

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4394624号  
(P4394624)

(45) 発行日 平成22年1月6日(2010.1.6)

(24) 登録日 平成21年10月23日(2009.10.23)

(51) Int. Cl. F I  
G O 6 F 9 / 4 6 ( 2 0 0 6 . 0 1 ) G O 6 F 9 / 4 6 3 5 0

請求項の数 10 (全 26 頁)

(21) 出願番号	特願2005-273400 (P2005-273400)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成17年9月21日(2005.9.21)	(74) 代理人	100075513 弁理士 後藤 政喜
(65) 公開番号	特開2007-86963 (P2007-86963A)	(74) 代理人	100084537 弁理士 松田 嘉夫
(43) 公開日	平成19年4月5日(2007.4.5)	(74) 代理人	100114236 弁理士 藤井 正弘
審査請求日	平成20年1月28日(2008.1.28)	(72) 発明者	垂井 俊明 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	保田 淑子 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】 計算機システム及びI/Oブリッジ

(57) 【特許請求の範囲】

【請求項1】

複数のCPUコアと、前記CPUコアに接続されたI/Oブリッジと、前記CPUコアまたはI/Oブリッジからアクセス可能な主記憶と、を備えたCPUモジュールと、

前記CPUモジュールのI/OブリッジとI/Oモジュールとを接続するI/Oスイッチと、を備えた計算機システムにおいて、

前記CPUモジュールは、

前記複数のCPUコアと主記憶とを複数の論理区画に分割するファームウェアを備え、

前記I/Oブリッジは、

前記論理区画とI/Oモジュールとの間で送受信されるI/Oアクセス情報を中継する際に、前記論理区画毎に設定された仮想経路情報と、当該I/OブリッジからI/Oモジュールまでの経路情報とを前記I/Oアクセス情報の経路情報に付加して、論理区画毎に前記I/Oモジュールとの間のI/Oアクセス情報を切り換える仮想スイッチを備えたことを特徴とする計算機システム。

【請求項2】

請求項1に記載の計算機システムにおいて、

前記仮想スイッチは、

前記論理区画からI/Oモジュールへ送信するI/Oアクセス情報の経路情報に、前記論理区画毎の仮想経路情報を前記仮想スイッチの経路情報として付加し、前記I/Oスイッチへ送信する経路情報設定部と、

10

20

前記 I / O モジュールから受信した I / O アクセス情報に含まれる経路情報から前記仮想スイッチの仮想経路情報を抽出する宛先抽出部と、

前記抽出した仮想経路情報に対応する論理区画を特定して前記 I / O アクセス情報を当該論理区画へ転送する宛先論理区画特定部と、  
を備えたことを特徴とする計算機システム。

【請求項 3】

請求項 1 に記載の計算機システムにおいて、

前記仮想スイッチは、

前記論理区画が割り当てられた主記憶上の領域のベースアドレスとサイズを、論理区画に対応付けて格納するアドレス管理部と、

前記論理区画毎にマルチキャストによる書き込みを制御するための宛先論理区画設定部と、

前記 I / O モジュールから受信した I / O アクセス情報に含まれるコマンドを抽出するコマンド抽出部と、

前記抽出したコマンドがマルチキャスト DMA のときには、前記論理区画に対応するベースアドレスとサイズを前記アドレス管理部から取得し、前記マルチキャスト DMA コマンドに含まれる DMA アドレスに前記ベースアドレスを論理区画毎に加算するアドレス変換部と、

前記加算結果が前記サイズで示される前記論理区画に割り当てられた主記憶上の領域を超えていなければ、前記加算結果が示す論理区画に対応する主記憶のアドレスにそれぞれ DMA アクセスを行う DMA 処理部と、

マルチキャストの宛先である全ての論理区画に対して上記書き込みを行う、複数回書き込み処理部と、

を備えたことを特徴とする計算機システム。

【請求項 4】

請求項 1 に記載の計算機システムにおいて、

前記仮想スイッチは、

前記論理区画毎にマルチキャストによるイベントを制御するための宛先論理区画設定部と、

前記 I / O モジュールから受信した I / O アクセス情報に含まれるコマンドを抽出するコマンド抽出部と、

前記抽出したコマンドがマルチキャストイベントのときには、宛先論理区画番号を前記宛先論理区画設定部に設定し、前記ファームウェアにマルチキャストイベントに含まれるイベントを通知するマルチキャストイベント通知部と、を有し、

前記ファームウェアは、前記通知を受けたイベントを前記宛先論理区画設定部の設定に基づく複数の論理区画に通知することを特徴とする計算機システム。

【請求項 5】

請求項 1 に記載の計算機システムにおいて、

前記仮想スイッチは、

前記論理区画毎に設定する仮想ポートを複数備え、

前記論理区画と仮想ポートとの対応関係と、論理区画の状態を管理する論理区画管理部を有し、

前記ファームウェアは、前記論理区画の生成時または削除時に前記論理区画管理部の前記状態を更新することを特徴とする計算機システム。

【請求項 6】

請求項 1 に記載の計算機システムにおいて、

前記 I / O スイッチは、複数の CPU ブレードと複数の I / O モジュールとを接続し、

前記 CPU モジュールに設定された前記複数の論理区画に対する前記 I / O モジュールの割り当て状態を管理するファブリック管理装置を備え、

前記ファブリック管理装置は、ひとつの I / O モジュールを複数の論理区画または CP

10

20

30

40

50

Uモジュールに割り当ててることを特徴とする計算機システム。

**【請求項 7】**

複数のCPUコアと主記憶とをファームウェアによって分割した複数の論理区画をI/Oスイッチに接続し、前記I/Oスイッチに接続されたI/Oモジュールと論理区画との間でI/Oアクセスを行うI/Oブリッジであって、

前記I/Oブリッジは、

前記論理区画とI/Oモジュールとの間で送受信されるI/Oアクセス情報を中継する際に、前記論理区画毎に設定された仮想経路情報を前記I/Oアクセス情報の経路情報に付加して、論理区画毎に前記I/Oモジュールとの間のI/Oアクセス情報を切り換える仮想スイッチを備えたことを特徴とするI/Oブリッジ。

10

**【請求項 8】**

請求項 7に記載のI/Oブリッジにおいて、

前記仮想スイッチは、

前記論理区画からI/Oアクセス情報を受信すると、このI/Oアクセス情報を発行した論理区画を特定する発信元特定部と、

論理区画に対応する前記仮想スイッチの仮想経路情報を前記I/Oアクセス情報に付加し、前記I/Oスイッチへ送信する経路情報設定部と、  
を有することを特徴とするI/Oブリッジ。

**【請求項 9】**

請求項 7に記載のI/Oブリッジにおいて、

前記仮想スイッチは、

前記I/Oモジュールから受信したI/Oアクセス情報に含まれる経路情報から前記仮想スイッチの仮想経路情報を抽出する宛先抽出部と、

前記抽出した仮想経路情報に対応する論理区画を特定して前記I/Oアクセス情報を当該論理区画へ転送する宛先論理区画特定部と、  
を備えたことを特徴とするI/Oブリッジ。

20

**【請求項 10】**

請求項 7に記載のI/Oブリッジにおいて、

前記仮想スイッチは、

前記論理区画が割り当てられた主記憶上の領域のベースアドレスとサイズを、論理区画に対応付けて格納するアドレス管理部と、

前記論理区画毎にマルチキャストによる書き込みを制御するための宛先論理区画設定部と、

前記I/Oモジュールから受信したI/Oアクセス情報に含まれるコマンドを抽出するコマンド抽出部と、

前記抽出したコマンドがマルチキャストDMAのときには、前記論理区画に対応するベースアドレスとサイズを前記アドレス管理部から取得し、前記マルチキャストDMAコマンドに含まれるDMAアドレスに前記ベースアドレスを論理区画毎に加算するアドレス変換部と、

前記加算結果が前記サイズで示される前記論理区画に割り当てられた主記憶上の領域を超えていなければ、前記加算結果が示す論理区画に対応する主記憶のアドレスにそれぞれDMAアクセスを行うDMA処理部と、

マルチキャストの宛先である全ての論理区画に対して上記書き込みを行う、複数回書き込み処理部と、

を備えたことを特徴とするI/Oブリッジ。

30

40

**【発明の詳細な説明】**

**【技術分野】**

**【0001】**

本発明は、仮想計算機システムに関し、複数の論理区画とI/Oデバイスの割当をI/Oスイッチで実現する技術に関する。

50

## 【背景技術】

## 【0002】

近年、サーバ台数の増加と共に運用に関する複雑さが増加し、運用コストの増大が問題化している。この運用コストを低減する技術として複数サーバを1台にまとめるサーバコンソリデーション（サーバ統合）が注目を集めている。

## 【0003】

サーバコンソリデーションを実現する技術のひとつとして、一つの計算機を任意の割合で論理的に分割する仮想計算機が知られており、ハイパバイザなどのファームウェア（またはミドルウェア）により、物理計算機を複数の論理区画（LPAR: Logical PARTition）に分割し、各LPARに対して計算機資源（CPU、主記憶、I/O）を割当て、各LPAR上でそれぞれOSを動作させるものが知られている（例えば、特許文献1）。

10

## 【0004】

あるいは、VMware（登録商標）のように、ホストOS上で複数のゲストOSを稼働させ、各ゲストOSを論理区画として提供する技術も知られている。

## 【0005】

また、近年、基板単位で計算機を変更可能なブレードサーバにより、ITシステムの構成をより柔軟に変更する技術が知られており、多数のサーバをブレードサーバに統合するサーバコンソリデーション（サーバ統合）が行われている。

## 【0006】

さらに近年では、AS（Advanced Switching）等に代表される、I/Oスイッチが提案されている（例えば、非特許文献1）。この種のI/Oスイッチは、複数のCPU（物理的なCPU）からI/Oカード（またはI/Oデバイス）の共有を実現しようとするものである。ASでは、現在普及しつつある、PCI-EXPRESS規格のI/Oデバイス（またはI/Oカード）を利用できるため、汎用性が高い。

