

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6492083号
(P6492083)

(45) 発行日 平成31年3月27日(2019.3.27)

(24) 登録日 平成31年3月8日(2019.3.8)

(51) Int.Cl.

F I

G O 6 F 13/38 (2006.01)

G O 6 F 13/38 3 2 0 Z

G O 6 F 13/28 (2006.01)

G O 6 F 13/28 3 1 0 A

G O 6 F 13/10 (2006.01)

G O 6 F 13/10 3 1 0 A

G O 6 F 13/00 (2006.01)

G O 6 F 13/00 3 5 1 A

H O 4 L 12/061 (2013.01)

G O 6 F 13/38 3 5 0

請求項の数 11 (全 29 頁) 最終頁に続く

(21) 出願番号 特願2016-536668 (P2016-536668)
 (86) (22) 出願日 平成26年10月21日(2014.10.21)
 (65) 公表番号 特表2017-501492 (P2017-501492A)
 (43) 公表日 平成29年1月12日(2017.1.12)
 (86) 国際出願番号 PCT/US2014/061640
 (87) 国際公開番号 W02015/084506
 (87) 国際公開日 平成27年6月11日(2015.6.11)
 審査請求日 平成29年10月11日(2017.10.11)
 (31) 優先権主張番号 14/096,987
 (32) 優先日 平成25年12月4日(2013.12.4)
 (33) 優先権主張国 米国(US)
 (31) 優先権主張番号 14/096,949
 (32) 優先日 平成25年12月4日(2013.12.4)
 (33) 優先権主張国 米国(US)

(73) 特許権者 502303739
 オラクル・インターナショナル・コーポレ
 イション
 アメリカ合衆国カリフォルニア州9406
 5レッドウッド・シティ、オラクル・パ
 ークウェイ500
 (74) 代理人 110001195
 特許業務法人深見特許事務所
 (72) 発明者 アガーワル、ウッタム
 アメリカ合衆国、94583 カリフォル
 ニア州、サン・ラモン、サマークリーク・
 レーン、420

審査官 田名網 忠雄

最終頁に続く

(54) 【発明の名称】 インフィニバンド (I B) 上で仮想ホストバスアダプタ (vHBA) を管理およびサポートする
 ためのシステムおよび方法、ならびに単一の外部メモリインターフェイスを用いてバッファの効

(57) 【特許請求の範囲】

【請求項 1】

コンピューティング環境において、入力/出力 (I/O) 仮想化をサポートするための
 システムであって、

ネットワークファブリック上のサーバに関連付けられたチップを備え、前記チップは、
 複数のパケットバッファを含む外部メモリに関連付けられており、

前記チップは、物理ホストバスアダプタ (HBA) から受信したディスク読取データを含
 む1つ以上のパケットの状態を保存するオンチップメモリを含み、

前記チップは、

前記外部メモリ上の前記複数のパケットバッファ内の前記1つ以上のパケットをキュー
 に入れ、

前記1つ以上のパケットの状態に基づいて、前記外部メモリから前記1つ以上のパケ
 ットを読み出し、および

前記1つ以上のパケットを前記ネットワークファブリック上の前記サーバに送信する
 ように動作する、システム。

【請求項 2】

前記ネットワークファブリックは、インフィニバンド (I B) ファブリックであり、

前記サーバは、1つ以上のキューペア (QP) に関連付けられる、請求項1に記載のシ
 ステム。

【請求項 3】

10

20

前記複数のパケットバッファは、1つ以上のバッファリストに格納され、
各バッファリストは、前記サーバに関連付けられたキューペア（QP）に対応する、請求項2に記載のシステム。

【請求項4】

前記チップは、1つ以上のIBヘッダとシーケンス番号とを前記物理HBAから受信した各パケットに追加するように動作する、請求項2または3に記載のシステム。

【請求項5】

少なくとも1つのIBヘッダは、IB命令を含む、請求項4に記載のシステム。

【請求項6】

前記IBファブリック内のキューペア（QP）をターゲティングする複数のパケットは、パケットシーケンス番号スペースを共有するように配置される、請求項2～5のいずれか1項に記載のシステム。

10

【請求項7】

前記サーバに関連付けられたQPのために、複数のディスク読取命令、リモートダイレクトメモリアクセス（RDMA）読取リクエストおよび送信命令を多重化することをサポートするように、複数のコンテキストが開かれる、請求項2～6のいずれか1項に記載のシステム。

【請求項8】

前記1つ以上のパケットの状態は、前記外部メモリに格納された前記1つ以上のパケットがすべてキューに入れられたことを示す、および/または、1つ以上の関連IB命令が更新されたことを示す、請求項1～7のいずれか1項に記載のシステム。

20

【請求項9】

前記オンチップメモリ内の各エントリは、2ビットの幅を有しており、

前記オンチップメモリ内の各エントリの一方のビットは、関連IB命令を更新する必要があるか否かを示し、

他方のビットは、チップがキューに入れられたパケットを前記外部メモリから読出することができるか否かを示す、請求項1～8のいずれか1項に記載のシステム。

【請求項10】

コンピューティング環境において、効率的なパケット処理をサポートするための方法であって、

30

前記コンピューティング環境は、ネットワークファブリック上のサーバに関連付けられたチップを備え、前記チップは、複数のパケットバッファを含む外部メモリに関連付けられており、

物理ホストバスアダプタ（HBA）から受信したディスク読取データを含む1つ以上のパケットの状態を、前記チップに設けられるオンチップメモリ上に保存するステップとを含む、

前記チップは、

前記外部メモリ上の前記複数のパケットバッファ内の前記1つ以上のパケットをキューに入れ、

前記1つ以上のパケットの状態に基づいて、前記外部メモリから前記1つ以上のパケットを読出し、および

40

前記1つ以上のパケットを前記ネットワークファブリック上の前記サーバに送信するように動作する、方法。

【請求項11】

命令を含む機械読取可能プログラムであって、前記命令は、実行されると、複数のステップをシステムに実行させ、

前記システムは、ネットワークファブリック上のサーバに関連付けられたチップを備え、前記チップは、複数のパケットバッファを含む外部メモリに関連付けられており、

前記複数のステップは、

物理ホストバスアダプタ（HBA）から受信したディスク読取データを含む1つ以上の

50

パケットの状態を、前記チップに設けられるオンチップメモリ上に保存するステップと、前記外部メモリ上の前記複数のパケットバッファ内の前記１つ以上のパケットをキューに入れ、前記１つ以上のパケットの状態に基づいて、前記外部メモリから前記１つ以上のパケットを読み出し、および前記１つ以上のパケットを前記ネットワークファブリック上の前記サーバに送信するように、前記チップを動作させるステップとを含む、機械読取可能プログラム。

【発明の詳細な説明】

【技術分野】

【０００１】

著作権表示

10

この特許文書の開示の一部には、著作権保護の対象となるものが含まれている。著作権者は、特許商標庁の特許ファイルまたは記録に掲載された特許文書または特許開示の複製に対しては異議を唱えないが、その他の場合、すべての著作権を留保する。

【０００２】

発明の分野

本発明は、一般に、コンピュータシステムに関連し、特に、コンピューティング環境において、入力／出力（Ｉ／Ｏ）仮想化のサポートに関連する。

【背景技術】

【０００３】

背景

20

インフィニバンド（ＩＢ）技術は、クラウドコンピューティングファブリックの基盤として、その応用が増加しつつある。より大きなクラウドコンピューティングアーキテクチャが導入されるため、従来のネットワークおよびストレージに関連する性能上および管理上のボトルネックが重大な問題となっている。このような問題を対処することは、本発明の実施形態の一般的な意図である。

【発明の概要】

【課題を解決するための手段】

【０００４】

概要

システムおよび方法は、コンピューティング環境において、入力／出力（Ｉ／Ｏ）仮想化をサポートすることができる。システムは、チップを含むことができる。このチップは、ネットワークファブリック上のサーバに関連付けられている。さらに、このチップは、複数のパケットバッファを含む外部メモリに関連付けられている。また、オンチップメモリは、物理ホストバスアダプタ（ＨＢＡ）から受信したディスク読取データを含む１つ以上のパケットの状態を保存する。さらに、チップは、外部メモリ上の複数のパケットバッファ内の１つ以上のパケットをキューに入れ、１つ以上のパケットの状態に基づいて、外部メモリから１つ以上のパケットを読み出し、１つ以上のパケットをネットワークファブリック上のサーバに送信するように動作する。

30

【０００５】

システムおよび方法は、コンピューティング環境において、入力／出力（Ｉ／Ｏ）仮想化をサポートすることができる。システムは、１つ以上の仮想ホストバスアダプタ（ｖＨＢＡ）に関連付けられた複数のパケットバッファを含む空きバッファプールを含み、ｖＨＢＡの各々は、１つ以上のパケットバッファに指向するバッファポイントの主要リンクリストを空きバッファプールに保存する。また、入力／出力（Ｉ／Ｏ）装置に関連付けられたオンチップメモリ上で、コンテキストテーブルを定義することができる。このコンテキストテーブルは、ディスク読取操作のために、空きバッファプールから割当てられた１つ以上のパケットバッファに指向するバッファポイントの一時リンクリストを保存する。Ｉ／Ｏ装置は、ディスク読取操作を実行する物理ホストバスアダプタ（ＨＢＡ）からディスク読取データを受信すると、コンテキストテーブルを開き、バッファポイントの一時リンクリストを更新し、およびコンテキストテーブルが閉じられると、バッファポイントの

40

50

一時リンクリストをバッファポイントの主要リンクリストに合併するように、動作する。

【0006】

本明細書に記載のシステムおよび方法は、コンピューティング環境において、入力/出力(I/O)仮想化をサポートすることができる。システムは、メモリにおいて、空きバッファプールを備えることができる。I/O装置は、空きバッファプールを用いて、物理ホストバスアダプタ(HBA)から受信したディスク読取データを保存するように動作する。空きバッファプールは、2次元リンクリストおよび1次元リンクリストを含むことができる。2次元リンクリストの各エントリは、連続したメモリ位置で複数のパケットバッファを含み、1次元リンクリストの各エントリは、単一のパケットバッファを含む。

【図面の簡単な説明】

10

【0007】

【図1】さまざまなメモリインターフェイスを用いて、入力/出力(I/O)仮想化をサポートすることを示す図である。

【図2】本発明の一実施形態に従って、単一のメモリインターフェイスを用いて、入力/出力(I/O)仮想化をサポートすることを示す図である。

【図3】本発明の一実施形態に従って、単一のメモリインターフェイスを用いて、入来トラフィックを処理するための体系をサポートすることを示す図である。

【図4】本発明の一実施形態に従って、I/O装置上でディスク読取操作を開始することを示す図である。

【図5】本発明の一実施形態に従って、I/O装置上で要求したIOCBを抽出することを示す図である。

20

【図6】本発明の一実施形態に従って、I/O装置上でディスク読取データを処理することを示す図である。

【図7】本発明の一実施形態に従って、I/O装置上でディスク読取操作を完了する処理を示す図である。

【図8】本発明の一実施形態に従って、複数の仮想ホストバスアダプタ(vHBA)を用いて、I/O仮想化をサポートすることを示す図である。

【図9】本発明の一実施形態に従って、仮想ホストバスアダプタ(vHBA)において、複数のコンテキストをサポートすることを示す図である。

【図10】本発明の一実施形態に従って、オンチップメモリを用いて、外部メモリ上でキューに入れられたパケットの状態を保存することを示す図である。

30

【図11】本発明の一実施形態に従って、単一のメモリインターフェイスを用いて、入力/出力(I/O)仮想化をサポートすることを示す例示的なフローチャートである。

【図12】本発明の一実施形態に従って、空きバッファプールを用いて、複数の仮想ホストバスアダプタ(vHBA)をサポートすることを示す図である。

【図13】本発明の一実施形態に従って、ハイブリッドリンクリスト構造を用いて、ディスク読取操作をサポートすることを示す図である。

【図14】本発明の一実施形態に従って、ハイブリッドのリンクリスト構造を用いて、ヘッドラインブロッキングを回避することを示す図である。

【図15】本発明の一実施形態に従って、ハイブリッドリンクリスト構造を用いて、ヘッドラインブロッキングを回避することを示す例示的なフローチャートである。

40

【図16】本発明の一実施形態に従って、I/O装置のために2次元リンクリスト構造をサポートすることを示す図である。

【図17】本発明の一実施形態に従って、I/O装置のためにメモリの効率的な使用をサポートすることを示す図である。

【図18】本発明の一実施形態に従って、コンピューティング環境において、効率的なパケット処理をサポートすることを示す例示的なフローチャートである。

【発明を実施するための形態】

【0008】

詳細な説明

50

本発明は、限定することなく、例示として、添付の図面に示される。図面において、同様の参照番号は、同様の要素を標記する。本開示において、「一実施形態」または「１つの実施形態」または「いくつかの実施形態」を言及する場合、必ずしも同様の実施形態に限定されず、少なくとも１つの実施形態を意味することに留意すべきである。

