

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5272265号
(P5272265)

(45) 発行日 平成25年8月28日(2013.8.28)

(24) 登録日 平成25年5月24日(2013.5.24)

(51) Int.Cl.

F I

G 0 6 F 13/10 (2006.01)

G 0 6 F 13/10 3 3 0 C

請求項の数 13 (全 30 頁)

(21) 出願番号	特願2008-250208 (P2008-250208)	(73) 特許権者	000005108
(22) 出願日	平成20年9月29日(2008.9.29)		株式会社日立製作所
(65) 公開番号	特開2010-79816 (P2010-79816A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成22年4月8日(2010.4.8)	(74) 代理人	100114236
審査請求日	平成23年6月16日(2011.6.16)		弁理士 藤井 正弘
		(74) 代理人	100075513
			弁理士 後藤 政喜
		(74) 代理人	100120260
			弁理士 飯田 雅昭
		(72) 発明者	馬場 貴成
			東京都国分寺市東恋ヶ窪一丁目280番地
			株式会社日立製作所 中央研究所内
		(72) 発明者	森木 俊臣
			東京都国分寺市東恋ヶ窪一丁目280番地
			株式会社日立製作所 中央研究所内
			最終頁に続く

(54) 【発明の名称】 P C I デバイス共有方法

(57) 【特許請求の範囲】

【請求項 1】

仮想化部を備えて1以上の仮想サーバを提供する複数の物理サーバと、前記複数の物理サーバで共有されるI/Oデバイスと、前記複数の物理サーバと前記I/Oデバイスを接続するスイッチと、前記スイッチを初期化する管理部とを有するサーバシステムであって、

前記I/Oデバイスが、

仮想的な機能であるVirtual Function(VF)を有し、

前記スイッチが、

前記物理サーバが認識する識別子である第1の識別子と、前記管理部が管理する識別子である第2の識別子と、の間の対応関係を示す識別子対応情報を保持し、

前記物理サーバからパケットを受信した場合に、

送信元の物理サーバを識別するサーバ識別子を含む前記パケットを前記スイッチ内で変換し、

前記パケットを前記I/Oデバイスに送信する際には、前記識別子対応情報に基づいて、前記パケットに含まれる前記第1の識別子を、当該第1の識別子に前記対応付けられた前記第2の識別子に変換し、かつ、前記サーバ識別子を削除した前記パケットを、前記I/Oデバイスに送信する、ことを特徴とするサーバシステム。

【請求項 2】

請求項1記載のサーバシステムであって、

10

20

前記物理サーバは、前記第1の識別子に基づいて、前記仮想サーバから前記VFへのアクセスを管理し、

前記I/Oデバイスは、前記第2の識別子に基づいて、前記仮想サーバから前記VFへのアクセスを管理する、ことを特徴とするサーバシステム。

【請求項3】

請求項1記載のサーバシステムであって、

前記スイッチは、

更に前記物理サーバが認識する前記VFのMMIOアドレス領域情報と、前記管理部が管理する前記VFのMMIOアドレス領域情報との間の対応関係を示すMMIOアドレス対応情報を保持し、

前記物理サーバからパケットを受信した場合に、更に、前記MMIOアドレス対応情報に従って前記パケットの宛先アドレスを変換する処理を行う、ことを特徴とするサーバシステム。

【請求項4】

請求項1記載のサーバシステムであって、

前記I/Oデバイスは、PCI sig規定のSingle Root I/O Virtualization and Sharing Specification (SR-IOV) に準拠、もしくはSR-IOVの上位互換のI/Oデバイスであることを特徴とするサーバシステム。

【請求項5】

仮想化部を備えて1以上の仮想サーバを提供する複数の物理サーバと、仮想的な機能であるVirtual Function (VF) を有して前記複数の物理サーバで共有されるI/Oデバイスと、を接続するスイッチング装置であって、

前記物理サーバと接続する第1のポートと、

前記I/Oデバイスと接続する第2のポートと、

前記物理サーバが認識する識別子である第1の識別子と、前記スイッチを初期化する管理部が管理する識別子である第2の識別子と、の間の対応関係を管理する識別子管理部と、

前記物理サーバにより送信された前記I/Oデバイス宛のパケットの前記I/Oデバイスへの転送を制御する転送制御部と、を備え、

前記第1のポートが前記物理サーバから前記I/Oデバイス宛のパケットを受信した場合に、

前記転送制御部は、送信元の物理サーバを識別するサーバ識別子を含む前記パケットを前記スイッチ内で変換し、

前記パケットを前記第2のポートから前記I/Oデバイスに送信する際には、前記識別子管理部が管理する前記対応関係に基づいて、前記パケットに含まれる前記第1の識別子を、当該第1の識別子に前記対応付けられた前記第2の識別子に変換し、かつ、前記サーバ識別子を削除した前記パケットを、前記I/Oデバイスに送信する、ことを特徴とするスイッチング装置。

【請求項6】

請求項5記載のスイッチング装置であって、

前記第1の識別子は、前記仮想サーバから前記VFへのアクセスの管理のために前記物理サーバが用いる識別子であり、

前記第2の識別子は、前記仮想サーバから前記VFへのアクセスの管理のために前記I/Oデバイスが用いる識別子である、ことを特徴とするスイッチング装置。

【請求項7】

請求項5記載のスイッチング装置であって、

前記物理サーバが認識する前記VFのMMIOアドレス領域情報と、前記管理部が管理する前記VFのMMIOアドレス領域情報と、の間の対応関係を管理するMMIOアドレス対応関係管理部を、更に有し、

前記第 1 のポートが前記物理サーバから前記 I / O デバイス宛のパケットを受信した場合に、

前記転送制御部は、更に、前記 M M I O アドレス対応管理部に従って前記パケットの宛先アドレスを変換する処理を行う、ことを特徴とするスイッチング装置。

【請求項 8】

仮想化部を備えて 1 以上の仮想サーバを提供する複数の物理サーバと、I / O デバイスと、前記複数の物理サーバと前記 I / O デバイスとを接続するスイッチと、前記スイッチを初期化する管理部と、を有するサーバシステムにおいて、前記 I / O デバイスを前記複数の物理サーバ間で共有する I / O デバイス管理方法であって、

前記 I / O デバイスによって、仮想的な機能である Virtual Function (V F) を複数生成し、

前記スイッチによって、

前記物理サーバが認識する前記 V F の識別子である第 1 の識別子と、前記管理部が管理する前記 V F の識別子である第 2 の識別子と、の間の対応関係を示す識別子対応情報を管理し、

前記スイッチが前記物理サーバからパケットを受信した場合に、送信元の物理サーバを識別するサーバ識別子 (= V H N) を含む前記パケットを前記スイッチ内で変換し、

前記パケットを前記 I / O デバイスに送信する際には、前記識別 x 子対応情報に基づいて、前記パケットに含まれる前記第 1 の識別子を、当該第 1 の識別子に前記対応付けられた前記第 2 の識別子に変換し、かつ、前記サーバ識別子を削除した前記パケットを、前記 I / O デバイスに送信する、ことを特徴とする I / O デバイス管理方法。

【請求項 9】

請求項 8 記載の I / O デバイス管理方法であって、

前記第 1 の識別子に基づいて、前記仮想サーバから前記 V F へのアクセスに関して前記物理サーバによる管理を行い、

前記第 2 の識別子に基づいて、前記仮想サーバから前記 V F へのアクセスに関して前記 I / O デバイスによる管理を行う、ことを特徴とする I / O デバイス管理方法。

【請求項 10】

請求項 8 記載の I / O デバイス管理方法であって、

前記スイッチによって、

更に前記物理サーバが認識する前記 V F の M M I O アドレス領域情報と、前記管理部が管理する前記 V F の M M I O アドレス領域情報と、の間の対応関係を示す M M I O アドレス対応情報を管理し、

前記物理サーバからパケットを受信した場合に、更に、前記 M M I O アドレス対応情報に従って前記パケットの宛先アドレスを変換する処理を行う、ことを特徴とした、I / O デバイス管理方法。

【請求項 11】

請求項 3 記載のサーバシステムであって、

前記物理サーバが認識する前記 V F の M M I O アドレス領域情報は、前記仮想サーバがアクセスする前記 V F の M M I O 空間のアドレスのオフセットであって、

前記管理部が管理する前記 V F の M M I O アドレス領域情報は、前記管理部がアクセスする前記 V F の M M I O 空間のアドレスのオフセットであって、

前記サーバ識別子は、送信元の物理サーバを識別する仮想階層ナンバーであることを特徴とするサーバシステム。

【請求項 12】

請求項 7 記載のスイッチング装置であって、

前記物理サーバが認識する前記 V F の M M I O アドレス領域情報は、前記仮想サーバがアクセスする前記 V F の M M I O 空間のアドレスのオフセットであって、

前記管理部が管理する前記 V F の M M I O アドレス領域情報は、前記管理部がアクセスする前記 V F の M M I O 空間のアドレスのオフセットであって、

前記サーバ識別子は、送信元の物理サーバを識別する仮想階層ナンバーであることを特徴とするスイッチング装置。

【請求項 13】

請求項 10 記載の I / O デバイス管理方法であって、

前記物理サーバが認識する前記 V F の M M I O アドレス領域情報は、前記仮想サーバがアクセスする前記 V F の M M I O 空間のアドレスのオフセットであって、

前記管理部が管理する前記 V F の M M I O アドレス領域情報は、前記管理部がアクセスする前記 V F の M M I O 空間のアドレスのオフセットであって、

前記サーバ識別子は、送信元の物理サーバを識別する仮想階層ナンバーであることを特徴とする I / O デバイス管理方法。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数の計算機を備えたブレードサーバに関し、特に、ひとつの I / O デバイスを複数の計算機で共有する技術に関する。

【背景技術】

【0002】

情報セキュリティやコンプライアンスへの意識の高まりから、サーバサイドでのウィルスチェックやメールフィルタなど、企業の情報システムに求められる処理要求が増大している。これらの処理要求の増大に対し、従来は処理内容ごとに個別にサーバを導入して対応してきた。しかし、サーバ台数の増大は運用コストの増大を招き、企業の IT 予算を圧迫して問題化している。

20

【0003】

この問題に対し、複数のサーバで実行されていた処理を 1 台の高性能サーバに集約して、情報システムを構成するサーバ台数を削減するサーバ統合が注目されている。サーバ統合ではサーバの台数に比例して継続して発生する消費電力やスペース(占有床面積)、ハードウェア障害時のメンテナンスコスト等を削減できる。

【0004】

サーバ統合を実現する手段として、C P U を高密度に集積したブレードサーバが台頭している(例えば、特許文献 1)。典型的なブレードサーバは、C P U (プロセッサ)、メモリ及び N I C (Network Interface Card) を搭載した複数のブレードと、ネットワークスイッチ、および拡張用 I / O スロットを搭載した I / O ドロワーが 1 つの筐体に格納された構成をとる。ブレードサーバを用いると、ネットワーク経由で処理を行うサーバに対して効果的にサーバ統合を実現できる。

30

【0005】

一方、ブレードサーバは I / O ドロワーの拡張用 I / O スロットは、ブレードと I / O スロットとの対応が固定的であった。このため、N I C 以外の I / O カードを稀にしか使用しない、もしくは全く使用しないブレードに対しても I / O スロットが予約されており、スペースに無駄が生じるという問題があった。また、ブレードの用途によっては多数の I / O カードを必要とする場合(例; 複数ブレード間での H o t - S t a n d b y 構成等)に、I / O ドロワー内の I / O スロットの数を超えて I / O スロットを割り当てることができない、という問題があった。

40

【0006】

ブレードサーバのように多数のブレードで限られた数の I / O スロット(または I / O デバイス)を利用する技術として、ひとつの I / O デバイスを複数のホストで利用する技術が提案されている。

【0007】

このような技術としては、P C I - S i g で標準化が行われているシングルルート I / O 仮想化技術(S R - I O V : Single Root IO Virtualization)やマルチルート I / O 仮想化技術(M R - I O V : Multi Root IO Virtualization)が知られている(例えば、

50

特許文献2、非特許文献1等)。

【0008】

この他、任意のサーバ間でストレージ装置を共有する技術としては、iSCSIが広く知られている。iSCSIは、サーバとストレージ装置のアクセスに使用されるSCSIコマンドを、ネットワーク通信のプロトコルであるTCP/IPパケットにカプセル化する技術である。これにより任意のサーバ間でiSCSI対応ターゲットデバイスを共有できる。

【特許文献1】特開特開2002-32153号

【特許文献2】米国特許第7058738号

【非特許文献1】「Single-Root I/O Virtualization and Sharing Specification, Revision 1.0, Section 1 - Architectural」、2008年発行、著者PCI-SIG、第15～34頁

10

【発明の開示】

【発明が解決しようとする課題】

【0009】

上記従来例のシングルルートI/O仮想化技術(以下、SR-IOV)では、1つのI/Oデバイスの物理機能(PF; Physical Function)が複数の仮想機能(VF; Virtual Function)を提供することができる。同一ブレード内の複数のOSインスタンス間で、複数の仮想機能をそれぞれ占有することでひとつのSR-IOV対応のI/Oデバイスを共有することができる。

