

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5869135号
(P5869135)

(45) 発行日 平成28年2月24日 (2016. 2. 24)

(24) 登録日 平成28年1月15日 (2016. 1. 15)

(51) Int. Cl.	F I
G 0 6 F 13/36 (2006. 01)	G 0 6 F 13/36 3 1 0 E
G 0 6 F 13/38 (2006. 01)	G 0 6 F 13/38 3 4 0 A
G 0 6 F 13/28 (2006. 01)	G 0 6 F 13/38 3 5 0
	G 0 6 F 13/28 3 1 0 Z

請求項の数 20 (全 16 頁)

(21) 出願番号	特願2014-533274 (P2014-533274)	(73) 特許権者	593096712
(86) (22) 出願日	平成23年9月30日 (2011. 9. 30)		インテル コーポレーション
(65) 公表番号	特表2014-531685 (P2014-531685A)		アメリカ合衆国 9 5 0 5 4 カリフォル
(43) 公表日	平成26年11月27日 (2014. 11. 27)		ニア州 サンタ クララ ミッション カ
(86) 国際出願番号	PCT/US2011/054394		レッジ ブールバード 2 2 0 0
(87) 国際公開番号	W02013/048477	(74) 代理人	100107766
(87) 国際公開日	平成25年4月4日 (2013. 4. 4)		弁理士 伊東 忠重
審査請求日	平成26年5月13日 (2014. 5. 13)	(74) 代理人	100070150
			弁理士 伊東 忠彦
		(74) 代理人	100091214
			弁理士 大貫 進介
		(72) 発明者	マーゴ, ウィリアム, アール.
			アメリカ合衆国 5 1 8 2 2 - 9 2 2 8
			イリノイ州 シャンペーン コブルフィー
			ルド ロード 1 6 0 3

最終頁に続く

(54) 【発明の名称】 コプロセッサのためのダイレクト I/O アクセス

(57) 【特許請求の範囲】

【請求項 1】

リモートダイレクトメモリアクセス (R D M A) デバイスと ;
中央処理装置 (C P U) とメモリーを有する周辺機器と ;
前記 R D M A デバイスと前記周辺機器に対して通信可能に接続されたホストコントローラと、を含み、

前記ホストコントローラーは、

前記ホストコントローラーの中に含まれ、前記周辺機器のメモリー又は前記 R D M A デバイスのメモリーのうちの一つに対してマップされた仮想アドレスにおけるアクセスリクエストメッセージに係る通知を受け取り ;

前記 R D M A デバイスの前記メモリーに対してマップされた仮想アドレスにおける前記アクセスリクエストメッセージの受け取りに応じて、前記周辺機器からの前記 R D M A デバイスのメモリーに対する前記アクセスリクエストメッセージに基づき、R D M A アクセスリクエストを前記 R D M A デバイスに対して伝送し、前記 R D M A アクセスリクエストメッセージは仮想アドレス R D M A パラメーターを含み ; かつ

前記周辺機器の前記メモリーに対してマップされた仮想アドレスにおける前記アクセスリクエストメッセージの受け取りに応じて、アクセスリクエストメッセージを前記周辺機器の前記メモリーに対して伝送する、

ように構成されており、

前記 R D M A デバイスは、リクエストの完了後、前記周辺機器のメモリーにマップされ

10

20

た仮想アドレスに対して完了メッセージを送付し、前記仮想アドレスは、前記完了メッセージを保管するための完了キューを表している、
ことを特徴とするシステム。

【請求項 2】

前記 R D M A アクセスリクエストは、R D M A 読み出しリクエストを有し、

前記ホストコントローラーは、さらに；

前記 R D M A デバイスからの前記周辺機器によりリクエストされたデータを受け取り、かつ、

前記周辺機器の前記メモリーに対して、前記データを含んでいるデータメッセージを
伝送する、

10

請求項 1 に記載のシステム。

【請求項 3】

前記 R D M A アクセスリクエストは、R D M A 書き込みリクエストを有し、

前記ホストコントローラーは、さらに；

前記周辺機器のメモリーからの前記 R D M A デバイスによりリクエストされたデータ
を受け取り、かつ、

前記 R D M A デバイスに対して、前記データを含んでいるデータメッセージを伝送する、

請求項 1 に記載のシステム。

【請求項 4】

20

前記周辺機器の前記 C P U は、複数のプロセッサコアのうちの一つ、または、前記周辺
機器の中に含まれているマルチコアプロセッサを有する、

請求項 1 に記載のシステム。

【請求項 5】

前記周辺機器と前記 R D M A デバイスは、周辺コンポーネント相互接続エクスプレス (P C I e) リンクを介して、前記ホストコントローラーに対して通信可能に接続されている、

請求項 1 に記載のシステム。

【請求項 6】

前記 R D M A デバイスと前記ホストコントローラーは、P C I e ルートポートを介して
通信可能に接続されている、

30

請求項 5 に記載のシステム。

【請求項 7】

前記周辺機器は、グラフィックサブシステムを有し、かつ

前記 C P U は、グラフィックプロセッサユニット (G P U) を有する、

請求項 1 に記載のシステム。

【請求項 8】

前記アクセスリクエストメッセージは、インフィニバンド仕様に準じている、

請求項 1 に記載のシステム。

【請求項 9】

40

前記アクセスリクエストメッセージは、R D M A コンソーシアム仕様に準じている、

請求項 1 に記載のシステム。

【請求項 10】

リモートダイレクトメモリーアクセス (R D M A) デバイスのメモリーと周辺機器のメモリーのうちの一つに対してマップされた仮想アドレスにおけるアクセスリクエストメッセージに係る通知を受け取るステップであり、前記周辺機器は、さらに、中央処理装置 (C P U) とメモリーを有するステップと；

前記 R D M A デバイスの前記メモリーに対してマップされた仮想アドレスにおける前記
アクセスリクエストメッセージの受け取りに応じて；

R D M A パラメーターとして前記 R D M A デバイスの前記メモリーに対してマップさ

50

れた前記仮想アドレスを利用するステップと；

前記周辺機器からの前記 R D M A デバイスのメモリーに対する前記アクセスリクエストメッセージに基づき、R D M A アクセスリクエストを前記 R D M A デバイスに対して伝送するステップであり、前記 R D M A アクセスリクエストメッセージは仮想アドレス R D M A パラメーターを含んでいるステップと；

前記周辺機器の前記メモリーに対してマップされた仮想アドレスにおける前記アクセスリクエストメッセージの受け取りに応じて、

アクセスリクエストメッセージを前記周辺機器の前記メモリーに対して伝送するステップと、を含み、

前記 R D M A デバイスは、リクエストの完了後、前記周辺機器のメモリーにマップされた仮想アドレスに対して完了メッセージを送付し、前記仮想アドレスは、前記完了メッセージを保管するための完了キューを表している、