20

【特許文献1】特開2002-304364号

【非特許文献1】「Advanced Switching 技術概要」、[online]、「平成17年8月16日検索」、インターネット<[http://www.asi-sig.org/education/ASI\\_AdvSwitch\\_TB\\_JP.pdf](http://www.asi-sig.org/education/ASI_AdvSwitch_TB_JP.pdf)>

## 【発明の開示】

## 【発明が解決しようとする課題】

30

## 【0007】

上記VMwareによる論理区画間のI/O共有では、I/Oデバイス側で論理区画を識別することができないため、各論理区画のI/O要求を取りまとめて物理的なI/Oデバイスにリクエストを出す必要がある。このため、ホストOSがI/O操作を集中管理することで、ゲストOSで構成される各論理区画間でI/Oデバイスの共有を実現することになる。この結果、各論理区画からのI/O要求は、一旦、ホストOSの処理を待ってから物理的なI/Oデバイスに対して操作が行われるため、ホストOSの処理がオーバーヘッドとなって、I/O操作のレスポンスが低下するという問題がある。

## 【0008】

上記ハイパバイザでは、論理区画を意識することができる特別なI/O装置（チャネル等）を用いることで、仮想計算機間でI/Oの共有を行っている。I/O装置がパーティションを意識してI/O操作を行い、各論理区画からI/Oデバイスを直接アクセスできる。しかし、特別なI/O装置を必要とするため製造コストが高く、また、汎用性も低いという問題がある。

40

## 【0009】

一方、上記従来例のハイパバイザにAdvanced Switching等のI/Oスイッチを用いて論理区画間でI/Oデバイスの共有を行う場合には、次のような問題が生じる。

## 【0010】

上述のAdvanced Switching（以下、単にASとする）では、物理的なノード（またはデバイス）しか識別できないため、ASではハイパバイザが提供する論理区画を識別できず

50

、上記ハイパバイザにASをそのまま適用することはできない。このため、上述のVMwareと同様にASのI/O操作を集中管理するソフトウェアが必要になってしまい、このソフトウェアの処理がそのままI/Oアクセスのオーバーヘッドになるという問題を生じていた。

【0011】

そこで本発明は、上記問題点に鑑みてなされたもので、AS等に代表される汎用性の高いI/Oスイッチを用いて、仮想計算機の論理区画間でI/Oの共有を実現する際のオーバーヘッドを低減することを目的とする。

【課題を解決するための手段】

【0012】

本発明は、複数のCPUコアと、前記CPUコアに接続されたI/Oブリッジと、前記CPUコアまたはI/Oブリッジからアクセス可能な主記憶と、を備えたCPUモジュールと、前記CPUモジュールのI/OブリッジとI/Oモジュールとを接続するI/Oスイッチと、を備えた計算機システムにおいて、前記CPUモジュールは、前記複数のCPUコアと主記憶とを複数の論理区画に分割するファームウェアを備え、前記I/Oブリッジは、前記論理区画とI/Oモジュールとの間で送受信されるI/Oアクセス情報を中継する際に、前記論理区画毎に設定された仮想経路情報と、当該I/OブリッジからI/Oモジュールまでの経路情報とを前記I/Oアクセス情報の経路情報に付加して、論理区画毎に前記I/Oモジュールとの間のI/Oアクセス情報を切り換える仮想スイッチを備える。

【発明の効果】

【0013】

したがって、本発明は、I/OスイッチでI/Oモジュールを複数の論理区画間で共有する際に、CPUモジュールに備えたI/Oブリッジを仮想スイッチとして機能させ、各論理区画毎に仮想経路情報を設定することで、ひとつのCPUモジュールで複数の論理区画を提供する仮想計算機システムにおいてI/Oスイッチを用いたI/O共有を実現することができる。

【0014】

さらに、I/Oブリッジのハードウェアレベルで論理区画間のI/O共有を行うようにすることで、論理区画を識別するためのソフトウェアは不要となってI/Oアクセスのオーバーヘッドを低減して高速化を図ることができる。

【発明を実施するための最良の形態】

【0015】

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。

【0016】

図1は第1の実施形態を示し、本発明を適用するブレードサーバシステム（物理計算機）のブロック図である。

【0017】

ブレードサーバシステムは、複数のCPUブレード#0~2と、各種I/Oインターフェースを備えたI/Oブレード（またはI/Oカード）#0~5と、CPUブレード#0~2とI/Oブレード#0~5を接続する複数のI/Oスイッチ#0~2を備えたスイッチブレード3が接続されて、図示しない筐体に格納されている。なお、スイッチブレード3にはI/Oスイッチ#0~2やCPUブレード#0~2に対するI/Oブレード#0~5の割り当てを管理するファブリック管理サーバ4が接続され、ファブリック管理サーバ4には管理者に対する入出力が可能なコンソール5が接続されている。なお、CPUブレード（CPUモジュール）#0~2、I/Oブレード（I/Oモジュール）#0~5、スイッチブレード（スイッチモジュール）3は図示しないバックプレーンを介して接続されている。

【0018】

本実施形態のブレードサーバシステムでは、後述するように、CPUブレード#0~2

10

20

30

40

50

の計算機資源を複数の論理区画（パーティション）に分割した仮想計算機が稼動する。

【0019】

以下に、ブレードサーバの各ブレードの概要を説明する。

【0020】

<CPUブレード>

CPUブレード#0～2はそれぞれ、複数のCPUコア#0、#1を備えたCPU（所謂、マルチコアCPU）を複数備えている。なお、CPUブレード#0～2は同一の構成であるので、CPUブレード#0の構成についてのみ以下に説明する。

【0021】

CPU#0は複数のCPUコア#0-0、#0-1を有し、CPU#1は複数のCPUコア#1-0、#1-1を有する。CPU#0、#1は、フロントサイドバス11を介してノースブリッジ12に接続される。

10

【0022】

ノースブリッジ12はメモリバスを介して主記憶13に接続され、CPU#0、#1のメモリアクセス要求に応じて主記憶13をアクセスする。また、ノースブリッジ12はI/Oバス14を介してI/OブリッジとしてのAS（Advanced Switching）ブリッジ15に接続され、CPU#0、#1のI/Oアクセス要求に応じてASブリッジ15にアクセスを行う。ここでは、I/Oバス14はPCI-EXPRESSで構成される。

【0023】

ASブリッジ15は、AS（Advanced Switching）の規格に準じて構成され、I/Oバス14からのPCI-EXPRESS（以下、PCIeとする）の packets（I/Oアクセス情報）に経路情報を付加してAS packetsとしてから、後述するスイッチ#0（SW0）～#2（SW2）に送信する。また、ASブリッジ15は、後述するスイッチ#0（SW0）～#2（SW2）から受信した packets を、PCIeの packets に変換し、CPUコア#0-0～CPUコア#1-1に割り当てられた論理区画を識別して配信する。このため、ASブリッジ15には、制御部16とメモリ17を備えて、スイッチブレード3に接続されている。

20

【0024】

また、ASブリッジ15は、後述するように、AS（スイッチブレード3）側からの packets を論理区画に割り振る仮想スイッチとして機能する。

30

【0025】

さらに、ノースブリッジ12とASブリッジ15は、CPUブレード#0上のハードウェアを監視するBMC（Baseboard Management Controller または Blade Management Controller）18に接続され、各ブリッジに接続されたハードウェアの監視が行われる。このBMC7は、基板上のハードウェアの電圧、温度、エラーなどを監視し、OSやハイパバイザ等に通知するものである。なお、BMC7は各CPUブレード毎に配置されて、各BMC7はネットワークを介して接続されており、各BMC7を管理するコンソール70を有する。

【0026】

本実施例では、コンソール70とファブリック管理サーバ4は別としたが、同一のサーバとして実装しても良い。

40

【0027】

<スイッチブレード>

スイッチブレード3には、複数のASスイッチ（I/Oスイッチ）#0（SW0）、#1（SW1）、#2（SW2）が配置される。ASスイッチ#0はCPUブレード#0～2に接続され、ASスイッチ#1、2はI/Oブレード#0～5と接続されて、ASスイッチ#0～2が相互に接続される。また、ASスイッチ#0～2はそれぞれ相互に接続される。

【0028】

各ASスイッチ#0～2は、複数のポートを有しており、例えば、8ポートで構成され

50

る。ASスイッチ#0は、ポート7をCPUブレード#0のASブリッジ15に接続し、ポート0をCPUブレード#1のASブリッジ(図示省略)に接続し、ポート1をCPUブレード#2のASブリッジ(図示省略)に接続し、ポート2をASスイッチ#1のポート7に接続し、ポート5をASスイッチ#2のポート1に接続している。

【0029】

ASスイッチ#1は、ポート3~5をI/Oブレード#3~5に接続し、ポート7をASスイッチ#0のポート2に接続し、ポート6をASスイッチ#1のポート2に接続している。

【0030】

ASスイッチ#2は、ポート3~5をI/Oブレード#0~2に接続し、ポート1をASスイッチ#0のポート5に接続し、ポート2をASスイッチ#1のポート6に接続している。

【0031】

ここで、ASにおけるパケット形式とルーティングについて述べる。ASスイッチ#0~2を通過する従来のAdvanced Switchingパケット(以下、ASパケットとする)は、図5~6で示すように、ヘッダー部P0とデータ部P1から構成されたPCIEのパケットに、ASヘッダーが付加されている。

【0032】

ASパケットには単一の宛先に送られるユニキャストパケットと、複数の宛先に送られるマルチキャストパケットがある。ユニキャスト/マルチキャストは、ASヘッダー内のプロトコルインタフェース(コマンド種別を表す)フィールド(図示しない)により区別される。以下では、各々の場合のルーティング方式を説明する。

(1)ユニキャスト(マルチキャストでない)パケットの場合

ユニキャストパケットでは、経路情報として各スイッチでの切り換え数を示すターンプール値TP1~n、パケットの転送方向(上り/下り)を表すビットDIR等の情報を持つ(図では、TP値、DIR以外のヘッダー情報は略す)。DIRは0が下り、1が上りをあらわす。このターンプール値TP1~nは、ASパケットが通過するASスイッチの数に応じてASブリッジ15またはI/Oブレード#0~5が付加する。