20

【0010】

一方、上記従来例のマルチルートI/O仮想化技術(以下、MR-IOV)では、1つのI/Oデバイスに複数の物理機能(PF; Physical Function)を搭載し、各物理機能上で複数の仮想機能(VF; Virtual Function)を提供することができる。MR-IOV対応のI/Oデバイスを利用し、仮想化ソフトウェアを導入することで、異なる物理サーバ上で動作する仮想計算機(VM; Virtual Machine)間で1つのI/Oデバイスを共有することができる。すなわち、MR-IOVでは、特定のブレード(物理計算機)に対して1つのPFを占有させ、各VMはPFが提供する複数のVFを占有する。

【0011】

ここで、SR-IOVに対応するI/Oデバイスでは、一つの物理機能を備えればよいので、MR-IOVに対応するI/Oデバイスに低コストで調達することができる。

30

【0012】

上記従来例のブレードサーバにおいて複数のブレードでI/Oデバイスを共有する場合、上記SR-IOVに対応するI/Oデバイスを用いれば、I/Oデバイスの調達コストを抑制できるものの、ひとつのブレードに対してひとつのI/Oデバイスを割り当てることになり、複数のブレード間でひとつのI/Oデバイスを共有できない、という問題があった。

【0013】

一方、従来例のブレードサーバにおいて複数のブレードでI/Oデバイスを共有する場合、上記MR-IOVに対応するI/Oデバイスを用いれば、複数のブレード間でひとつのI/Oデバイスを共有することができる。しかしながら、MR-IOVに対応するI/OデバイスはSR-IOVに対応するI/Oデバイスに対して調達コストが高いため導入コストが増大する、という問題がある。

40

【0014】

なお、上記iSCSIをブレードサーバで利用した場合、ストレージ装置の共有は可能であるが、I/OデバイスがNIC等の場合は利用することができない。

【0015】

そこで本発明は、上記問題点に鑑みてなされたもので、複数のブレードを備えたブレードサーバにおいて、複数のブレード間で、I/Oドローに装着されるI/Oデバイスを安価に共有することを目的とする。

50

【課題を解決するための手段】

【0016】

本発明は、仮想化部を備えて1以上の仮想サーバを提供する複数の物理サーバと、前記複数の物理サーバで共有されるI/Oデバイスと、前記複数の物理サーバと前記I/Oデバイスを接続するスイッチと、前記スイッチを初期化する管理部とを有するサーバシステムであって、前記I/Oデバイスが、仮想的な機能であるVirtual Function(VF)を有し、前記スイッチが、前記仮想サーバが認識する識別子である第1の識別子と、前記管理部が管理する識別子である第2の識別子と、の間の対応関係を示す識別子対応情報を保持し、前記仮想サーバからパケットを受信した場合に、送信元の物理サーバを識別するサーバ識別子を含む前記パケットを前記スイッチ内で変換し、前記パケットを前記I/Oデバイスに送信する際には、前記識別子対応情報に基づいて、前記パケットに含まれる前記第1の識別子を、当該第1の識別子に前記対応付けられた前記第2の識別子に変換し、かつ、前記サーバ識別子を削除した前記パケットを、前記I/Oデバイスに送信する、ことを特徴とするサーバシステム。

10

【発明の効果】

【0017】

したがって、本発明は、ひとつのI/Oデバイス、特にSR-IOV対応のI/Oデバイスを複数のサーバで共有することが可能となる。

【発明を実施するための最良の形態】

【0018】

20

以下、本発明の一実施形態を添付図面に基づいて説明する。

【0019】

図1は、第1の実施形態を示し、本発明を適用したブレードサーバ(複合型計算機システム)のブロック図を示す。図1において、ブレードサーバ1は、物理計算機として機能するn台のブレード10-1~10-nと、シングルルートI/O仮想化技術(以下、SR-IOV)に対応したI/Oデバイス50と、マルチルートI/O仮想化技術(以下、MR-IOV)でトランザクション(パケット)を処理するブレード10-1~10-nとI/Oデバイス50を接続するPCI-eスイッチ40と、I/Oデバイス50のブレード10-1~10-nへの割当てを管理するPCI管理サーバ(PCI管理計算機)20と、ユーザや管理者がブレードサーバ1を制御するための管理用端末30と、ブレードサーバ1の各部に電力を供給する電源供給装置60から構成されている。なお、PCI-eスイッチ40、I/Oデバイス50は、上記PCI-SIG(<http://www.pcisig.com/>)が策定したPCI-expressの規格に準拠したインターフェースで接続される。なお、図1においては、PCI-eスイッチ40が1つ、I/Oデバイス50が1つの例を示したが、複数のPCI-eスイッチ40と複数のI/Oデバイス50を備えることができる。

30

【0020】

ブレード10-1は1以上のCPU(プロセッサ)11-1, 11-2と、1以上のメモリ12-1と、1以上のチップセット13-1と、を含むハードウェアで構成される。なお、ブレード10-nとPCI管理サーバ20もブレード10-1と同一のハードウェアで構成される。

40

【0021】

また、ブレード10-1~10-nとPCI管理サーバ20の間はネットワーク70を介して接続される。ネットワークとしては、LAN(Local Area Network)やI²C(Integrated Circuit)を用いることができる。各ブレード10-1~10-nにはBMC(Baseboard Management Controller)を搭載してもよく、BMCによって各ブレード10-1~10-nの構成情報(コンフィギュレーション)や電源状態を収集し、PCI管理サーバ20に通知することができる。

【0022】

なお、各ブレード10-1~10-n及びPCI管理サーバ20のチップセット13-

50

1 ~ 13 - nは、それぞれPCI - expressに準拠したI/Oポート131 - 1 ~ 131 - nと、イーサネット（登録商標）等のLANの規格に準拠したネットワークインターフェース132 - 1 ~ 132 - nを備え、チップセット13 - 1 ~ 13 - nのネットワークインターフェース132 - 1 ~ 132 - nがネットワーク70に接続され、チップセット13 - 1 ~ 13 - nのI/Oポート131 - 1 ~ 131 - nがPCI - eスイッチ40に接続される。

【0023】

PCI - eスイッチ40は、ブレード10 - 1 ~ 10 - n及びPCI管理サーバ20のチップセット13と接続されるアップストリームポート41 - 0 ~ 41 - nと、I/Oデバイス50が接続されるダウンストリームポート42 - 1 ~ 42 - nを備える。

10

【0024】

図1の例では、ブレード10 - 1のチップセット13 - 1のI/Oポート131 - 1がPCI - eスイッチ40のアップストリームポート41 - 1に接続され、ブレード10 - nのチップセット13 - nのI/Oポート131 - nがPCI - eスイッチ40のアップストリームポート41 - nに接続され、PCI管理サーバ20のチップセット13 - MのI/Oポート131 - MがPCI - eスイッチ40のアップストリームポート41 - 0に接続され、ダウンストリームポート42 - 1にI/Oデバイス50が接続された例を示す。

【0025】

なお、ブレード10 - 1 ~ 10 - nとPCI - eスイッチ40の接続及びI/Oデバイス50とPCI - eスイッチ40の接続は、ブレードサーバ1のバックプレーなどを利用することができる。

20

【0026】

PCI - eスイッチ40は、アップストリームポート41 - 1 ~ 41 - nがMR - IOVに対応しており、ダウンストリームポート42 - 1 ~ 42 - nがSR - IOVに対応し、アップストリームポート41 - 1 ~ 41 - nはブレード10 - 1 ~ 10 - nとの間でマルチルートI/O仮想化技術のトランザクション（パケット）を送受信し、ダウンストリームポート42 - 1 ~ 42 - nとI/Oデバイス50の間ではシングルルートI/O仮想化技術のパケットを送受信する。このため、PCI - eスイッチ40は、マルチルートI/O仮想化技術のパケットとシングルルートI/O仮想化技術のパケットを相互に変換する処理を行う。

30

【0027】

PCI管理サーバ20には、入出力装置（図示省略）を備えた管理端末30が接続され、管理者などの操作により後述するPCIマネージャ202に指令を行うことができる。

【0028】

図2は、ブレードサーバ1の機能要素を示すブロック図である。ブレード10 - 1 ~ 10 - n上ではソフトウェア構成要素である仮想マシンモニタ（VMM；Virtual Machine Monitor）100 - 1 ~ 100 - nが稼動し、各仮想マシンモニタ100 - 1 ~ 100 - n上では複数の仮想計算機101 - 0 ~ 101 - k - 1が提供される。ブレード10 - 1では、仮想マシンモニタ100 - 1上に2つの仮想計算機101 - 0、101 - 1が生成され、各各仮想計算機101 - 0、101 - 1ではOS102 - 0とOS - 102 - 1がそれぞれ実行される。仮想計算機101 - 0、101 - 1のOS102 - 0とOS - 102 - 1にはI/Oデバイス50の仮想機能（VF；Virtual Function）を利用するためのVFドライバ103がそれぞれロードされている。他のブレード10 - nもブレード10 - 1と同様に構成され、それぞれ複数の仮想計算機101 - k - 1でOS102を実行し、VFドライバ103によってI/Oデバイス50の仮想機能を利用することができる。また、各仮想計算機101 - 1 ~ 101 - k - 1のOS102 - 0 ~ 102 - K - 1では、任意のアプリケーションまたはサービスが実行される。

40

【0029】

一方、PCI管理サーバ20では、OS201上で各ブレード10 - 1 ~ 10 - nとP

50

PCI-eスイッチ40及びI/Oデバイス50を管理するPCIマネージャ202が動作する。PCI管理サーバ20のOS201にはI/Oデバイス50の物理機能(PF; Physical Function)を利用するためのPFドライバ203がロードされている。

【0030】

ブレード10-1~10-nの各仮想計算機101-1~101-k-1と、PCI管理サーバ20からアクセスされるI/Oデバイス50は、ひとつの物理機能(PF)501と複数の仮想機能(VF)502を提供するSR-IOVに準拠したPCI-expressのI/Oデバイスである。物理機能501と仮想機能502は、I/Oデバイス50の制御部510が提供する機能であり、仮想機能502の数等はPCI管理サーバ20からの要求によって決定される。なお、図2では、I/Oデバイス50が仮想機能502としてVF1~VFkまでのk個の仮想機能502を提供する例を示す。

10

【0031】

ここで、ブレード10-1~10-nチップセット13-1~13-nと各仮想計算機101-1~101-k-1のVFドライバ103は、MR-IOVに対応したPCI-expressのパケットを送受信する。一方、I/Oデバイス50はSR-IOVに対応したPCI-expressのパケットを送受信する。

【0032】

本発明の特徴部であるPCI-eスイッチ40は、後述するように、MR-IOVのアップストリームポート41-1~41-nとSR-IOVのダウンストリームポート42-1~42-nとの間で、パケットの変換を行って、複数のブレード10-1~10-nでSR-IOVのI/Oデバイス50を共有する。

20

【0033】

各ブレード10-1~10-n及びPCI管理サーバ20を構成するチップセット13-1~n、13-Mは、PCI-Expressのプロトコル階層としてルートコンプレックス(RC; Root Complex)を含み、ルートコンプレックスはエンドポイントとしてのI/Oデバイス50までのローカルPCIツリーを管理する。

【0034】

PCI管理サーバ20のPCIマネージャ202は、後述するように、ブレード10-1~10-nに割り当てるPCIツリーの初期化と、PCI-eスイッチ40からI/Oデバイス50のトポロジーと、ブレード10-1~10-nへ割り当てるPCIツリーの対応関係を管理し、各ブレード10-1~10-nとPCI-eスイッチ40間のMR-IOVの設定を行う。

30

【0035】

また、PCIマネージャ202は、ブレード10-1~10-nの管理部としても機能し、各ブレード10-1~10-nの仮想マシンモニタ100-1~100-nが生成する仮想計算機101-1~101-k-1や起動するOS102-0~102-k-1を管理する。

【0036】

図3は、PCI-eスイッチ40の構成を示すブロック図である。PCI-eスイッチ40は、各ブレード10-1~10-nに接続されてMR-IOVのパケットを送受信するアップストリームポート41-1~41-nがマルチルートスイッチ(以下、MRS)論理410に接続され、I/Oデバイス50に接続されるダウンストリームポート42-1~42-nは、SR-IOVのパケットとMR-IOVのパケットを相互に変換するマルチルート-シングルルート(以下、MR-SR)変換論理430に接続される。なお、MR-SR変換論理430は、ダウンストリームポート42-1に対応するものだけを図示したが、実際には、各ダウンストリームポート42-1~42-n毎にMR-SR変換論理430が設けられる。

40

【0037】

MRS(マルチルートスイッチ)論理410は、アップストリームポート41-1~4

50

1 - n に接続されるポート 4 1 1 - 0 ~ 4 1 1 - n と、ポート 4 1 1 - 0 ~ 4 1 1 - n にそれぞれ接続されたアップストリームポートブリッジ 4 1 2 と、アップストリームポートブリッジ 4 1 2 に接続されたダウンストリームポートブリッジ 4 1 3 と、ダウンストリームポートブリッジ 4 1 2 に接続されたマルチルートポート 4 1 4 と、ポート 4 1 1 - 0 ~ 4 1 1 - n とマルチルートポート 4 1 4 と、ポート 4 1 1 - 0 ~ 4 1 1 - n とマルチルートポート 4 1 4 の接続関係等の M R S (マルチルートスイッチ) 構成情報 (コンフィギュレーション) 4 1 5 とを備える。M R S 構成情報 4 1 5 は、ポート 4 1 1 - 0 ~ 4 1 1 - n とマルチルートポート 4 1 4 のルーティングを管理する。この M R S 構成情報 4 1 5 は、コンフィギュレーション空間アクセスパス 4 2 1 によりアップストリームポートブリッジ 4 1 2 に接続されてアップストリームポート 4 1 - 1 ~ 4 1 - n からアクセス可能となっており、例えば、アップストリームポート 4 1 - 0 に接続された P C I 管理サーバ 2 0 からアクセスすることができる。さらに、コンフィギュレーション空間アクセスパス 4 2 1 は、M R - S R 変換論理 4 3 0 にも接続され、P C I 管理サーバ 2 0 は M R S 論理 4 1 0 からコンフィギュレーション空間アクセスパス 4 2 1 を介して M R - S R 変換論理 4 3 0 にもアクセス可能となっている。