10

ことを特徴とする方法。

【請求項 1 1】

前記 R D M A アクセスリクエストは、R D M A 読み出しリクエストを有し、

前記方法は、さらに；

前記 R D M A デバイスからの前記周辺機器によりリクエストされたデータを受け取るステップと、

前記周辺機器の前記メモリーに対して、前記データを含んでいるデータメッセージを伝送するステップと、を含む、

20

請求項 1 0 に記載の方法。

【請求項 1 2】

前記 R D M A アクセスリクエストは、R D M A 書き込みリクエストを有し、

前記方法は、さらに；

前記周辺機器のメモリーからの前記 R D M A デバイスによりリクエストされたデータを受け取るステップと、

前記 R D M A デバイスに対して、前記データを含んでいるデータメッセージを伝送するステップと、を含む、

請求項 1 0 に記載の方法。

【請求項 1 3】

30

前記周辺機器の前記 C P U は、複数のプロセッサコアのうちの一つ、または、前記周辺機器の中に含まれているマルチコアプロセッサを有する、

請求項 1 0 に記載の方法。

【請求項 1 4】

前記周辺機器と前記 R D M A デバイスは、周辺コンポーネント相互接続エクスプレス (P C I e) リンクを介して、ホストコントローラーに対して通信可能に接続されている、

請求項 1 0 に記載の方法。

【請求項 1 5】

前記 R D M A デバイスと前記ホストコントローラーは、P C I e ルートポートを介して通信可能に接続されている、

40

請求項 1 4 に記載の方法。

【請求項 1 6】

前記周辺機器は、グラフィックサブシステムを有し、かつ

前記 C P U は、グラフィックプロセッサユニット (G P U) を有する、

請求項 1 0 に記載の方法。

【請求項 1 7】

リモートダイレクトメモリアccess (R D M A) デバイスに対して通信可能に接続された第 1 の相互接続リンクと；

中央処理装置 (C P U) とメモリーを有する周辺機器に対して通信可能に接続された第 2 の相互接続リンクと；

50

ホストコントローラーと、を含み、
前記ホストコントローラーは、

前記ホストコントローラーの中に含まれ、前記周辺機器のメモリー又は前記 R D M A デバイスのメモリーのうちの一つに対してマップされた仮想アドレスにおけるアクセスリクエストメッセージに係る通知を受け取り；

前記 R D M A デバイスの前記メモリーに対してマップされた仮想アドレスにおける前記アクセスリクエストメッセージの受け取りに応じて、前記周辺機器からの前記 R D M A デバイスのメモリーに対する前記アクセスリクエストメッセージに基づき、R D M A アクセスリクエストを前記 R D M A デバイスに対して伝送し、前記 R D M A アクセスリクエストメッセージは仮想アドレス R D M A パラメーターを含み；かつ

前記周辺機器の前記メモリーに対してマップされた仮想アドレスにおける前記アクセスリクエストメッセージの受け取りに応じて、アクセスリクエストメッセージを前記周辺機器の前記メモリーに対して伝送し、

前記 R D M A デバイスは、リクエストの完了後、前記周辺機器のメモリーにマップされた仮想アドレスに対して完了メッセージを送付し、前記仮想アドレスは、前記完了メッセージを保管するための完了キューを表している、

ことを特徴とする装置。

【請求項 18】

前記 R D M A アクセスリクエストは、R D M A 読み出しリクエストを有し、

前記ホストコントローラーは、さらに；

前記 R D M A デバイスからの前記周辺機器によりリクエストされたデータを受け取り、かつ、

前記周辺機器の前記メモリーに対して、前記データを含んでいるデータメッセージを伝送する、

請求項 17 に記載の装置。

【請求項 19】

前記 R D M A アクセスリクエストは、R D M A 書き込みリクエストを有し、

前記ホストコントローラーは、さらに；

前記周辺機器のメモリーからの前記 R D M A デバイスによりリクエストされたデータを受け取り、かつ、

前記 R D M A デバイスに対して、前記データを含んでいるデータメッセージを伝送する、

請求項 17 に記載の装置。

【請求項 20】

前記第 1 および第 2 の相互接続リンクのそれぞれは、周辺コンポーネント相互接続エクスプレス (P C I e) リンクを有しており、かつ、

前記装置は、さらに、前記 R D M A デバイスに対して通信可能に接続されたルートポートを含んでいる、

請求項 17 に記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施例は、一般的に、コンピューターデバイスに適用される。より特定的には、周辺機器のためのメモリーアクセス管理に関する。

【背景技術】

【0002】

コンピューターシステムは、ネットワーク又はバス構成といった相互接続構造を介してお互いに接続されている種々のデバイスを含んでいる。これらのデバイスは、典型的にはローカルメモリーを有しており、コンピューター環境の中で処理速度と適応性を提供するために複数のデバイスが並行して動作される。

10

20

30

40

50

【 0 0 0 3 】

リモートダイレクトメモリアクセス (R D M A) は、ネットワークインターフェイスカード (N I C) の機能であり、コンピューターデバイスが、別のコンピューターデバイスのメモリの情報にアクセスできるようにする。特に、R D M A 技術を介して、コンピューターデバイスは、ホストオペレーティングシステム (O S) を巻き込むことなく、別のコンピューターデバイスのメモリから情報を読み出すことができる。別のコンピューターデバイスのメモリに情報を書き込むことも同様である。

【 0 0 0 4 】

図 1 は、C P U とメモリの複合物を有している周辺機器を含んだ従来技術に係るシステムを図示している。システム 1 0 0 は、システム C P U 1 1 0、システムメモリ 1 2 0、周辺機器コントローラー 1 3 0、周辺機器 1 4 0、および、R D M A デバイス 1 5 0 を含んでいる。周辺機器 1 4 0 は、プロセッサ 1 4 1 及びメモリ 1 4 2 を含んでいる。周辺機器 1 4 0 及び R D M A デバイス 1 5 0 は、「ピア (" p e e r ") 」デバイスとして参照されてよい。

10

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 5 】

周辺機器 1 4 0 は、R D M A 1 5 0 の中に保管されているデータにアクセスする必要があり、その逆もまた同様である。現在、インフィニバンド (I n f i n i B a n d) 仕様又は R D M A コンソーシアム仕様、等といった多くの相互接続構造規格の実施は、ピアデバイスが、別のピアデバイスのアドレス空間において保管されているデータに直接的にアクセスすることをできなくしている。