【０００９】

本発明の以下の説明において、高性能ネットワークの一例として、インフィニバンド（ＩＢ）ネットワークを使用する。限定することなく、他の種類の高性能ネットワークを使用できることは、当業者には明らかであろう。また、本発明の以下の説明において、ストレージネットワークの一例として、ファイバチャネル（ＦＣ）ストレージネットワークを使用する。限定することなく、他の種類のストレージネットワークを使用できることは、10

【００１０】

本明細書に記載のシステムおよび方法は、１つ以上の仮想ホストバスアダプタ（ｖＨＢＡ）を用いて、入力／出力（Ｉ／Ｏ）仮想化をサポートすることができる。

【００１１】

入力／出力（Ｉ／Ｏ）仮想化

２つの異なるメモリアンターフェイスに基づいて、ＩＢファブリック上で、Ｉ／Ｏ仮想化をサポートすることができる。

【００１２】

図１は、異なるメモリアンターフェイスを用いて、入力／出力（Ｉ／Ｏ）仮想化をサポートすることを示す図である。図１に示すように、Ｉ／Ｏ装置１００は、ファイバチャネル（ＦＣ）ドメイン１０１およびインフィニバンド（ＩＢ）ドメイン１０２を用いて、入来トラフィック、たとえばストレージネットワーク１０５からＩＢファブリック１０４へのディスク読取データを処理することができる。20

【００１３】

図１に示すように、ファイバチャネル（ＦＣ）ドメイン１０１は、物理ホストバスアダプタ（ＨＢＡ）１０３に接続することができる。物理ＨＢＡ１０３は、たとえばＦＣ命令を用いて、ディスク読取操作を実行することができ、たとえば周辺機器相互接続エクスプレス（ＰＣＩエクスプレスまたはＰＣＩｅ）命令を用いて、データおよびコンテキストをＦＣドメイン１０１に送信することができる。30

【００１４】

ＦＣドメイン１０１は、ＦＣコンテキストリスト１２１を保存することができる。ＦＣコンテキストリスト１２１は、さまざまな仮想ホストバスアダプタ（ｖＨＢＡ）に関連する情報およびコンテキストを含むことができる。また、ＦＣドメイン１０１は、受信したディスク読取データおよび／またはコンテキストを外部メモリ、たとえば、シンクロナスダイナミックランダムアクセスメモリ（ＳＤＲＡＭ）１１１に記憶することができる。

【００１５】

図１に示すように、ＦＣドメイン１０１とＩＢドメイン１０２とは、シリアル相互接続を介して直接に接続されている。ＩＢドメイン１０２は、ＦＣドメイン１０１からＦＣデータおよびコンテキストを受信することができ、シーケンスの順序付けおよびコンテキストの管理のために、受信したＦＣデータおよびコンテキストをＩＢコンテキストリスト１２２内の異なるキューペア（ＱＰ）にマッピングすることができる。また、ＩＢドメイン１０２は、受信したディスク読取データおよびコンテキストを外部メモリ、たとえばＳＤＲＡＭ１１２に保存することができる。受信したディスク読取データおよびコンテキストは、ＩＢフォーマットであってもよい。次に、ＩＢドメイン１０２は、これらの情報をＩＢファブリック１０４に転送することができる。40

【００１６】

このように、システムは、複数の異なるメモリアンターフェイスを用いて、ストレージネットワーク１０５からＩＢファブリック１０４への入来トラフィックを処理することができる。50

【 0 0 1 7 】

単一のメモリアンターフェイス

本発明の一実施形態によれば、システムは、単一のメモリアンターフェイスを用いて、I / O 仮想化をサポートする、たとえば入来トラフィックおよび送出トラフィックの両方のために I B ファブリック上に作成された異なる仮想 H B A 用の並列 F C コンテキストを管理することができる。

【 0 0 1 8 】

図 2 は、本発明の一実施形態に従って、単一のメモリアンターフェイスを用いて、入力 / 出力 (I / O) 仮想化をサポートすることを示す図である。図 2 に示すように、I / O 装置 2 0 0 は、単一の F C / I B ドメイン 2 0 1 を表すチップを用いて、ストレージネットワーク 2 0 5 から I B ファブリック 2 0 4 への入来トラフィック、たとえばディスク読取データを処理することができる。

10

【 0 0 1 9 】

単一の F C / I B ドメイン 2 0 1 は、物理ホストバスアダプタ (H B A) 2 0 3 に直接に接続することができる。物理 H B A 2 0 3 は、F C 命令を用いて、ディスク読取操作を実行することができる。物理 H B A 2 0 3 は、P C I e 命令を用いて、ディスク読取データおよびコンテキストを F C / I B ドメイン 2 0 1 に送信することもできる。その後、F C / I B ドメイン 2 0 1 は、I B プロトコルを用いて、受信したディスク読取データおよびコンテキストを I B ファブリック 2 0 4 に送信することができる。

20

【 0 0 2 0 】

図 2 に示すように、F C / I B ドメイン 2 0 1 は、v H B A / Q P 情報リスト 2 2 0 を保存することができる。v H B A / Q P 情報リスト 2 2 0 は、受信した F C データおよびコンテキストを I B コンテキストリスト内の異なるキューペア (Q P) にマッピングすることができる。また、F C / I B ドメイン 2 0 1 は、たとえば I B フォーマットのディスク読取データおよびコンテキストを外部メモリ、たとえば S D R A M 2 1 0 に保存することができる。

【 0 0 2 1 】

本発明の一実施形態によれば、F C コンテキストリストを I B 信頼性のある接続 (R C) キューペア (Q P) リストに合併することによって、2 つの異なるメモリアンターフェイスの代わりに、単一のメモリアンターフェイスを使用することができる。たとえば、システムは、一時コンテキストリストを I B ドメインにマッピングする前に、外部メモリバッファのために、この一時コンテキストリストの動的リストを保有することができる。この手法は、2 つの異なる外部メモリを使用することを回避することができる、I B ドメインからバックプレッシャーメッセージ (back pressure message) を F C ドメインに送信することを回避することができる。したがって、システムは、同一のデータおよび / またはコンテキストを複数回格納することを回避することができ、遅延を改善する。さらに、2 つの異なるチップおよびメモリアンターフェイスの代わりに、単一のチップおよびメモリアンターフェイスを使用することは、システムのコストを低減することができる。

30

【 0 0 2 2 】

また、システムは、2 つの異なるドメイン間に通信を行うために、外部 (たとえば、ベンダ固有) インターフェイスに依存しない。単一のメモリアンターフェイスが使用されるため、F C / I B ドメイン 2 0 1 は、バッファサイズを知り、外部メモリ、たとえば S D R A M 2 1 0 のバッファを超過することを回避することができる。単一のメモリアンターフェイス手法によって、v H B A がダウンした場合、より良いフラッシュ操作を行うことができる。I B ドメインと F C ドメインとの間にメッセージの送受信がないため、フラッシュ操作を迅速かつ奇麗に行うことができる。

40

【 0 0 2 3 】

図 3 は、本発明の一実施形態に従って、単一のメモリアンターフェイスを用いて、入来トラフィックを処理するための体系 3 0 0 をサポートすることを示す図である。図 3 に示すように、単一のメモリアンターフェイスに関連付けられている F C / I B ドメイン 3 2

50

0を用いて、ストレージネットワークに接続している物理ホストバスアダプタ(HBA)330からIBファブリック上のサーバ310への入来トラフィックを処理することができる。

【0024】

ステップ301において、サーバ310は、たとえば、初期化ブロックをRC送信メッセージとしてFC/IBドメイン320に送信することによって、ディスク読取操作を開始することができる。次に、ステップ302において、FC/IBドメイン320は、メッセージを受信したことをサーバ310に知らせることができる。

【0025】

続いて、ステップ303において、サーバ310は、記述子リングの書込インデックスを更新することができ、1つ以上の新たな入力/出力制御ブロック(IOCB)が存在していることをFC/IBドメイン320に知らせることができる。次に、ステップ304において、FC/IBドメイン320は、メッセージを受信したことをサーバ310に知らせることができる。

10

【0026】

また、FC/IBドメイン320は、受信した書込インデックス値を読取インデックス値と比較することができる。値が異なる場合、FC/IBドメイン320は、ステップ305において、RDMA読取命令を用いて、サーバ310から1つ以上のIOCBを取得しようとする。よって、ステップ306において、サーバ310は、1つ以上のIOCBをRDMA読取応答データとしてFC/IBドメイン320に送信することができる。

20

【0027】

FC/IBドメイン320は、サーバ310からIOCBを受信すると、利用可能なコンテキストが存在する場合に、このコンテキストを開くことができる。本明細書において、コンテキストは、オンチップメモリを用いてチップ上に保存された特定命令の状態を示すものである。その後、FC/IBドメイン320は、このIOCB命令を物理HBA330にプッシュすることができる。

【0028】

たとえば、ステップ307において、FC/IBドメイン320は、ポインタ、たとえば応答書込インデックスを更新することができる。この応答書込インデックスは、要求したIOCBが利用可能であることをHBA330に表示する。次に、ステップ308において、HBA330は、IOCBリクエストの読取を試行することができ、ステップ309において、FC/IBドメイン320は、IOCBリクエスト読取データをHBA330に送信することができる。それに応じて、HBA330は、ディスク読取操作を実行することができる。

30

【0029】

本発明の一実施形態によれば、並列サーバのIOCB命令を処理するために、上記のステップ301~309は、同時に行うことができる。すなわち、FC/IBドメイン320は、複数の並列コンテキストを同時に保存および処理することができる。

【0030】

また、HBA330は、ステップ311~319において、ディスク読取データをFC/IBドメイン320に送信することができる。これに応じて、FC/IBドメイン320は、ステップ321~329において、RDMA書込操作を実行することによって、ディスク読取データをIBファブリック上のサーバ310に送信することができる。

40

【0031】

本発明の一実施形態によれば、システムは、ディスク読取データが完全にサーバ310またはホストに転送されたことを確保することができる。ステップ331において、サーバ310は、ディスク読取データの受信を確認するために、メッセージをFC/IBドメイン320に送信することができる。

【0032】

さらに、ステップ332において、ディスク読取データが完全に送信された場合、物理

50

H B A 3 3 0 は、対応する I O C B リクエストが完全に処理されたことを示す I O C B 応答を F C / I B ドメイン 3 2 0 に送信することができる。これに応じて、ステップ 3 3 3 において、F C / I B ドメイン 3 2 0 は、R C 送信メッセージを用いて、サーバ 3 1 0 に I O C B 応答を送信することができる。

【 0 0 3 3 】

最後に、ステップ 3 3 4 において、サーバは、I O C B 応答の受信を確認することができ、ステップ 3 3 5 において、F C / I B ドメイン 3 2 0 は、ポインタ、たとえば I O C B 応答がサーバ 3 1 0 に送信されたことを H B A 3 3 0 に示す応答読取インデックスを更新することができる。

【 0 0 3 4 】

本発明の一実施形態によれば、F C / I B ドメイン 3 2 0 は、データパス内の異なる種類の入来トラフィック、たとえば、v H B A 上のコンテキストのための R D M A 読取リクエスト、物理 H B A からのディスク読取データ、および v H B A 上のコンテキストのために、物理 H B A から受信した I O C B 応答を処理することができる。本実施形態において、ディスク書込データを取得するための R D M A 読取リクエストは、F C / I B ドメイン 3 2 0 によって内部で生成することができ、ディスク読取データおよび I O C B 応答は、P C I e バス (P C I - E x p r e s s b u s) を介して、物理 H B A から受信することができる。

【 0 0 3 5 】

図 4 は、本発明の一実施形態に従って、I / O 装置上でディスク読取操作を開始することを示す図である。図 4 に示すように、I / O 装置 4 0 0、たとえば F C / I B ドメイン 4 0 1 を表すチップは、I B ファブリック上のサーバ 4 0 2 から書込インデックス 4 1 2 を取得することができる。

【 0 0 3 6 】

F C / I B ドメイン 4 0 1 は、取得した書込インデックス 4 1 2 の値を読取インデックス値のコピーと比較することができる。値が異なる場合、F C / I B ドメイン 4 0 1 は、R D M A 読取命令 4 1 3 を用いて、サーバ 4 0 2 から 1 つ以上の要求した I O C B 4 1 1 を取得することができる。これらの R D M A 読取命令 4 1 3 は、I B フォーマットに変換することができ、F C / I B ドメイン 4 0 1 に関連付けられた外部入来メモリ 4 1 0 内の空きバッファプール 4 2 0 に保存することもできる。