10

【 0 0 3 8 】

M R S 論理 4 1 0 のマルチルートポート 4 1 4 は、内部マルチルートリンク 4 2 0 を介して M R - S R 変換論理 4 3 0 に接続される。

【 0 0 3 9 】

M R - S R 変換論理 4 3 0 は、内部マルチルートリンク 4 2 0 から受信した M R - I O V のパケットを S R - I O V のパケットに変換し、ダウンストリームポート 4 2 - 1 から I / O デバイス 5 0 へ送信するユニット 4 3 1 ~ 4 3 4 と、ダウンストリームポート 4 2 - 1 から受信した I / O デバイス 5 0 からの S R - I O V のパケットを M R - I O V のパケットに変換して、内部マルチルートリンク 4 2 0 に送信するユニット 4 3 5 ~ 4 3 8 と、M R - I O V のパケットと S R - I O V のパケットの変換を行うための情報を格納する T L P (Transaction Layer Packet) 変換情報 4 4 0 を主体に構成される。

20

【 0 0 4 0 】

ここで、本発明のパケットの構成について説明する。図 1 4 は、マルチルート T L P 及び P C I e (P C I - e x p r e s s) ベース T L P のフォーマットを示す説明図である。

30

【 0 0 4 1 】

マルチルート T L P (図中 M R T L P) 1 3 0 0 は、P C I - e x p r e s s のパケットである P C I e ベース T L P 1 2 0 0 のヘッダーの前に、マルチルート T L P プリフィックスヘッダー 1 3 1 0 を付加したものである。

【 0 0 4 2 】

P C I e ベース T L P 1 2 0 0 は、スタートフレーム (図中 S T P) と、シーケンス番号と、T L P ヘッダーと、E C R C (End to End Cyclic Redundancy Check) と、L C R C (Link Cyclic Redundancy Check) と、エンドフレームから構成される。

【 0 0 4 3 】

マルチルート T L P 1 3 0 0 は、P C I e ベース T L P 1 2 0 のシーケンス番号と T L P ヘッダーの間に、マルチルート T L P プリフィックスヘッダ 1 3 1 0 を挿入したものである。

40

【 0 0 4 4 】

マルチルート T L P プリフィックスヘッダ 1 3 1 0 は、M R - I O V において P C I e パケットの発行元のブレード 1 0 - 1 ~ 1 0 - n を特定するための仮想階層ナンバー (V H n ; Virtual Hierarchy Number) 1 3 1 1 を含む。なお、本実施例における仮想階層ナンバー 1 3 1 1 は P C I e スイッチ 4 0 の内部で付与される識別番号であり、M R S 構成情報 4 1 5 に格納されている。

【 0 0 4 5 】

本発明のブレードサーバ 1 では、仮想計算機 1 0 1 - 1 ~ 1 0 1 - k - 1、仮想マシン

50

モニタ100-1~100-n、チップセット13-1~13-n及びアップストリームポート41-1~41-n間では、MR-IOVのパケットであるマルチルートTLP1300で送受信を行い、ダウンストリームポート42-1とI/Oデバイス50の間では、マルチルートTLP1300からマルチルートTLPプリフィックスヘッダ1310を削除して、PCIeベースTLP1200で送受信を行う。

【0046】

ここで、本発明のブレードサーバ1では、ブレード10-1~10-nからI/Oデバイス50への下り方向(Outbound)のパケットではマルチルートTLP1300は、PCI-eスイッチ40のMR-SR変換論理430が、マルチルートTLPプリフィックスヘッダ1310を削除して、送信元をPCI管理サーバ20としたPCIeベースTLP1200に変換してI/Oデバイス50に送信する。

10

【0047】

逆に、I/Oデバイス50からブレード10-1~10-nへの上り方向(Inbound)のパケットでは、PCI-eスイッチ40のMR-SR変換論理430が、I/Oデバイス50から受信したPCIeベースTLP1200に、下り方向のパケットに付与されていた仮想階層ナンバー1311を含むマルチルートTLPプリフィックスヘッダ1310を付加してブレード10-1~10-nに送信する。

【0048】

以上のような構成により、PCI-eスイッチ40のブレード10-1~10-n側ではMR-IOVのパケットで通信を行い、PCI-eスイッチ40のI/Oデバイス50側ではSR-IOVのパケットで通信を行うことで、複数のブレード10-1~10-nでひとつのSR-IOVに準拠したI/Oデバイス50を共有できる。

20

【0049】

また、PCI-eスイッチ40は、発行元のブレード10-1~10-nをI/Oデバイス50に対して隠蔽するのに加え、各ブレード10-1~10-n間で異なるMMIO(Memory-mapped I/O)空間を吸収する。

【0050】

このため、PCI管理サーバ20は、仮想マシンモニタ100-1~100-nの起動時に、各仮想計算機101-1~101-k-1のMMIO空間のアドレスと、I/Oデバイス50へのPCIeベースTLP1200の発行元となるPCI管理サーバ20のMMIO空間の差分(オフセット)を、PCI-eスイッチ40のTLP変換情報440に格納しておき、上り方向のマルチルートTLP1300の宛先アドレスを、実際の宛先となる仮想マシンモニタ100-1~100-nのMMIO空間に変換する。

30

【0051】

以下に、PCI-eスイッチ40が行う仮想階層ナンバー1311の削除及び付加と宛先アドレスの変換について詳述する。

【0052】

図4は、ブレード10-1(Blade#1)~10-n(Blade#n)及びPCI管理サーバ20(PCI-M)が管理するメモリ12-1~n、12-M上に設定したMMIO空間を示す。例えば、図示の例では、ブレード10-1に2つの仮想機能502=VF1、VF2を割り当て、ブレード10-nにも2つの仮想機能502=VFk-1、VF2k割り当てた例を示す。

40

【0053】

ブレード10-1~10-nの仮想マシンモニタ100-1~nはローカルPCIツリーの後にI/Oデバイス50の仮想機能502(VF)へアクセスするためのMMIO空間を設定する。各ブレードのMMIO空間は、搭載しているメモリ12-1~nの容量等の違いにより、アドレスは相違する。各ブレードのMMIO空間は、各ブレード10-1~10-nが使用する仮想機能502(VF)毎にPCI管理サーバ20が設定するメモリ12-M上のMMIO空間に割り当てられる。

【0054】

50

各ブレード10-1~10-nが使用するMMIO空間のアドレスの違いは、PCIマネージャ202が設定する宛先アドレス修飾情報4401の「Offset」にその差が設定されることで、PCI-eスイッチ40は、I/Oデバイス50とブレード10-1~10-n間のパケットの宛先アドレスを変更することができる。

【0055】

つまり、ブレード10-1~10-nからI/Oデバイス50に対する読み込み要求等では、各ブレードが送信するパケットの発行元アドレスは、PCI-eスイッチ40が、宛先アドレス修飾情報4401の「Offset」によって、PCI管理サーバ20のMMIO空間が発行元アドレスのパケットに書き換えてI/Oデバイス50に送信する。

【0056】

逆に、I/Oデバイス50からブレード10-1~10-nに対する読み込み要求等への応答では、I/Oデバイス50が送信するパケットの宛先アドレスは、PCI管理サーバ20のMMIO空間となっているので、PCI-eスイッチ40は宛先アドレス修飾情報4401の「Offset」によって、PCI管理サーバ20のMMIO空間をブレード10-1~10-n毎のMMIO空間に書き換えて返信することができる。

【0057】

以上の処理により、PCI-eスイッチ40は、MMIO空間のアドレスを差し替えることで、複数のブレード10-1~10-nでSR-IOVに対応するI/Oデバイス50の共有を実現する。

【0058】

図5は、ブレード10-1(Blade#1)~10-n(Blade#n)及びPCI管理サーバ20(PCI-M)のチップセット13-1~n、Mが管理するルーティングIDの関係を示す。各チップセット13-1~n、Mは、それぞれ、各計算機内のデバイスについて認識したローカルPCIツリー用のルーティングIDと、PCI-eスイッチ40のMRS論理410から取得したMRS論理用ルーティングIDと、各チップセット13-1~n、MからのI/Oデバイス50のアクセス用ルーティングIDとを備える。

【0059】

ブレード10-1~10-nのチップセット13-1~nでは、接続しているPCI-eスイッチ40のアップストリームポート41-1~nが異なるため、チップセット13-1~nが認識するI/Oデバイス50のバス番号は異なる。例えば、図示の例では、ブレード10-1に2つの仮想機能502=VF1、VF2を割り当て、ブレード10-nにも2つの仮想機能502=VFk-1、VF2k割り当てた例を示す。

【0060】

チップセット13-1からI/Oデバイス50の仮想機能502=VF1へのルーティングIDは、11:0:1であり、バス番号は11であるのに対し、チップセット13-nからI/Oデバイス50の仮想機能502=VFk-1へのルーティングIDは、13:x:y-1でありバス番号は13である。MRS論理用ルーティングIDも同様であり、ブレード10-1~10-n毎に異なる。なお、ルーティングIDは、バス番号(Bus#):デバイス番号(Dev#):ファンクション番号(Fun#)の順に設定される。

【0061】

一方、I/Oデバイス50の物理機能501にアクセス可能なPCIマネージャ202は、起動時に物理機能501と仮想機能502のルーティングIDを初期化してチップセット13-Mで管理しており、図示の例ではルーティングID=10:0:0~10:x:y-1をI/Oデバイス50のアクセス用ルーティングIDに設定する。

【0062】

そして、PCIマネージャ202は、各ブレード10-1~10-nが起動する度に、図7に示すTLP変換情報440の仮想機能ID修飾情報4402に、各ブレード10-1~10-nのチップセット13-1~nからI/Oデバイス50の仮想機能502まで

10

20

30

40

50

のルーティングIDと、PCI管理サーバ20のチップセット13-Mが管理するI/Oデバイス50のルーティングIDを、各ブレード10-1~10-n毎の仮想階層ナンバー1311と共に設定する。

【0063】

ルーティングIDで各ブレード10-1~10-nとI/Oデバイス50がアクセスを行う場合には、PCI-eスイッチ40が仮想機能ID修飾情報4402を参照し、パケット中のルーティングIDを付け替えて通信を行う。

【0064】

すなわち、PCI-eスイッチ40は、ブレード10-1~10-nからI/Oデバイス50へのパケットがルーティングIDでアクセスする場合、各ブレード10-1~10-nのルーティングIDをPCI管理サーバ20のルーティングIDに差し替えてI/Oデバイス50に送信する。

10

【0065】

逆に、I/Oデバイス50からブレード10-1~10-nへのパケットがルーティングIDでアクセスする場合、PCI-eスイッチ40は、PCI管理サーバ20のルーティングIDを各ブレード10-1~10-nのルーティングID差し替えてI/Oデバイス50に送信する。

【0066】

以上の処理により、PCI-eスイッチ40は、ルーティングIDを差し替えることで、複数のブレード10-1~10-nでSR-IOVに対応するI/Oデバイス50の共有を実現する。

20

【0067】

ここで、ブレードサーバ1で使用されるPCI-expressのパケットのフォーマットを図15に示す。図15はPCI-eスイッチ40で転送されるPCIeベースTLP1200のうちTLPヘッダをアクセス形態毎に示す説明図である。

【0068】

図15は、PCIeベースTLP1200のTLPヘッダの詳細を示し、TLPヘッダ1200Aは、MMIO空間のアドレスでブレード10-1~10-nとI/Oデバイス50のアクセスを行うパケットを示し、TLPヘッダ1200B、1200CはルーティングIDでアクセスを行うパケットを示しており、TLPヘッダ1200Bは要求(書き込み要求など)に対する完了通知を通知するパケットの例を示し、TLPヘッダ1200Cは対象のデバイスの構成情報(コンフィギュレーション)の設定するパケットを示す。

30

【0069】

各TLPヘッダ1200A~Cは、バイト0~15の16バイトで構成される。

【0070】

MMIO空間のアドレスでアクセスするTLPヘッダ1200Aは、メモリリード要求やメモリライト要求の際に使用される。

【0071】

TLPヘッダ1200Aは、バイト4、5にパケットの発行元(リクエスタ)IDに要求元のルーティングIDを格納し、バイト8~15に読み込む対象のMMIO空間のアドレスを格納し、バイト0の0~4ビットに要求の種別を格納する。

40

【0072】

PCI-eスイッチ40は、ブレード10-1~10-nからI/Oデバイス50へ向かう下りパケットの時には、リクエスタIDに格納されたブレードのルーティングIDをPCI管理サーバ20のルーティングIDに書き換えて、要求元がPCI管理サーバ20であることを設定する。つまり、SR-IOVのI/Oデバイス50は、ひとつの計算機のみと接続されるため、I/Oデバイス50へのパケットはPCI管理サーバ20発に差し替える。