20

【 0 0 0 6 】

現在のソリューションは、ピアデバイスが、リクエストされたデータを共通に利用可能なメモリに対して書き込むことを要求する。この説明図においては、システムメモリ 1 2 0 であり、相互接続性構造に接続されたあらゆる周辺機器によってアクセス可能である。しかしながら、そうしたデータ転送のために共通のシステムメモリを使用することは時間の浪費であり、処理のオーバーヘッドを増加させるものである。さらに、共通システムのメモリを使用することは、周辺機器の処理動作を遅くしてしまう。

【 課題を解決するための手段 】

30

【 0 0 0 7 】

本発明の実施例は、ホストと、C P U 及びメモリの複合体を含む周辺機器 (代替的には、ここにおいてプロセッサアドインカードとして言及されるもの) との間で、リモートダイレクトメモリアクセス (R D M A) デバイスハードウェアを共有することができるシステム、装置、および、方法を説明する。本発明の実施例は、プロセッサアドインカード (a d d - i n c a r d) と R D M A デバイスとの間でのピアツ - ピア (p e e r - t o - p e e r) のデータ転送のための周辺コンポーネント相互接続エクスプレス (P C I e) といった相互接続ハードウェアを利用する。ホストシステムは、メモリ及びレジスターを R D M A デバイスに対して、及び/又は、R D M A デバイスからマップするためのモジュール又はロジックを含んでよい。これにより、ホストシステムの I / O オペレーティングシステムと同時に、アドインカードのプロセッサ上のユーザーモードアプリケーションに対して、または、アプリケーションから直接的に実行されるべき I / O をできるようにしている。

40

【 図面の簡単な説明 】

【 0 0 0 8 】

以降の記述は、本発明の具体化に係る実施例によって与えられる例解を有する図面の説明を含んでいる。図面は、例示として理解されるべきであり、本発明を限定するものとして理解されるべきではない。ここにおいて使用されるように、一つまたはそれ以上の「実施例」は、本発明に係る少なくとも一つの実施に含まれる所定の機能、構成、または、特性を表しているものとして理解されるべきである。このように、ここにおいて表れている

50

「一つの実施例において」または「代替的な実施例において」といったフレーズは、本発明に係る種々の具体化及び実施を説明するものであり、必ずしも全てが同一の実施例を参照する必要はない。しかしながら、それらはお互いに排他的であることも要しない。

【図１】図１は、ＣＰＵ及びメモリーの複合体を有する周辺機器を含む従来技術のシステムを示している。

【図２】図２は、本発明の一つの実施例に従ったブロックダイアグラムである。

【図３】図３は、本発明の一つの実施例に従ったシステムコンポーネントのブロックダイアグラムである。

【図４】図４は、本発明の一つの実施例に従ったプロセスのフローチャートである。

【図５】図５は、本発明の一つの実施例に従ったプロセスのフローチャートである。

【図６】図６は、本発明の一つの実施例に従ったホストとアドインカードモジュールのブロックダイアグラムである。

【図７】図７は、本発明の一つの実施例を利用するシステムのブロックダイアグラムである。

【０００９】

所定の詳細及び実施例に係る説明が後に続く。図面の説明を含むものであり、図面は以降に記述される実施例のいくつか又は全てを表している。他の可能性のある具体化、または、ここにおいて表される本発明の概念に係る実施例の説明も同様に後に続く。本発明の実施例の概要が以下に提供され、図面に関してより詳細な説明が続く。

【発明を実施するための形態】

【００１０】

本発明の具体例を利用するプロセッサアドインカードの実施例は、グラフィックプロセッサユニット（ＧＰＵ）を有するグラフィックプロセッササブシステムと、高度に並列なアプリケーションのパフォーマンスを改善するための複数の小さく、低消費電力なプロセッサコア又はマルチコアプロセッサを有するプロセッサアドインカードを含んでいる。異種混在（heterogeneous）が加速されたコンピューター環境においては、しかしながら、共通で、規格ベースのプログラミング及びコミュニケーションモデルを提供することは難しいと証明されている。このことは、クラスター（cluster）において特に真実であり、全てのプロセッサからの効率的なコミュニケーションメカニズムを有することが望ましい。プロセッサが、プライマリーシステムのＣＰＵであるか、アドインカードのコプロセッサコアであるかにかかわらずである。共通で、規格ベースのプログラミングモデルは、アプリケーションの開発とメンテナンスを簡素化するだけでなく、システムを使用するためのより大きな適応性を与えてシステムのパフォーマンスの全ての利点を得ることができる。

【００１１】

クラスターは、一般的に、お互いに密接に動作するようにリンク又は相互接続されたコンピューターシステムのグループを参照するものである。それらは、多くの観点で一つのコンピューターを形成するようにリンク等されている。クラスターは、一般的に、単独のコンピューターによって提供されるものよりも大幅に改善されたパフォーマンス及び/又は可用性を提供する。クラスターは、また、典型的には、同等な速度と可用性のある単独のコンピューターよりも、コスト効率が良い。

【００１２】

大規模なクラスターシステムの構築に対する重要な見地は、相互接続である。相互接続は、システムの全てを一緒に接続する「ファブリック（"fabric"）」を含んでよい。システムをファブリックに対してインターフェイスするホストアダプターも同様である。クラスターは、インフィニバンド仕様又はＲＤＭＡコンソーシアム仕様と調和した相互接続を利用してよい。インフィニバンドは、本来高性能コンピューターにおいて使用されたスイッチファブリックコミュニケーションリンク（switched fabric communication link）であり、スケーラビリティと同様に、サービス品質とフェイルオーバー機能を提供する。インフィニバンド相互接続は、一般的に、よ

10

20

30

40

50

り小さなレイテンシー (l a t e n c y)、より高いバンド幅、および、改善された信頼性を提供する。

【 0 0 1 3 】

インフィニバンドといった R D M A アーキテクチャーは、メッセージ受渡しオペレーティングシステムのレイテンシーを削減し、かつ、バンド幅を増加させることによって、高パフォーマンスコンピューティング (H P C) クラスターアプリケーションのパフォーマンス改善において、非常に成功してきた。R D M A アーキテクチャーは、カーネルバイパス (k e r n e l - b y p a s s)、ダイレクトデータ配置を通じて、ネットワークインターフェイスをアプリケーションに対して大変近くに移動することによってパフォーマンスを改善し、かつ、アプリケーションの要求に合せるように I / O オペレーティングシステムをより広くコントロールすることができる。

