

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4660362号
(P4660362)

(45) 発行日 平成23年3月30日(2011.3.30)

(24) 登録日 平成23年1月7日(2011.1.7)

(51) Int.Cl.

F I

G06F 13/14 (2006.01)

G06F 13/14 310H

請求項の数 10 (全 29 頁)

(21) 出願番号	特願2005-340088 (P2005-340088)	(73) 特許権者	000005108
(22) 出願日	平成17年11月25日(2005.11.25)		株式会社日立製作所
(65) 公開番号	特開2007-148621 (P2007-148621A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成19年6月14日(2007.6.14)	(74) 代理人	100075513
審査請求日	平成20年9月19日(2008.9.19)		弁理士 後藤 政喜
		(74) 代理人	100084537
			弁理士 松田 嘉夫
		(74) 代理人	100114236
			弁理士 藤井 正弘
		(72) 発明者	對馬 雄次
			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】 計算機システム

(57) 【特許請求の範囲】

【請求項1】

複数のサーバと、

該複数のサーバとI/Oカードを接続するスイッチと、を備え、

前記サーバが、前記I/Oカードから書き込み可能なメモリ空間を備えた計算機システムにおいて、

前記スイッチとI/Oカードの間に配置されて、前記複数のサーバと前記I/Oカードとの間のアクセス要求信号及び応答信号を操作するI/Oカード共有部と、

前記I/OカードとI/Oカード共有部にアクセス可能であって、前記I/Oカード共有部から書き込み可能なメモリと、前記I/Oカード共有部からの割り込みを受け付けるプロセッサとを備えて、前記複数のサーバに対する前記I/Oカードの割り当てを管理するI/Oプロセッサと、

を備え、

前記スイッチは、前記サーバからI/Oカードへの指令と前記メモリ空間を指定するベースアドレスとを含んだアクセス要求信号に、当該アクセス要求信号の宛先と要求元の経路情報を前記アクセス要求信号のヘッダ情報に付加して前記I/Oカード共有部へ送信し、前記I/Oカードからサーバへの応答と前記メモリ空間を指定するベースアドレスとを含んだ応答信号を、当該応答信号のヘッダ情報に含まれる宛先のサーバに転送するヘッダ処理部を有し、

前記I/Oカード共有部は、

10

20

前記スイッチから受信したアクセス要求信号に含まれる指令が、予め設定した第1の指令のときには前記I/Oカード共有部のメモリに当該アクセス要求信号を書き込む要求信号書込部と、

前記スイッチから受信した前記アクセス要求信号に含まれる指令が、予め設定した第2の指令のときには前記I/Oプロセッサへ割り込みをかける割り込み発生部と、

前記I/Oカードからサーバへの応答信号を受信したときには、当該応答信号に含まれるベースアドレスから要求元の経路情報を抽出して、当該抽出した要求元の経路情報を応答信号のヘッダ情報の宛先に設定するヘッダ修正部と、

前記応答信号のベースアドレスに埋め込まれた要求元の経路情報を削除するベースアドレス修正部と、

当該応答信号を前記スイッチへ送信する送信部と、を有し、

前記I/Oプロセッサは、

前記割り込みがあったときに、前記メモリに書き込まれたアクセス要求信号に含まれるベースアドレスの所定の上位ビットに、前記アクセス要求信号のヘッダ情報から前記要求元の経路情報を設定するアドレス変換部と、

前記割り込みに基づいて前記要求元のヘッダ情報を上位ビットに設定した前記ベースアドレスを前記I/Oカードの応答先として設定する応答アドレス設定部と、

前記割り込みに基づいて前記アクセス要求信号をI/Oカードへ送信し、前記I/Oカードの動作を起動するI/Oカード起動部と、を有し、

前記I/Oカードは、前記アクセス要求信号に応じた処理の応答信号を、前記I/Oカード共有部に返信することを特徴とする計算機システム。

【請求項2】

前記I/Oカードは、

前記アクセス要求信号に含まれる指令を設定するコマンドレジスタと、

前記アクセス要求信号に対して応答するサーバのメモリ空間を指定するアドレスレジスタと、を有し、

前記要求信号書込部は、前記アクセス要求信号に含まれる指令が前記コマンドレジスタを除くI/Oカードへの書き込み処理を含む第1の指令のときに、前記スイッチから受信したアクセス要求信号を前記I/Oカード共有部のメモリに書き込み、