【0033】

そして、本発明のASブリッジ15は、図5、図6で示すように、ターンプール値TP1~nに加えて、論理区画を識別するための仮想スイッチターンプール値TP0を設定する。

【0034】

ASスイッチ#0~2は、図5に示す下りのASパケットと図6に示す上りのASパケットのターンプール値に基づいて転送先を決定する。なお、本実施形態では、CPUブレードからI/Oブレードへ向かうパケットを下りのASパケットとし、逆に、I/OブレードからCPUブレードへ向かうパケットを上りのASパケットとする。

【0035】

ASスイッチ#0~2は、下りASパケットの場合は、ASパケットを受信した隣のポートから時計回りに0からターンプール値まで数えたポートを転送先として決定し、ASパケットを転送する。下りのパケット(DIR=0)の場合、ASスイッチを一段渡る毎に、ターンプールを(TP1~TPnの順番に)左から右に使用する。

【0036】

例えば、図4において、ASスイッチ#0がポート7に接続されたCPUブレード#0から受信した下りASパケットのターンプール値が2であれば、時計回りに0-1-2とターンプール値の数だけ0からカウントしてポート2をASパケットの転送先とする。

【0037】

それに対して上りパケット(DIR=1)の場合、ターンプールの値は、反時計回りに数えられ、ターンプールは(TPn~TP1の順番に)右から左に使用される。

【0038】

10

20

30

40

50

A Sスイッチ# 0が、ポート2に接続されたA Sスイッチ# 1から受信した上りA Sパケットのターンプール値が2であれば、反時計回りに1 - 0 - 7とターンプール値の数だけ0からカウントしてポート7をA Sパケットの転送先とする。

【0039】

以上により、要求(下り)パケットとそれに対する応答(上り)パケットのターンプール値は同一とすることができ、D I Rビットを(0から1に)反転するだけで応答パケットのルーティング情報を作成することができる。これにより、ブリッジにおける応答パケットの生成を容易に行うことができる。

【0040】

なお、ターンプール値T P 0 ~ nの詳細については、後述する。

10

【0041】

なお、P C I eパケットのデータ部P 1には、データやコマンド(I/Oアクセスコマンド)等が格納される。

(2) マルチキャストパケットの場合

A Sでは複数の宛先に同一のデータやイベントを送付するマルチキャストを行うことができる。マルチキャストの宛先はM G I D (マルチキャストグループID)と呼ばれるシステム内で一意な番号により管理される。同一のM G I Dを持つパケットは、常に同一の宛先にマルチキャストされる。それを実現するために、あらかじめ各A Sスイッチに、マルチキャストグループ毎に、どのポートにデータを出力する必要があるかがテーブルで設定される。図17にマルチキャストパケットの形式、図18にマルチキャストを行うA S

20

スイッチを示す。

【0042】

図17に示すように、マルチキャストパケットでは、パケットのA Sヘッダーでは、マルチキャストグループID (M G I D) をルーティング情報として持つ(M G I D以外のヘッダー情報は省略する)。A Sスイッチでは、各M G I Dにおいて、A Sスイッチのどのポートにパケットを送付すれば良いかを示すマルチキャストテーブルを持つ。図18にマルチキャストテーブル180の構成とA Sスイッチの動作を示す。マルチキャストテーブル180では、各M G I Dにおいて、A Sスイッチの各物理ポートにパケットへパケットを出力する/しないを、ビットマップで表す。例えば図18では、M G I Dが25のパケットはマルチキャストテーブル180の該当するエントリが1であるポート3、4、5

30

に出力される。

【0043】

本発明では、各A Sスイッチに加え、ブリッジにおいても、どの論理区画にパケットをマルチキャストするかを示す仮想スイッチのマルチキャストテーブルを持つことにより、論理区画に向けたマルチキャストを可能にする。仮想スイッチのマルチキャストテーブルの詳細は後述する。

【0044】

マルチキャストテーブル180はA Sスイッチ# 0 ~ 2 毎におかれ、システム初期化時や構成変更時にファブリック管理サーバ4により初期化される。

【0045】

40

なお、I/Oブレード側からのマルチキャストは、例えば、ネットワークインターフェース(図1のE t h e r I / Fを備えたI/Oブレード# 0、# 5等)からの受信データやI/Oブレード側からの割り込みであり、I/Oブレード側で宛先の論理区画を判断できない場合に使用される。

【0046】

< I/Oブレード >

A Sスイッチ# 0 ~ 2 を介してC P Uブレード# 0 ~ 2 からアクセスされるI/Oブレード# 0 ~ 5 には、各種I/Oインターフェースと、A Sスイッチ# 0 ~ 2 と各I/Oインターフェースを接続するためのL E A Fブリッジが設けられる。

【0047】

50



I/Oブレード#0~5のLEAFブリッジは、上記図4で示したように、I/Oブレードが受信した下りASパケットのターンブル値TP0~nを一旦記憶し、この下りASパケットの返信となる上りASパケットに、受信したASパケットのターンブル値TP0~nを設定するとともに、上述したようにDIRビットを反転する。

【0048】

そして、各I/Oブレード#0~5の各LEAFブリッジは、ASパケットからASヘッダを削除したPCIeパケットを各I/Oインターフェースに渡し、各I/OインターフェースはPCIeのデバイスとして機能する。I/Oインターフェースとしては、図1に示すように、FC(Fibre Channel)インターフェース、SAS(Serial Attached SCSI)インターフェース、イーサネット(登録商標)インターフェース、SATA(Serial ATA)インターフェース、SCSIインターフェースなどで構成される。

10

【0049】

<ファブリック管理サーバ>

スイッチブレード3に接続されたファブリック管理サーバ4は、CPU、メモリ、インターフェース(図示省略)を備え、各CPUブレード#0~2上で稼動する論理区画が使用するI/Oブレード#0~5の割り当ての管理、マルチキャストの宛先の管理、ASスイッチ#0~2やASブリッジの設定を行う。

【0050】

ファブリック管理サーバは、ASファブリック経由で管理を行うIn-bound管理方式、ASファブリック以外の管理用ネットワーク(イーサネット(登録商標)等)(図示しない)経由で管理を行うOut-of-bound管理の2方式がある。

20

【0051】

<仮想計算機システム>

次に、物理計算機を構成するCPUブレード#0~2上で稼動するソフトウェアについて、図2を参照しながら詳述する。

【0052】

各CPUブレード#0~2上ではハイパバイザ200がファームウェアとして稼動しており、ハイパバイザ200は物理計算機100を2つ以上の論理区画(LP AR: Logical Partition)LP AR0(210)~LP ARm(21m)に分割し、計算機資源の割り当てを管理する。

30

【0053】

論理区画LP AR0~LP ARmのそれぞれでOS0(220)~OSm(22m)を動作させ、各OS上でそれぞれアプリケーション0(230)~アプリケーションm(23m)を動作させる。

【0054】

ハイパバイザは、各論理区画LP AR(210~21m)に対して各CPUブレード#0~2のCPUコア、主記憶13、ASブリッジ15に接続されたI/Oブレード#0~5の割り当てのリソース(計算機資源)を割り当てる。

【0055】

なお、CPU#0、#1は複数のCPUコアを持つ場合を示したが、ひとつのCPUがひとつのCPUコアを備えるものであっても良い。また複数のCPUがSMPを構成する場合を示したが、単一CPUでも良い。各論理区画のCPUへの割り当ては、必ずしもコア毎ではなくて良く、単一のコアを複数の論理区画が時分割で共有しても良い。

40

【0056】

<本願発明の概要>

次に、本発明の要部となるASブリッジとASパケットの概要について、以下に説明する。

【0057】

各CPUブレード#0~2に配置されたASブリッジ15は、図7で示すように、仮想のASスイッチとして機能する。物理的には、図1、図4で示したように、ASブリッジ

50

15とASスイッチ#0がASにより接続され、ASブリッジ15とノースブリッジ12がPCIeにより接続されている。

【0058】

図4において、ターンプール値TP1～nに基づいてAS側（ASスイッチ#0～2）からCPUブレード#0を見たとき、ターンプール値はひとつのASスイッチ#0のポートを指し示すだけであるので、ASスイッチ#0からCPUブレード#0上の複数の論理区画を識別することはできない。これが、前記従来例の課題であり、ひとつのCPUブレード上で稼動する複数の論理区画（OS）からI/Oアクセスを行うと、ASスイッチ側はI/Oアクセス元の論理区画に返信を行うことはできなかった。

【0059】

そこで、図7のように、ASブリッジ15を仮想的なASスイッチ（以下、単に仮想スイッチとする）SWv1として機能させ、且つ、図5、図6のように実際のASスイッチ#0～2のターンプール値TP1～TPnに加えて、論理区画毎に仮想的に接続された仮想スイッチSWv1内のターンプール値TP0（仮想経路情報）を付加することで、CPUブレード#0では受信したASパケットがどの論理区画に向けて送信されたものであるかを識別できるのである。

【0060】

図7において、CPUブレード#0からI/Oブレード#5に向かう下りのASパケットは、1つの仮想スイッチSWv1と2つの物理的なASスイッチ#0、#1を通過するので、図5で示すように、ASパケットには第1、第2、第3のターンプール値TP0、1、2を設定する。なお、ターンプール値は通過するASスイッチの数（n）だけASブリッジ15が付加するもので、図5の例では、3つのターンプール値TP0～TP2を付加する。

【0061】

CPUブレード#0のASブリッジ15は、ハイパバイザ200が生成した論理区画の数に応じて仮想スイッチSWv1のポートを生成する。図7の例では、2つの論理区画#0、#1を生成したのでポート1、2がASブリッジ15により生成される。

【0062】

CPU#0、1（または各CPUコア）で稼動する複数の論理区画からI/Oアクセスがあると、ASブリッジ15は第1のターンプール値TP0にアクセス元となる論理区画からASスイッチへ向かう仮想スイッチSWv1のターンプール値TP0を設定する。例えば、図7の論理区画#0からのアクセスの場合にはTP0=0とし、論理区画#1からのアクセスの場合にはTP0=1とする。このターンプール値の設定は、上述の図4と同様である。