【0073】

50

そして、P C I - e スイッチ 4 0 は、ブレードの M M I O 空間のアドレスを宛先アドレス修飾情報 4 4 0 1 の「o f f s e t」で P C I 管理サーバ 2 0 (P C I マネージャ 2 0 2) の M M I O 空間のアドレスに変更し、I / O デバイス 5 0 が認識可能なアドレス空間に変更する。

【 0 0 7 4 】

ルーティング I D でアクセスする場合の、T L P ヘッダ 1 2 0 0 B は、要求 (書き込み要求等) に対する完了通知の際に使用される。T L P ヘッダ 1 2 0 0 B は、バイト 4 , 5 のコンプリータ I D に完了通知 (コンプリーション) を発行するデバイスのルーティング I D を格納し、バイト 8、9 のリクエスト I D に要求元 (リクエスト) のルーティング I D を格納し、バイト 0 の 0 ~ 4 ビットにコンプリーションを示す値を格納する。

10

【 0 0 7 5 】

P C I - e スイッチ 4 0 は、仮想機能 I D 修飾情報 4 4 0 2 を参照して、I / O デバイス 5 0 へ向かう上りパケットの時には、リクエスト I D に格納された P C I 管理サーバ 2 0 のルーティング I D を P C I 管理サーバ 2 0 のルーティング I D から送信先のブレードのルーティング I D に差し替えて送信する。

【 0 0 7 6 】

ルーティング I D でアクセスする T L P ヘッダ 1 2 0 0 C は、デバイスに対する構成情報の設定を要求するパケットで使用される。T L P ヘッダ 1 2 0 0 C は、バイト 4 , 5 のリクエスト I D に構成情報の設定を要求するデバイスのルーティング I D を格納し、バイト 8、9 には構成情報の設定対象のバス番号とデバイス番号とファンクション番号からなるルーティング I D を格納し、バイト 0 の 0 ~ 4 ビットに構成情報の設定要求 (コンフィギュレーション) を示す値を格納する。

20

【 0 0 7 7 】

P C I - e スイッチ 4 0 は、ブレード 1 0 - 1 ~ 1 0 - n から I / O デバイス 5 0 へ向かう下りパケットの時には、リクエスト I D に格納されたブレードのルーティング I D を P C I 管理サーバ 2 0 のルーティング I D に書き換えて、要求元が P C I 管理サーバ 2 0 であることを設定する。また、P C I - e スイッチ 4 0 は、仮想機能 I D 修飾情報 4 4 0 2 を参照して、ブレードが設定した I / O デバイス 5 0 のルーティング I D (バス番号、デバイス番号、ファンクション番号を、) P C I 管理サーバ 2 0 が認識するルーティング I D に差し替えて、I / O デバイス 5 0 には P C I 管理サーバ 2 0 からの要求であることに変更する。

30

【 0 0 7 8 】

図 1 6 は、ブレードサーバ 1 の起動手順を示すフローチャートである。この処理は、ブレードサーバ 1 の管理者またはユーザのスイッチ操作に基づいて開始される。

【 0 0 7 9 】

ステップ S 1 では、管理者 (またはユーザ) が P C I - e スイッチ 4 0 の電源を投入する。P C I - e スイッチ 4 0 は、起動すると P C I - e リンクの初期化処理と各種レジスタの初期化処理を実行する (S 2)。すなわち、図 3 の M R S 構成情報 4 1 5 及び T L P 変換情報 4 4 0 の初期化が行われる。また、管理者などにより I / O デバイス 5 0 の電源が投入される。なお、I / O デバイス 5 0 の電源は、P C I - e スイッチ 4 0 と連動していても良い。

40

【 0 0 8 0 】

ステップ S 3 では、管理者 (またはユーザ) が、P C I 管理サーバ 2 0 の電源を投入する。P C I 管理サーバ 2 0 では、後述する図 1 7 のように各種初期化が行われる。P C I 管理サーバ 2 0 の起動が完了すると、ステップ S 4 では P C I マネージャ 2 0 2 が管理端末 3 0 からのブレードに対する起動の指示を待ち受ける。

【 0 0 8 1 】

P C I マネージャ 2 0 2 は管理端末 3 0 から起動するブレード 1 0 - 1 ~ 1 0 - n に対する指令を受け付けると、ステップ S 5 に進んで指令を受け付けたブレード 1 0 - 1 ~ 1 0 - n の電源を投入し、その後、ステップ S 6 では、P C I マネージャ 2 0 2 が、起動が

50

完了したブレードに対して仮想マシンモタ100-1~100-nを起動する。仮想マシンモタ100-1~100-nの起動については後述の図18にて詳述する。ステップS7では、仮想マシンモタの起動が完了した後に、受け付けた指令に基づいて仮想計算機101-1~101-k-1を生成し、各仮想計算機101-1~101-k-1でOS102-0~102-K-1を起動させる。仮想計算機の起動が完了すると再びステップS4に戻って次のブレード10-1~10-nの起動の指令を待つ。

【0082】

上記処理により、PCI-eスイッチ40、PCI管理サーバ20、I/Oデバイス50の順で電源が投入されて初期化が行われた後、管理端末30から指令されたブレード10-1~10-nが起動される。

10

【0083】

図17は、図16に示したステップS3のPCI管理サーバ20の起動の際に行われる処理の一例を示すフローチャートである。

【0084】

PCI管理サーバ20は、電源を投入されるとステップS11でBIOS(Basic I/O System)またはEFI(Extensible Firmware Interface)が起動する。BIOS(またはEFI)の起動時には、ステップS12でチップセット13-MがPCI管理サーバ20内のデバイスについて、ルートコンプレックスの下にローカルPCIツリーを構成し、PCIローカルツリーの初期化を行う。すなわち、チップセット13-Mは、図5で示したような、ローカルPCIツリーのルーティングIDをチップセット13-Mのレジスタ等の所定の領域に設定する。この例では、PCI管理サーバ20内のローカルルーティングIDとして、0:0:0~7:1F:7が各エントリに設定される。

20

【0085】

ステップS13では、BIOS(またはEFI)の起動が完了するとOS201が起動し、その後に、PCIマネージャ202(PCI-M)が起動する。このとき、OS201はPFドライバ203を読み込んで、PCIマネージャ202がI/Oデバイス50の物理機能PF501を利用する準備を行う。なお、I/Oデバイス50の物理機能を使用するのはPCI-eスイッチ40及びI/Oデバイス50を管理するPCIマネージャ202のみであり、他のブレード10-1~10-nは仮想機能を利用する。

【0086】

ステップS14では、OS201がネットワークインターフェース132-Mの初期化を実施する。ネットワークインターフェース132-Mの初期化は、ネットワーク70への接続とIPアドレスの付与及びブレード10-1~10-nの仮想計算機からの接続要求に対する待ち受けの開始が含まれる。

30

【0087】

次に、ステップS15では、PCIマネージャ202がI/Oポート131-MからPCI-eスイッチ40にアクセスして、マルチルートスイッチ論理410のアップストリームポートブリッジ412からアップストリームポート41-0に接続されたポート411-0に設定された構成情報(例えば、マルチルート化可能ポインタ: Capability pointer)を取得し、当該ポート411-0のバス番号(例えば、9)を取得する。ここで、構成情報としてのマルチルート化可能ポインタは、図19で示すように、4バイト目の0~15ビットのVSEC_IDがマルチルート化可能ポインタを示す。PCIマネージャ202は、PCI-eスイッチ40から仮想階層ナンバー1311を取得する。図4の例では、PCI管理サーバ20の仮想階層ナンバー1311(VH)は「0」となる。

40

【0088】

ステップS16では、PCIマネージャ202が、I/Oポート131-MからPCI-eスイッチ40にアクセスして、MRS論理410のダウンストリームポートブリッジ413の構成情報(例えば、コンフィギュレーションヘッダ)を取得する。

【0089】

次に、ステップS17では、PCIマネージャ202は、I/Oポート131-M及び

50

アップストリームポート 4 1 - 0 から I / O デバイス 5 0 にアクセスし、I / O デバイス 5 0 の S R - I O V に関する構成情報（例えば、コンフィギュレーションヘッダ）を取得する。P C I マネージャ 2 0 2 は、I / O デバイス 5 0 の構成情報から仮想機能（V F）の個数や M M I O 空間などを取得する。また、P C I マネージャ 2 0 2 は I / O デバイス 5 0 からバス番号及びデバイス番号を取得する。例えば、P C I 管理サーバ 2 0 のチップセット 1 3 - M のルートコンプレックスから見た I / O デバイス 5 0 のアクセス用ルーティング I D は 0 : 2 : 0 となる。

【 0 0 9 0 】

ステップ S 1 8 では、上記ステップ S 1 5 ~ S 1 7 で取得したマルチルートスイッチ論理 4 1 0 のアップストリームポートブリッジ 4 1 2 とダウンストリームポートブリッジ 4 1 3 の構成情報と、I / O デバイス 5 0 の構成情報から、図 5 で示すように、M R S 論理用ルーティング I D と、I / O デバイス 5 0 をアクセスするためのルーティング I D（D E V アクセス用 Routing I D）をチップセット 1 3 - M のレジスタ等の所定の領域に設定する。また、P C I マネージャ 2 0 2 は、メモリ 1 2 - M の所定の領域に物理機能 5 0 1 にアクセスするための M M I O と、仮想機能 5 0 1 にアクセスするための M M I O を k 個設定する。

【 0 0 9 1 】

図 5 の例では、P C I マネージャ 2 0 2 が P C I 管理サーバ 2 0 の I / O ポート 1 3 1 - M から P C I - e スwitch 4 0 のアップストリームポート 4 1 - 0 からダウンストリームポートブリッジ 4 1 3 までの M R - I O V のルーティング I D として、ポート 4 1 1 - 0 がバス番号 = 9、M R S 構成情報 4 1 5 がバス番号 = 8、I / O デバイス 5 0 のバス番号 = 1 0 となり、M R S 論理 4 1 0 用のルーティング I D は、8 : 0 : 0、9 : 0 : 0 となる。また、I / O デバイス 5 0 のデバイスアクセス用ルーティング I D は、1 0 : 0 : 0 が I / O デバイス 5 0 の物理機能にアクセスするためのルーティング I D に設定され、1 0 : 0 : 1 ~ 1 0 : 0 : V F k までの k 個のルーティング I D が I / O デバイス 5 0 のアクセス用に設定される。そして、P C I マネージャ 2 0 2 が P F ドライバ 2 0 3 で I / O デバイス 5 0 の物理機能 P F 5 0 1 と仮想機能 5 0 2 を初期化する。

【 0 0 9 2 】

上記処理によって、I / O デバイス 5 0 の物理機能 P F 5 0 1 と仮想機能 5 0 2 を初期化が P C I マネージャ 2 0 2 によって行われる。

【 0 0 9 3 】

図 1 8 は、図 1 6 のステップ S 6、S 7 で行われるブレード 1 0 - 1 ~ 1 0 - n の起動と仮想計算機（ゲスト V M）の起動処理の一例を示すフローチャートである。この例では、管理端末 3 0 から図 1 のブレード 1 0 - 1（B L A D E # 1）を起動し、ブレード 1 0 - 1 の仮想マシンモニタ 1 0 0 - 1 上で仮想計算機 1 0 1 - 1（V M # 0）、1 0 1 - 2（V M # 1）を生成する指令があった例を示す。

【 0 0 9 4 】

P C I マネージャ 2 0 2 は、既に起動してステップ S 3 1 にて仮想マシンモニタ 1 0 0 - 1 ~ 1 0 0 - n からの通信を待ち受ける状態となっている（図 1 7 の S 1 4）。

【 0 0 9 5 】

P C I 管理サーバ 2 0 は管理端末 3 0 から指令を受けたブレード 1 0 - 1（B R A D E # 1）に起動の指令を出力し、ブレード 1 0 - 1 を起動する。この処理は、例えば、ブレード 1 0 - 1 ~ 1 0 - n が B M C を備えている場合には、該当するブレードの B M C に対して起動を指令し、ネットワークインターフェース 1 3 2 - 1 ~ n へ起動を指令する場合にはマジックパケットを送信することで行われる。

【 0 0 9 6 】

P C I 管理サーバ 2 0 からの指令で起動したブレード 1 0 - 1 は、電源を投入されるとステップ S 2 1 で B I O S または E F I が起動する。B I O S（または E F I）の起動時には、ステップ S 2 2 でチップセット 1 3 - 1 がブレード 1 0 - 1 内のデバイスについてルートコンプレックスの下にローカル P C I ツリーを構成し、P C I ローカルツリーの初

10

20

30

40

50

期化を行う。すなわち、チップセット13-1は、図5で示したBlade#1のように、ローカルPCIツリーのルーティングIDをチップセット13-1のレジスタ等の所定の領域に設定する。この例では、ブレード10-1内のローカルルーティングIDとして、0:0:0~8:1F:7が各エントリに設定される。

【0097】

ステップS23では、BIOS(またはEFI)の起動が完了すると仮想マシンモニタ100-1が起動する。仮想マシンモニタ100-1がネットワークインターフェース132-1の初期化を実施する。ネットワークインターフェース132-1の初期化は、ネットワーク70への接続とIPアドレスの付与が含まれる。