前記割り込み発生部は、前記スイッチから受信したアクセス要求信号に含まれる指令が前記コマンドレジスタへの書き込み処理を含む第2の指令のときに、前記I/Oプロセッサへ割り込みをかけ、

前記応答アドレス設定部は、前記割り込みに基づいて、前記要求元のヘッダ情報を上位ビットに設定した前記ベースアドレスを前記I/Oカードのアドレスレジスタに書き込むことを特徴とする請求項1に記載の計算機システム。

【請求項3】

前記要求信号書込部は、前記アクセス要求信号に含まれる指令がアドレスレジスタへの書き込みを要求するDMA初期化処理を含むときに、前記スイッチから受信したアクセス要求信号を前記I/Oカード共有部のメモリに書き込み、

前記割り込み発生部は、前記スイッチから受信したアクセス要求信号にコマンドレジスタへの書き込みを要求するDMA開始指令が含まれるときに、前記I/Oプロセッサへ割り込みをかけることを特徴とする請求項2に記載の計算機システム。

【請求項4】

前記I/Oプロセッサは、

前記各サーバが前記I/Oカードを利用可能か否かを識別し、かつ前記I/Oカードが共有されているか否かを識別する属性を予め設定したI/Oカード共有設定管理部を有し、

前記I/Oカード共有部は、

前記I/Oカード共有設定管理部を参照して、前記スイッチから受信したアクセス要求信号の宛先のI/Oカードと要求元のサーバの前記属性が利用可能かつ占有のときには、

10

20

30

40

50

前記サーバからの要求信号をそのまま I / O カードに転送し、当該 I / O カードから前記サーバへの応答信号をそのまま転送することを特徴とする請求項 1 に記載の計算機システム。

【請求項 5】

前記 I / O プロセッサは、

前記各サーバが前記 I / O カードを利用可能か否かを識別し、かつ前記 I / O カードが共有されているか否かを識別する属性を予め設定した I / O カード共有設定管理部を有し、

前記 I / O カード共有部は、

前記 I / O カード共有設定管理部を参照して、前記スイッチから受信したアクセス要求信号の宛先に設定されている I / O カードと要求元のサーバの前記属性が利用不可のときには、前記サーバからの要求信号に対してマスタポートまたは全ビットが 1 のデータの何れか一方を返信することを特徴とする請求項 1 に記載の計算機システム。

10

【請求項 6】

前記 I / O プロセッサは、

前記各サーバが前記 I / O カードを利用可能か否かを識別し、かつ前記 I / O カードが共有されているか否かを識別する属性を予め設定した I / O カード共有設定管理部と、

前記 I / O プロセッサに接続された入力部を介して前記 I / O カード共有設定部の属性を設定する属性設定部と、

を有することを特徴とする請求項 1 に記載の計算機システム。

20

【請求項 7】

前記 I / O カード共有部は、

前記要求信号書込部または割り込み発生部で用いるアドレス情報を、前記スイッチから受信したアクセス要求信号に含まれるヘッダ情報と予め設定した情報とを比較して出力するアドレス情報判定部を有し、

当該アドレス情報判定部は、前記サーバの起動時に実行される I / O カードへの書き込み指令に基づいて前記情報を設定することを特徴とする請求項 1 に記載の計算機システム。

【請求項 8】

前記 I / O カードは、

前記アクセス要求信号に含まれる指令を設定するコマンドレジスタと、

前記アクセス要求信号に対して応答するサーバのメモリ空間を指定するアドレスレジスタと、

前記サーバのメモリに格納されたデータ構造体を読み込んでコマンドチェイン処理を行うコマンドチェイン処理部と、

を有し、

前記第 1 及び第 2 の指令が一つの起動要求で構成されて、当該起動要求を前記 I / O カードのコマンドレジスタに書き込むことで行われ、

前記要求信号書込部は、前記要求信号を I / O プロセッサのメモリへ書き込み、また、前記サーバのメモリに設定されたデータ構造体を I / O プロセッサのメモリへ書き込み、