10

【 0 0 1 4 】

R D M A アーキテクチャーは、ハードウェアにおいて、プロセス分離、保護、および、アドレス変換をすることができる。このアーキテクチャーは、ホストとコプロセッサアプリケーションが分離されたアドレスドメインにおいて実行されるコプロセッサコア環境に対して相性がいい。しかしながら、従来技術においては、取り付けられたプロセッサ (つまり、プロセッサアドインカード) に対して R D M A の利益は利用可能ではなかった。本発明の実施例は、R D M A アーキテクチャーの利益を取り付けられたプロセッサに対して直接的に提供し、図 1 に示すようなソリューションに対する必要性を取り除く。図 1 では、取り付けられたプロセッサに係る内外のコミュニケーションは、ホストメモリーの中への追加のデータコピーを招くことを要し、実質的に、メッセージレイテンシーと達成可能なバンド幅の両方に影響を与えている。

20

【 0 0 1 5 】

図 2 は、本発明の一つの実施例に従ったシステムのブロックダイアグラムである。システム 2 0 0 は、システム C P U 2 1 0、システムメモリー 2 2 0、プラットフォームコントローラーハブ (P C H) 2 3 0、周辺機器 2 4 0、および、R D M A デバイス 2 5 0 を含んでいる。この実施例において、周辺機器 2 4 0 は、プロセッサアドインカードであり、プロセッサ 2 4 1 とメモリー 2 4 2 を含んでいる。周辺機器 2 4 0 と R D M A デバイス 2 5 0 は、「ピア (" p e e r ")」デバイスとして参照されてよい。

30

【 0 0 1 6 】

周辺機器 2 4 0 は、R D M A デバイス 2 5 0 の中に保管されているデータにアクセスをリクエストし得る。周辺機器が、P C I e リンクを介して P C H 2 3 0 と通信可能に接続されているものと示されている一方で、R D M A デバイスは、P C I e ルートポート 2 3 1 を介して P C H 1 3 0 に対して動作可能に接続されているものと示されている。この実施例において、周辺機器 2 4 0 は、P C H を含んでおらず、従って、ホストチャンネルアダプター/ネットワークインターフェイスコントローラー (H C A / N I C) カードを R D M A デバイス 2 5 0 に専念させる能力を有していない (しかしながら、P C H 又はオンボードの R D M A デバイスを伴う内部 P C I e バスを含んでいる周辺機器も、また、以下に説明される本発明の実施例を利用し得ることが理解されるべきである)。P C I e 相互接続は、単独では、ピアデバイスが、別のピアデバイスのアドレス空間の中に保管されているデータにアクセスするようにはできない。

40

【 0 0 1 7 】

本発明の実施例は、システム 2 0 0、P C H 2 3 0、および、周辺機器 2 4 0 に含まれているロジック及び/又はモジュールについて説明する。周辺機器が、R D M A デバイス 2 5 0 の中に含まれているデータに対して直接的にアクセスできるようにするものである。つまり、システムメモリー 2 2 0 は、図 1 のシステム 1 0 0 といった従来技術のソリューションにおいて必要とされたようには必要とされない。リクエストされたデータをホストメモリーの中に受け取るための、ホストメモリーに対する必要性を取り除くことによって、本発明の実施例は、著しく、メッセージレイテンシーを減少させ、達成可能なバンド幅を増加させる。

50

【 0 0 1 8 】

図 3 は、本発明の一つの実施例に従ったシステムコンポーネントのブロックダイアグラムである。この実施例において、PCH 310、周辺機器 320、および、RDMA 330 は、システムバス 390（ここにおいては、代替的にイントラノード（intra-node）バスとして参照される）を介してお互いに通信可能に接続されている。周辺機器 320 は、図 2 に係る機器 240 と類似のものであり、RDMA デバイス 330 を専念させる能力を有していない。

【 0 0 1 9 】

この実施例において、周辺機器 320 のメモリー 321 はリクエストキュー（queue）322 を含んでいる。リクエストキューは、RDMA デバイス 330 のメモリー 331 に向けて指示されたリクエストを保管している。周辺機器の CPU 323 は、リクエストキューがメモリー 331 へのアクセスに対する未解決のリクエストを有していることを、RDMA デバイス 330 に通知する。一つの実施例において、RDMA デバイス 330 は、キュー 322 の中に未解決のリクエストがいくらかでもある場合に通知される。他の実施例において、RDMA デバイス 330 は、キュー 322 の中の未解決のリクエストの数量が閾値を超えた場合に通知される。

【 0 0 2 0 】

CPU 323 は、PCH 310 の中に含まれる仮想アドレス 311 に対して書き込むことによって、未解決のリクエストについて RDMA デバイス 350 に通知する。仮想アドレスは、RDMA デバイスのメモリー 331 に対してマップされたメモリーである。未解決のリクエストは、実行されるべきオペレーションを記述することができる。つまり、send、receive、write、read、atomic compare/exchange、atomic fetch/add、等である。このように、本発明の実施例は、周辺機器 320 からホストシステムの PCH 310 への「プロキシ（"proxy"）」リクエストに対して記述されてよい。

【 0 0 2 1 】

図 4 は、本発明の一つの実施例に従ったプロセスのフローチャートである。ここにおいて説明されるように、フローチャートは、種々のプロセス動作のシーケンスに係る実施例を提供する。所定のシーケンス又は順序において示されてはいるが、別に指定がなければ、動作の順序は変更することができる。従って、図示された実施例は、単なる例示として理解されるべきであり、図示されたプロセスは、異なる順序で実行することができ、いくつかの動作は並行に実行され得る。加えて、一つまたはそれ以上の動作は、本発明に種々の実施例において除外することができる。従って、全ての動作が、それぞれの実施例において必要とされるわけではない。他のプロセスフローも可能である。

【 0 0 2 2 】

プロセス 400 は、プロセッサコアとメモリーを有する周辺機器による、RDMA デバイスのメモリーへのアクセスに対するリクエストを受け取るオペレーションを含んでいる、410。そのリクエストは、周辺機器のメモリーにおけるキューの中に保管されてよい。いくつかの実施例において、周辺機器は複数のプロセッサコアを含んでおり、それぞれのプロセッサコアに対して分離されたキューを維持している。

【 0 0 2 3 】

周辺機器は、PCH の中に含まれている仮想アドレスに対してデータを送付することによって、キューの中のアクセスリクエストメッセージを RDMA デバイスに通知する、420。その仮想アドレスは、PCH の中に含まれており、RDMA デバイスのメモリーにマップされている。PCH は、RDMA デバイスにマップされた仮想アドレスを利用する、430。そして、RDMA デバイスのメモリーに対して RDMA アクセスリクエストメッセージを伝送する、440。その RDMA アクセスリクエストメッセージは、周辺機器の中にキューされたアクセスリクエストメッセージに基づくものである。