【0098】

次に、ステップS24では、仮想マシンモニタ100-1がI/Oポート131-1からPCI-eスイッチ40にアクセスして、マルチルートスイッチ論理410のアップストリームポートブリッジ412からアップストリームポート41-0に接続されたポート411-1に設定された構成情報(例えば、マルチルート化可能ポインタ:Capability pointer)を取得し、当該ポート411-1のバス番号(例えば、9)を取得する。構成情報としてのマルチルート化可能ポインタは、図19で示したように、4バイト目の0~15ビットのVSEC_IDがマルチルート化可能ポインタを示す。仮想マシンモニタ100-1は、PCI-eスイッチ40から仮想階層ナンバー1311を取得する。図5の例では、ブレード10-1の仮想階層ナンバー1311(VH)は「1」となる。

【0099】

ステップS25では、仮想マシンモニタ100-1が、I/Oポート131-1からPCI-eスイッチ40にアクセスして、MRS論理410のダウストリームポートブリッジ413の構成情報(例えば、コンフィギュレーションヘッダ)を取得する。仮想マシンモニタ100-1は、ルートコンプレックスから見たI/Oデバイス50のルーティングIDのデバイス番号として、MRS論理410のダウストリームポートブリッジ413のデバイスを割り当てる。この例では、ルートコンプレックスから見たI/Oデバイス50のルーティングIDは、0:6:0となる。

【0100】

次に、ステップS26では、仮想マシンモニタ100-1が必要とする仮想機能502の数を含む確保要求をPCIマネージャ202にVF確保要求データとしてネットワーク70から送信する。この確保要求は、図20で示すように、仮想マシンモニタ100-1が必要とする仮想機能(VF)502の数(=2)と、仮想マシンモニタ100-1がひとつの仮想機能502に対してメモリ12-1に設定可能なMMIOのサイズの最大値と、仮想機能502を設定するMMIOのベースアドレス及び仮想階層ナンバー1311(VH=1)と、I/Oデバイス50からアクセスするためのルーティングID(11:0:1及び11:0:2)とを含む。なお、ブレード10-1から見たI/Oデバイス50のバス番号は11となり、PCI管理サーバ20から見たバス番号とは異なる。また、ブレード10-1は、確保要求にルートコンプレックスが取得したI/Oデバイス50のアクセス用ルーティングID(0:6:0)も通知する。なお、I/Oデバイス50からアクセスするためのルーティングIDは、各ブレード10-1~10-nの起動時にBIOSが決定した値を用いることができる。

【0101】

仮想機能502を確保する要求を受け付けたPCIマネージャ202は、ブレード10-1から見たPCI-eスイッチ40のMRS論理410及びI/Oデバイス50のルーティングIDと、ブレード10-1のMMIO空間を認識する。

【0102】

PCIマネージャ202は、ブレード10-1から要求された仮想機能502に対するI/Oデバイス50のアクセス用ルーティングIDから図5に示したVF1及びVF2をブレード10-1に割り当て、図4に示したMMIO空間に確保したVF1, VF2のMMIOをブレード10-1に割り当てる。

10

20

30

40

50

【 0 1 0 3 】

そして、P C I マネージャ 2 0 2 は、ブレード 1 0 - 1 から受信した確保要求に含まれるブレード 1 0 - 1 の I / O デバイス 5 0 のアクセス用ルーティング I D と、仮想階層ナンバー 1 3 1 1 と、M M I O のベースアドレスと、仮想階層ナンバー 1 3 1 1 を、M R - S R 変換論理 5 3 0 の T L P 変換情報 4 4 0 に書き込む。このとき、P C I マネージャ 2 0 2 は、確保要求で要求された仮想機能 5 0 2 の数に応じて P C I マネージャ 2 0 2 が管理する仮想機能 5 0 2 の M M I O とブレード 1 0 - 1 が管理する M M I O のオフセットを求める。

【 0 1 0 4 】

P C I マネージャ 2 0 2 は、M R - S R 変換論理 5 3 0 の T L P 変換情報 4 4 0 (図 7) を構成する宛先アドレス変換論理 4 4 0 1 へ、仮想機能 5 0 2 毎に仮想階層ナンバー 1 3 1 1、ブレード 1 0 - 1 の M M I O のベースアドレスと、M M I O のサイズと、P C I マネージャ 2 0 2 が管理する仮想機能 5 0 2 の M M I O とブレード 1 0 - 1 が管理する M M I O のオフセットを書き込む。P C I マネージャ 2 0 2 は、書き込みを行ったエントリの V a l i d に有効であることを示す「 1 」をセットする。

10

【 0 1 0 5 】

次に、P C I マネージャ 2 0 2 は、図 5 のルーティング I D からブレード 1 0 - 1 に割り当てた仮想機能 5 0 2 (V F 1 と V F 2) のルーティング I D (1 0 : 0 : 1 と 1 0 : 0 : 2) と、ブレード 1 0 - 1 から通知を受けたルーティング I D (1 1 : 0 : 1 と 1 1 : 0 : 2) を対応づける。

20

【 0 1 0 6 】

そして、P C I マネージャ 2 0 2 は、M R - S R 変換論理 5 3 0 の T L P 変換情報を構成する仮想機能 I D 修飾情報 4 4 0 2 (図 7) へ、各ブレードから I / O デバイス 5 0 のアクセス用ルーティング I D 毎に、ブレード 1 0 - 1 ~ 1 0 - n の仮想階層ナンバー 1 3 1 1 (V H) を V H に書き込み、I / O デバイス 5 0 のアクセス用ルーティング I D を V H x _ R o u t i n g I D に書き込み、ブレードのルーティング I D に対応する P C I マネージャ 2 0 2 が対応付けたルーティング I D を V H 0 _ R o u t i n g I D に書き込む。P C I マネージャ 2 0 2 は、書き込みを行ったエントリの V a l i d に有効であることを示す「 1 」をセットする。

【 0 1 0 7 】

30

次に、P C I マネージャ 2 0 2 は、図 7 のルートポート I D 情報 4 4 0 3 に、ブレード 1 0 - 1 のルートコンプレックスが認識した I / O デバイス 5 0 のルーティング I D (0 : 6 : 0) を仮想階層ナンバー 1 3 1 1 毎に書き込む。P C I マネージャ 2 0 2 は、書き込みを行ったエントリの V a l i d に有効であることを示す「 1 」をセットする。

【 0 1 0 8 】

以上により、M R - S R 変換論理 4 3 0 の T L P 変換情報 4 4 0 にブレード 1 0 - 1 から受信した V F 確保要求データの値を設定し、ブレード 1 0 - 1 が I / O デバイス 5 0 を共有することを設定する。

【 0 1 0 9 】

P C I マネージャ 2 0 2 は、M R - S R 変換論理 4 3 0 の T L P 変換情報 4 4 0 にブレード 1 0 - 1 から受信した V F 確保要求データの設定が完了すると、T L P 変換情報 4 4 0 に設定した情報を V F 確保完了データとしてネットワーク 7 0 からブレード 1 0 - 1 に送信する。V F 確保完了データの一例としては図 2 1 のようになる。図 2 1 において、V F 確保完了データは、P C I マネージャ 2 0 2 が T L P 変換情報 4 4 0 に確保した仮想機能 5 0 2 の数 (= 2) と、仮想機能 5 0 2 当たりの M M I O のサイズと、I / O デバイス 5 0 の物理機能 P F 5 0 1 の構成情報の空間サイズと、I / O デバイス 5 0 の物理機能 5 0 1 の構成空間データと、I / O デバイス 5 0 の仮想機能 5 0 2 の構成空間サイズと、I / O デバイス 5 0 の仮想機能 5 0 2 の構成空間データが含まれる。

40

【 0 1 1 0 】

V F 確保完了データを受信した仮想マシンモニタ 1 0 0 - 1 は、ステップ S 2 7 で管理

50

端末 30 から指令された仮想計算機 101 - 1、101 - 2 (ゲスト VM) を起動し、PCI マネージャ 202 から受信した I/O デバイス 50 の仮想機能 502 (VF1, VF2) をそれぞれ割り当てる。

【0111】

上記の処理により、仮想マシンモニタ 100 - 1 が必要とする I/O デバイス 50 の仮想機能 502 を PCI マネージャ 202 が確保し、仮想マシンモニタ 100 - 1 は PCI マネージャ 202 が確保した仮想機能 502 を起動する仮想計算機 101 - 1、101 - 2 に割り当てることができる。

【0112】

なお、上記の処理は他のブレード 10 - 2 ~ 10 - n を起動するときにも同様に実行され、TLP 変換情報 440 には I/O デバイス 50 を共有するブレード 10 - 1 ~ 10 - n の情報が加えられることになる。

【0113】

図 21 は、上記図 18 のステップ S27 で行われる仮想計算機 101 - 0 (101 - 1) の起動処理の一例を示すフローチャートである。

【0114】

仮想マシンモニタ 100 - 1 は、ステップ S41 で仮想計算機 101 - 0 を起動し、ステップ S42 で仮想計算機 101 - 0 に割り当てた仮想 BIOS を起動する。ステップ S42 では、仮想 BIOS によって仮想 PCI ツリーの初期化を実施する。

【0115】

ステップ S44 では、仮想 BIOS が仮想マシンモニタ 100 - 1 によって提供される仮想コンフィギュレーションにアクセスを開始する。ステップ S45 では、仮想マシンモニタ 100 - 1 は仮想 BIOS のアクセス対象が、I/O デバイス 50 の仮想機能 502 に対するものであるか否かを判定する。仮想機能 502 以外に対して仮想 BIOS がアクセスした場合には、ステップ S46 で仮想マシンモニタ (VMM) 100 - 1 がアクセス対象のエミュレーションを実施する。一方、仮想機能 502 に対して仮想 BIOS がアクセスした場合には、ステップ S47 で仮想マシンモニタ 100 - 1 が、PCI マネージャ 202 から仮想機能確保応答情報により取得した I/O デバイス 50 の物理機能 PF501 の構成情報の空間サイズと、I/O デバイス 50 の物理機能 501 の構成空間データと、I/O デバイス 50 の仮想機能 502 の構成空間サイズと、I/O デバイス 50 の仮想機能 502 の構成情報の空間データを応答する。

【0116】

ステップ S48 では、仮想マシンモニタ 100 - 1 が仮想 BIOS による仮想 PCI ツリーの初期化が完了したか否かを判定する。仮想 PCI ツリーの初期化が完了していなければステップ S44 に戻り上記初期化を繰り返す。一方、仮想 PCI ツリーの初期化が完了すると、ステップ S49 で仮想マシンモニタ 100 - 1 が OS 100 - 0 (ゲスト OS) の起動を開始する。ステップ S50 で OS 100 - 0 の起動が完了し、仮想計算機 101 - 0 でアプリケーションやサービスの提供が可能となる。

【0117】

上記により、仮想 BIOS や OS 100 - 0 (ゲスト OS) が I/O デバイス 50 の仮想機能 502 の構成情報を取得する場合には、仮想マシンモニタ 100 - 1 が PCI マネージャ 202 から取得した物理機能 501 及び仮想機能 502 の情報を提供することで、OS 100 - 0 等が VF ドライバ 103 を組み込んで仮想機能 502 を利用することが可能となる。

【0118】

図 23 は、上記実施形態の変形例を示し、図 3 に示した MRS 論理 410 のうち、アップストリームポートブリッジ 412 に接続されるポート 411 - 1 ~ 411 - n を MRS 論理 410 から分離して、アップストリームポート 41 - 0 ~ 41 - n の増設を可能にしたものである。

【0119】

10

20

30

40

50

P C I - e スイッチ 4 0 のアップストリームポート 4 1 - 0 ~ 4 1 - n は、所定数（例えば、2 つ）毎に M R S 論理 4 1 1 0 - 1 ~ 4 1 1 0 - n に接続される。各 M R S 論理 4 1 1 0 - 1 ~ n には、2 つのポート 4 1 1 - 0、4 1 1 - 1 等がそれぞれ設けられて、アップストリームポート 4 1 - 0 ~ 4 1 - n に接続される。各 M R S 論理 4 1 1 0 - 1 ~ 4 1 1 0 - n の M R S 論理 4 1 0 側にはひとつのポート 4 1 1 1 - 1 ~ 4 1 1 1 - n が設けられており、それぞれのポート 4 1 1 1 - 1 ~ 4 1 1 1 - n は M R S 論理 4 1 0 に設けたポート 4 1 1 2 - 1 ~ 4 1 1 2 - n に接続される。

【 0 1 2 0 】

この変形例では、M R S 論理 4 1 0 のアップストリームポート 4 1 - 0 ~ 4 1 - n 側を多段化することによって、ブレードサーバ 1 のアップストリームポート 4 1 - 0 ~ 4 1 - n の構成を柔軟に行うことが可能となる。

10

【 0 1 2 1 】

なお、上記実施形態では、チップセット 1 3 - 1 ~ n とプロセッサ 1 1 - 1、2 が独立した構成を示したが、プロセッサ 1 1 - 1、2 にチップセット 1 3 - 1 ~ n が組み込まれていても良い。

【 0 1 2 2 】

また、上記実施形態では、P C I 管理サーバ 2 0 で P C I マネージャ 2 0 2 を稼働させる例を示したが、ブレード 1 0 - 1 ~ 1 0 - n の何れかで P C I マネージャ 2 0 2 を稼働させてもよい。