【 0 0 3 7 】

本実施形態において、R D M A 読取命令 4 2 1 を入来 D R A M 4 1 0 に格納する前に、キューロジックは、R D M A 読取リクエストに利用可能なバッファが外部入来メモリ 4 1 0 に存在することを確認することができる。その後、F C / I B ドメイン 4 0 1 は、物理 H B A 4 0 3 からサーバ 4 0 2 への入来トラフィックを処理することができる。

【 0 0 3 8 】

図 5 は、本発明の一実施形態に従って、I / O 装置上で要求した I O C B を抽出することを示す図である。図 5 に示すように、I / O 装置 5 0 0、たとえば F C / I B ドメイン 5 0 1 を表すチップは、I B ファブリック上のサーバ 5 0 2 から R D M A 読取応答データを受信することができる。

【 0 0 3 9 】

I B プロトコルを用いて、予想通りにサーバ 5 0 2 から R D M A 応答読取データを完全に受信すると、F C / I B ドメイン 5 0 1 は、外部入来メモリ 5 1 0 内の空きバッファプール 5 2 0 に保存された R D M A 読取リクエスト 5 2 1 をキューから除外することができる。その後、F C / I B ドメイン 5 0 1 は、受信した R D M A 読取応答データ 5 1 2 を保存された R D M A 読取リクエスト 5 2 1 と比較することができる。

【 0 0 4 0 】

また、F C / I B ドメイン 5 0 1 は、受信した R D M A 読取応答データ 5 1 2 を解析することができ、I O C B リクエストを H B A 5 0 3 に転送する前に、R D M A 読取応答データ 5 1 2 に含まれた I O C B リクエスト 5 1 1 を抽出することができる。

【 0 0 4 1 】

10

20

30

40

50

図6は、本発明の一実施形態に従って、I/O装置上でディスク読取データを処理することを示す図である。図6に示すように、I/O装置600、たとえばFC/IBドメイン601を表すチップは、IOCBリクエスト613をHBA603に転送する前に、IOCBリクエスト613用のコンテキスト612を開くことができる。

【0042】

本発明の一実施形態によれば、FC/IBドメイン601は、IOCBリクエスト613用のコンテキスト612を開く前に、IOCBリクエスト命令613用のHBA603からのディスク読取データ611を格納するのに十分なスペース（たとえば、外部入来メモリ610に保留したDRAMスペース621）を有することを確認することができる。したがって、システムは、（たとえば、IOCBリクエスト命令613内の）ディスク読取命令が発行されると、FC/IBドメイン601が物理HBA603にバックプレッシャーを与えないことを確認することができる。

10

【0043】

FC/IBドメイン601から、IOCBリクエスト命令613内のディスク読取命令を受信した後、HBA603は、（たとえば、FCプロトコルを用いて）ストレージ上で、実際のディスク読取操作を実行することができる。HBA603は、PCI/PCIe書込トランザクションを用いて、ディスク読取データ611をFC/IBドメイン601に返送することができる。

【0044】

FC/IBドメイン601は、コンテキストを開くときに、ディスク読取命令のために外部入来メモリ610の空きバッファプール620にスペース621を保留しているため、受信したディスク読取データ611を外部入来メモリ610内のパケットバッファに書込む操作を開始することができる。また、FC/IBドメイン601は、ディスク読取データ611を外部入来メモリ610内のパケットバッファに書込む前に、IBヘッダおよびシーケンス番号をディスク読取データ611のために受信したパケットに追加することができる。よって、ディスク読取データ611のために受信した格納パケットは、IBフォーマットであってもよい。

20

【0045】

また、FC/IBドメイン601は、完全メッセージ（たとえば、RDMA読取リクエスト）またはIB最大伝送ユニット（MTU）パケット（たとえば、RDMA書込専用パケット）が利用可能な場合、格納されたディスク読取データ611を外部入来メモリ610に読出すことができる。その後、FC/IBドメイン601は、外部入来メモリ610内の空きバッファプール620から読出されたIBパケットをディスク読取データ631としてIBファブリック上のサーバ602に送信することができる。

30

【0046】

図7は、本発明の一実施形態に従って、I/O装置上でディスク読取操作を完了する処理を示す図である。図7に示すように、たとえば、I/O装置700、たとえばFC/IBドメイン701を表すチップを用いて、IBファブリック上の物理HBA703からサーバ702への入来トラフィックを処理することができる。

【0047】

40

ディスク読取データが完全に送信された場合、HBA703は、コンテキスト712に関連付けられた対応のIOCBリクエストが完全に処理されたことを示すIOCB応答711をFC/IBドメイン701に送信することができる。その後、FC/IBドメイン701は、IBヘッダおよびシーケンス番号をIOCB応答711に追加することができ、IOCB応答721を外部入来メモリ710内の空きバッファプール720に格納することができる。

【0048】

メッセージまたはパケットを送信する準備ができると、FC/IBドメイン701は、IBプロトコルを用いて、ホスト/サーバ702にIOCB応答721を送信することができる。IOCB713応答を受信すると、ホスト/サーバ702は、ディスク読取IO

50

ＣＢリクエスト命令７３１がハードウェアによって完全に処理されたことを確認することができる。

【００４９】

また、ＦＣ／ＩＢドメイン７０１は、コンテキスト用のＩＯＣＢ応答７２１を送信した後、関連するコンテキスト７１２を閉じることができる（すなわち、状態メモリを消去し、外部入来メモリ７１０内の保留スペースを削除することができる）。

【００５０】

複数のコンテキスト

図８は、本発明の一実施形態に従って、複数の仮想ホストバスアダプタ（ｖＨＢＡ）を用いて、Ｉ／Ｏ仮想化をサポートすることを示す図である。図８に示すように、Ｉ／Ｏ装置８００、たとえばＦＣ／ＩＢドメイン８０１を表すチップを用いて、入来トラフィック８３０を処理することができる。入来トラフィック８３０は、物理ＨＢＡ８０３からＩＢファブリック上のサーバ８０２に転送された複数のパケット、たとえばパケット８３１～８３９を含むことができる。

10

【００５１】

また、ＦＣ／ＩＢドメイン８０１は、１つ以上のｖＨＢＡ、たとえば、ｖＨＢＡ Ａ８５１、ｖＨＢＡ Ｂ８５２およびｖＨＢＡ Ｃ８５３をサポートすることができる。ｖＨＢＡ Ａ８５１、ｖＨＢＡ Ｂ８５２およびｖＨＢＡ Ｃ８５３は、ＩＢサーバ８０２に関連付けられたキューペア（ＱＰ）、たとえばＱＰ Ａ８４１、ＱＰ Ｂ８４２およびＱＰ Ｃ８４３にそれぞれ対応することができる。

20

【００５２】

さらに、ＦＣ／ＩＢドメイン８０１は、外部入来メモリ８１０を用いて、１つ以上の受信パケットを格納することができる。ＦＣ／ＩＢドメイン８０１は、単一のメモリインターフェイスの使用をサポートするために、ＦＣコンテキスト情報、たとえばｖＨＢＡ Ａ８５１、ｖＨＢＡ Ｂ８５２およびｖＨＢＡ Ｃ８５３をＩＢコンテキストリスト、たとえばＱＰ Ａ８４１、ＱＰ Ｂ８４２およびＱＰ Ｃ８４３に合併することができる。

【００５３】

図８に示すように、外部入来メモリ８１０は、空きバッファプール８２０を提供することができる。空きバッファプール８２０は、１つ以上のバッファリスト、たとえばバッファリストＡ８２１、バッファリストＢ８２２およびバッファリストＣ８２３を含む。バッファリストＡ８２１、バッファリストＢ８２２およびバッファリストＣ８２３の各々を用いて、特定のＱＰ（またはｖＨＢＡ）をターゲットする１つ以上の受信パケットを格納することができる。

30

【００５４】

たとえば、ＦＣ／ＩＢドメイン８０１は、ｖＨＢＡ Ａ８５１に関連付けられたバッファリストＡ８２１内のＱＰ Ａ８４１をターゲットするパケット８３２および８３９をキューに入れることができる。同様に、ＦＣ／ＩＢドメイン８０１は、ｖＨＢＡ Ｂ８５２に関連付けられたバッファリストＢ８２２内のＱＰ Ｂ８４２をターゲットするパケット８３３および８３８をキューに入れることができ、ｖＨＢＡ Ｃ８５３に関連付けられたバッファリストＣ８２３内のＱＰ Ｃ８４３をターゲットするパケット８３１をキューに入れることができる。

40

【００５５】

また、ＦＣ／ＩＢドメイン８０１は、受信した複数のパケット８３１～８３９の状態を保存することができる制御構造８１１を含むことができる。さらに、ＦＣ／ＩＢドメイン８０１は、読取ロジック８１２を用いて、格納されたパケット８３１～８３９のうち１つ以上を読出すことができる。

【００５６】

本発明の一実施形態によれば、ＦＣ／ＩＢドメイン８０１は、ＩＢドメイン内のＱＰのために、複数のディスク読取命令、ＲＤＭＡ読取リクエストおよびＲＣ送信命令の多重化をサポートするために、ｖＨＢＡ Ａ８５１～ｖＨＢＡ Ｃ８５３内の複数のコンテキスト

50

を開くことができる。

【0057】

図9は、本発明の一実施形態に従って、仮想ホストバスアダプタ（vHBA）において、複数のコンテキストをサポートすることを示す図である。図9に示すように、I/O装置、たとえばFC/IBドメイン900を表すチップは、たとえば、QP904のために物理HBA903上で複数のディスク読取命令を実行するために、単一のvHBA901内の複数のコンテキスト、たとえばコンテキストI 910およびコンテキストII 920を開くことができる。

【0058】

たとえば、コンテキストI 910は、物理HBA903から受信したいくつかのパケット、たとえばC1D1 911、C1D2 912およびC1D3 913を含むことができる。C1D1 911は、コンテキストI 910用のディスク読取データD1を含み、C1D2 912は、コンテキストI 910用のディスク読取データD2を含み、C1D3 913は、コンテキストI 910用のディスク読取データD3を含むことができる。

10

【0059】

また、コンテキストII 920は、物理HBA903から受信したいくつかのパケットC2D1 921およびC2D2 922を含むことができる。C2D1 921は、コンテキストII 920用のディスク読取データD1を含み、C2D2 922は、コンテキストII 920用のディスク読取データD2を含むことができる。

20

【0060】

また、FC/IBドメイン900は、物理HBA903から受信したパケットをIBファブリック上のQP904に送信する前に、このパケットに対応のシーケンス番号（PSN）および異なるIBヘッダを追加することができる。

【0061】

本発明の一実施形態によれば、同一のQPたとえば（vHBA901に関連付けられた）QP904をターゲットするすべてのパケットは、IBドメイン内で単一PSNスペース902を共有することができる。図9に示すように、PSNスペース902において、パケットをP0、P1、・・・、P（N）の順序に編成することができる。ここでは、 $P1 = P0 + 1$ 、 $P2 = P1 + 1$ 、・・・、 $P（N） = P（N - 1） + 1$ 。

30

【0062】

一方、IBドメイン内でPSNスペース902を共有する場合、動作中に、IBドメイン内のPSN番号割当体系を用いて発信パケットの順序を変更することができないため、異なるコンテキストにおいて、IBヘッダおよびシーケンス番号を単一のメモリアンターフェイスに基づいてHBA803から受信したパケットに追加することを複雑化する可能性がある。

【0063】

図9に示すように、コンテキストI 910用のディスク読取データが完全に処理される前に、コンテキストII 920用のディスク読取データが入来すると、vHBA901にヘッドラインブロッキングという問題が生じる可能性がある。たとえば、システムが別のディスク書込操作を行っている最中に、FC/IBドメイン900がディスク書込操作のためにRDMA読取リクエストをスケジュールしようとするときに、このような問題が生じる。

40

【0064】

図10は、本発明の一実施形態に従って、オンチップメモリを用いて、外部メモリ上でキューに入れられたパケットの状態を保存することを示す図である。図10に示すように、I/O装置1000、たとえばFC/IBドメイン1000を表すチップは、単一のvHBA/QP内の複数のコンテキスト、たとえばコンテキストI 1010およびコンテキストII 1020を開くことができる。各コンテキストは、1つ以上のパケットを含むことができる。たとえば、コンテキストI 1010は、パケットC1D1 1011