【0063】

そして、ASブリッジ15は、I/Oアクセス先を読み込んで、CPUブレード#0からアクセス先のI/Oブレードまでのターンプール値TP1～nを設定する。

【0064】

例えば、図4～図7では、CPUブレード#0の論理区画#0がI/Oブレード#5にアクセスする場合を示す。

【0065】

図7において、ASブリッジ15は、仮想のポート0～2を生成し、仮想ポート0をASスイッチ#0のポート7に接続し、仮想ポート1を論理区画#1に接続し、仮想ポート2を論理区画#0に接続する。

【0066】

論理区画から見て1番目のスイッチとなる仮想スイッチSWv1では、論理区画#0の隣のポート0に転送するので、第1のターンプール値である仮想スイッチターンプール値TP0に「0」を設定する。AS側では、ASスイッチ#0からASスイッチ#1のポート3からI/Oブレード#5にASパケットを転送することになる。

【0067】

10

20

30

40

50

このため、2番目のスイッチとなる物理的なASスイッチ#0ではポート7からポート2へASパケットを転送するので、ASブリッジ15は、時計回りに0からカウントした「2」を第2のターンプール値TP1に設定する。

【0068】

次に、3番目のスイッチとなる物理的なASスイッチ#1ではポート7からポート3へASパケットを転送するので、ASブリッジ15は、時計回りに0からカウントした「3」を第3のターンプール値TP1に設定する。

【0069】

こうして、ASブリッジ15は図5で示すように、下りASパケットのターンプール値TP0~TP2を「0、2、3」と設定する。さらに下りパケットであることを表すために、方向ビットDIRを1に設定する。

10

【0070】

CPUブレード#0のASブリッジ15は、図1、図5で示すように、ASスイッチ#0のポート番号7に接続されている。論理区画#0から見て2番目のスイッチとなるASスイッチ#0は、第2のターンプール値TP1を読み込んで、図4、図7の時計回りに0からターンプール値TP1の2まで数え、ポート番号2のポートを当該ASパケットの転送先に設定し、ポート番号2へASパケットを転送する。

【0071】

ASスイッチ#0のポート2は、ASスイッチ#1のポート7に接続されているので、ASパケットはASスイッチ#1に転送される。

20

【0072】

ASスイッチ#1は、論理区画#0から見て3番目のスイッチであるので、第3のターンプール値TP2を読み込んで、図4、図7の時計回りに0からターンプール値TP2の3まで数え、ポート番号3のポートを当該ASパケットの転送先に設定し、ポート番号3へASパケットを転送する。ASスイッチ#1のポート3は、I/Oブレード#5が接続されており、ASパケットはCPUブレード#0からI/Oブレード#5に送られる。こうして、下りのASパケットは、各ASスイッチがターンプール値を順次読み込んで時計回りにカウントしたポートに順次転送する。

【0073】

逆に、I/Oブレード#5からCPUブレード#0へのASパケットは、図6で示す上りのASパケットのように、下りのASパケットのターンプール値TP0~nを、方向ビットDIRを反転(0から1)して設定する。方向ビットDIRが1の上りパケットのターンプールは右から左に解釈され、さらに各ASスイッチでターンプール値は反時計回りに数えられる。

30

【0074】

つまり、図4、図7の場合の上りASパケットのターンプール値は、図6のように、第1のターンプール値(方向ビットDIRが1の上りパケットの場合一番右の)TP2が「3」、第2のターンプール値TP1が「2」、第3の値には仮想スイッチターンプール値TP0が「0」に設定される。なお、上りASパケットのターンプール値TP0~nの設定は各I/OブレードのLEAFブリッジが行う。

40

【0075】

そして、I/Oブレード#5が送信したASパケットは、I/Oブレード#5から見て1番目のASスイッチとなるASスイッチ#1が、第1のターンプール値TP2=3を読み込んで、図4、図7の反時計回りに0からターンプール値TP1の3まで数え、ポート番号7のポートを当該ASパケットの転送先に設定し、ポート番号7へASパケットを転送する。

【0076】

ASスイッチ#1のポート7は、ASスイッチ#0のポート2に接続されているので、ASパケットはASスイッチ#0に転送される。

【0077】

50

A Sスイッチ# 0は、I/Oブレード# 5から見て2番目のA Sスイッチであるので、第2のターンプール値T P 1を読み込んで、図4の反時計回りに0からターンプール値T P 1の2まで数え、ポート番号7のポートを当該A Sパケットの転送先に設定し、ポート番号7へA Sパケットを転送する。

【0078】

図7で示すように、A Sスイッチ# 0のポート7は仮想スイッチS W v 1に接続されており、仮想スイッチS W v 1は、I/Oブレード# 5から見て3番目のA Sスイッチであるので、第3のターンプール値T P 0を読み込んで、図4の反時計回りに0からターンプール値T P 0の0まで数え、ポート番号2を当該A Sパケットの転送先に設定し、ポート2と仮想的に接続された論理区画# 0へA Sパケットを転送する。

10

【0079】

こうして、上りのA Sパケットは、各A Sスイッチがターンプール値を反時計回りにカウントしたポートに順次転送することで、最後に仮想スイッチターンプール値T P 0とA Sブリッジ15の仮想スイッチS W v 1を用いることで、A Sスイッチ側からI/Oアクセス元の論理区画# 0を識別してA Sパケットを確実に返信することができるのである。

【0080】

なお、A Sスイッチを管理するファブリック管理サーバ4は、A Sスイッチ# 0～# 2に加えて、各C P UブレードのA Sブリッジ15が提供する仮想スイッチS W v 1も、A Sスイッチ群のトポロジとして管理する。

【0081】

20

以上は宛先が単一のユニキャストパケットであるが、以下ではマルチキャストの場合について述べる。マルチキャストでは、各パケットが、マルチキャストの宛先を一意に表すM G I D (マルチキャストグループI D)を持つ。各スイッチは、あらかじめ決められているマルチキャストテーブル180を引き、M G I Dに対応する出力ポート全てにパケットを同時に出力する。

【0082】

本発明では、各C P Uブレードに置かれた、A Sブリッジ15にも仮想スイッチS W v 1のマルチキャストテーブルを新たに置き、同一C P Uブレード上の複数の論理区画へのマルチキャストをサポートする。各ブリッジは、マルチキャストパケットを受け取ると、パケットのM G I Dに対応する仮想スイッチのマルチキャストテーブルエントリを読み、どの論理区画にマルチキャストをする必要があるかを定める。その後、該当する論理区画の持つ主記憶空間へのデータの書き込み、該当する論理区画へのイベント配信が行われる。

30

【0083】

なお、A Sスイッチを管理するファブリック管理サーバ4は、A Sスイッチ# 0～# 2に加えて、各C P UブレードのA Sブリッジ15が提供する仮想スイッチS W v 1のマルチキャストテーブルも、A Sスイッチ群のマルチキャストテーブルとして管理する。

【0084】

< A Sブリッジの詳細 >

次に、図3を参照しながら、仮想スイッチとして機能するA Sブリッジ15の詳細について説明する。なお、図3はC P Uブレード# 0のA Sブリッジ15を主体にした機能ブロック図であり、他のC P Uブレードも同様に構成される。また、図3においてA Sブリッジ15とC P U# 0、# 1または主記憶13との間で行われる通信は、図1で示したようにノースブリッジ12を介して行うものである。

40

【0085】

A Sブリッジ15には、図1で示した制御部16とメモリ17を備え、A Sスイッチ# 0～2側から見て複数の論理区画を複数のノードとして見せる仮想A Sスイッチとして機能する仮想スイッチを提供する。このため、A Sブリッジ15のメモリ17には、ハイバイザ200上で稼動する論理区画と仮想スイッチの仮想ポート番号の対応関係を設定するパーティション管理表155と、各論理区画の主記憶13上のアドレスを管理するアド

50

レス管理表 160 と、論理区画からの I/O アクセスコマンドを受け付ける I/O アクセスコマンドレジスタ 153 と、I/O ブレード側からのイベント（割り込みやエラー通知など）を受け付けるイベントレジスタ 156 と、I/O ブレードからのイベント転送があった場合に、転送先の論理区画を決定するための宛先 LPAR 番号レジスタ 158 が BMC 18 によって生成される。宛先 LPAR 番号レジスタ 158 は、イベントのマルチキャスト転送をサポートするために、ビットマップで実装される。マルチキャストの場合には、複数のビットが 1 となる。さらに、マルチキャストの場合に、MGID と宛先論理区画の関係をあらわす、仮想スイッチマルチキャストテーブル 170 がファブリック管理サーバにより作成される（この処理は新しいマルチキャストグループが定義されるたびに行われる）。仮想スイッチマルチキャストテーブル 170 は、図 19 で示すように構成され、マルチキャストグループ ID（MGID）毎に、各論理区画（図中出力 LPAR）に対応したビットマップで構成される。この仮想スイッチマルチキャストテーブル 170 は、各論理区画にマルチキャスト転送を行うか否かを示すビットが設定される。このビットは 1 のときにマルチキャスト転送を許可し、0 のときにマルチキャスト転送を禁止する。図 19 において、MGID = 1 のときには論理区画 # 1 へのマルチキャスト転送を許可し、論理区画 # 0 へのマルチキャスト転送を禁止する。

10

## 【0086】

さらに、ファブリック管理サーバ 4 が各 CPU ブレードにおかれたハイパバイザと In-bound で（AS のファブリックを介して）通信する場合は、宛先 LPAR 番号レジスタ 158 には、各論理区画に対応するビットに加えて、ハイパバイザに対応するビットを持つ。この場合、ハイパバイザ宛にイベントを送る場合は、ハイパバイザに該当するビットを ON にする。

20

## 【0087】

上記パーティション管理表 155、アドレス管理表 160、仮想スイッチマルチキャストテーブル 170 及び各レジスタを除く AS ブリッジ 15 の構成要素は、上記図 1 の制御部 16 に相当する。

## 【0088】

まず、パーティション管理表 155 は、図 8 で示すように、論理区画（パーティション）の番号に割り当てられた仮想スイッチの仮想ポート番号と、この論理区画が有効か無効かの状態を示すフィールドから構成されている。このパーティション管理表 155 は、ハイパバイザ 200 が論理区画を生成、削除するときに更新される。