【 0 0 2 4 】

RDMA デバイスは、RDMA アクセスリクエストメッセージの受け取りに応じて、周

10

20

30

40

50

辺機器からのリクエストを完了する、450。いくつかの実施例において、RDMAデバイスは、PCHの中に含まれ、かつ、周辺機器のメモリーにマップされている仮想アドレスに対して完了メッセージを送付する。その仮想アドレスは、完了メッセージを保管するための完了キューを表してよい。例えば、RDMAアクセスリクエストメッセージがRDMA読み出しリクエストを有する場合、周辺機器によってリクエストされたデータを含んでいるデータメッセージが、PCHを介して周辺機器に対して送付される。RDMAアクセスリクエストメッセージがRDMA書き込みリクエストを有する場合、PCHは、書き込まれるべきデータを含んでいるデータメッセージをRDMAデバイスに対して伝送する。いくつかの実施例において、RDMA動作の完了は、RDMAデバイスによって完了メッセージがいつ送付されるかを決定するものではない。例えば、インフィニバンド仕様は、リクエストは、それらが投稿された順序で実行されることを要求する。先の全てのRDMA動作が完了するまで、周辺機器は完了メッセージを受け取らないことを意味するものである。

10

【0025】

図5は、本発明の一つの実施例に従ったプロセスのフローチャートである。プロセス500は、周辺機器（プロセッサコアも含んでいる）のメモリーへのアクセスに対するリクエストを受け取るオペレーションを含んでいる、510。そのリクエストは、RDMAデバイスのメモリーにおけるキューの中に保管されてよい。

【0026】

RDMAデバイスは、PCHの中に含まれている仮想アドレスに対してデータを送付することによって、キューの中のアクセスリクエストメッセージをホストシステムに通知する、520。その仮想アドレスは、周辺機器のメモリーにマップされている。PCHは、アクセスリクエストメッセージパラメーターとして、RDMAデバイスにマップされた仮想アドレスを利用する、530。そして、周辺機器のメモリーに対してアクセスリクエストメッセージを伝送する、540。そのアクセスリクエストメッセージは、RDMAデバイスの中にキューされたアクセスリクエストメッセージに基づいて、周辺機器に対して伝送される。

20

【0027】

周辺機器は、PCHからのアクセスリクエストメッセージの受け取りに応じて、RDMAデバイスからのリクエストを完了する。例えば、アクセスリクエストメッセージが読み出しリクエストを有する場合、RDMAデバイスによってリクエストされたデータを含んでいるデータメッセージが、PCHを介してRDMAデバイスに対して送付される。アクセスリクエストメッセージが書き込みリクエストを有する場合、RDMAデバイスは、PCHを介して、書き込まれるべきデータを含んでいるデータメッセージを周辺機器に対して伝送する。RDMAデバイスは、RDMA動作及び/又はパラメーターのタイプに応じて、PCHの中に含まれ、かつ、周辺機器のメモリーにマップされている仮想アドレスに対して完了メッセージを送付する、550。その仮想アドレスは、完了メッセージを保管するための完了キューを表してよい。

30

【0028】

図6は、本発明の一つの実施例に従ったホストとプロセッサアドインカードモジュールのブロックダイヤグラムである。上述のように、本発明の実施例は、インフィニバンド仕様（例えば、規格リリース1.0.a、2001年7月19日発行）に準じた相互接続を利用してよい。インフィニバンドは、本来高性能コンピューターにおいて使用されるスイッチファブリックコミュニケーションリンクであり、スケーラビリティと同様に、サービス品質とフェイルオーバー機能を提供する。インフィニバンド相互接続は、一般的に、より小さなレイテンシー、より高いバンド幅、および、改善された信頼性を提供する。インフィニバンドは、コンピューターシステムのコンポーネントの中、および、コンピューターの中で情報を移動する方法を提供する。インフィニバンドは、コンピューターのCPUが、非常に高いパフォーマンスでI/Oデバイス及び他のCPUと直接的にコミュニケーションできるようにする。インフィニバンド技術は、あらゆるネットワークに係るデータ

40

50

センターのバックエンド (back end) を対象としている。ネットワークインフラストラクチャーのフロントエンド及びミドルエンドは、典型的には、従来のイーサネット (登録商標) (Ethernet) 技術を含んでいる。別の言葉で言えば、インフィニバンド及びイーサネット技術は、両方とも同一のホストによって使用され得るものである。

【0029】

ホスト610と周辺機器カード630は、種々のインフィニバンドモジュール (以下に説明される) を含むように示されており、PCIe相互接続660を介してRDMAデバイス650と通信可能に接続されている。以下に説明するように、RDMAデバイス650によって受け取られるRDMAメッセージは、RDMAコンソーシアム仕様 (例えば、RDMAプロトコル規格 (バージョン1.0)、2002年10月21日発行) に準じてよい。RDMAコンソーシアムのRDMAプロトコルは、TCP/IPプロトコルにおけるTCP層の上で規定されている。従って、RDMA動作は、プロトコルスタックのトップから送信器側のボトムまで進み、次に、プロトコルスタックを越えて受信器側のトップまで進む。本発明の実施例によって利用されるRDMAプロトコルは、基本のTCP/IP処理ハードウェアを構築することを含んでよい。パケットを受け取り、TCP/IPを停止し、TCP/IPを通してアプリケーションに対してパケットを処理し、そして、データ及びメモリーに書き込むためのアプリケーション層でのアドレスを抽出することである。RDMAプロトコルは、データ転送をより効率的にする (特により大きなデータペイロードに対して) ことによって、パケットをメモリーにコピーすることを防ぐことができる (そして、その後データペイロードをメモリーにコピーする)。

【0030】

この実施例において、ホスト610上のモジュールと周辺機器630はお互いにコミュニケーションし、PCIe相互接続660をわたりRDMAデバイス650に対して直接的なアクセスを有している。モジュールは、RDMAデバイスのリソースを管理するためにPCIe相互接続660をわたりプロキシオペレーションに対するスプリットドライバー (split-driver) モデルを使用する。

【0031】

ホスト610は、インフィニバンドベースのソフトウェアコミュニケーションスタックを含むように図示されている。メッセージパッシングインターフェイス (MPI) アプリケーション611、RDMA APIであるユーザーモードダイレクトアクセスプロバイダーライブラリー (uDAPL) 612、IBバーブ (verb) (つまり、機能) ライブラリー613、ベンダーライブラリー614、IBユーバーブ (uverb) 635、および、IBコア636を含んでいる。ホストと周辺機器は、さらに、以下に説明されるモジュールとして実行される本発明の実施例を利用する。