50

、C 1 D 2 1 0 1 2 および C 1 D 3 1 0 1 3 を含み、コンテキスト I I 1 0 2 0 は、パケット C 2 D 1 1 0 2 1 および C 2 D 2 1 0 2 2 を含むことができる。

【 0 0 6 5 】

本発明の一実施形態によれば、F C は / I B ドメイン 1 0 0 0 は、シーケンス番号およびさまざまな I B ヘッダを物理 H B A 1 0 0 5 から受信した各パケットに追加することができる。I B ヘッダは、特定の packets を I B ファブリック 1 0 0 4 に送信するときに、この特定の packets に適用することができる I B 命令を含むことができる。

【 0 0 6 6 】

たとえば、コンテキスト用のディスク読取データのサイズが I B 最大伝送ユニット (M T U) のサイズと同一である場合、I B ヘッダ内のキューに入れられた I B 命令は、コン
10 テキストメモリに指定された仮想アドレス (V A) を有する「R D M A 書込専用」命令に
することができる。逆に、所定のコンテキスト用のディスク読取データのサイズが I B
M T U のサイズよりも大きい場合、ディスク読取データを複数の packets に分割すること
ができる。各 packets の I B ヘッダ内のキューに入れられた命令は、コンテキストメモリ
により指定された V A を有する「最初の R D M A 書込」命令、「中間の R D M A 書込」命
令および「最後の R D M A 書込」命令のいずれか 1 つであってもよい。ディスク読取デ
ータのサイズに応じて、「中間の R D M A 書込」命令を有するように packets をキューに入
れなくてもよく、または「中間の R D M A 書込」命令を有するように複数の packets をキ
ューに入れてもよい。

【 0 0 6 7 】

図 1 0 に示すように、初期 P S N は、P 0 であり、ディスク読取データは、I B M T
U よりも大きい。第 1 の packets (C 1 D 1 1 0 1 1) を受信すると、システムは、P
S N としての P 0 および「最初の R D M A 書込」命令を packets に追加することができる
。次の packets (C 1 D 2 1 0 1 2) を受信すると、システムは、P S N としての P 1
(すなわち、P 0 + 1) および「中間の R D M A 書込」命令を packets に追加すること
ができる。第 3 の packets (C 2 D 1 1 0 2 1) を受信すると、システムは、P S N とし
ての P 2 (すなわち、P 1 + 1) および「最初の R D M A 書込」命令を packets に追加す
ることができる。さらに、packets (C 2 D 2 1 0 2 2) を受信すると、システムは、
P S N としての P 3 (すなわち、P 2 + 1) および「最後の R D M A 書込」命令を packets
20 30 に追加することができる。

【 0 0 6 8 】

しかしながら、上記一連の操作には、整合性問題が存在している。packets C 1 D 2
1 0 1 2 に追加された I B 命令は、「中間の R D M A 書込」命令である。システムは、デ
ィスク読取データのサイズに応じて、次の packets に追加された I B 命令が「中間の R D
M A 書込」命令または「最後の R D M A 書込」命令のいずれかであると预期している。し
かしながら、図 1 0 に示すように、次の packets C 2 D 1 1 0 2 1 が異なるコンテキス
ト I I 1 0 2 0 からのものであるため、F C / I B ドメイン 1 0 0 0 は、新たな命令、
たとえば、「最初の R D M A 書込」命令または「R D M A 書込専用」命令を packets (そ
の P S N が正確であっても) に追加することができる。また、F C / I B ドメイン 1 0 0
0 が、別のコンテキストのために記述子の取得または R C 送信 I O C B 応答を行うため、
R D M A 読取リクエストをキューに入れようとする際に、同様の問題が生じる可能性もあ
る。

【 0 0 6 9 】

この問題を解決するために、システムは、既にキューに入れられた packets に関連付け
られた I B 命令を更新することができる。たとえば、F C / I B ドメイン 1 0 0 0 がパケ
ット C 2 D 1 1 0 2 1 を受信した後、システムは、C 1 D 2 1 0 1 2 の I B 命令、す
なわち「中間の R D M A 書込」命令 1 0 0 7 を「最後の R D M A 書込」命令 1 0 0 8 に変
更することができる。

【 0 0 7 0 】

本発明の一実施形態によれば、整合性を確保するために、F C / I B ドメイン 1 0 0 0

10

20

30

40

50

は、オンチップメモリ 1 0 0 2 を用いて、外部メモリ 1 0 0 1 上でキューに入れられたパケット 1 0 0 6 の状態 1 0 0 9 を格納することができる。

【 0 0 7 1 】

このオンチップメモリ 1 0 0 2 は、異なるパースペクティブ (perspective) に有益であり得る。まず、パケットがキューに入れられ、必要に応じてパケットに関連付けられた I B 命令が更新された場合に限り、外部メモリからパケットを読み取りおよびパケットをホストに送信するように、読取ロジックを確保することができる。次に、複数のコンテキストをサポートするために、既にキューに入れられたパケットに関連付けられた I B 命令を必要に応じて更新することができる。

【 0 0 7 2 】

たとえば、オンチップメモリ 1 0 0 2 は、2 ビットの幅 (および 6 4 K の深さ) を有することができる。オンチップメモリ 1 0 0 2 内のエントリの第 1 ビット、たとえばビット 0 は、I B 命令を変更または更新する必要があるか否かを示し、第 2 ビット、たとえばビット 1 は、F C / I B ドメイン 1 0 0 0 の読取ロジックが外部メモリ 1 0 0 1 からキューに入れられたパケットを取出すことができるか否かを示すことができる。

【 0 0 7 3 】

以下の表 1 は、一連のパケットが到着する場合に、例示のオンチップメモリに格納された異なるパケットの状態を示している。

【 0 0 7 4 】

【表 1】

物理 H B A からのパ ケット	外部メモリ上で、 パケットをキュー に入れる	現在のパケットをキュー に入れるときのオン チップメモリ	次のパケットをキュー に入れるときのオン チップメモリ
C2D1	最初の RDMA 書込 PSN P0, C2D1	ビット 0 : 0 ビット 1 : 0	ビット 0 : 0 ビット 1 : 1
C2D2	中間の RDMA 書込 PSN P1, C2D2	ビット 0 : 0 ビット 1 : 0	ビット 0 : 1 ビット 1 : 1
C1D1	最初の RDMA 書込 PSN P2, C1D1	ビット 0 : 0 ビット 1 : 0	ビット 0 : 0 ビット 1 : 1
C1D2	最後の RDMA 書込 PSN P3, C1D2	ビット 0 : 0 ビット 1 : 1	ビット 0 : 0 ビット 1 : 1

【 0 0 7 5 】

上記の表 1 に示すように、第 1 のパケット、すなわち C 2 D 1 をキューに入れるときに、オンチップステートメモリは、2 b 0 0 である。このとき、読取ロジックは、このパケットを読み取ることができない。その理由は、パケットが他のコンテキストから来る場合、システムがこのパケットの命令を変更する必要があるからである。

【 0 0 7 6 】

次のパケット C 2 D 2 が来的时候に、前のパケット、すなわち C 2 D 1 のオンチップ状態は、2 b 1 0 に変更されている。このとき、パケットは、正常にキューに入れられ、読取ロジックは、このパケットを読み出すことができる。この場合、C 2 D 2 が同様のコンテキスト I I (C 2) から来たため、命令を変更する必要がない。

【 0 0 7 7 】

さらに、第 3 のパケット、すなわち C 1 D 1 が来的时候に、C 2 D 2 のオンチップメモリの状態は、2 b 1 1 に変更されている。このとき、パケットは、キューに入れられ読取側で命令を変更する必要がある。読取ロジックは、このパケットを読み出すことができ、このパケットを送信する前に、命令を「中間の R D M A 書込」命令から「最後の R D M A 書込」命令に変更することができる。

【 0 0 7 8 】

図 1 1 は、本発明の一実施形態に従って、単一のメモリインターフェイスを用いて、入

10

20

30

40

50

力／出力（Ｉ／Ｏ）仮想化をサポートすることを示す例示的なフローチャートである。図 11 に示すように、ステップ 1101 において、システムは、ネットワークファブリック上のサーバに関連付けられ、且つ、複数のパケットバッファを含む外部メモリに関連付けられているチップを提供することができる。また、ステップ 1102 において、システムは、物理ホストバスアダプタ（ＨＢＡ）から受信したディスク読取データを含む 1 つ以上のパケットの状態をオンチップメモリに保存することができる。また、ステップ 1103 において、システムは、外部メモリ上の複数のパケットバッファ内の 1 つ以上のパケットをキューに入れ、1 つ以上のパケットの状態に基づいて、外部メモリから 1 つ以上のパケットを読み出し、およびネットワークファブリック上のサーバに 1 つ以上のパケットを送信するように、チップを動作させることができる。

10

【0079】

ハイブリッドリンクリスト構造

本発明の一実施形態によれば、システムは、ハイブリッドリンクリスト構造を用いて、仮想ホストバスアダプタ（ｖＨＢＡ）内の複数のコンテキストに関連付けられた入来トラフィックを処理することができる。このハイブリッドリンクリスト構造は、主要リンクリストおよび一時リンクリストを含むことができる。

【0080】

図 12 は、本発明の一実施形態に従って、空きバッファプールを用いて、複数の仮想ホストバスアダプタ（ｖＨＢＡ）をサポートすることを示す図である。図 12 に示すように、Ｉ／Ｏ装置 1200、たとえば ＦＣ／ＩＢドメイン 1204 を表すチップは、たとえば、空きバッファプール 1210 を用いて、異なる仮想ホストバスアダプタ（ｖＨＢＡ）、たとえば ｖＨＢＡ Ⅰ 1201 および ｖＨＢＡ Ⅱ 1202 に関連付けられ得る入来トラフィック 1203 を処理することができる。

20

【0081】

また、各 ｖＨＢＡ は、入来トラフィック 1203 から受信した各種パケットを空きバッファプール 1210 内のパケットバッファに格納するために、バッファポインタの 1 つ以上のリンクリストを保存することができる。たとえば、ｖＨＢＡ Ⅰ 1201 は、主要リンクリスト Ⅰ 1211 を保存することができ、ｖＨＢＡ Ⅱ 1202 は、主要リンクリスト Ⅱ 1212 を保存することができる。

【0082】

図 13 は、本発明の一実施形態に従って、ハイブリッドリンクリスト構造を用いて、さまざまなディスク読取操作をサポートすることを示す図である。図 13 に示すように、Ｉ／Ｏ装置 1300、たとえば ＦＣ／ＩＢドメインを表すチップは、複数のディスク読取操作を並列で実行するために、単一の ｖＨＢＡ 1303 内の複数のコンテキスト、たとえば コンテキスト Ａ 1301 および コンテキスト Ｂ 1302 を開くことができる。

30

【0083】

本発明の一実施形態によれば、ＦＣ／ＩＢドメインは、各ディスク読取操作のために、オンチップメモリ上で、コンテキストテーブルを定義することができる。たとえば、ＦＣ／ＩＢドメインは、コンテキスト Ａ 1301 のディスク読取データを受信した場合、オンチップメモリ 1310 上で、コンテキストテーブル Ａ 1311 を開くことができる。コンテキストテーブル Ａ 1311 は、空きバッファプール 1320 から割当てられた 1 つ以上のパケットバッファに指向するバッファポインタの一時リンクリストを保存することができる。また、コンテキストテーブル Ａ 1311 は、所定のトランザクション用の仮想アドレスを保存することもできる。

40

【0084】

コンテキスト Ａ 1301 の入来データが来るときに、ディスクデータを ＩＢヘッダおよび ＰＳＮ番号とともに外部 ＳＤＲＡＭメモリ内の空きバッファプール 1320 に書込むことができる。図 13 に示すように、コンテキストメモリ内のバッファポインタの一時リンクリスト 1321 が更新されるが、所定の ｖＨＢＡ用のバッファポインタの主要リンクリスト 1330 が変更されない。

50

【0085】

本発明の一実施形態によれば、異なるディスク読取操作のために、複数のコンテキストを開くことができる。新たに開いたコンテキストのディスク読取データが来るときに、システムは、外部S D R A Mメモリから以前に開かれたコンテキストに格納されたデータを読み出すことができ、必要に応じて、コンテキストリスト末尾のI Bヘッダ内の命令を更新することができる。たとえば、所定のコンテキストのためのディスク読取データがI B M T Uよりも大きいときに、I Bヘッダにキューに入れられた命令が「最初のR D M A書込」命令である場合、その命令を「R D M A書込専用」命令に変更することができ、I Bヘッダにキューに入れられた命令が「中間のR D M A書込」命令である場合、その命令を「最後のR D M A書込」命令に変更することができる。