30

## 【0089】

さらに、In-bound 管理を行う場合には、ハイパバイザに対応する仮想スイッチの情報も示される。

## 【0090】

次に、アドレス管理表 160 は、図 9 で示すように、主記憶 13 上の各パーティションの位置を示すため、論理アドレス空間上の論理区画のベースアドレスと、各論理区画のサイズを論理区画番号毎に設定したものである。このアドレス管理表 160 も、論理区画の生成や削除の際にハイパバイザ 200 によって更新される。なお、主記憶 13 の論理アドレス空間は、例えば、図 10 で示すように構成され、各論理区画（LPAR）のベースアドレスとサイズは、ハイパバイザ 200 によって管理される。なお、各論理区画のアドレスマップ上には、ハイパバイザ 200 によって I/O レジスタがそれぞれマッピングされる。

40

## 【0091】

I/O アクセスコマンドレジスタ 153 は、各論理区画上で稼動する OS からの I/O アクセスコマンドを受け付ける。I/O アクセスコマンドは、例えば I/O レジスタの PIO 要求（I/O デバイスに持つコントロール/ステータスレジスタへの読み書き）もしくは DMA 転送（READ/WRITE）要求や、ターゲットとなる I/O ブレードの情報が含まれる。

## 【0092】

50

イベントレジスタ156には、ディスクアクセスの終了や、イーサネット（登録商標）インターフェース等のネットワークからの通知（データ/コマンド）に起因する、I/Oブレードからの割り込み、もしくはエラー通知を受け付ける。そして、ハイパバイザ200を經由して各論理区画に転送される。イベントはユニキャスト/マルチキャストの場合がある。ユニキャストの場合は、仮想スイッチのターンプール情報より求められた宛先論理区画の番号は、宛先L P A R番号レジスタ158に設定される。マルチキャストの場合は、パケットのM G I Dにより、仮想スイッチのマルチキャストテーブル170を引き、宛先の論理区画（L P A R）のビットマップが求められ、宛先L P A R番号レジスタ158に設定される。図11に示す宛先L P A R番号レジスタ158に設定された論理区画（L P A R）毎のビットマップに基づいてハイパバイザ200が転送を行う。仮想スイッチマルチキャストテーブル170や、図11の宛先L P A R番号レジスタ158は、各論理区画毎にビットを有し、論理区画毎に対応するビットが1であればマルチキャストを行い、0であればマルチキャストを禁止することを示す。この仮想スイッチのマルチキャストテーブル170は、ファブリック管理サーバ4によって管理される。

10

**【0093】**

次に、制御部16を構成する各部について説明する。まず、論理区画からI/Oアクセスがあると、P C I eのパケットがA Sブリッジ15に送信される。A Sブリッジ15は、受信したP C I eのパケットからI/Oアクセスコマンド及びターゲットのI/Oブレードの情報を表すI/Oデバイスのアドレスを、I/Oアクセスコマンドレジスタ153に格納する。また、P C I eのパケットに含まれるI/OレジスタアドレスやD M Aを行う主記憶のアドレスを抽出してアドレスデコード部151に入力する。これにより、I/Oアクセスを発行した論理区画を特定することができる。

20

**【0094】**

アドレスデコード部151は、アドレス管理表160を参照して、抽出されたアドレスがどの論理区画のものあるかを判定する。つまり、抽出されたアドレスがアドレス管理表160のベースアドレス+サイズの範囲内となる論理区画を、I/Oアクセスを行う論理区画として判定し、この論理区画番号を仮想スイッチターンプール情報付加部152へ送る。

**【0095】**

また、I n - b o u n d管理を行う場合には、アドレス管理表160には、ハイパバイザのアドレス範囲の情報をもち、ハイパバイザに該当する仮想スイッチターンプール情報を付加する。

30

**【0096】**

なお、アドレスデコード部151で行う論理区画の判定は、I/OレジスタのアドレスまたはD M Aアドレスのアドレスとアドレス管理表160に設定されたアドレスのうち、予め設定した上位ビット（例えば、上位8ビット）について比較を行えばよい。

**【0097】**

仮想スイッチターンプール情報付加部152は、受信した論理区画番号からパーティション管理表155を参照し、論理区画番号に対応する仮想スイッチの仮想ポート番号を決定し、これを上記図5に示した仮想スイッチターンプール値T P 0としてパケット組み立て部154に送る。

40

**【0098】**

パケット組み立て部154は、I/Oアクセスコマンドレジスタ153からターゲットのI/Oブレードの情報を読み込んで、図5に示したターンプール値T P 1 ~ T P nを経路情報として決定する。A Sブリッジ15からターゲットのI/Oブレードまでの経路情報は、パケット組み立て部154が予めファブリック管理サーバ4から取得した経路情報に基づいて決定される。つまり、A Sブリッジ15は、I/Oバス14から受信したI/Oアクセス要求に含まれるアクセス先（ターゲット）のI/Oブレード（デバイス）の識別子に基づいて、ファブリック管理サーバ4からあらかじめ取得していた経路情報に基づき、C P UブレードからI/Oブレードまでの経路情報を決定する。

50

## 【0099】

パケット組み立て部154は、ASブリッジ15から目的のI/Oブレードまでのターンプール値TP1~TPnに、仮想スイッチターンプール値TP0を読み込んで、図5で示すようにターンプール値TP1~TPnの先頭に付加し、ASパケットのヘッダー部を生成する。

## 【0100】

そして、パケット組み立て部154は、受信したPCIeのパケットの先頭に、仮想スイッチターンプール値TP0、ターンプール値TP1~TPn、DIRビット0を、ASパケットのヘッダー部としてASパケットを生成してから、ASスイッチ#0~2へ送信する。

10

## 【0101】

ASブリッジ15は、上記アドレスデコード部151、仮想スイッチターンプール情報付加部152、I/Oアクセスコマンドレジスタ153、パケット組み立て部154により下りASパケットを生成し、ASブリッジ15に接続されたASスイッチ#0にASパケットを送る。

## 【0102】

次に、ASブリッジ15がASスイッチ#0からASパケットを受信した場合について説明する。

## 【0103】

ASブリッジ15は、ASスイッチ#0から上りASパケットを受信すると、コマンド解析部157で上りASパケットを受け付けて、ASパケット内のPCIeパケットに含まれるコマンドを解析する。コマンド解析部157は、コマンドの種類がDMA転送、割り込みやエラー通知等のイベント処理の何れであるかを判定する。さらにマルチキャストが必要かどうか判断する。

20

## 【0104】

そして、コマンド解析部157は、受信したASパケットのコマンドがDMA転送またはマルチキャストDMAであれば、ASパケットから抽出したPCIeパケットをDMA処理部162へ送る。受信したASパケットが割り込みやエラー通知等のイベント処理もしくはマルチキャストイベントの場合には、ASパケットから抽出したPCIeパケットをイベントレジスタ156に送る。

30

## 【0105】

宛先LPAR番号抽出部159は、ユニキャストとマルチキャストの場合で異なる動作をする。

## 【0106】

ユニキャストの場合、上りASパケットのヘッダー部から図6で示したように、最後のターンプール値に格納された仮想スイッチターンプール値TP0を抽出する。そして、宛先LPAR番号抽出部159は、パーティション管理表155から仮想スイッチターンプール値TP0に対応する論理区画番号を取得し、受信したASパケットの宛先となる論理区画を特定する。宛先LPAR番号抽出部159は、特定した論理区画の番号を表すビットマップ(この場合は何れか1ビットのみがセットされている)をアドレス管理表160と宛先LPAR番号レジスタ158に送る。

40

## 【0107】

マルチキャストの場合、宛先LPAR番号抽出部159は、ASパケットのヘッダー部からMGIDを求め、仮想スイッチのマルチキャストテーブル170を引き、受信したパケットの宛先となる論理区画の番号(複数)を表すビットマップを、アドレス管理表160と宛先LPAR番号レジスタ158に送る。ここで、上流のスイッチでのみマルチキャストが行われ、当該CPUでは宛先のLPARが一つのみの場合があるため、仮想スイッチのマルチキャストテーブルのビットパターンで1ビットのみがセットされている場合もある。

## 【0108】

50

DMA処理部162は、コマンドがDMA転送の場合には、宛先LPA R番号抽出部159が決定した論理区画番号からアドレス管理表160に基づいてベースアドレスを求め、PCIeパケット内のDMAアドレスから実際にアクセスする主記憶13上のアドレスをアドレス変換部161で変換する。そして、ASブリッジ15は、アドレス変換部161が決定した主記憶13上のアドレスに対してDMA転送(READまたはWRITE)を実施する。

【0109】

このとき、アドレス変換部161では、DMA転送を行う主記憶13上のアドレスが、転送を行う論理区画のベースアドレスとサイズから求められるアドレス範囲内であるかを判定する。そして、判定結果が該当する論理区画内であればDMA転送を行い、DMA転送を行う主記憶13上のアドレスが該当する論理区画を超えていれば、DMA転送を中止し、異なる論理区画に操作が行われるのを防いで仮想計算機の信頼性を確保する。

10

【0110】

In-bound管理を行う場合に、ファブリック管理サーバ4がハイパバイザに対して出したDMAについては、上記の論理区画のアドレス範囲チェックは行わない。

【0111】

次に、ASパケットのコマンドがマルチキャストDMAの場合は、DMA処理部162から複数回書込処理部163へPCIeパケットが転送される。複数回書込処理部163は、図11に示した宛先LPA R番号レジスタ158を参照し、論理区画番号に対応するビットが1となっている論理区画について、アドレス変換部161に転送を行う論理区画のI/Oレジスタの主記憶13上のアドレスを問い合わせ、複数の論理区画に対して複数回の書込を実施する。このマルチキャストの場合も、アドレス変換部161では書込の対象となる論理区画のアドレスが正当であるか否かを上記DMA転送と同様に判定し、不正な書込が行われるのを防止する。