【産業上の利用可能性】

20

【 0 1 2 3 】

以上のように、本発明は、複数の物理計算機を P C I - e x p r e s s のスイッチで I / O デバイスに接続する計算機システムに適用することができ、特に、物理計算機が M R - I O V でアクセスを行い、I / O デバイスが S R - I O V でアクセスする計算機システム及び P C I スイッチに好適である。

【図面の簡単な説明】

【 0 1 2 4 】

【図 1】本発明を適用したブレードサーバのハードウェア構成を示すブロック図。

【図 2】同じくブレードサーバの機能要素を示すブロック図。

【図 3】P C I - e スイッチの内部を示すブロック図。

30

【図 4】各ブレードと P C I 管理サーバの M M I O 空間の関係を示す説明図。

【図 5】各ブレードと P C I 管理サーバのルーティング I D の関係を示す説明図。

【図 6】M R - S R 変換論理の宛先情報変換回路の構成を示すブロック図。

【図 7】T L P 情報の構成を示す説明図。

【図 8】宛先情報変換回路を構成する宛先アドレス修飾部の構成を示すブロック図。

【図 9】宛先情報変換回路を構成する宛先 I D 修飾部の構成を示すブロック図。

【図 1 0】発行元 I D 変換回路 1 の構成を示すブロック図。

【図 1 1】仮想階層ナンバー付加回路の構成を示すブロック図。

【図 1 2】発行元 I D 変換回路 2 の構成を示すブロック図。

【図 1 3】宛先 I D 変換回路 2 の構成を示すブロック図。

40

【図 1 4】マルチルート T L P 及び P C I e ベース T L P のフォーマットを示す説明図。

【図 1 5】マルチルート T L P 及び P C I e ベース T L P の T L P ヘッダのフォーマットを示す説明図。

【図 1 6】ブレードサーバの電源投入の手順を示すフローチャート。

【図 1 7】P C I 管理サーバ 2 0 で行われる起動処理を示すフローチャート。

【図 1 8】ブレードで行われる仮想マシンモニタの起動処理を示すフローチャート。

【図 1 9】マルチルート化可能フォーマットを示す説明図。

【図 2 0】仮想機能確保要求情報のフォーマットを示す説明図。

【図 2 1】仮想機能確保完了情報のフォーマットを示す説明図。

【図 2 2】仮想計算機の起動処理を示すフローチャート。

50

【図 23】PCI-e スイッチのMRS 論理の変形例を示すブロック図。

【符号の説明】

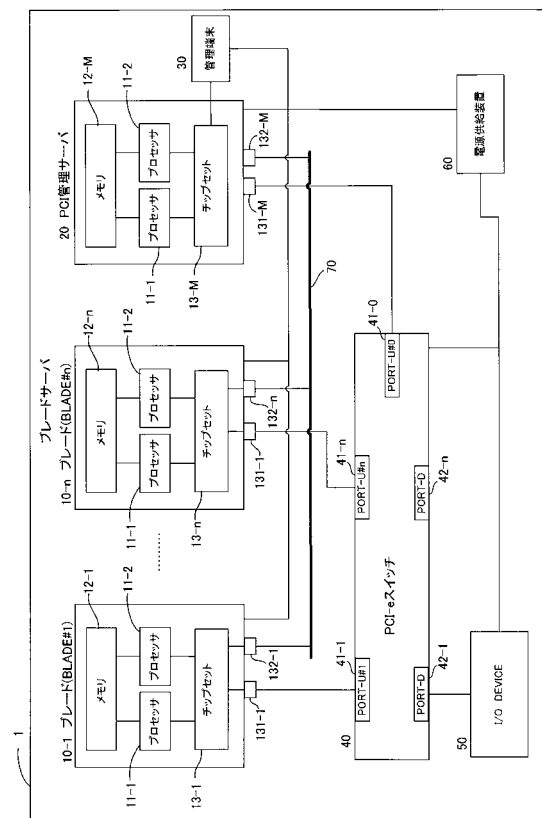
【0125】

- 10 - 1 ~ 10 - n ブレード
- 13 - 1 ~ 13 - n、13 - M チップセット
- 20 PCI 管理サーバ 20
- 40 PCI-e スイッチ
- 50 I/O デバイス 50
- 41 - 0 ~ 41 - n アップストリームポート
- 42 - 1 ダウンストリームポート
- 410 MRS 論理
- 430 MR - SR 変換論理
- 431 宛先情報変換回路
- 432 発行元 ID 変換回路 (1)
- 433 Vhn 消去部
- 434 LCRC 再生部
- 435 Vhn 付加部
- 436 発行元 ID 変換回路 (1)
- 437 宛先 ID 変換回路 (2)
- 438 LCRC 再生部

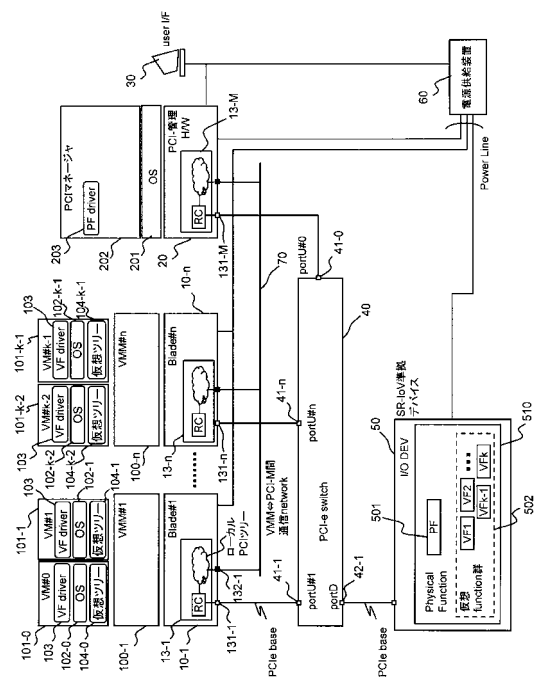
10

20

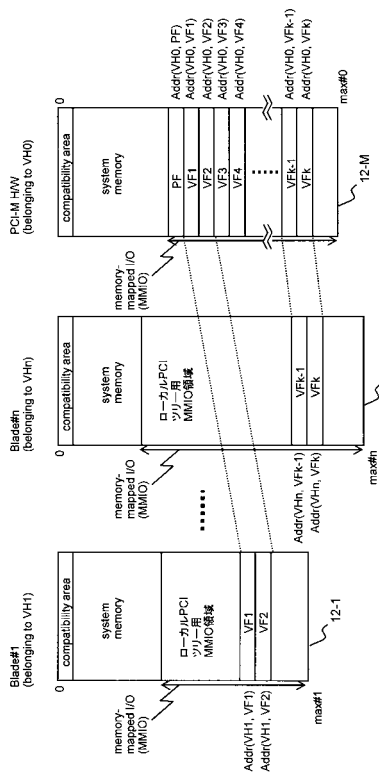
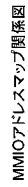
【図 1】



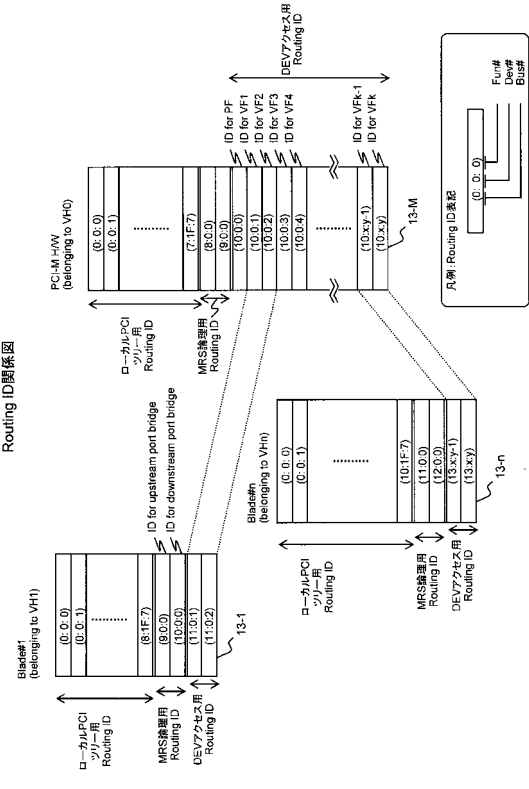
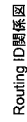
【図 2】