10

【0086】

図13に示すように、F C / I Bドメインが異なるコンテキストB 1302からデータを受信する場合、一時リンクリスト1321をv H B A 1303の主要リンクリスト1330に合併することができる。たとえば、v H B A 1303の主要リンクリスト1330のテールポインタを一時リンクリスト1321のヘッドポインタに変更することができ、一時リンクリスト1321のテールポインタを主要リンクリスト1330の新たなテールポインタになる。よって、新たなコンテキストのデータは、新たなP S Nを有する新たなコンテキストメモリに書込まれることができる。それに応じて、そのコンテキストの一時ポインタは、更新されることができる。

【0087】

20

同様に、システムは、新たなコンテキスト内の「R D M A書込専用」命令、「最後のR D M A書込」命令、「送信専用」命令、およびR D M A読取リクエストなどの命令を実行する必要がある場合、以前に開かれたコンテキストを閉じ、一時リンクリスト1321を主要リンクリスト1330に合併することができる。

【0088】

図14は、本発明の一実施形態に従って、ハイブリッドのリンクリスト構造を用いて、ヘッドラインブロッキングを回避することを示す図である。図14に示すように、コンテキスト1401を閉じ、v H B Aに1403の主要リンクリスト1430を更新した後、I / O装置1400、たとえばF C / I Bドメインを表すチップは、コンテキストB 1402のために、オンチップメモリ1410内の新たなコンテキストテーブルB 1412を開くことができる。コンテキストテーブルB 1412は、空きバッファプール1420から割当てられたパケットバッファに指向するバッファポインタを含む新たな一時リンクリストB 1422を保存することができる。

30

【0089】

図14に示すように、コンテキストB 1402内の受信したディスク読取データ（またはR D M A読取リクエスト）に対する処理は、コンテキストA 1401内の受信したディスク読取データの処理によってブロックされるため、システムは、ヘッドラインブロッキングを回避することができる。よって、所定のv H B A内の異なるコンテキストのディスク読取データを並列に処理することができる。

【0090】

40

図15は、本発明の一実施形態に従って、ハイブリッドリンクリスト構造を用いて、ヘッドラインブロッキングを回避することを示す例示的なフローチャートである。図15に示すように、ステップ1501において、システムは、ネットワーク環境内の1つ以上の仮想ホストバスアダプタ（v H B A）に関連付けられた複数のパケットバッファを含む空きバッファプールを提供することができる。各々のv H B Aは、空きバッファプール内の1つ以上のパケットバッファに指向するバッファポインタの主要リンクリストを保存する。ステップ1502において、システムは、入力/出力（I / O）装置に関連付けられたオンチップメモリ上で、コンテキストテーブルを定義することができる。コンテキストテーブルは、ディスク読取操作のために空きバッファプールから割当てられた1つ以上のパケットバッファに指向するバッファポインタの一時リンクリストを保存する。ステップ1

50

503において、システムは、I/O装置がディスク読取操作を実行する物理ホストバスアダプタ(HBA)からディスク読取データを受信すると、コンテキストテーブルを開き、一時リンクリストを更新し、およびコンテキストテーブルを閉じたと、バッファポインタの一時リンクリストをバッファポインタの主要リンクリストに合併するように、I/O装置を動作させることができる。

【0091】

一体化メモリ構造

本発明の一実施形態によれば、ディスク読取データがHBAチップから送信されるときに、システムは、ディスク読取データを一体化メモリ構造内のさまざまなデータバッファに格納することができる。

10

【0092】

図16は、本発明の一実施形態に従って、I/O装置のために2次元リンクリスト構造をサポートすることを示す図である。図16に示すように、システムは、空きバッファプール1600内の2次元リンクリスト1610を用いて、入来パケットバッファを管理することができる。2次元リンクリスト1610は、複数のエントリを含むことができる。各エントリは、スーパーブロック(たとえば、スーパーブロック1601~1609)であってもよい。

【0093】

本発明の一実施形態によれば、スーパーブロック1601~1609は、連続したメモリ位置に格納された複数のパケットを表すことができる。また、スーパーブロック1601~1609の各々は、バッファ管理のために、内部でパケットバッファリストに指向することができる。したがって、オンチップリソース使用の観点から、2次元リンクリスト構造は、非常に効率的である。システムは、バッファされるパケットの数を最大にしなから、オンチップメモリ上のリンクリストのサイズを最小にすることができる。

20

【0094】

たとえば、さまざまなサイズの(オーバーヘッドを含む)IBパケットを収容するために、2次元リンクリスト1610は、8Kのスーパーブロックを含むことができる。また、各スーパーブロックは、(各々が8KBのサイズを有する)8つのパケットを保存することができる64KB(512Kb)のサイズを有することができる。図16に示すように、スーパーブロック1601は、8つのパケットバッファ、すなわち、パケットバッファ1611~1618を含むことができる。

30

【0095】

2次元リンクリスト1610によって、FC/IBドメインは、異なるQPをターゲットする読取データディスクをIBドメインに格納することができる。図16に示すように、FC/IBドメインは、異なるポインタを用いて、2次元リンクリスト1610内のさまざまなスーパーブロックのリンクリストにアクセスすることができる。たとえば、FC/IBドメインは、スーパーブロック1602、スーパーブロック1604およびスーパーブロック1608を含むスーパーブロックのリンクリストに指向するQP Aヘッドポインタ1621(および/またはQP Aテールポインタ1622)を保存することができる。また、FC/IBドメインは、スーパーブロック1606、スーパーブロック1605およびスーパーブロック1609を含むスーパーブロックのリンクリストに指向するQP Bヘッドポインタ1623(および/またはQP Bテールポインタ1624)を保存することができる。

40

【0096】

本発明の一実施形態によれば、システムは、2次元ハイブリッドリンクリスト1610を所定のインフィニバンド(IB)RC QP接続の1次元リンクリストと動的に合併することによって、外部DRAMメモリの効率的な使用をサポートすることができる。したがって、システムは、小さいサイズのパケットを固定サイズのスーパーブロックに格納することによるメモリスペースの浪費を避けることができる。

【0097】

50

たとえば、FC/IBドメインは、ディスク読取リクエストを実行するために、空きバッファプール1600にバッファの有無を照会することができる。空きバッファプール1600に十分なパケットバッファがある場合に、FC/IBドメインは、物理HBAにディスク読取IOCBリクエストを発行することができる。ディスク読取リクエストに要求されたバッファは、空きバッファプール1600に保留され、現在のコンテキストがFC/IBドメインによって解放されるまで、他の後続のリクエストに使用されない。

【0098】

また、システムは、RDMA読取リクエストを保存するバッファ（たとえば、4Kバッファ）のリストを定義することができる。システムは、RDMA読取リクエストが発行されるたびに、RDMA読取リクエストに利用可能なスペースが外部メモリに保留されており、RDMA読取リクエストがRDMA書込操作にブロックされないことを保証することができる。

10

【0099】

2次元リンクリスト1610のみを使用する場合、システムは、RDMA読取リクエストに128個のキューペア（またはvHBA）に共有される4Kのパケットバッファを提供するために、64K（スーパーブロックのサイズ）×4K×128バイトのスペースをメモリに確保する必要がある。このスペースが8K（パケットバッファのサイズ）×4K×128バイトであるパケットバッファのメモリ使用量よりも実質的に多いため、この手法は、メモリを浪費する。

【0100】

20

図17は、本発明の一実施形態に従って、I/O装置のためにメモリの効率的な使用をサポートすることを示す図である。図17に示すように、I/O装置1700、たとえばFC/IBドメインを表すチップは、空きバッファプール1701を用いて、さまざまなパケットのキュー入れ（1730）をサポートすることができる。空きバッファプール1701は、スーパーブロック1711～1719を包含する2次元リンクリスト1710、およびパケットバッファ1721～1729を包含する1次元リンクリスト1720を含むことができる。この2次元リンクリスト1710は、図16に示された2次元リンクリスト1610と類似してもよい。

【0101】

本発明の一実施形態によれば、空きバッファプール1701において、異なる種類のトランザクションをキューに入れることができる。たとえば、これらのトランザクションを用いて、RDMA書込命令1742およびRDMA読取リクエスト1741を実行することができる。

30

【0102】

パケットをキューに入れた（1730）場合、トランザクションの種類に基づいて、2次元リンクリスト1710からまたは単一の1次元リンクリスト1720から、空きバッファを割当てることができる。

【0103】

また、システムは、さまざまなバッファされたパケットの状態を保存するために、チップ上でリンクリストの制御構造1740を保有することができる。制御構造1740は、たとえば、メモリスーパーブロック位置のヘッドポインタ（たとえば、13ビットのSBLKHEAD）、スーパーブロック内のパケットオフセット位置のヘッドポインタ（たとえば、3ビットのPKTHEAD）、メモリスーパーブロック位置のテールポインタ（たとえば、13ビットのSBLKTAIL）、スーパーブロック内のパケットオフセット位置のテールポインタ（たとえば、3ビットのPKTTAIL）、およびパケットバッファが2次元リンクリストから割当てられているかまたは1次元リンクリストから割当てられているかを示すフラグ（たとえば、1ビットのLISTTYPE）を格納することができる。また、制御構造1740は、QP/vHBAの数に基づいて深さ情報を格納することができ、必要な制御情報に基づいて幅情報を格納することができる。

40

【0104】

50

本発明の一実施形態によれば、システムは、異なるキュー入れシナリオをサポートすることができる。

【0105】

キューに入れられたトランザクションがR D M A 書込命令用のものである場合、システムは、2次元リンクリスト1710からバッファまたはスーパーブロックを取得することができる。

【0106】

逆に、キューに入れられたトランザクションがR D M A 読取命令用のものである場合、以前にキューに入れられたトランザクションに割当てられたスーパーブロックにパケットバッファが残されていないときに、システムは、1次元リンクリスト1720からバッファを取得することができる。

【0107】

一方、R D M A 書込操作が進行中に、R D M A 読取リクエスト用のトランザクションがキューに入れられる場合もある。利用可能なパケットバッファが存在する場合、システムは、R D M A 書込動作に割当てられたスーパーブロック内の現在のパケット位置で、R D M A 読取リクエストをキューに入れることができる。

【0108】

また、特定のQ P / v H B A のために、単一リンクリスト1720からのパケットバッファを保留することができる。単一リンクリスト1720から保留されたバッファは、R D M A 書込パケットまたはR D M A 読取リクエストパケットのいずれかに使用されることができる。さらに、システムは、制御メモリ内のL I S T T Y P E フィールドにフラグを付けることができる。よって、デキュー（dequeue）ロジックおよび/または読取ロジックは、1つのパケットが単一リンクリスト1720のキューに入れられたことを知ることができる。

【0109】

このように、システムは、効率的なパケット処理を実現することができ、外部メモリの浪費を回避することができる。

【0110】

図18は、本発明の一実施形態に従って、コンピューティング環境において、効率的なパケット処理をサポートすることを示す例示的なフローチャートである。図18に示すように、ステップ1801において、システムは、2次元リンクリストおよび1次元リンクリストを含む空きバッファプールをメモリに提供することができる。また、ステップ1802において、システムは、2次元リンクリストの各エントリが連続したメモリ位置に複数のパケットバッファを有すること、および1次元リンクリストの各エントリが単一のパケットバッファを有することを可能にする。次に、1803において、I / O 装置は、物理ホストバスアダプタ（H B A ）から受信したディスク読取データを空きバッファプールに格納することができる。

【0111】

本発明の一実施形態は、コンピューティング環境において、I / O 仮想化をサポートするためのシステムを提供する。このシステムは、コンピューティング環境内の1つ以上の仮想ホストバスアダプタ（v H B A ）に関連付けられた複数のパケットバッファを含む空きバッファプールを含み、v H B A の各々は、1つ以上のパケットバッファに指向するバッファポインタの主要リンクリストを空きバッファプールに保存し、入力/出力（I / O ）装置に関連付けられたオンチップメモリ上で定義されたコンテキストテーブルを含み、コンテキストテーブルは、ディスク読取操作のために、空きバッファプールから割当てられた1つ以上のパケットバッファに指向するバッファポインタの一時リンクリストを保存する。I / O 装置は、ディスク読取操作を実行する物理ホストバスアダプタ（H B A ）からディスク読取データを受信すると、コンテキストテーブルを開き、バッファポインタの一時リンクリストを更新し、およびコンテキストテーブルが閉じられると、バッファポインタの一時リンクリストをバッファポインタの主要リンクリストに合併するように動作

