信バッファーへデータをコピーすること及びその逆を行うことができる。

【0065】

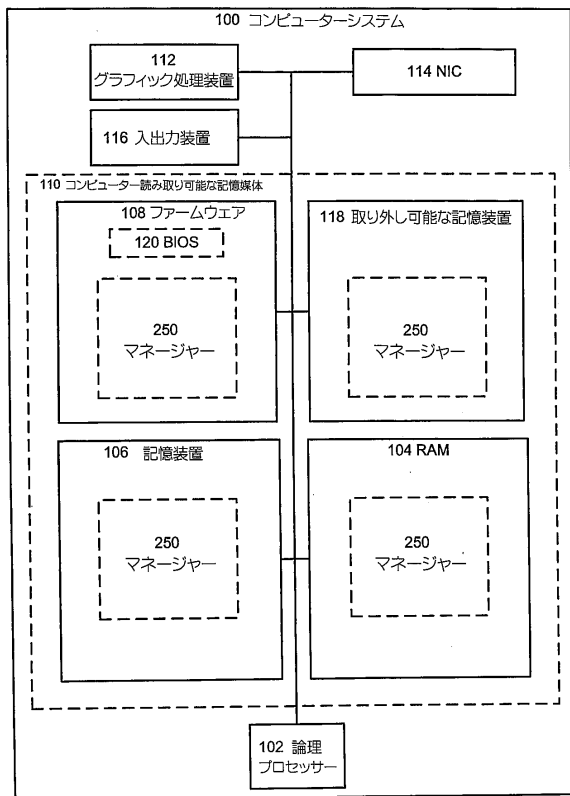
[0075] 前述の詳細な説明においては、例及び / 又は動作図によってシステム及び / 又はプロセスの様々な実施例を記載してきた。そのようなブロック図及び / 又は実施例が 1 つ又は複数の機能及び / 又は動作を含む限りにおいて、当業者には、そのようなブロック図や例における各々の機能及び / 又は動作を、ハードウェア、ソフトウェア、ファームウェア、又はそれらの実質的に任意の組み合わせの広い範囲によって、個別に及び / 又は集合的に実施することができることが理解されよう。

【0066】

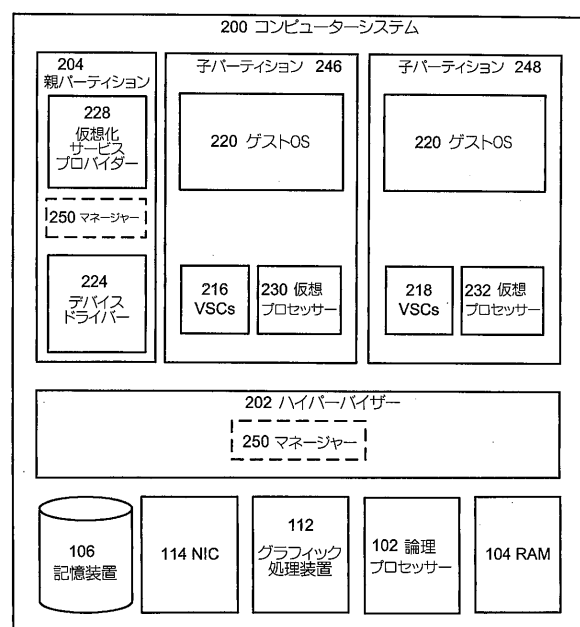
50

[0076]本明細書に記載された本願の主題の特定の態様が示され、記載されたが、本明細書の教示に基づいて、本明細書に記載された主題及びそのより広い態様を逸脱することなく変更及び修正をなすことができ、したがって、添付の特許請求の範囲は、本明細書に記載された主題の真の趣旨及び範囲内にあるようにすべてそのような変更及び修正をその範囲内に包含するものであることは、当業者にとって明らかであろう。