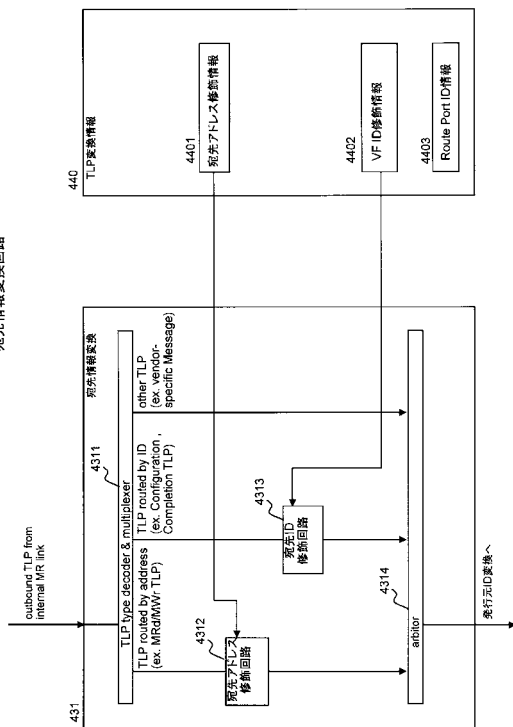
【 図 4 】



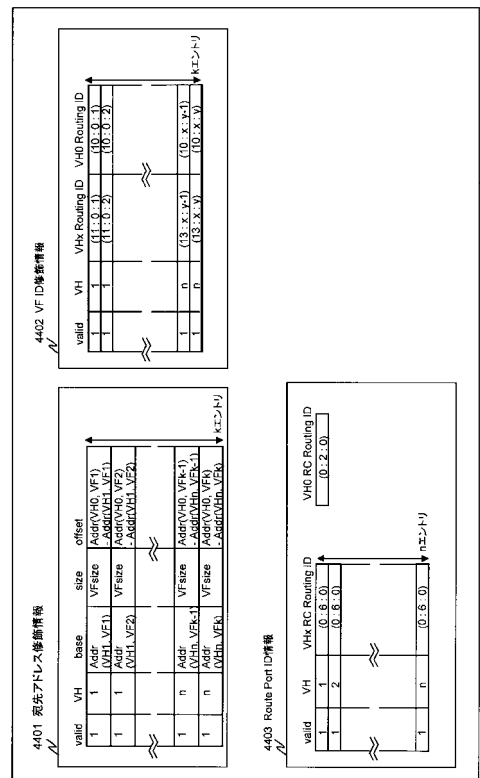
【 図 5 】



【 図 6 】

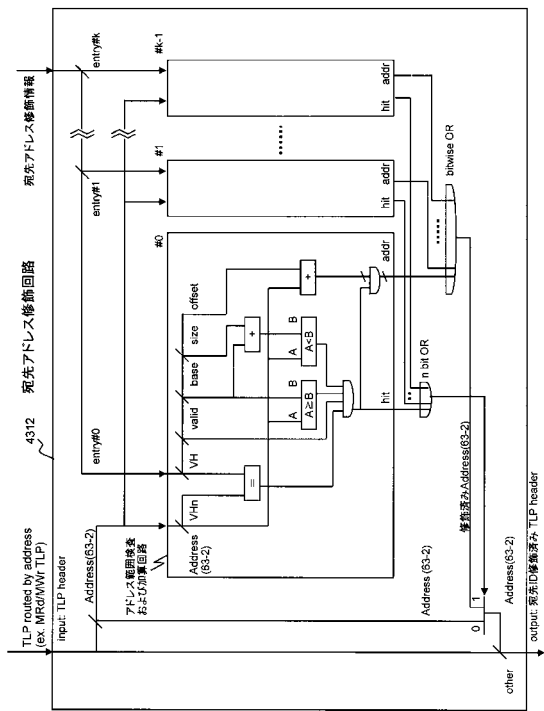


【圖 7】

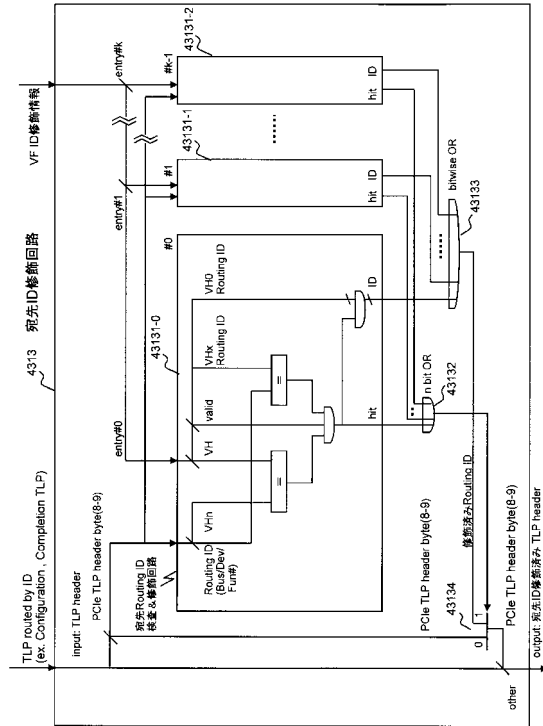


440 TLP变换情报

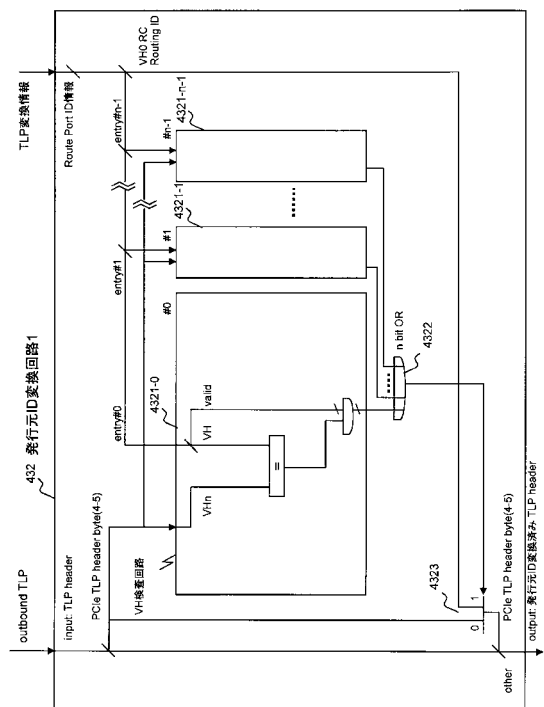
【図 8】



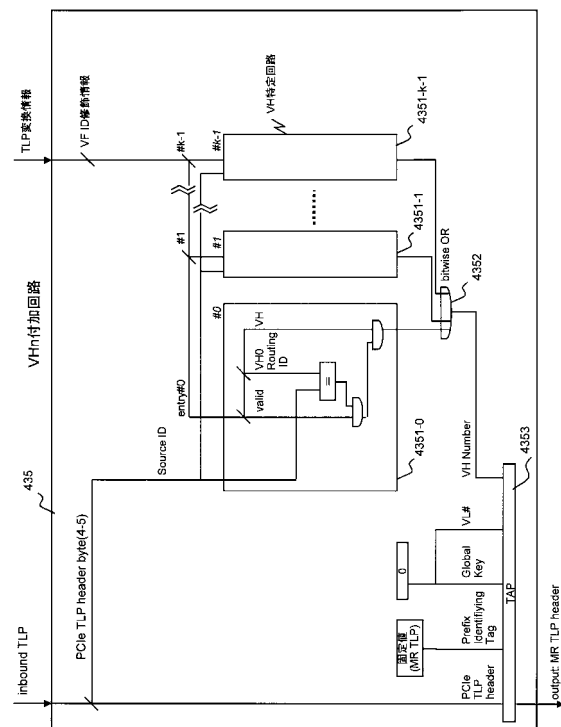
【図 9】



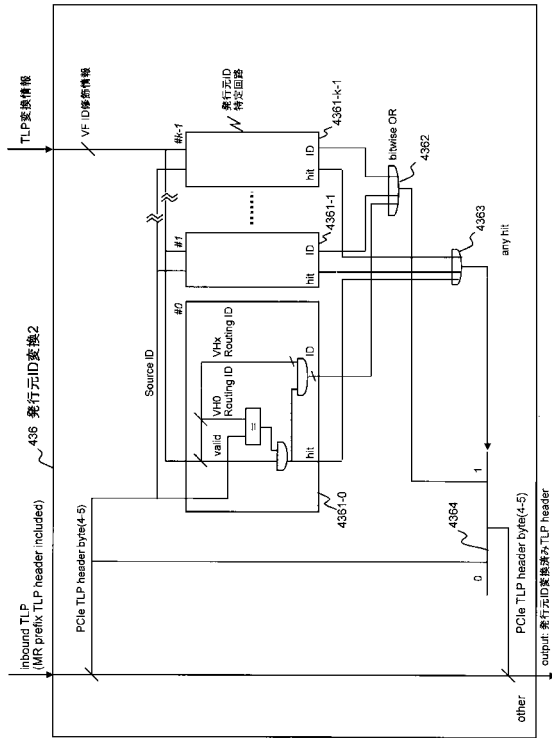
【図 10】



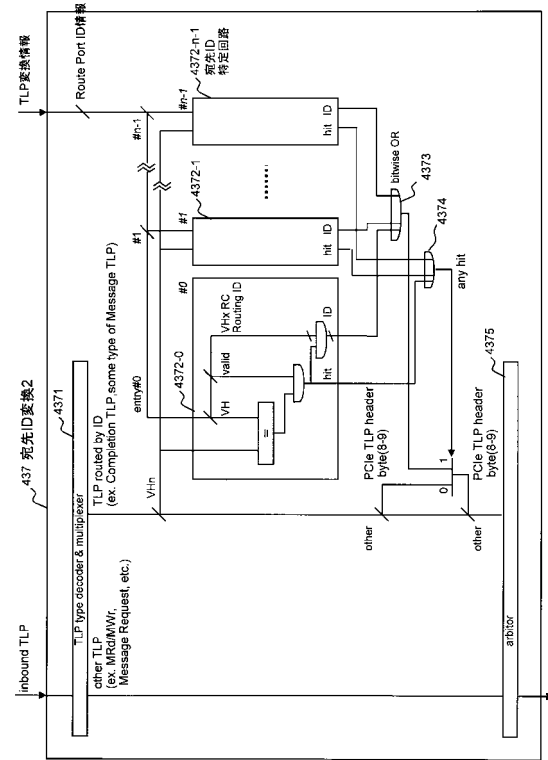
【図 11】



【図 1 2】

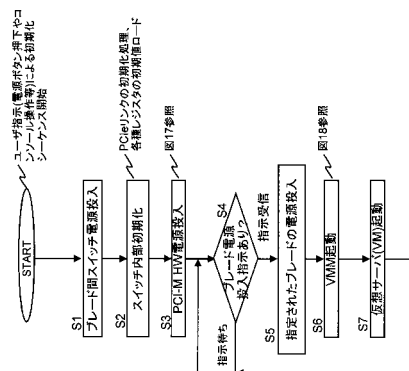


【図 1 3】



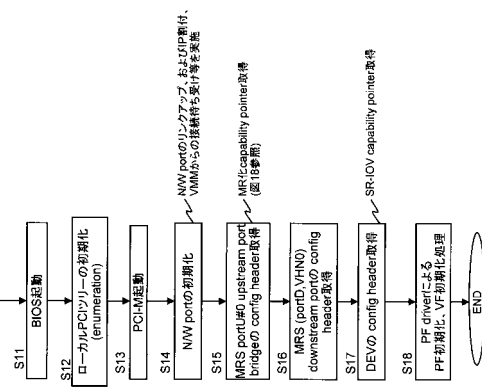
【図 1 6】

システム全体の電源投入シーケンス

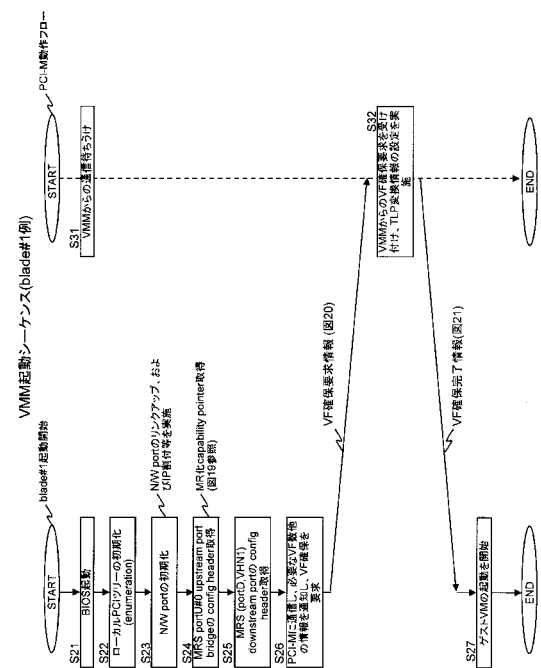


【図 1 7】

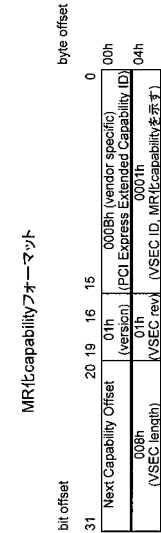
PCI-M HW電源投入シーケンス



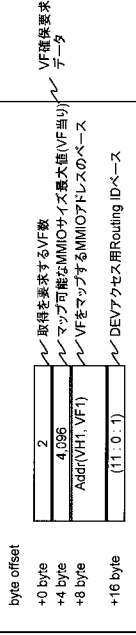
【図 18】



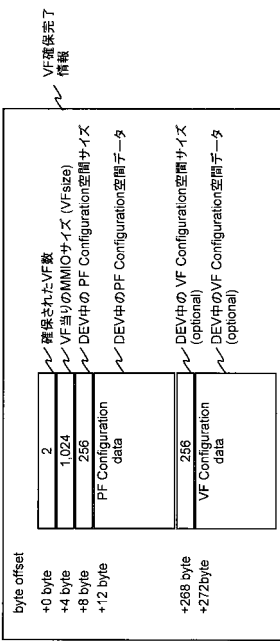
【図 19】

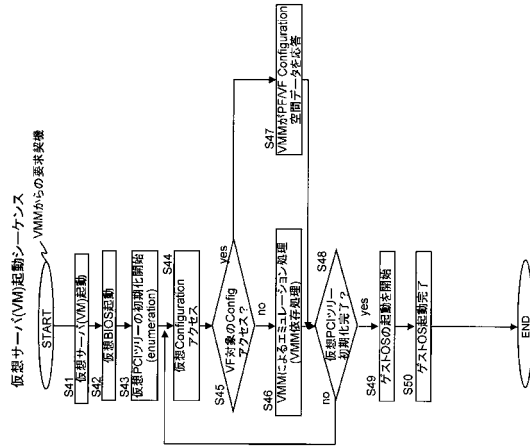


【図 20】

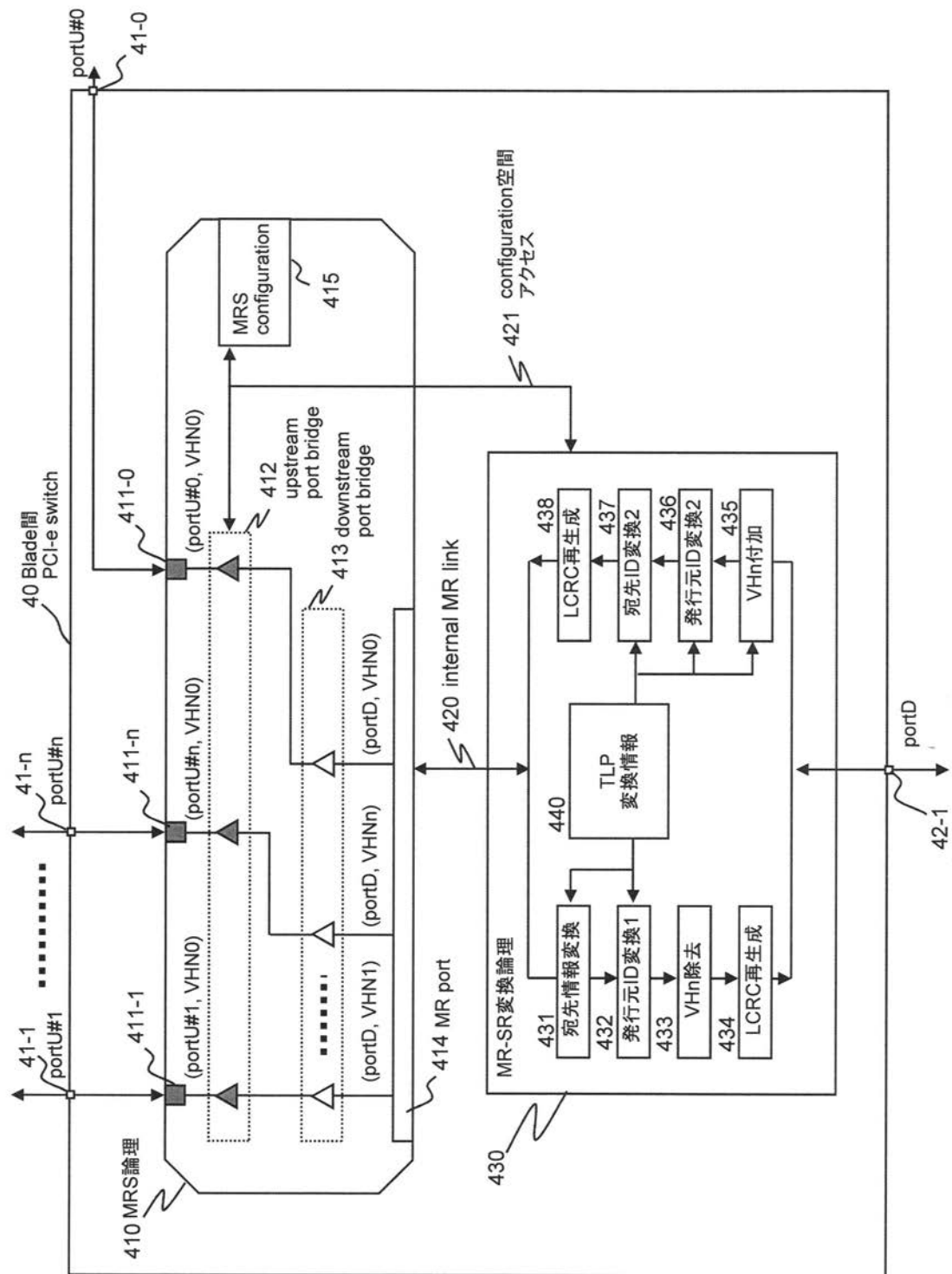