20

【0112】

次に、ASパケットのコマンドが割り込みやエラー通知等のイベント処理(マルチキャスト含む)の場合は、コマンド解析部157は上りASパケットから抽出したPCIeパケットをイベントレジスタ156に送り、イベントレジスタ156からハイパバイザ200へ通知を行う。ハイパバイザは宛先LPA R番号レジスタ158を参照し、ビットマップが1である論理区画にイベントを転送する。

30

【0113】

以上のように、ASブリッジ15は、I/Oバス14からI/Oアクセス要求を受け取ると、目的のI/Oブレードの識別子に基づいてターンプール値TP0~nと、I/Oアクセスを要求した論理区画に基づく仮想スイッチターンプール値TP0を生成する。そして、ASブリッジ15は、I/Oバス14から受信したI/Oアクセスのパケットに仮想スイッチターンプール値TP0とターンプール値TP0~nを付加して下りASパケットを生成することで、下りASパケットの仮想スイッチとして機能する。

【0114】

ASブリッジ15は、接続されたASスイッチ#0から上りASパケットを受信すると、ターンプール値の最後に付加された仮想スイッチターンプール値TP0を抽出して、この上りASパケットの宛先となる論理区画を特定し、ASパケットのヘッダー部を除去したものをI/Oバス14に送信することで、上りASパケットの仮想スイッチとして機能する。

40

【0115】

<ASブリッジの初期化・更新>

次に、ASブリッジ15は、CPUブレードの起動時や、ハイパバイザ200がパーティションを変更したとき、I/Oブレードの変更があったときなどに、BMC7やハイパバイザ200またはファブリック管理サーバ4によりメモリ17上の表やレジスタの初期化や更新が行われる。

【0116】

50



A Sブリッジ15が配置されたCPUブレードが起動する際の初期化について図12を参照しながら説明する。図12は、CPUブレードの起動時にBMC7で行われる初期化処理のフローチャートである。

【0117】

BMC7はCPUブレードが起動すると、図1に示したコンソール70でCPUブレードに生成可能な論理区(パーティション)の数を設定する(S1)。なお、CPUブレードで生成可能な論理区画の最大数の設定は、コンソール70からの入力に加え、予め設定したファイルをBMC7が読み込むことで行っても良い。

【0118】

次に、BMC7はASブリッジ15のメモリ17にアクセスを行い、メモリ17の所定のアドレスにパーティション管理表155とアドレス管理表160を、上記S1で設定した論理区画の最大数に応じて作成する(S2)。

【0119】

上記S1、S2の処理が終了した時点では、ASブリッジ15の各表がメモリ17に生成されただけであり、内容については未設定である。すなわち、上記図12の処理では、CPUブレードに設定可能な論理区画の最大数に応じた各表の大きさがASブリッジ15のメモリ17に確保される。

【0120】

次に、ハイパバイザ200が論理区画の作成または変更(削除を含む)を行う際のASブリッジ15に対する処理について、図13を参照しながら説明する。図13は、ハイパバイザ200が論理区画の設定を行う際に行う処理の一例を示すフローチャートである。

【0121】

まず、S11ではハイパバイザ200が作成(または変更)する論理区画の主記憶13上のベースアドレスとサイズ及び論理区画番号を取得する。

【0122】

次に、S12では変更を行う論理区画番号について、図9に示したアドレス管理表160に、ハイパバイザ200がベースアドレスとサイズを書き込む。

【0123】

S13では、ハイパバイザ200が、アドレス管理表160に書き込みを行った論理区画番号について、図8のパーティション管理表155のパーティション番号を検索し、該当する論理区画番号に対応するパーティション管理表155のエントリ(仮想スイッチポート番号)について、有効/無効の欄を更新する。つまり、論理区画を作成するときにはパーティション管理表155の該当エントリを有効に設定し、論理区画を削除するときには、パーティション管理表155の該当エントリを無効に設定する。

【0124】

なお、アドレス管理表160に書き込みを行った論理区画番号が図8のパーティション管理表155のパーティション番号にない場合には、パーティション番号が空となっているパーティション管理表155の仮想ポート番号に論理区画番号を書き込み、有効/無効の欄を更新する。

【0125】

S14では、ハイパバイザ200が、パーティション管理表155に書き込みを行った仮想スイッチポート番号を、ファブリック管理サーバ4に通知する。ファブリック管理サーバ4では、ASスイッチのトポロジに、ASブリッジ15が提供する仮想スイッチSWv1の情報を取得し、作成または更新された仮想スイッチ上の論理区画を管理することができる。

【0126】

さらに、新しいマルチキャストグループが定義された場合、既存のマルチキャストグループが削除された場合には、ファブリック管理サーバ4は、該当するマルチキャストグループがどのLPARにパケットを送出する必要があるかを調べ、ASブリッジ15の持つ仮想スイッチマルチキャストテーブル170の、パケットを出力する必要がある論理区画

10

20

30

40

50

に該当するビットに 1 を、パケットを出力する必要の無い論理区画に該当するビットには 0 を書き込む。

【 0 1 2 7 】

仮想スイッチでは、通常の A S スイッチのように、入力ポートにマルチキャスト結果を戻す処理（リフレクション）は行われなため、仮想スイッチマルチキャストテーブルには、ブリッジの入力を表すビットは設けられない。

【 0 1 2 8 】

< A S ブリッジのデータ書き込み、イベント通知処理 >

論理区画が要求した I / O アクセスに対する応答は、図 1 4 で示すように、A S スイッチの機能により、要求を送信した I / O ブレードに対して上り A S パケット（返答パケット）が送られ、さらに、A S ブリッジ 1 5 が提供する仮想スイッチにより C P U ブレード内の複数の論理区画に対して I / O アクセス要求に対する上り A S パケットが送られる。

【 0 1 2 9 】

ここで、A S スイッチ側から A S ブリッジ 1 5 を見た場合、ファブリック管理サーバ 4 では A S ブリッジ 1 5 を仮想スイッチ S W v 1 として識別するので、各 C P U ブレードの論理区画（または物理的な区画）は、仮想スイッチ S W v 1 の下位のノードとして見ることができる。

【 0 1 3 0 】

上り A S パケットを受信した A S ブリッジ 1 5 は、パーティション管理表 1 5 5 から宛先の論理区画が有効であることを確認した後、図 3 で示すように、A S パケットの内容（コマンド）に応じて各論理区画のメモリ空間への D M A による読み書きや、ハイパバイザ 2 0 0 を介した各論理区画に対するイベントの送付を行う。なお、パーティション管理表 1 5 5 の該当する論理区画が無効であれば、A S パケットを破棄する。

【 0 1 3 1 】

上りパケットは I / O 発の D M A 要求、イベント処理要求がある。I / O 発の D M A では、指定された論理区画の主記憶にアクセスが行われる。イベント処理では、指定された論理区画に対して割り込み等のイベントが送られる。これに対してマルチキャスト（マルチキャスト D M A、マルチキャストイベント）の場合は、該当する I / O ブレードから複数の論理区画に対して同一の D M A またはイベントを送ることができる。すなわち、マルチキャストの場合には、仮想スイッチマルチキャストテーブル 1 7 0 に予め設定された複数の論理区画に同一 D M A データが書き込まれ、同一イベントが送られる。このマルチキャストによる書き込みは、図 3 のように、マルチキャスト D M A の場合には、A S ブリッジ 1 5 が主記憶 1 3 上の各論理区画に対して複数回の書き込みを発行し、マルチキャストイベントの場合には、ハイパバイザ 2 0 0 に対して通知が行われ、ハイパバイザ 2 0 0 が宛先となる複数の論理区画にイベントを通知する。

【 0 1 3 2 】

図 1 5 は、上り A S パケットがマルチキャスト D M A のときに A S ブリッジ 1 5 が実行する処理の一例を示すフローチャートである。なお、この処理は図 3 における D M A 処理部 1 6 2 と複数回書込部 1 6 3 の機能に相当する。

【 0 1 3 3 】

まず、S 2 1 では、上り A S パケットから D M A の書き込み先アドレス A と書き込みデータ D を抽出する。S 2 2 では、宛先 L P A R 番号レジスタ 1 5 8 から D M A の宛先となる論理区画（パーティション番号）の番号を表すビットマップを取得するとともに、パーティション管理表 1 5 5 の論理区画のうち、有効となっている論理区画番号を取得する。そして、S 2 3 では A S ブリッジ 1 5 は、宛先 L P A R 番号レジスタ 1 5 8 のビットマップが 1 であり（かつ有効な）論理区画の番号を一つ取得し、マルチキャスト D M A 書き込みを行う論理区画を決定する。この書き込み時には、上り A S パケットの D M A アドレスに各論理区画のベースアドレスを加算し、この加算結果が各論理区画に割り当てられた主記憶 1 3 上の領域（ベースアドレス～ベースアドレス+サイズの範囲）を超えていないことを確認してから D M A によって各論理区画のアドレス A 毎にデータ D を書き込む。

10

20

30

40

50

## 【 0 1 3 4 】

S 2 4では、宛先 L P A R 番号レジスタ 1 5 8 のビットが 1 となっている全ての論理区画に対してマルチキャストの書き込みが完了したか否かを判定する。完了していれば処理を終了し、未了であれば、S 2 3 に戻って残りの論理区画へマルチキャスト D M A の書き込みを行う。

## 【 0 1 3 5 】

以上のように、宛先 L P A R 番号レジスタ 1 5 8 のビットマップ、及び、パーティション管理表 1 5 5 に設定された論理区画毎の有効 / 無効の欄から、マルチキャストの宛先となっている論理区画の中から現在使用している論理区画を判定し、さらに、A S ブリッジ 1 5 はマルチキャスト D M A を実施する論理区画を判定することができる。これにより、A S ブリッジ 1 5 は、マルチキャスト D M A が必要な論理区画に対してのみ同一のデータ D を複数の論理区画に対して書き込むことが可能となる。また、L P A R の有効 / 無効をチェックすることにより、L P A R が O N / O F F を繰り返す場合に、マルチキャストテーブルを変更する手間を削減することができる（特に全ビットが 1 のブロードキャストアドレスの場合に有効である）。