10

20

30

40

50

する。

【 0 1 1 2 】

上記に提供されたシステムにおいて、I/O装置は、ディスク読取操作を開始するように、インフィニバンド（IB）ファブリック上のサーバを動作させる。

【 0 1 1 3 】

上記に提供されたシステムにおいて、I/O装置は、IBヘッダとシーケンス番号とを物理HBAから受信した各パケットに追加する。

【 0 1 1 4 】

上記に提供されたシステムにおいて、I/O装置は、完全メッセージまたはIB最大伝送ユニット（MTU）パケットを受信すると、外部メモリに格納されたディスク読取データを読出すように動作する。

10

【 0 1 1 5 】

上記に提供されたシステムにおいて、各仮想HBAは、IBドメインにおいて、異なるパケットシーケンス番号（PSN）スペースを保留する。

【 0 1 1 6 】

上記に提供されたシステムにおいて、I/O装置は、vHBAに関連付けられた異なるディスク読取操作のために、異なるコンテキストテーブルを保存する。

【 0 1 1 7 】

上記に提供されたシステムにおいて、コンテキストテーブルは、I/O装置が仮想HBAに関連付けられた別のコンテキストテーブルを開くと、閉じられる。

20

【 0 1 1 8 】

上記に提供されたシステムにおいて、別のコンテキストテーブルは、空きバッファプールから割当てられた1つ以上のパケットバッファに指向するバッファポインタの新たな一時リンクリストを保存する。

【 0 1 1 9 】

上記に提供されたシステムにおいて、I/O装置は、仮想HBAが別のディスク読取操作からデータを受信すると、別のコンテキストテーブルを開く。

【 0 1 2 0 】

上記に提供されたシステムにおいて、I/O装置は、仮想HBAが書込専用命令、RDMA書込最終命令、送信専用命令およびRDMA読取リクエスト命令のうち1つの命令を受信すると、別のコンテキストテーブルを開く。

30

【 0 1 2 1 】

本発明の一実施形態は、コンピューティング環境において、効率的なパケット処理をサポートするための方法を提供する。この方法は、コンピューティング環境内の1つ以上の仮想ホストバスアダプタ（vHBA）に関連付けられた複数のパケットバッファを含む空きバッファプールを提供するステップを含み、vHBAの各々は、1つ以上のパケットバッファに指向するバッファポインタの主要リンクリストを空きバッファプールに保存し、入力/出力（I/O）装置に関連付けられたオンチップメモリにおいて、コンテキストテーブルを定義するステップを含み、コンテキストテーブルは、ディスク読取操作のために、空きバッファプールから割当てられた1つ以上のパケットバッファに指向するバッファポインタの一時リンクリストを保存し、ディスク読取操作を実行する物理ホストバスアダプタ（HBA）からディスク読取データを受信すると、コンテキストテーブルを開き、バッファポインタの一時リンクリストを更新し、およびコンテキストテーブルが閉じられると、バッファポインタの一時リンクリストをバッファポインタの主要リンクリストに合併するように、I/O装置を動作させるステップを含む。

40

【 0 1 2 2 】

上記に提供された方法は、ディスク読取操作を開始するように、インフィニバンド（IB）ファブリック上のサーバを動作させる。

【 0 1 2 3 】

上記に提供された方法は、IBヘッダとシーケンス番号とをHBAから受信した各パケ

50

ットに追加するステップをさらに含む。

【0124】

上記に提供された方法は、完全メッセージまたはIB最大伝送ユニット(MTU)パケットを受信すると、外部メモリに格納されたディスク読取データを読み出すステップをさらに含む。

【0125】

上記に提供された方法は、IBドメインにおいて、異なるパケットシーケンス番号(PSN)スペースを保留するように、各仮想HBAを構成するステップをさらに含む。

【0126】

上記に提供された方法は、vHBAに関連付けられた異なるディスク読取操作のために、異なるコンテキストテーブルを保存するステップをさらに含む。

10

【0127】

上記に提供された方法は、I/O装置が仮想HBAに関連付けられた別のコンテキストテーブルを開くと、コンテキストテーブルを閉じるステップをさらに含む。

【0128】

上記に提供された方法は、空きバッファプールから割当てられた1つ以上のパケットバッファに指向するバッファポインタの新たな一時リンクリストを保存するように、別のコンテキストテーブルを構成するステップをさらに含む。

【0129】

上記に提供された方法は、仮想HBAが別のディスク読取操作からデータ、または書込専用命令、RDMA書込最終命令、送信専用命令およびRDMA読取リクエスト命令のうち1つの命令を受信すると、別のコンテキストテーブルを開くステップをさらに含む。

20

【0130】

本発明の一実施形態は、命令を格納する非一時的な機械読取可能記憶媒体を提供する。これらの命令は、実行されると、以下のステップをシステムに実行させ、当該以下のステップは、ネットワーク環境内のつ以上の仮想ホストバスアダプタ(vHBA)に関連付けられた複数のパケットバッファを含む空きバッファプールを提供するステップを含み、vHBAの各々は、1つ以上のパケットバッファに指向するバッファポインタの主要リンクリストを空きバッファプールに保存し、入力/出力(I/O)装置に関連付けられたオンチップメモリにおいて、コンテキストテーブルを定義するステップを含み、コンテキストテーブルは、ディスク読取操作のために、空きバッファプールから割当てられた1つ以上のパケットバッファに指向するバッファポインタの一時リンクリストを保存し、ディスク読取操作を実行する物理ホストバスアダプタ(HBA)からディスク読取データを受信すると、コンテキストテーブルを開き、バッファポインタの一時リンクリストを更新し、およびコンテキストテーブルが閉じられると、バッファポインタの一時リンクリストをバッファポインタの主要リンクリストに合併するように、I/O装置を動作させるステップを含む。

30

【0131】

本発明の一実施形態は、コンピューティング環境において、入力/出力(I/O)仮想化をサポートするためのシステムを提供する。このシステムは、メモリ内の空きバッファプールを含み、空きバッファプールは、2次元リンクリストおよび1次元リンクリストを備え、2次元リンクリストの各エントリは、連続したメモリ位置で複数のパケットバッファを含み、1次元リンクリストの各エントリは、単一のパケットバッファを含み、I/O装置は、空きバッファプールを用いて、物理ホストバスアダプタ(HBA)から受信したディスク読取データを保存するように動作する。

40

【0132】

上記に提供されたシステムにおいて、I/O装置は、ディスク読取操作を開始するように、インフィニバンド(IB)ファブリック上のサーバを動作させる。

【0133】

上記に提供されたシステムにおいて、I/O装置は、1つ以上のIBヘッダとシーケン

50

ス番号とを物理 H B A から受信した各パケットに追加する。

【 0 1 3 4 】

上記に提供されたシステムにおいて、I / O 装置は、完全メッセージまたは I B 最大伝送ユニット (M T U) パケットを受信すると、外部メモリに格納されたディスク読取データを読み出すように動作する。

【 0 1 3 5 】

上記に提供されたシステムにおいて、I / O 装置は、1 つ以上の仮想ホストバスアダプタ (v H B A) をサポートしており、各 v H B A は、I B ドメインにおいて、異なるパケットシーケンス番号 (P S N) スペースを保留する。

【 0 1 3 6 】

上記に提供されたシステムにおいて、I / O 装置は、パケットがリモートダイレクトメモリアクセス (R D M A) 書込トランザクションまたは R D M A 読取リクエストトランザクションのいずれかを実行する場合、このパケットを外部メモリにキュー入れするように動作する。

【 0 1 3 7 】

上記に提供されたシステムにおいて、I / O 装置は、キューに入れられたパケットが R D M A 書込トランザクション用のものである場合、2 次元リンクリストからスーパーブロックを割当てるように動作する。

【 0 1 3 8 】

上記に提供されたシステムにおいて、I / O 装置は、I / O 装置は、キューに入れられたパケットが R D M A 読取リクエストトランザクション用のものであり、且つ、スーパーブロックに 1 つ以上のパケットバッファが残された場合、1 次元リンクリストからパケットバッファを割当てるように動作する。

【 0 1 3 9 】

上記に提供されたシステムにおいて、I / O 装置は、キューに入れられたパケットが R D M A 読取リクエストトランザクション用のものであり、且つ、スーパーブロックに 1 つ以上のパケットバッファが残された場合、2 次元リンクリストからスーパーブロックを割当てるように動作する。

【 0 1 4 0 】

上記に提供されたシステムにおいて、I / O 装置は、v H B A のために、単純リンクリストにおいてパケットバッファを保留するように動作する。

【 0 1 4 1 】

本発明の一実施形態は、ネットワーク環境において、効率的なパケット処理をサポートするための方法を提供する。この方法は、メモリに空きバッファプールを提供するステップを含み、空きバッファプールは、2 次元リンクリストおよび 1 次元リンクリストを含み、2 次元リンクリスト各エントリを連続したメモリ位置で複数のパケットバッファを含み、および 1 次元リンクリストの各エントリを単一のパケットバッファを含むようにするステップを含み、I / O 装置を介して、空きバッファプールを用いて、物理ホストバスアダプタ (H B A) から受信したデータ読取ディスクを保存するステップを含む。

【 0 1 4 2 】

上記に提供された方法は、ディスク読取操作を開始するように、インフィニバンド (I B) ファブリック上のサーバを動作させるステップをさらに含む。

【 0 1 4 3 】

上記に提供された方法は、1 つ以上の I B ヘッドとシーケンス番号とを物理 H B A から受信した各パケットに追加するステップをさらに含む。

【 0 1 4 4 】

上記に提供された方法は、完全メッセージまたは I B 最大伝送ユニット (M T U) パケットを受信すると、外部メモリに格納されたディスク読取データを読み出すステップをさらに含む。

【 0 1 4 5 】

10

20

30

40

50

上記に提供された方法は、1つ以上の仮想ホストバスアダプタ（vHBA）をサポートするステップをさらに含み、各vHBAは、IBドメインにおいて、異なるパケットシーケンス番号（PSN）スペースを保留する。

【0146】

上記に提供された方法は、パケットがリモートダイレクトメモリアクセス（RDMA）書込トランザクションまたはRDMA読取リクエストトランザクションのいずれかを実行する場合、このパケットを外部メモリにキュー入れするステップをさらに含む。

【0147】

上記に提供された方法は、キューに入れられたパケットがRDMA書込トランザクション用のものである場合、2次元リンクリストからスーパーブロックを割当てするステップをさらに含む。

10

【0148】

上記に提供された方法は、キューに入れられたパケットがRDMA読取リクエストトランザクション用のものであり、且つ、スーパーブロックに1つ以上のパケットバッファが残された場合、1次元リンクリストからパケットバッファを割当てするステップをさらに含む。

【0149】

上記に提供された方法は、キューに入れられたパケットがRDMA読取リクエストトランザクション用のものであり、且つ、スーパーブロックに1つ以上のパケットバッファが残された場合、2次元リンクリストからスーパーブロックを割当てするステップと、

20

vHBAのために、単純リンクリストにおいてパケットバッファを保留するステップとをさらに含む。

【0150】

本発明の一実施形態は、命令を格納する非一時的な機械読取可能記憶媒体を提供する。命令は、実行されると、以下のステップをシステムに実行させ、当該以下のステップは、メモリに空きバッファプールを提供するステップを含み、空きバッファプールは、2次元リンクリストおよび1次元リンクリストを含み、2次元リンクリスト各エントリを連続したメモリ位置で複数のパケットバッファを含み、および1次元リンクリストの各エントリを単一のパケットバッファを含むようにするステップを含み、I/O装置を介して、空きバッファプールを用いて、物理ホストバスアダプタ（HBA）から受信したデータ読取ディスクを保存するステップを含む。

30

【0151】

本発明の多くの特徴は、ハードウェア、ソフトウェア、ファームウェア、またはそれらの組合せの内部で、またはそれらを用いて、またはそれらの援助をもって実現されることができる。したがって、本発明の特徴は、（たとえば、1つ以上のプロセッサを含む）処理システムを用いて実施されることができる。

【0152】

本発明の特徴は、コンピュータプログラム製品の内部で、またはそれを用いて、またはその援助をもって実現されることができる。コンピュータプログラム製品は、本明細書に記載の特徴のいずれかを実現するように、処理システムをプログラムさせるために使用することができる命令をその上／中に格納する記憶媒体またはコンピュータ読取可能媒体である。記憶媒体は、フロッピーディスク（登録商標）、光ディスク、DVD、CD-ROM、マイクロドライブおよび光磁気ディスクを含む任意の種類のディスク、ROM、RAM、EPROM、EEPROM、DRAM、VRAM、フラッシュメモリデバイス、磁気または光カード、（分子メモリICを含む）ナノシステム、または指令および／またはデータの格納に適した任意の種類の媒体またはデバイスを含むことができるが、これらに限定されない。