【図 21】

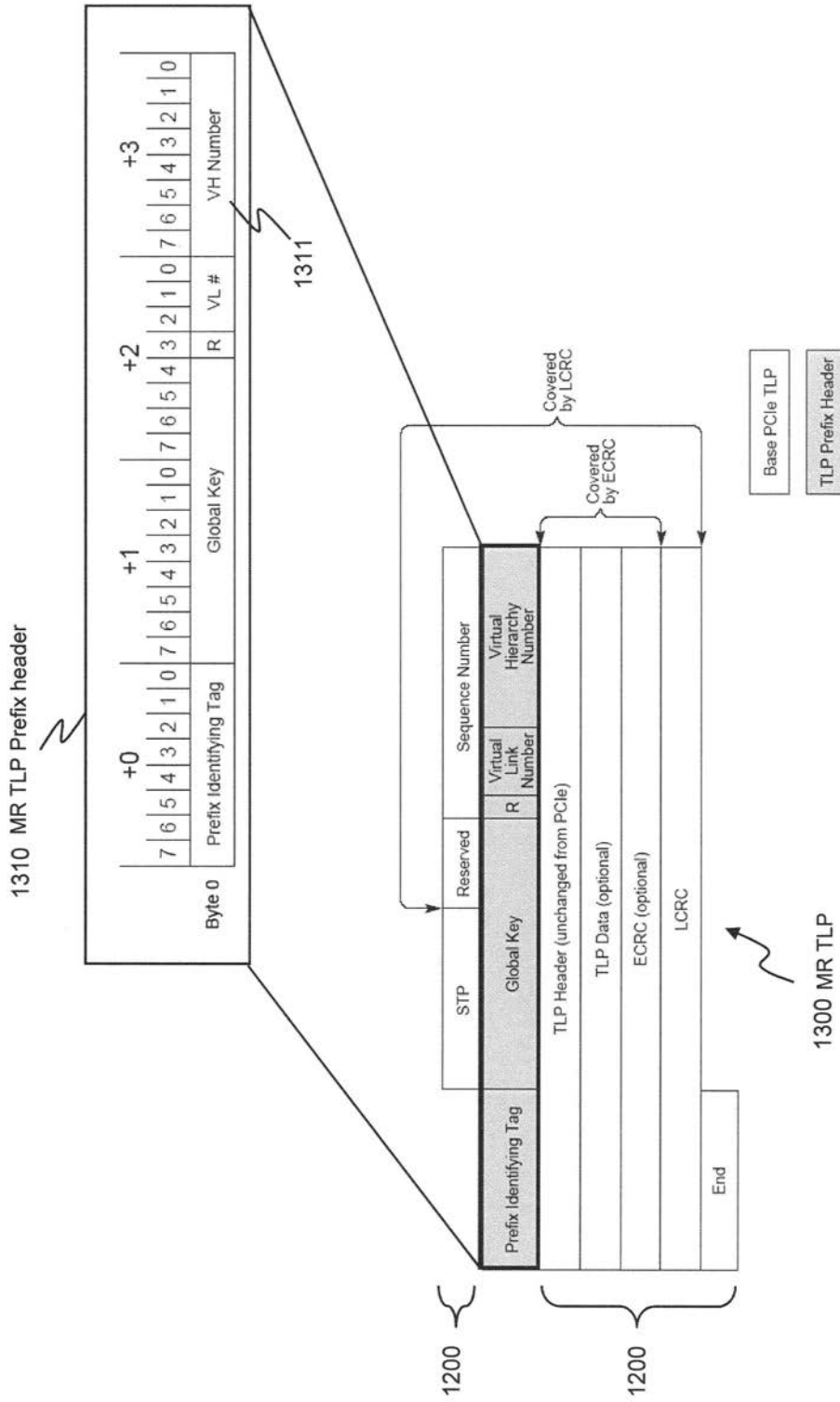




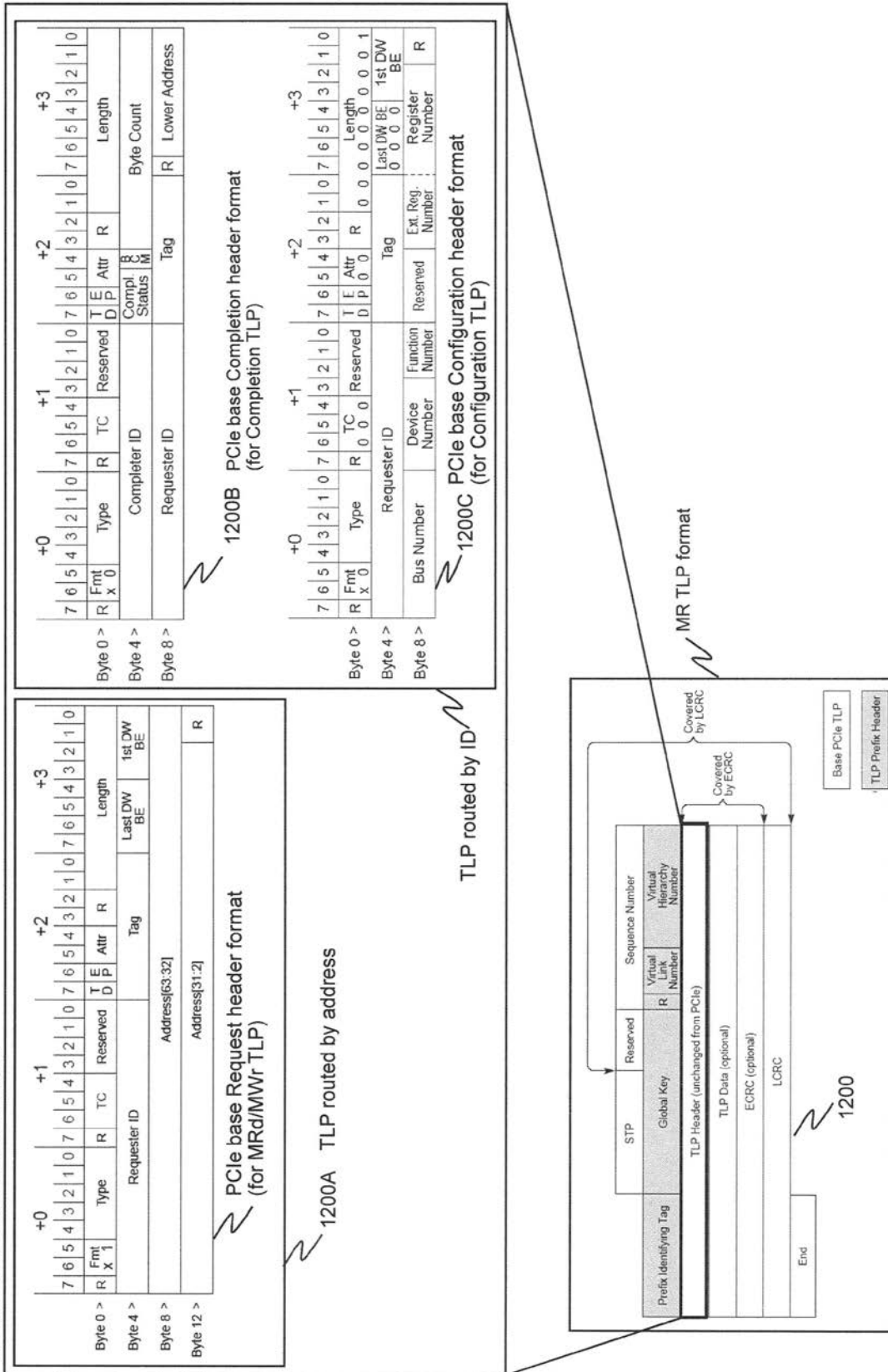
【図 3】



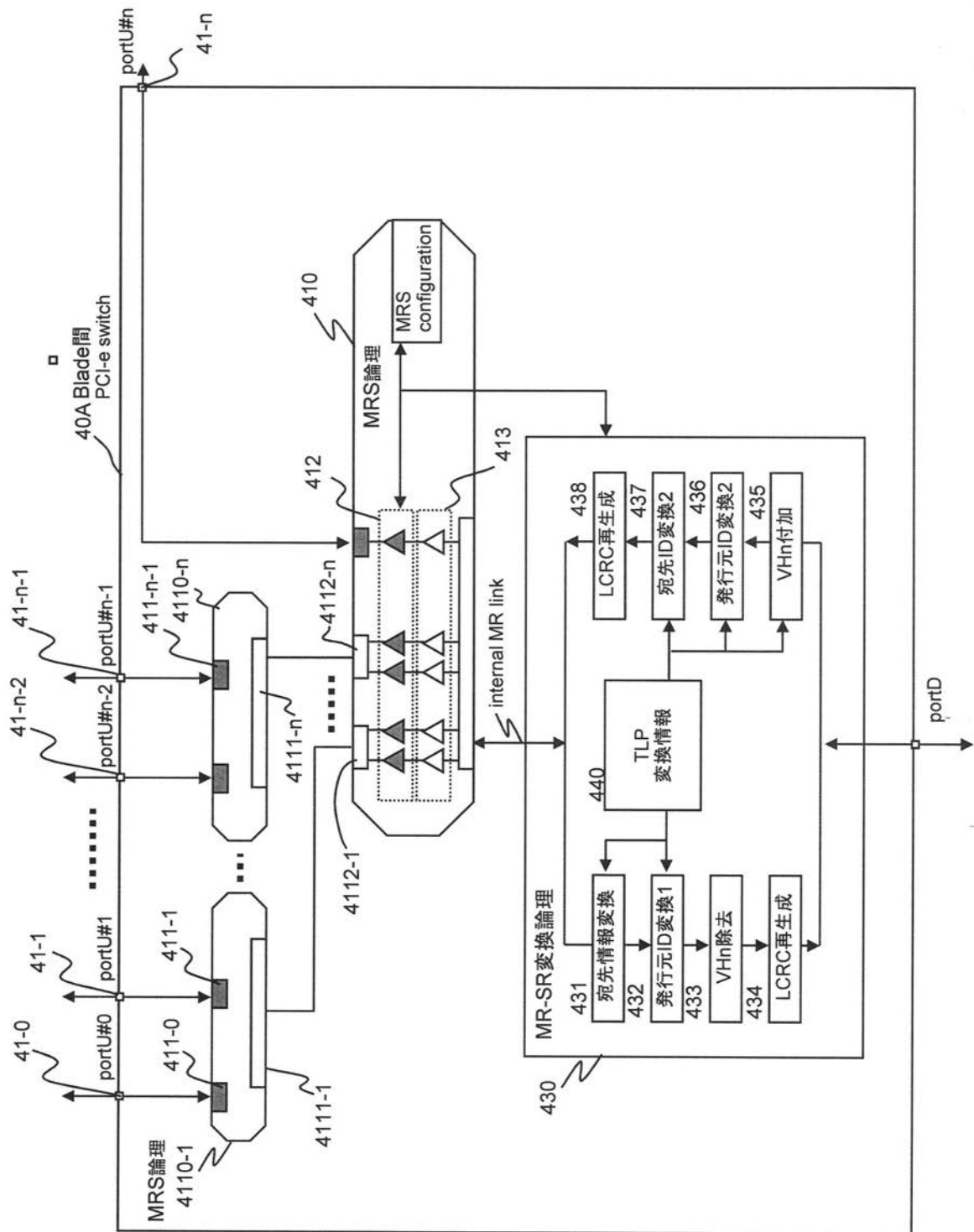
MR TLPフォーマット、およびPCIe base TLPフォーマット



PCle base header format example MR TLPフォーマット、およびPCle base TLPフォーマット



【図23】



フロントページの続き

(72)発明者 上原 敬太郎

東京都国分寺市東恋ヶ窪一丁目 2 8 0 番地 株式会社日立製作所 中央研究所内

審査官 横山 佳弘

(56)参考文献 特開 2 0 0 8 - 0 7 8 8 8 7 (J P , A)

特開 2 0 0 9 - 3 0 1 1 6 2 (J P , A)

国際公開第 2 0 0 8 / 0 1 8 4 8 5 (W O , A 1)

特開 2 0 0 8 - 0 2 1 2 5 2 (J P , A)

特開 2 0 0 7 - 2 1 9 8 7 3 (J P , A)

鈴木 順 Jun Suzuki, E x p E t h e r (エクスプレスイーサ) による単一ホスト仮想化対
応 I / O のマルチホスト同時共有, 第 7 0 回 (平成 2 0 年) 全国大会講演論文集 (1) アーキ
テクチャ ソフトウェア科学・工学 データベースとメディア, 日本, 社団法人情報処理学会,
2 0 0 8 年 3 月 1 3 日, p . 1 - 2 3 ~ p . 1 - 2 4

(58)調査した分野(Int.Cl., D B 名)

G 0 6 F 1 3 / 1 0

G 0 6 F 1 3 / 1 4

G 0 6 F 9 / 4 6