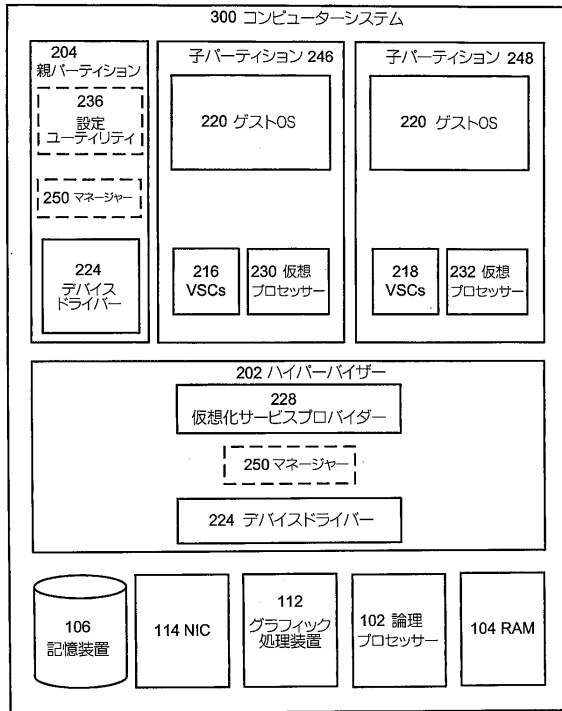
【図 1】



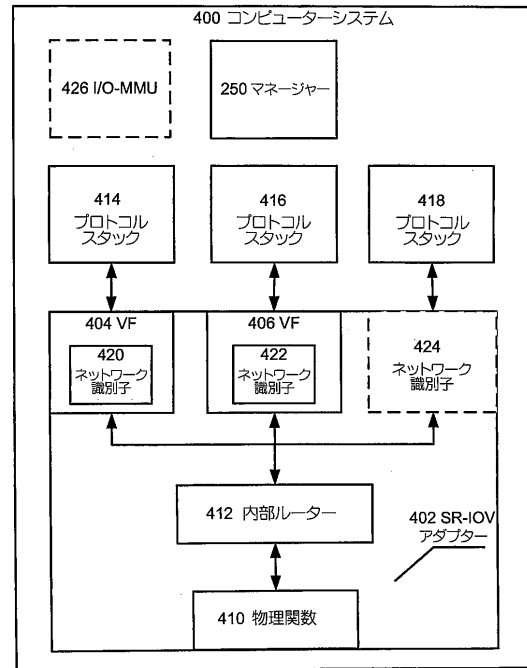
【図 2】



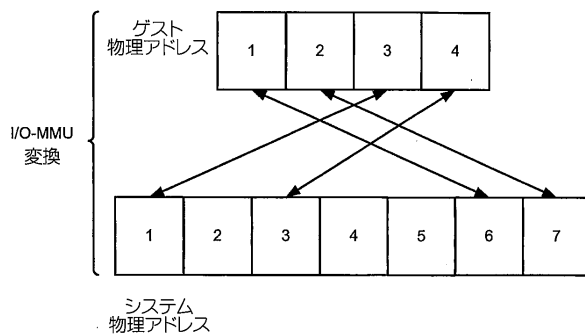
【図 3】



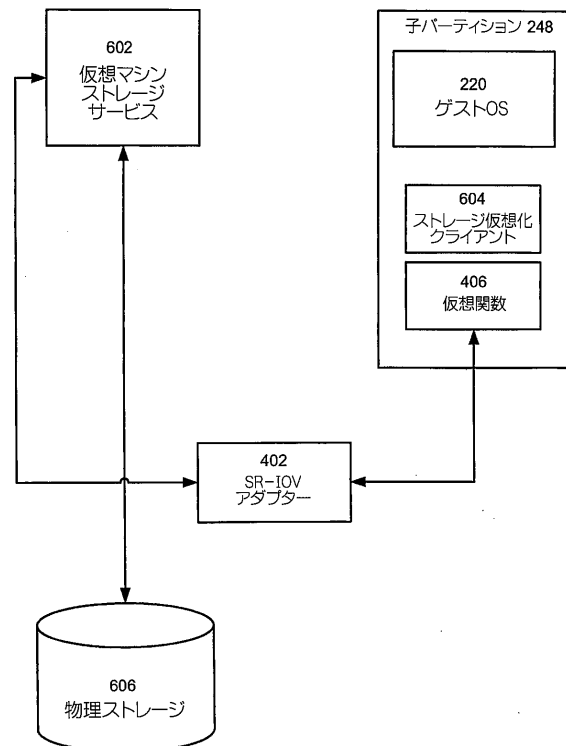
【図 4】



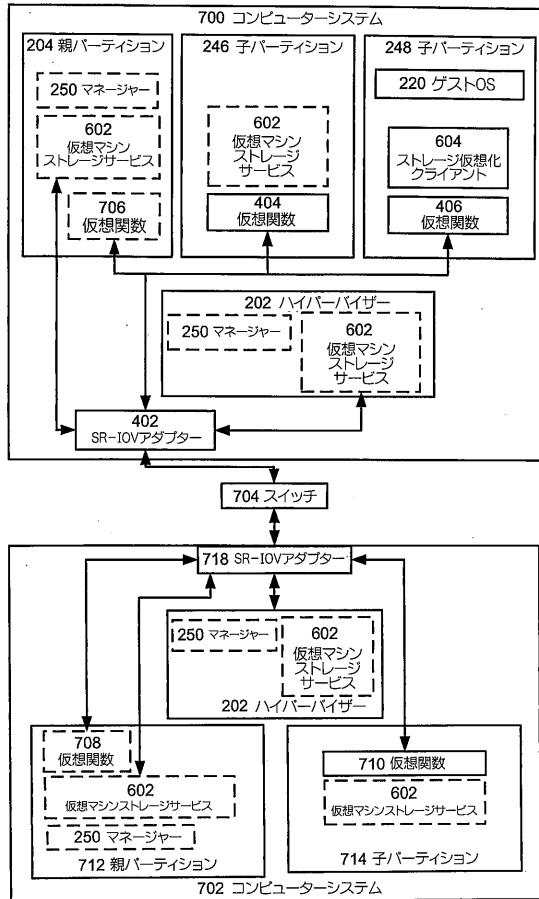
【図 5】



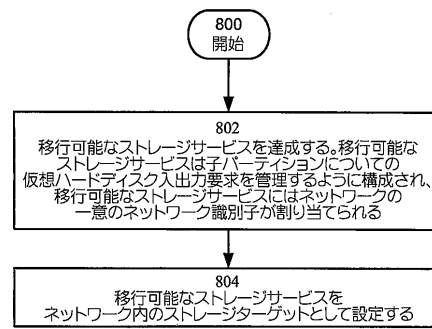
【図 6】



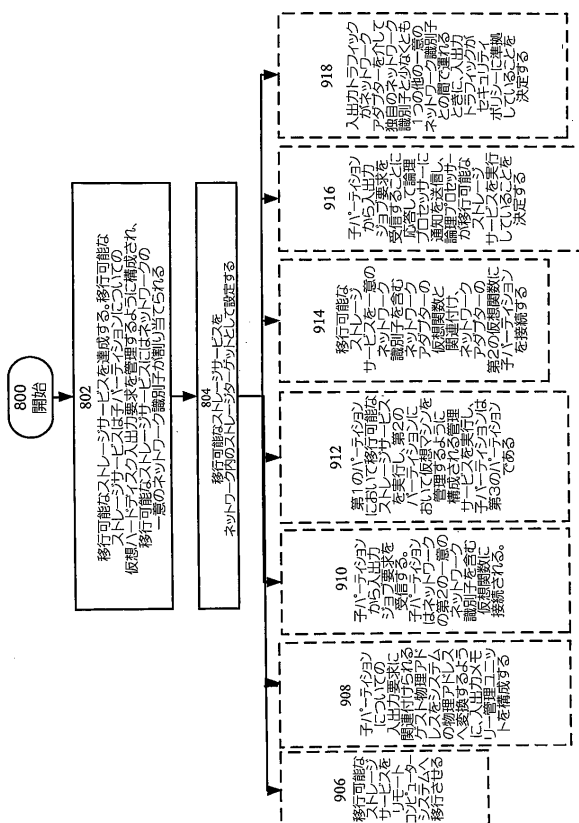
【図 7】



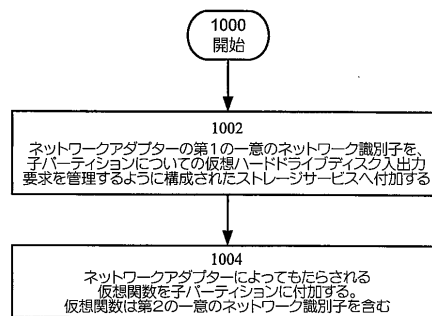
【図 8】



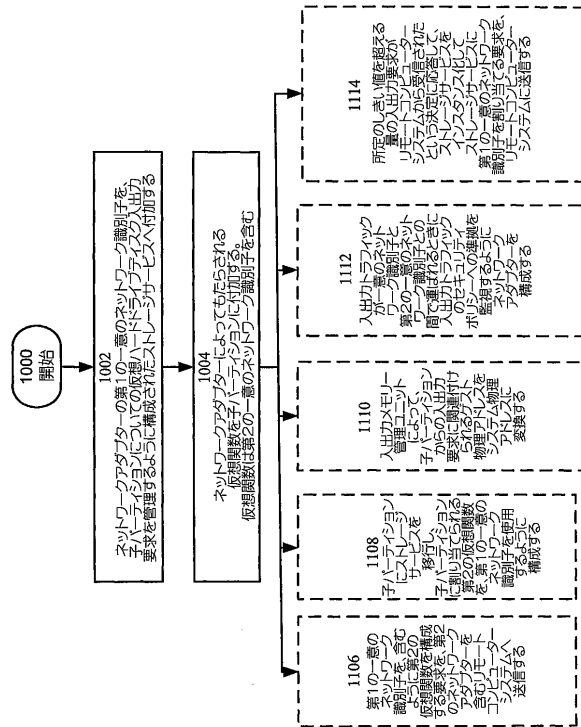
【図 9】



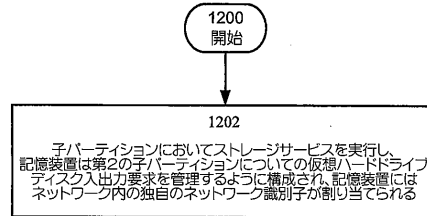
【図 10】



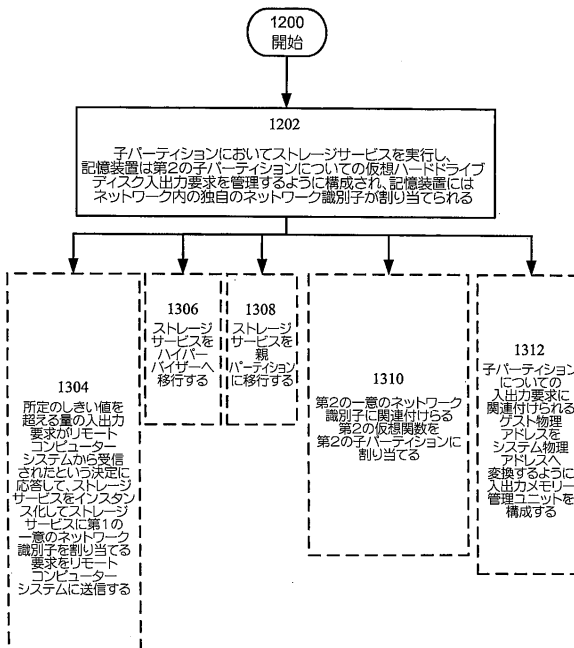