【0032】

ホスト610は、IBプロキシデーモン (daemon) 618を含んでいる。IBプロキシデーモンは、ホストユーザーモードのアプリケーションであり、基底にあるベンダーライバー617に対するコール (call) のためにIBプロキシサーバー619 (以下に説明される) に対してユーザーモードプロセスコンテキスト (context) を提供する。ユーザーモードプロセスコンテキストは、ベンダーライバー617を変更することなくRDMAデバイス650のメモリーの仮想アドレスマッピングを実行するために使用され得る。

【0033】

ホスト610は、さらに、IBプロキシサーバー619を含んでおり、IBプロキシサーバーはホストカーネルモジュールを有している。この実施例において、そのIBプロキシサーバーは、コミュニケーションと周辺機器630 (以下に説明される) のIBプロキシクライアント638のためのコマンドサービスを提供する。この実施例において、IBプロキシサーバー619は、クライアント接続をリスン (listen) し、RDMAデバイスの追加、除去、および、イベント通知メッセージをリレーする。IBプロキシサーバー619は、さらに、IBプロキシクライアント638のためにIBコ

アレイヤー 6 1 6 に対するカーネルモード I B バープコールを開始して、その結果を戻すことができる。

【 0 0 3 4 】

周辺機器 6 3 0 は、I B プロキシークライアント 6 3 8 を含んでおり、I B プロキシークライアントはカーネルモジュールを有している。その I B プロキシークライアントは、ホスト 6 1 0 上でカーネルモード I B バープを実行するために、ベンダープロキシードライバー 6 3 7 (以下に説明される) に対してプログラミングインターフェイスを提供する。インターフェイスは、さらに、コマンドのフォーマット及びコミュニケーションの実行に係る詳細を抽出し得る。I B プロキシークライアント 6 3 8 は、所定のデバイスの追加、除去、および、ベンダープロキシードライバー 6 3 7 へのイベント通知に対するコールバック (c a l l b a c k) を呼び出す。

10

【 0 0 3 5 】

周辺機器 6 3 0 は、さらに、ベンダープロキシードライバー 6 3 7 を含んでおり、ベンダープロキシードライバーはカーネルモジュールを有している。所定の R D M A デバイスをサポートするために異なるベンダープロキシードライバーが使用されてよい。そのベンダープロキシードライバーそれぞれは、R D M A のデバイスの追加、除去、および、所定の P C I e ドライバーからのイベント通知を I B プロキシークライアント 6 3 8 に登録し得る。ベンダープロキシードライバー 6 3 7 は、カーネルモード I B バープコールを実行するために、I B プロキシークライアント 6 3 8 によって提供されるプログラミングインターフェイスを使用してよい。そのベンダープロキシードライバーは、さらに、ベンダーライブラリー 6 3 4 とホスト 6 1 0 上のベンダードライバー 6 1 7 との間で共有されるあらゆるプライベートデータの解釈及び変換を取り扱う。

20

【 0 0 3 6 】

この実施例において、ホスト 6 1 0 と周辺機器 6 3 0 の両方は、シンメトリックコミュニケーションインターフェイス (S C I F) モジュール 6 2 0 と 6 4 0 を、それぞれに含んでいる。その S C I F モジュールは、シングルプラットフォームの中でイントラノード (i n t r a - n o d e) コミュニケーションのためのメカニズムを提供する。S C I F は、ホスト 6 1 0 と周辺機器 6 3 0 との間で対称な A P I を提供する一方で、P C I e にわたるコミュニケーション (および、関連する周辺機器ハードウェアのコントロール) の詳細を抽出する。

30

【 0 0 3 7 】

上述のモジュールに加えて、本発明の実施例は、I B コアレイヤー 6 1 6 におけるコールを利用して、ベンダーライブラリー 6 1 4 及び 6 3 4 とベンダードライバー 6 1 7 との間でプライベートデータを転送する。メモリーを R D M A デバイス 6 5 0 に対してマッピングすることも同様である。

【 0 0 3 8 】

上記の「バープ (" v e r b ") 」 (つまり、機能) は、R D M A デバイス 6 5 0 に向けた P D M A オペレーションを実行する。バープは、特権クラスと非特権クラスに分類されてよい。特権バープは、典型的に R D M A ハードウェアのリソースを割り当て、管理するために使用され、ベンダードライバー 6 1 7 によって実施される。周辺機器 6 3 0 上で稼働しているアプリケーションのために、これらの特権バープが、ベンダープロキシードライバー 6 3 7 を通じてホスト 6 1 0 上のベンダードライバー 6 1 7 に対してフォワードされる。一旦、ハードウェアが割り当てられ開始されると、非特権バープは、カーネルをバイパスし、リソース割り当ての最中にアプリケーションアドレス空間の中にマップされたメモリーを使用して、ユーザーモードからハードウェアに対する直接的なアクセスを許可する。同様に、R D M A デバイスは、キューにアクセスすることができ、プロセスアドレス空間へ、または、プロセスアドレス空間から直接的にデータ転送を実行することができる。このように、本発明の実施例は、周辺機器 6 3 0 上のクライアントプロセスを、まるでホスト 6 1 0 上の別の「ユーザーモード」プロセスであるかのようにする。

40

【 0 0 3 9 】

50

従って、上記のモジュールにより、ホスト 610 は、RDMA デバイス 650 のメモリーにアクセスするための、周辺機器 630 のプロセッサコアからのアクセスリクエストメッセージに係る通知を受け取ることができる。その通知は、ホスト 610 及び RDMA デバイス 650 のメモリーにマップされたメモリー（例えば、図 3 に示されるように）の中に含まれる仮想アドレスにおいて受け取られる。ホストデバイス 610 は、そのアクセスリクエストメッセージに基づいて、RDMA デバイスに対して RDMA アクセスリクエストを転送する。ここで、リクエストは仮想アドレス RDMA パラメーターを含んでいる。従って、データリクエストのタイプが何であれ、アクセスリクエストメッセージは、例えば、send、receive、write、read、atomic compare/exchange、atomic fetch/add、等を含んでいる。データリクエストは、RDMA デバイス 650 に対して、周辺機器 630 からというよりむしろ、まるでホスト 610 上の「ユーザーモード」から生じたものであるように見える。

10

【0040】

図 7 は、本発明の一つの実施例を利用するシステムのブロックダイアグラムである。システム 700 は、サーバプラットフォームを記述しているが、例えば、以下のものに含まれてのよい。デスクトップコンピューター、ラップトップコンピューター、タブレットコンピューター、ネットブック、ノートブックコンピューター、パーソナルデジタルアシスタント (PDA)、サーバー、ワークステーション、携帯電話、モバイルコンピューター機器、インターネット機器、MP3 又はメディアプレーヤー、または、あらゆる他のタイプのコンピューターデバイス、である。