40

【0153】

機械読取可能媒体のいずれかに格納された本発明の特徴は、処理システムのハードウェアを制御するため、および本発明の結果を利用する他の機構と対話できる処理システムを

50

可能にするために、ソフトウェアおよび／またはファームウェアに組込むことができる。このようなソフトウェアまたはファームウェアは、アプリケーションコード、デバイスドライバ、オペレーティングシステム、実行環境／コンテナ含むことができるが、これらに限定されない。

【0154】

また、本発明の特徴は、たとえば特定用途向け集積回路（ASIC）のようなハードウェア部品を用いて、ハードウェアにおいて実現されてもよい。本明細書に記載の機能を実行するようにハードウェア状態マシンを実装することは、当業者には明らかであろう。

【0155】

さらに、本発明は、1つ以上の従来の汎用または専用デジタルコンピュータ、コンピューティング装置、コンピューティング機械、または1つ以上のプロセッサ、メモリおよび／または本開示の教示に従ってプログラムされたコンピュータ読取可能記憶媒体を含むマイクロプロセッサを用いて、簡便に実施することができる。ソフトウェア分野の当業者には明らかなように、本開示の教示に基づいて、熟練したプログラマは、適切なソフトウェアコーディングを容易に用意することができる。

10

【0156】

上記で本発明のさまざまな実施形態を説明したが、これらの実施形態は、限定の目的ではなく、例示として提示されていることが理解すべきである。本発明の精神および範囲から逸脱することなく、本発明に形式上および詳細上のさまざまな変更を行うことができることは、当業者には明らかであろう。

20

【0157】

本発明は、特定の機能およびそれらの関係を示す機能的構造ブロックを用いて説明しました。説明の便宜のために、これらの機能的構造ブロックは、多くの場合、本明細書において任意に定義されている。特定の機能およびそれらの関係が適切に実行される限り、代替的な構造ブロックを定義することができる。任意の代替的な構造ブロックは、本発明の範囲および精神に含まれる。

【0158】

本発明の上記説明は、例示および説明のために提供されている。本発明を網羅的であることにまたは開示された形態に厳密に限定することを意図するものではない。本発明の幅および範囲は、上述した例示的な実施形態のいずれかに限定されない。多くの修正および変更は、当業者にとって明らかであろう。修正および変更は、開示された特徴の任意の適切な組合せを含む。実施形態は、本発明の原理およびその実際の応用を最善に説明するために選択され説明された。よって、当業者は、さまざまな実施形態により本発明を理解し、考えられる特定の用途に適したさまざまな修正を行うことができる。なお、本発明の範囲が添付の特許請求の範囲およびその等価物によって定義されることが意図される。

30

【図 1】

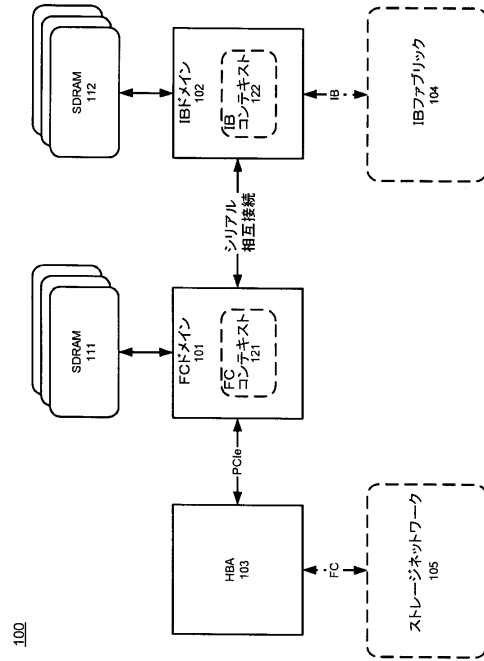


FIGURE 1

【図 2】

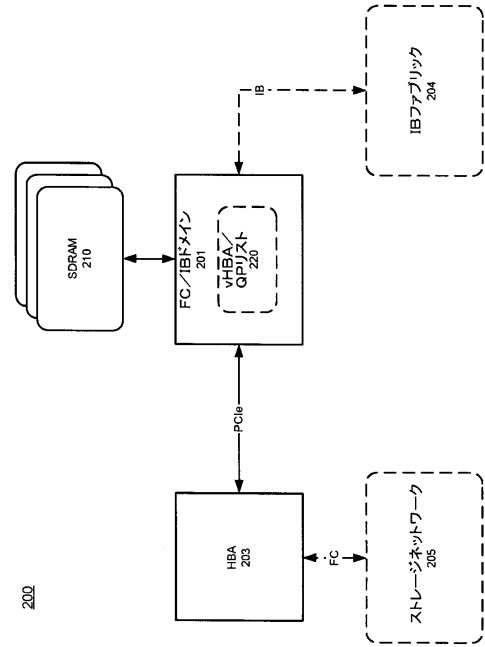


FIGURE 2

【図 3】

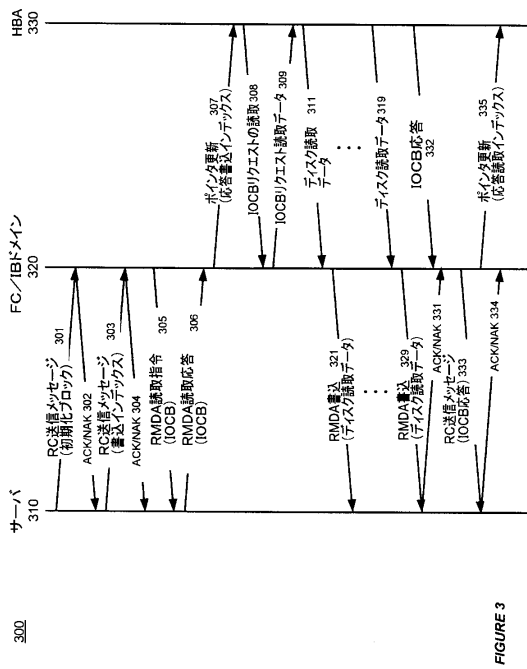


FIGURE 3

【図 4】

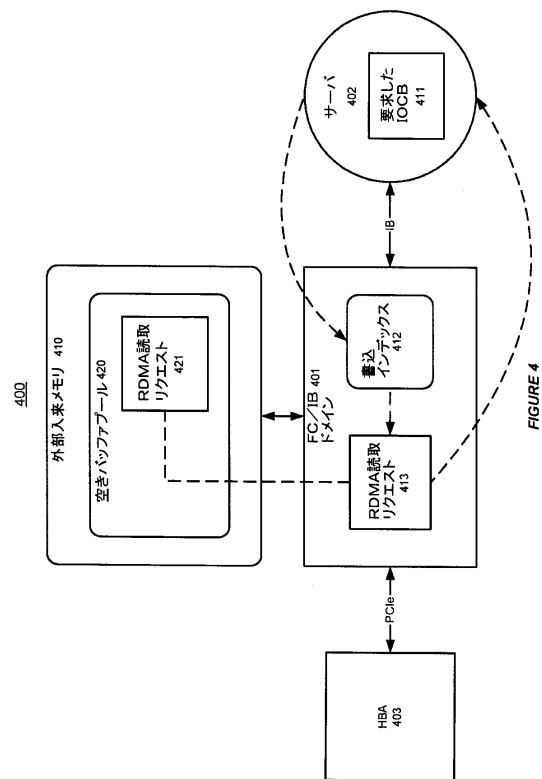
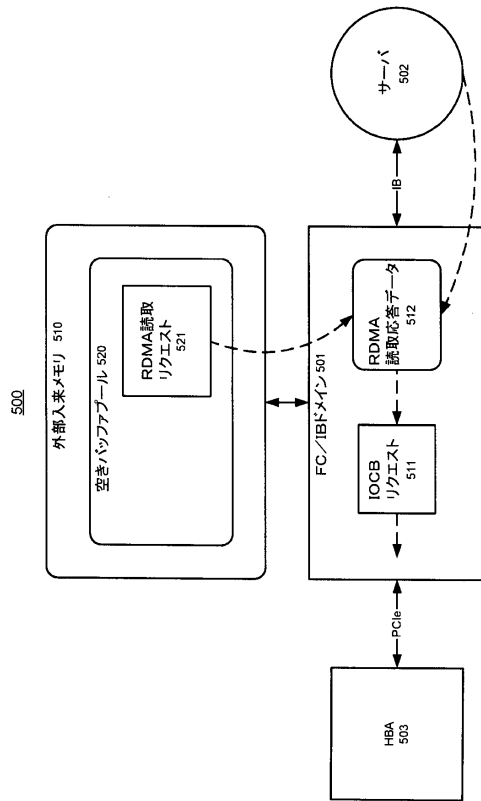
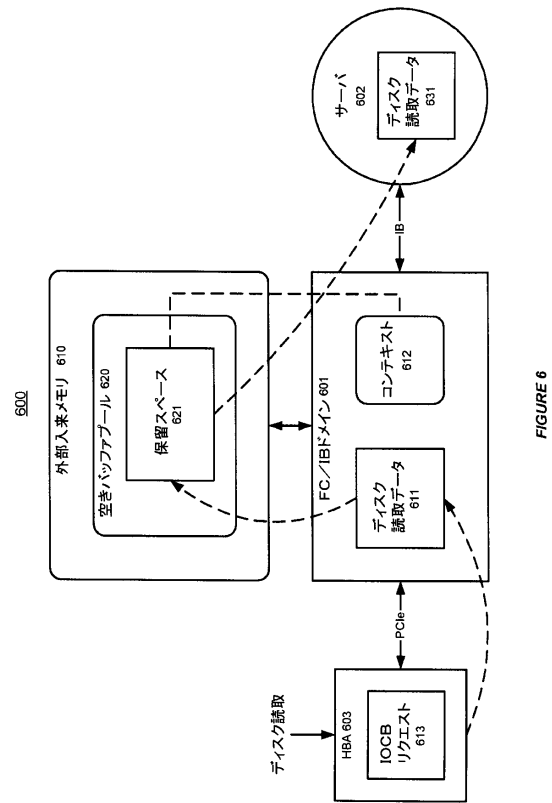