【図 1 1】



【図 1 2】



【図 1 3】



フロントページの続き

- (74)代理人 100153028
弁理士 上田 忠
- (74)代理人 100120112
弁理士 中西 基晴
- (74)代理人 100196508
弁理士 松尾 淳一
- (74)代理人 100147991
弁理士 鳥居 健一
- (74)代理人 100119781
弁理士 中村 彰吾
- (74)代理人 100162846
弁理士 大牧 綾子
- (74)代理人 100173565
弁理士 末松 亮太
- (74)代理人 100138759
弁理士 大房 直樹
- (72)発明者 オシンズ, ジェイコブ
アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9 , レッドモンド, ワン・マイクロソフト・ウェイ, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテント
- (72)発明者 グリーン, ダスティン・エル
アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9 , レッドモンド, ワン・マイクロソフト・ウェイ, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテント

審査官 井上 宏一

- (56)参考文献 特開 2 0 0 9 - 2 6 2 9 5 (J P , A)
特開 2 0 0 9 - 1 2 3 2 1 7 (J P , A)
米国特許出願公開第 2 0 0 9 / 0 0 3 7 9 4 1 (U S , A 1)
特開 2 0 0 9 - 2 5 9 1 0 8 (J P , A)
国際公開第 2 0 0 9 / 0 2 5 3 8 1 (W O , A 1)
特開 2 0 0 7 - 1 0 2 6 3 3 (J P , A)
特開 2 0 0 7 - 3 2 8 6 1 1 (J P , A)
特開 2 0 0 7 - 1 2 2 4 3 2 (J P , A)
特開 2 0 0 4 - 1 3 4 5 4 (J P , A)
特開 2 0 1 0 - 9 3 9 6 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 9 / 4 6 - 9 / 5 4