20

【0041】

システム 700 は、システムバス 720 を介して、データを交換するためのプロセッサ 710、ユーザーインターフェイス 760、システムメモリー 730、周辺機器コントローラー 740、および、ネットワークコネクタ 750 を含んでいる。その周辺機器コントローラーは、周辺機器及び RDMA デバイスとコミュニケーション可能に接続され、上記の本発明の実施例のいずれかに従ってデバイス関連情報の I/O リクエストを管理する。

【0042】

システム 700 は、さらに、システム 700 の種々のエレメントによって処理されるべき信号を送信及び受信するためのアンテナと RF 回路 770 を含んでいる。上記のアンテナは、指向性アンテナ又は無指向性アンテナであってよい。ここにおいて使用されるように、無指向性アンテナという用語は、少なくとも一平面において実質的に均一なパターンを有するあらゆるアンテナを参照するものである。例えば、いくつかの実施例において、アンテナは、ダイポール (dipole) アンテナ、または、四分の一波長 (quarter wave) アンテナといった、無指向性アンテナであってよい。例えば、いくつかの実施例において、アンテナは、パラボラアンテナ、パッチアンテナ、または、八木アンテナといった、指向性アンテナであってもよい。いくつかの実施例において、システム 700 は、複数の物理的なアンテナを含んでよい。

30

【0043】

ネットワークコネクタ 750 から離れているように示されているが、他の実施例において、アンテナ及び RF 回路 770 は、無線インターフェイスを含んでよいことが理解されるべきである。無線インターフェイスは、これらに限定されるわけではないが、IEEE 802.11 規格及び関連ファミリー規格、Home Plug AV (HPAV)、ウルトラワイドバンド (UWB)、Bluetooth (登録商標)、WiMax、または、無線通信プロトコルの他の形式に従って動作するものである。

40

【0044】

ここにおいて、プロセス、サーバー、または、ツールとして説明され、上記に参照された種々のコンポーネントは、説明された機能を実行するための手段であってよい。ここにおいて説明されたそれぞれのコンポーネントは、ソフトウェア又はハードウェア、または、これらの組み合わせを含んでいる。それぞれ及び全てのコンポーネントは、ソフトウェ

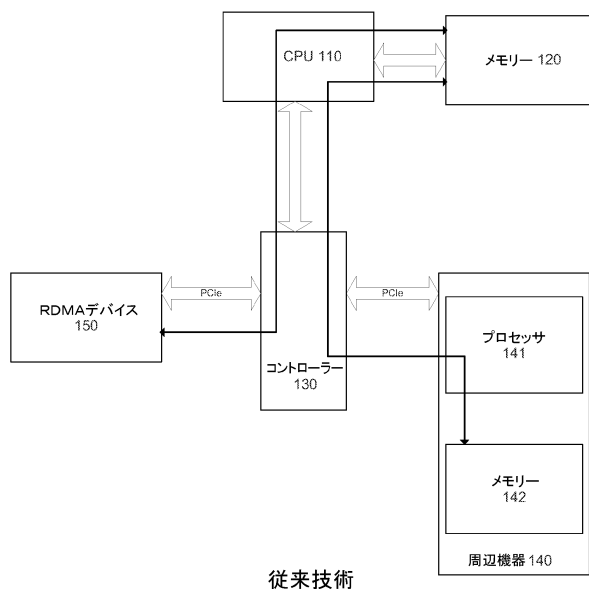
50

アモジュール、ハードウェアモジュール、特定用途ハードウェア（例えば、アプリケーション所定のハードウェア、ASIC、DSP、等）、エンベッドコントローラー、ハードワイヤード回路、ハードウェアロジック、等として実施されてよい。ソフトウェアコンテンツ（例えば、データ、インストラクション、コンフィグレーション）は、固定で有形なコンピューター又はマシンで読取り可能な媒体を含む製品を介して提供されてよい。実行され得るインストラクションを表すコンテンツを提供するものである。コンテンツは、結果として、ここにおいて説明された種々の機能／動作のコンピューターによる実行を生じる。

【0045】

コンピューターで読取り可能な固定記録媒体は、コンピューター（例えば、計算機器、電子システム、等）によってアクセス可能な形式で情報を提供（つまり、保管及び/又は送信）するあらゆるメカニズムを含んでいる。記録可能／記録不能媒体（例えば、読み出し専用メモリー（ROM）、ランダムアクセスメモリー（RAM）、磁気ディスク記録媒体、光記録媒体、フラッシュメモリーデバイス、等）といったものである。コンテンツは、直接的に実行可能な（「オブジェクト」又は「実行可能」形式）ソースコード、または、異なるコード（「デルタ」又は「パッチ」コード）であってよい。コンピューターで読取り可能な固定記録媒体は、また、ストレージ又はデータベースを含んでよく、そこからコンテンツをダウンロードすることができる。コンピューターで読取り可能な媒体は、また、販売時又は引き渡し時に、媒体上にコンテンツが保管されているデバイス又はプロダクトを含んでよい。従って、コンテンツが保管されたデバイスを引き渡すこと、または、コミュニケーション媒体にわたるダウンロードのためにコンテンツを提供することは、ここにおいて説明されたようなコンテンツを伴う製品を提供することとして理解されてよい。