10

## 【 0 1 3 6 】

D M A をアドレス範囲が複数の場合は、あるアドレスのデータを対象とする全ての論理区画に書いてから、次のアドレスのデータを書き込む。これにより、パケット全体を一旦記憶する必要が無く、バッファ領域に必要なメモリを削減できる他、全体が到着する前に書き込みを開始できるため、処理を高速化できる。

20

## 【 0 1 3 7 】

なお、マルチキャストイベントの上り A S パケットの場合は、A S ブリッジ 1 5 はイベントをハイパバイザ 2 0 0 に伝達する。マルチキャストイベントの通知を受けたハイパバイザ 2 0 0 は、宛先 L P A R 番号レジスタ 1 5 8 から対象となる複数の論理区画を求め、イベントを複数の論理区画に通知する。

## 【 0 1 3 8 】

上記は、各論理区画宛の D M A、イベント送付の手順について述べた。それに加え、I n - b o u n d 管理を行う際には、ハイパバイザに宛てた D M A、イベントもサポートする。その場合には、宛先 L P A R 番号レジスタ 1 5 8、パーティション管理表 1 5 5、アドレス管理表 1 6 0 はハイパバイザに該当するエントリを持つことにより、仮想スイッチがハイパバイザに対応するポート、ターンブル情報をサポートすることを可能にする。

30

## 【 0 1 3 9 】

< I / O コンフィグレーション >

次に、ファブリック管理サーバ 4 で行われる I / O コンフィグレーションについて、図 1 6 を参照しながら説明する。

## 【 0 1 4 0 】

A S スイッチを管理するファブリック管理サーバ 4 に接続されたコンソール 5 には、図 1 6 に示すようなブレードサーバ管理画面が表示される。ファブリック管理サーバ 4 は、A S スイッチに接続された I / O ブレードと、C P U ブレード上の論理区画の関係を管理する。すなわち、ファブリック管理サーバ 4 では、C P U ブレード上の論理区画の状態を把握して、論理区画に割り当てる I / O ブレードを設定することができる。

40

## 【 0 1 4 1 】

図 1 6 において、C P U ブレード上 # 0 ~ 2 が示されており、各 C P U ブレード上には、図 8 のパーティション管理表 1 5 5 に設定されている論理区画の数に応じて複数のボックス 4 1 が表示され、各ボックス 4 1 の中には論理区画の状態を示す「State」欄が表示される。この「State」欄は、パーティション管理表 1 5 5 で有効となっている論理区画が「ON」として表示され、無効となっている論理区画が「off」として表示される。

## 【 0 1 4 2 】

図 1 6 では、C P U ブレード # 0 には論理区画はなく物理的な区画が有効となっている状態を示し、C P U ブレード # 1 には、2 つの論理区画が生成されて一方が有効、他方が

50

初期化中となっている状態を示す。CPUブレード#2には、3つの論理区画が生成され、2つの論理区画が有効で、ひとつのパーティションが無効となっている状態を示している。

【0143】

これらCPUブレードは、3つのI/Oブレード#0~2が接続されたASスイッチに接続される。図16においては、I/Oブレード#0がHBA1という名称でコンソール5に表示され、I/Oブレード#1がHBA2、I/Oブレード#3がEther0という名称でそれぞれ表示される。

【0144】

各I/Oブレードには、I/O共有の状態を示す「State」欄が表示され、複数の論理区画または物理区画に接続されているI/Oブレードには「共有」と表示され、単一の論理区画または物理区画に接続されているI/Oブレードには「占有」と表示される。

10

【0145】

そして、ASスイッチはひとつのボックスで表示され、I/Oブレードと論理区画（または物理区画）の接続状態を示す破線30が表示される。破線30は、現在割り当てられている論理区画（物理区画）とI/Oブレードを結ぶ線である。図16において、I/Oブレード#0（HBA1）は、物理区画のみのCPUブレード#0と、CPUブレード#1の論理区画0で共有されている状態を示している。

【0146】

論理区画（または物理区画）に対するI/Oブレードの割り当てを変更するには、論理区画のボックス41または物理区画をカーソルでクリックする。このクリックにより、コンソール5の該当する論理区画（または物理区画）には階層型のメニュー40が表示される。

20

【0147】

階層型メニュー40には、I/Oブレードの設定の種類（割り当て（Attach）または割り当て解除（Detach））と、割り当て（解除）可能なI/Oブレードの名称が表示される。例えば、図示のように、現在初期化中のCPUブレード#1の論理区画1に新たなI/Oブレードを割り当てた場合は、「Attach」をクリックすると、割り当て可能なI/Oブレードの名称が表示されるので、この名称のいずれかひとつをクリックすると、このパーティション1に新たなI/Oブレードが割り当てられ、選択しI/Oブレードとこの論理区画1の間に破線30が表示される。逆に割り当てを解除した場合には、破線30が消去される。

30

【0148】

以上のように、ファブリック管理サーバ4はASスイッチを介して接続されたI/Oブレードの論理区画（または物理区画）への割り当てを管理することができる。なお、ファブリック管理サーバ4には、図示しないCPU、記憶装置が設けられ、上記ASスイッチを介して接続されたI/Oブレードの論理区画（または物理区画）への割り当てを管理するテーブルが、記憶装置に格納される。このテーブルは、例えば、I/Oブレード毎に、割り当てられているCPUブレードの番号と、論理区画番号を格納するフィールドを備えていけばよい。

40

【0149】

また、ファブリック管理サーバ4は、論理区画へのI/Oブレードの割り当てを変更すると、ハイバイザ200に変更した内容を通知する。この際、In-boundで管理を行う場合には、ファブリック管理サーバ4は、上記I/Oブレードからのイベントや割り込みと同様にして、ハイバイザ200にI/Oコンフィグレーションの変更を、ASファブリック経由で通知すればよい。Out-boundな管理を行う場合には、I/Oコンフィグレーションの変更は、管理用ネットワーク（図示しない）により転送される。

【0150】

<まとめ>

50

以上のように、本発明によればA SスイッチなどのI / Oスイッチを用いてI / Oブレード(デバイス)を複数の論理区画間で共有する際に、C P Uブレードに備えたI / Oブリッジを仮想スイッチとして機能させ、各論理区画毎に仮想経路情報(仮想スイッチテーブル情報)を設定することで、ひとつのC P Uブレードで複数の論理区画を提供する仮想計算機システムにおいてI / Oスイッチを用いたI / O共有を実現することができる。

**【0151】**

これにより、従来では仮想計算機間でI / Oの共有を行うためには、論理区画を識別するソフトウェアが必要になって、I / Oアクセスのオーバーヘッドが過大になっていたが、本発明のように、汎用のI / OバスであるP C I - E X P R E S Sと、P C I - E X P R E S SをスイッチングするA Sを用いてハードウェアレベルで仮想計算機(論理区画)間のI / O共有を行うようにすることで、論理区画を識別するためのソフトウェアは不要となってI / Oアクセスの高速化を図ることができる。

10

**【0152】**

また、従来では仮想計算機間でI / Oの共有を行うために仮想計算機システムに固有の特別なI / O装置を必要としていたため、システムの価格が高価になっていたが、本発明のように、汎用のI / OバスであるP C I - E X P R E S Sと、P C I - E X P R E S SをスイッチングするA Sを用いて仮想計算機間のI / O共有を行うようにすることで、システムの価格を大幅に低減できるのである。

20

**【0153】**

特に、複数サーバを1台にまとめるサーバコンソリデーション(サーバ統合)を実現する際には、I / Oブレードを複数の論理区画で共有できるため、従来ではC P Uブレード毎に設けていたN I C (Network Interface Card) やF C A (Fibre Channel Adaptor) 等のI / Oデバイスを、必要な数のI / Oブレードとして導入すればよいので、I / Oデバイスを大幅に低減することで計算機システムの導入コストを低減できる。また、C P UブレードにはI / Oスロットが不要となるので、C P Uブレードの製造コストも低減でき、サーバコンソリデーションを効果的に実現できるのである。

**【0154】**

なお、上記実施形態では、1つのスイッチブレード3に複数のA Sスイッチを備えた例を示したが、1つのスイッチブレード3に1つのA Sスイッチを備え、複数のスイッチブレード3でA Sスイッチ群を構成しても良い。

30

**【0155】**

また、上記実施形態では、アドレスデコード部151で抽出した論理区画のアドレス情報のうち上位ビットで論理区画を探索する例を示したが、論理区画のI Dを明示的に設定するレジスタをA Sブリッジ15のメモリ17に設定するようにしてもよい。あるいは、このレジスタを論理区画の数に応じて複数持つようにしてもよい。

**【産業上の利用可能性】****【0156】**

以上のように、本発明はI / Oデバイスと複数のC P Uの間でスイッチングを行うI / Oスイッチを仮想計算機システムに適用することができる。

40

**【図面の簡単な説明】****【0157】**

【図1】本発明を適用する計算機システムのブロック図。

【図2】本発明を適用する仮想計算機システムのソフトウェアの機能ブロック図。

【図3】A Sブリッジを中心とする機能ブロック図。

【図4】A Sスイッチの機能を説明する説明図。

【図5】下りA Sパケットの一例を示す説明図。

【図6】上りA Sパケットの一例を示す説明図。

【図7】A Sブリッジの仮想スイッチ機能を示す説明図。

【図8】パーティション管理表の一例を示す説明図。

50

【図9】アドレス管理表の一例を示す説明図。

【図10】CPUブレードの主記憶の内容を示すアドレスマップ。

【図11】宛先LPAR番号レジスタの一例を示す説明図。

【図12】CPUブレードの初期化時のBMCの制御の一例を示すフローチャート。

【図13】論理区画作成または更新時のハイバイザの制御の一例を示すフローチャート

【図14】下りASパケットと上りASパケットの様子を示す説明図。

【図15】マルチキャスト時のASブリッジの制御の一例を示すフローチャート。

【図16】ファブリック管理サーバの管理画面の一例を示す説明図。

【図17】マルチキャストパケットの一例を示す説明図。

【図18】マルチキャストを行うスイッチの一例を示す説明図。

【図19】仮想スイッチマルチキャストテーブルの一例を示す説明図。

【符号の説明】

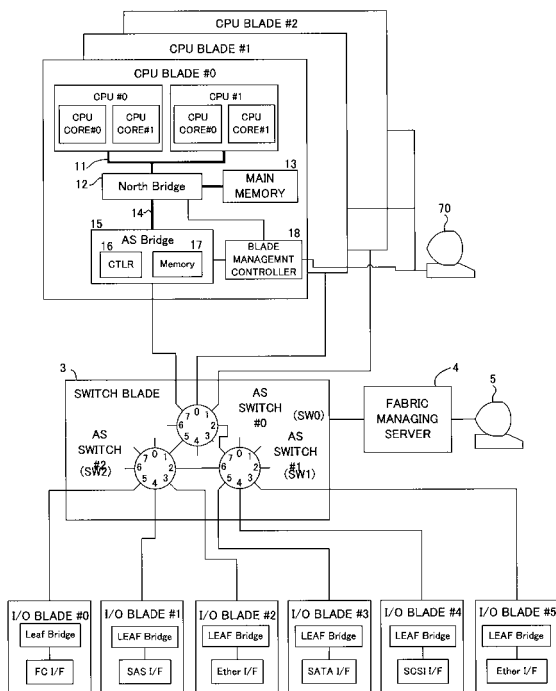
【0158】

- # 0、# 1 CPUブレード
- 3 スイッチブレード
- 4 ファブリック管理サーバ
- 5 コンソール
- 13 主記憶
- 15 ASブリッジ
- 18 BMC
- # 0 ~ # 5 I/Oブレード