FIGURE 4

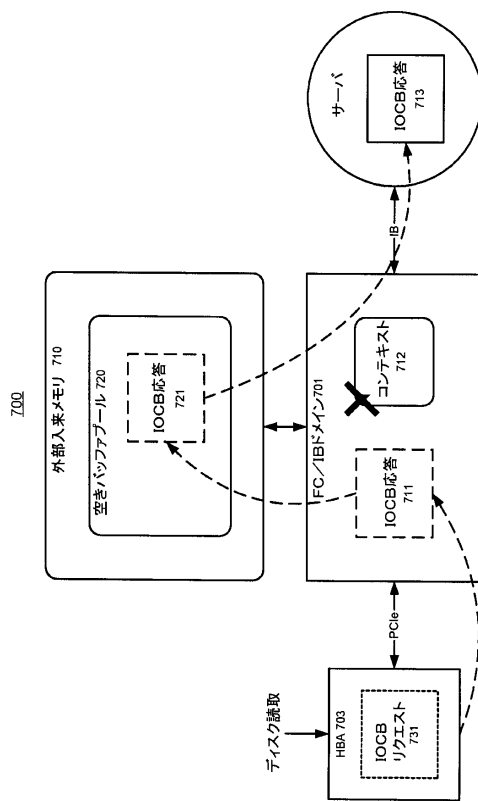
【図 5】



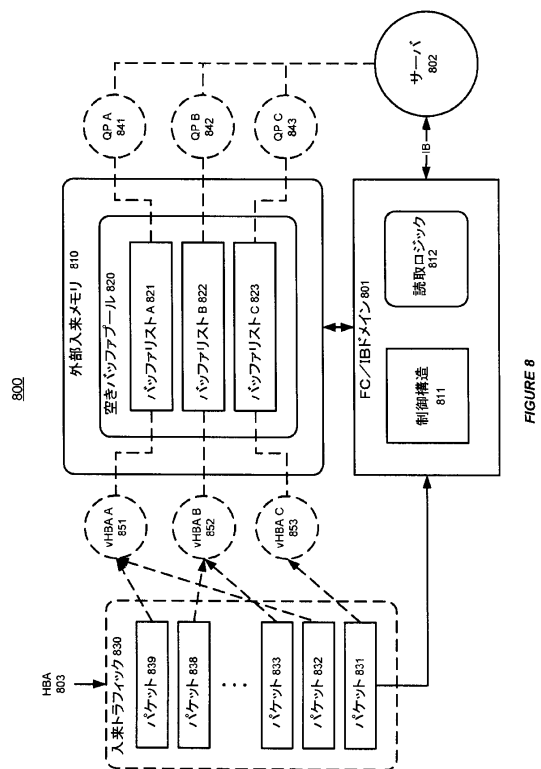
【図 6】



【図 7】



【図 8】



【図 9】

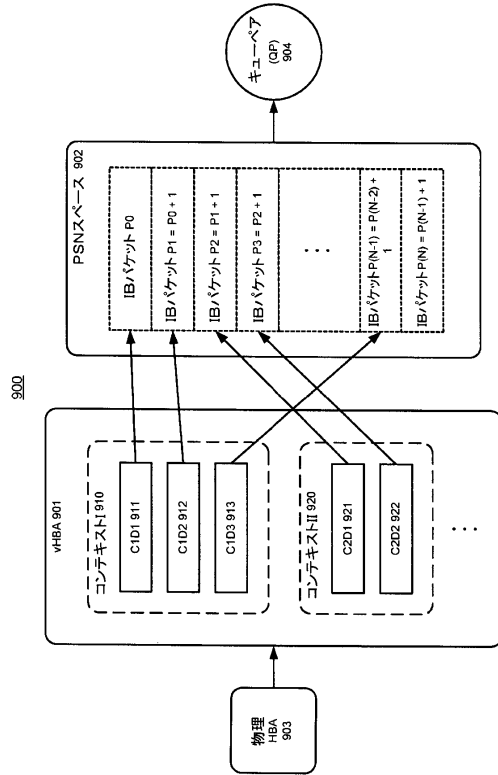


FIGURE 9

【図 10】

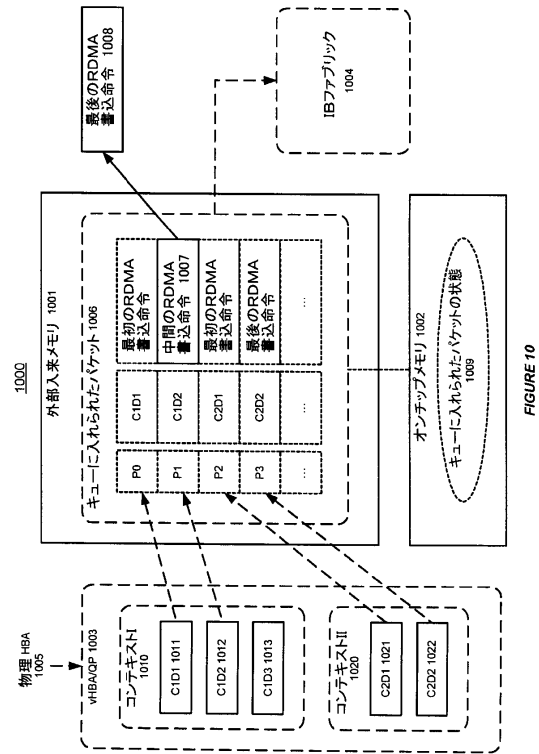


FIGURE 10

【図 11】

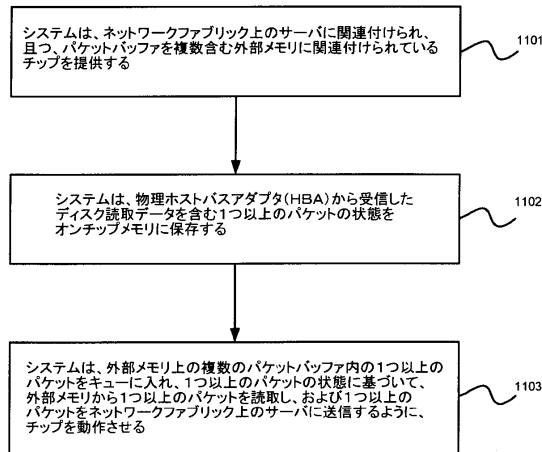


FIGURE 11

【図 12】

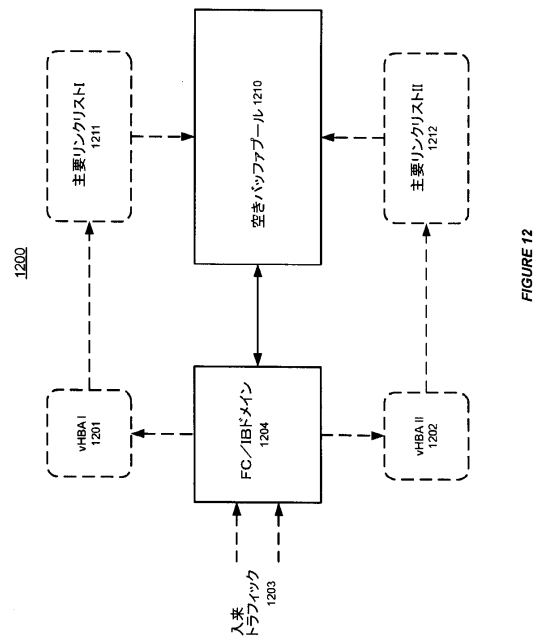


FIGURE 12

【図 13】

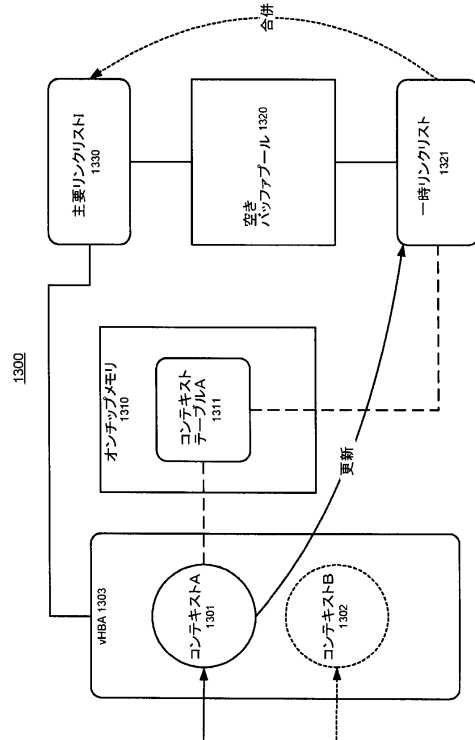


FIGURE 13

【図 14】

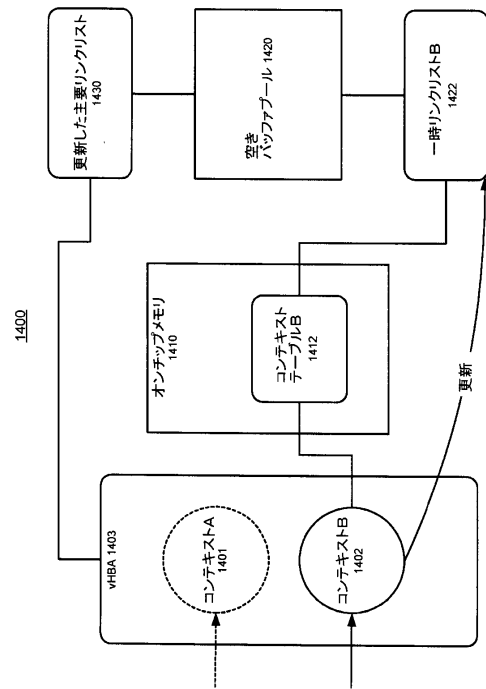


FIGURE 14

【図 15】

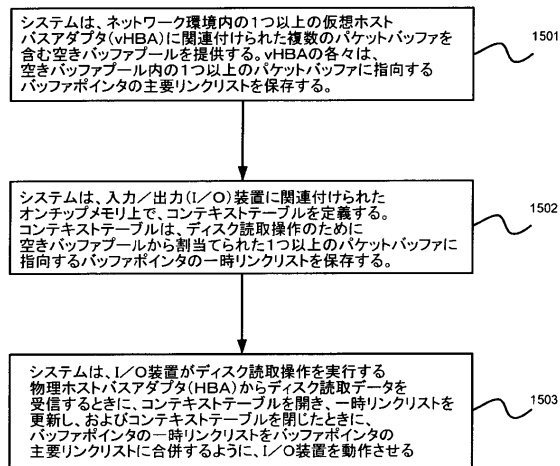


FIGURE 15

【図 16】

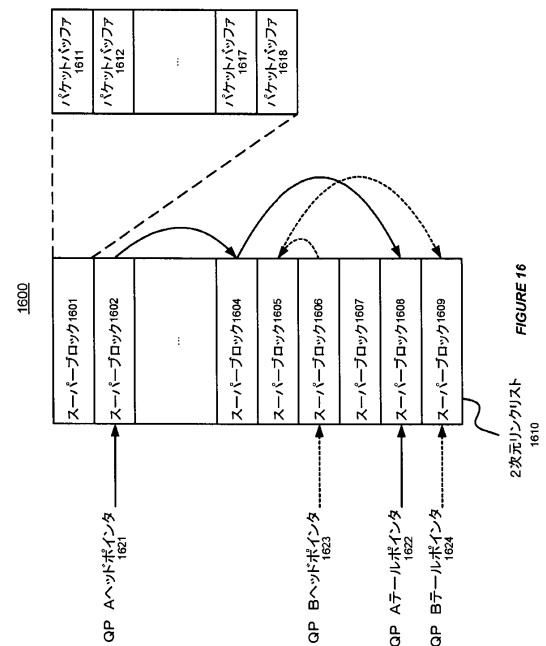


FIGURE 16

【図 17】

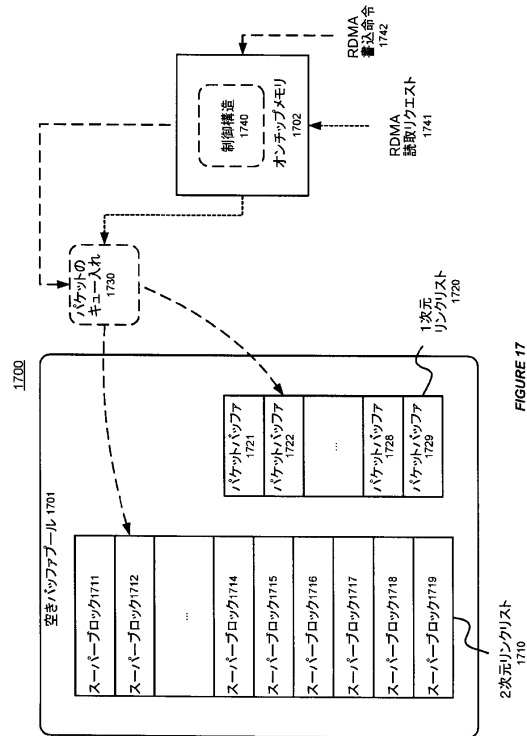


FIGURE 17

【図 18】

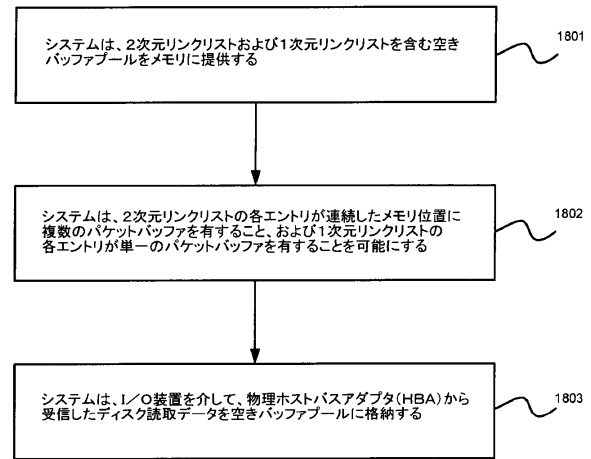


FIGURE 18

フロントページの続き

(51)Int.Cl. F I
H 0 4 L 12/861

(31)優先権主張番号 14/097,009

(32)優先日 平成25年12月4日(2013.12.4)

(33)優先権主張国 米国(US)

(56)参考文献 米国特許出願公開第2013/0138836(US,A1)
米国特許第08458306(US,B1)
特開2005-216283(JP,A)

(58)調査した分野(Int.Cl.,DB名)

G 0 6 F 1 3 / 0 0 - 1 3 / 1 4
G 0 6 F 1 3 / 2 0 - 1 3 / 4 2
G 0 6 F 1 5 / 1 6 - 1 5 / 1 7 7
H 0 4 L 1 2 / 8 6 1

(54)【発明の名称】インフィニバンド(I B)上で仮想ホストバスアダプタ(v H B A)を管理およびサポートするためのシステムおよび方法、ならびに単一の外部メモリアンターフェイスを用いてバッファの効率的な使用をサポートするためのシステムおよび方法