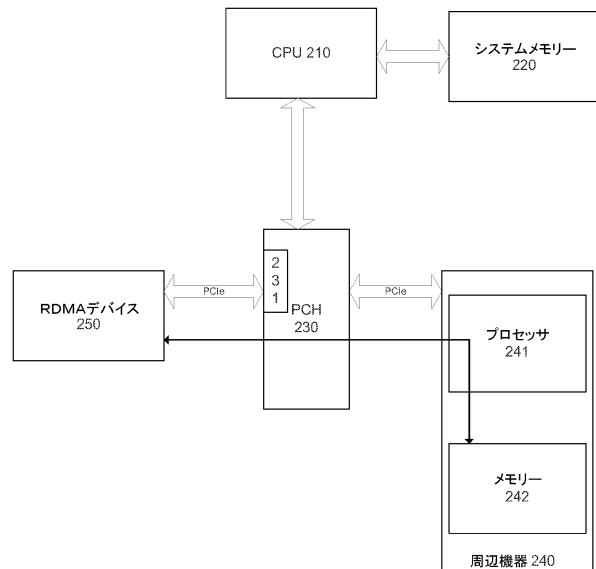
【図1】



従来技術

100

【図2】

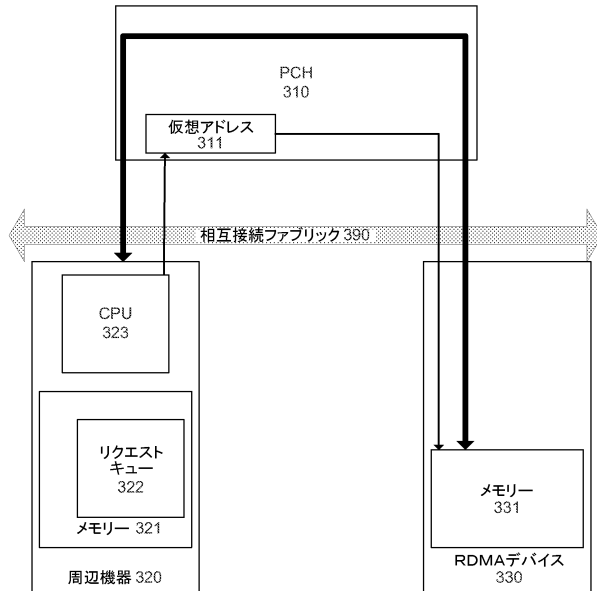


200

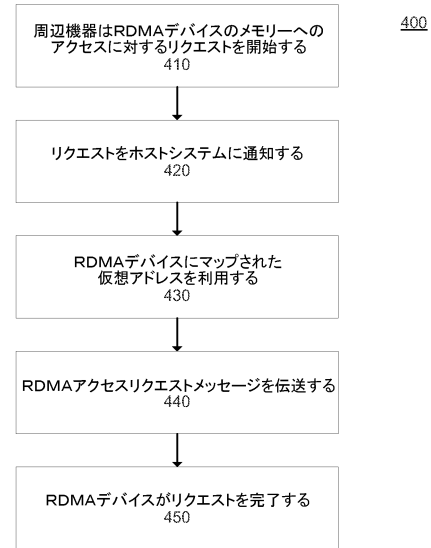
10

20

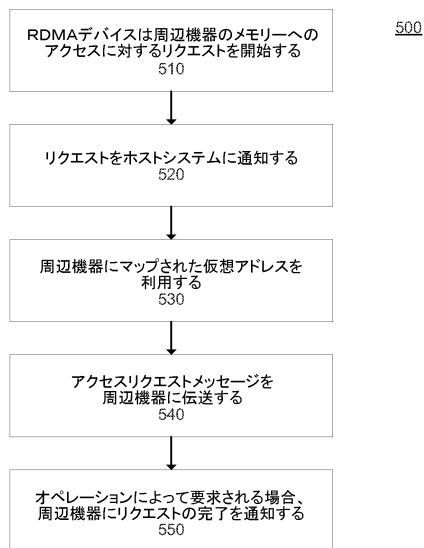
【図 3】



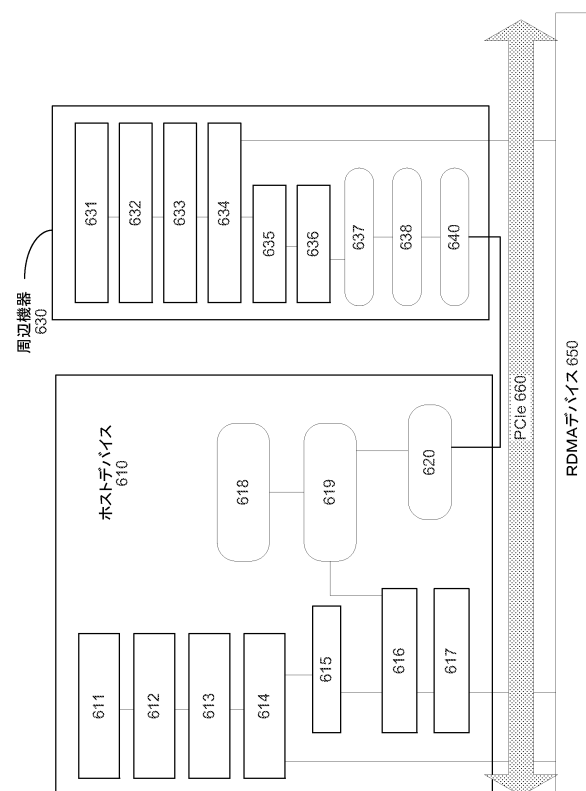
【図 4】



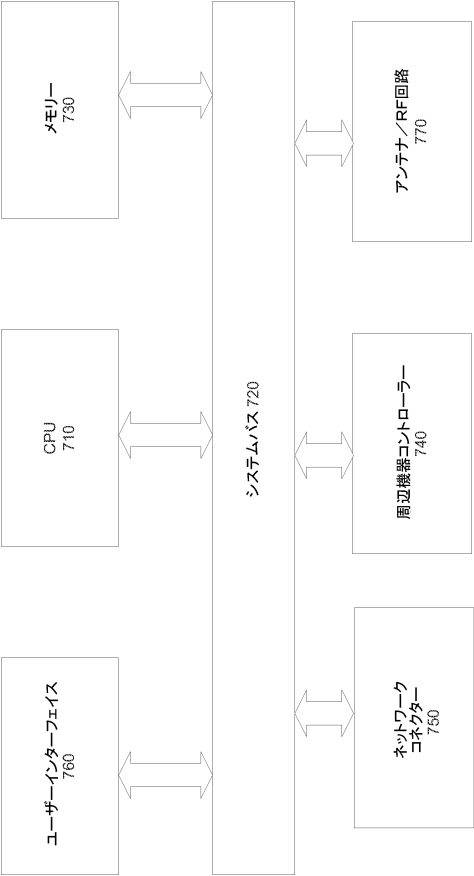
【図 5】



【図 6】



【図 7】



700

フロントページの続き

- (72)発明者 ウッドラフ, ロバート, ジェイ.
アメリカ合衆国 97106 オレゴン州 バンクス ノースウエスト グランドビュー 139
83
- (72)発明者 リー, デイヴィッド, エム.
アメリカ合衆国 97229 オレゴン州 ポートランド ノースウエスト クレストビュー ウ
エイ 2203
- (72)発明者 デイヴィス, アーリン, アール.
アメリカ合衆国 97148 オレゴン州 ヤムヒル ノースイースト コーヴ オーチャード
ロード 23190
- (72)発明者 ヘフティ, マーク, ショーン
アメリカ合衆国 97007 オレゴン州 アロハ サウスウエスト ハート ロード 1856
0
- (72)発明者 コフマン, ジェリー, エル.
アメリカ合衆国 97124 オレゴン州 ヒルズボロ ノースイースト 11ス コート 31
53

審査官 寺谷 大亮

- (56)参考文献 米国特許第07711793 (US, B1)
米国特許出願公開第2004/0034725 (US, A1)

- (58)調査した分野(Int.Cl., DB名)
- | | |
|------|-------|
| G06F | 13/36 |
| G06F | 13/28 |
| G06F | 13/38 |