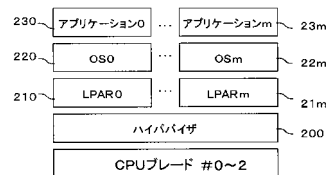
10

20

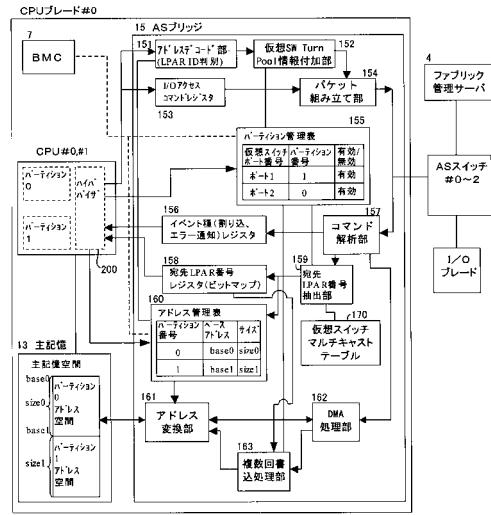
【図1】



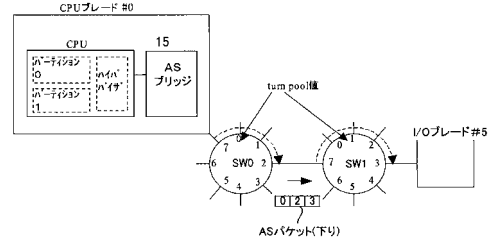
【図2】



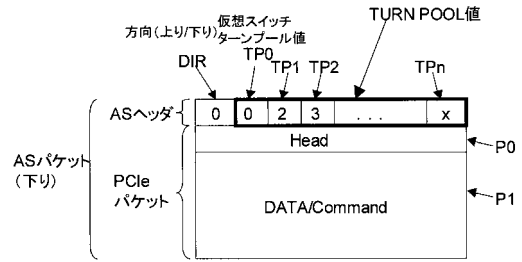
【図3】



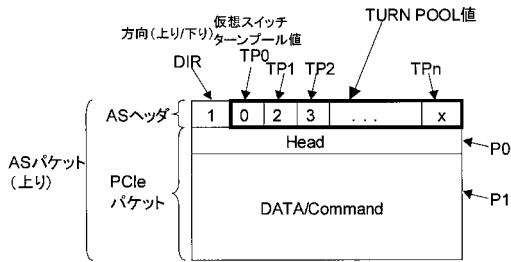
【図4】



【図5】



【図6】



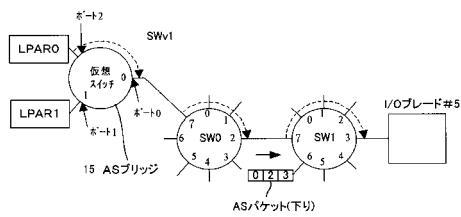
【図8】

仮想SWポートNo	パーティションNo	有効/無効
1	1	有効
2	0	有効

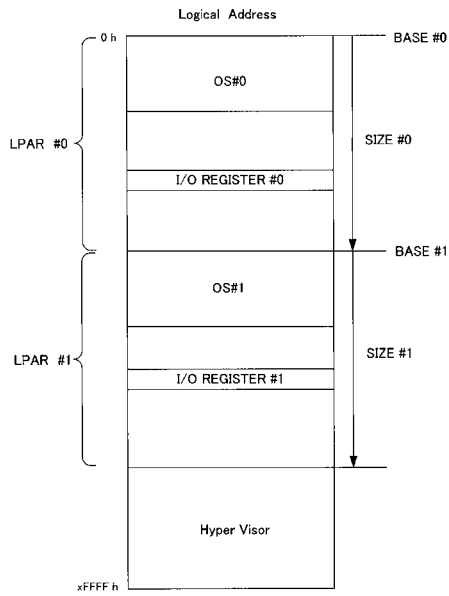
【図9】

パーティション番号	ベースアドレス	サイズ
0	Base0	Size0
1	Base1	Size1

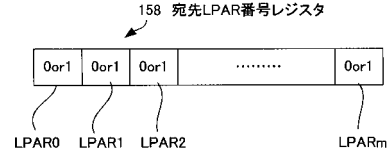
【図7】



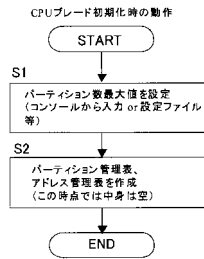
【図10】



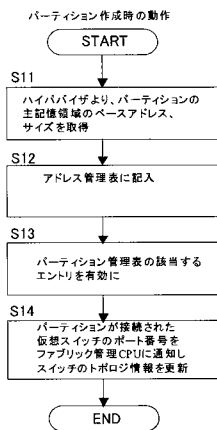
【図11】



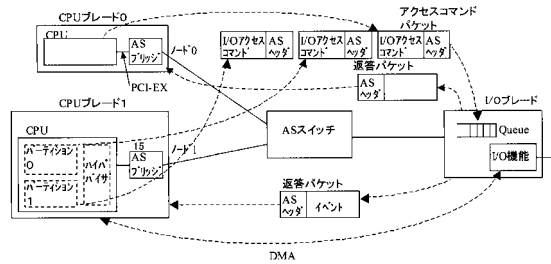
【図12】



【図13】

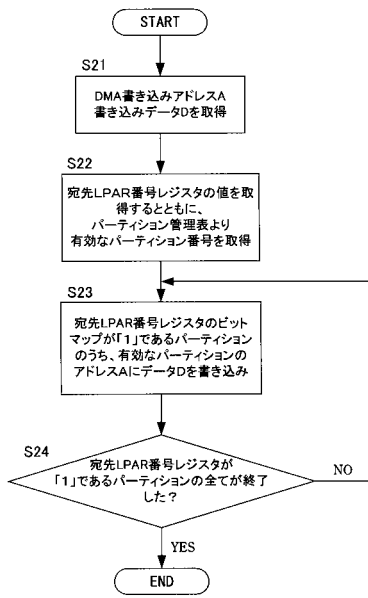


【図14】

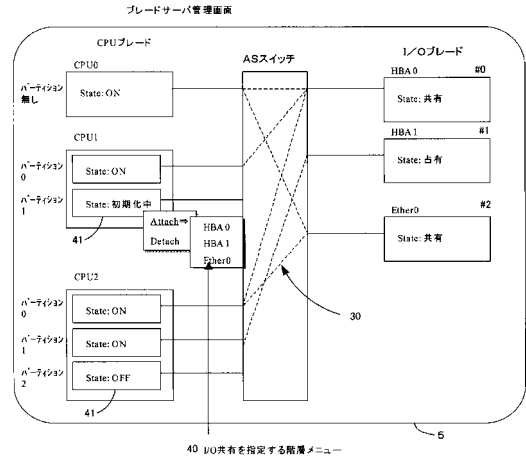




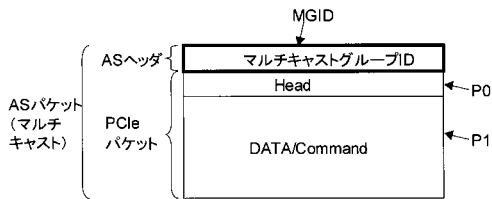
【図15】



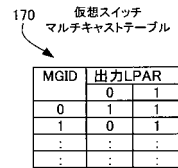
【図16】



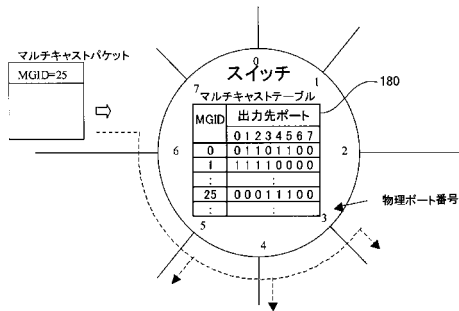
【図17】



【図19】



【図18】



---

フロントページの続き

審査官 鈴木 修治

- (56)参考文献 特開2002-318701(JP,A)  
特開2003-316752(JP,A)  
特開2002-222110(JP,A)  
特開2005-122640(JP,A)  
MAYHEW D, PCI Express and Advanced Switching, 2003 P  
ROCEEDINGS. 11TH SYMPOSIUM ON HIGH PERFORMANCE INTERCONNECTS, 米国, IEEE, 2003年  
8月20日, 頁21-29  
WONG W, ADVANCED SWITCHING FOR PCI EXPRESS, ELECTRONIC DESIGN, 米国, PENTON MEDIA, 2  
003年 6月23日, 頁36

(58)調査した分野(Int.Cl., DB名)

G06F 9/46 - 9/54

G06F 13/10