前記 I / O カードは、当該 I / O カードの起動後に、前記 I / O プロセッサのメモリに設定されたデータ構造体を読み込んで処理を行うことを特徴とする請求項 2 に記載の計算機システム。

40

【請求項 9】

複数のサーバと、

該複数のサーバと I / O カードを接続するスイッチと、を備え、

前記サーバが、前記 I / O カードから書き込み可能なメモリ空間を備えた計算機システムにおいて、

前記スイッチと I / O カードの間に配置されて、前記複数のサーバと前記 I / O カードとの間のアクセス要求信号及び応答信号を操作する I / O カード共有部と、

前記 I / O カードと I / O カード共有部にアクセス可能であって、前記 I / O カード共

50

有部から書き込み可能なメモリと、前記 I / O カード共有部からの割り込みを受け付けるプロセッサとを備えて、前記複数のサーバに対する前記 I / O カードの割り当てを管理する I / O プロセッサと、
を備え、

前記スイッチは、

前記サーバから I / O カードへの指令と前記メモリ空間を指定するベースアドレスとを含んだアクセス要求信号に、当該アクセス要求信号の宛先と要求元の経路情報を前記アクセス要求信号のヘッダ情報に付加して前記 I / O カード共有部へ送信し、前記 I / O カードからサーバへの応答と前記メモリ空間を指定するベースアドレスとを含んだ応答信号を、当該応答信号のヘッダ情報に含まれる宛先のサーバに転送するヘッダ処理部と、前記 I / O カード共有部を有し、

10

前記 I / O カード共有部は、

前記ヘッダ処理部から受信したアクセス要求信号に含まれる指令が、予め設定した第 1 の指令のときには前記 I / O カード共有部のメモリに当該アクセス要求信号を書き込む要求信号書込部と、

前記ヘッダ処理部から受信した前記アクセス要求信号に含まれる指令が、予め設定した第 2 の指令のときには前記 I / O プロセッサへ割り込みをかける割り込み発生部と、

前記 I / O カードからサーバへの応答信号を受信したときには、当該応答信号に含まれるベースアドレスから要求元の経路情報を抽出して、当該抽出した要求元の経路情報を応答信号のヘッダ情報の宛先に設定するヘッダ修正部と、

20

前記応答信号のベースアドレスに埋め込まれた要求元の経路情報を削除するベースアドレス修正部と、

当該応答信号を前記スイッチへ送信する送信部と、を有し、

前記 I / O プロセッサは、

前記割り込みがあったときに、前記メモリに書き込まれたアクセス要求信号に含まれるベースアドレスの所定の上位ビットに、前記アクセス要求信号のヘッダ情報から前記要求元の経路情報を設定するアドレス変換部と、

前記割り込みに基づいて前記要求元のヘッダ情報を上位ビットに設定した前記ベースアドレスを前記 I / O カードの応答先として設定する応答アドレス設定部と、

前記割り込みに基づいて前記アクセス要求信号を I / O カードへ送信し、前記 I / O カードの動作を起動する I / O カード起動部と、を有し、

30

前記 I / O カードは、前記アクセス要求信号に応じた処理の応答信号を、前記 I / O カード共有部に返信することを特徴とする計算機システム。

【請求項 10】

前記サーバとスイッチの間と、前記スイッチと I / O カード共有部及び前記 I / O カード共有部と I / O カードの間を、PCI または PCI - EXPRESS で接続したことを特徴とする請求項 1 または請求項 9 に記載の計算機システム。

【発明の詳細な説明】

【技術分野】

【0001】

40

本発明は、I / O デバイスの共有を行う計算機システムに関し、特に、サーバ統合を実現可能なブレードサーバシステムに関する。

【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】

サーバ統合を実現する技術として、バックプレーンに複数のサーバブレード及び複数の I / O ブレードを装着可能なブレードサーバシステムが知られており、統合するサーバに

50

応じてサーバブレードを装着し、一つの筐体内に複数のサーバを統合している。つまり、一つの筐体に複数のサーバを集約した高性能なサーバを提供する。

【0004】

サーバ統合では、複数の筐体で独立して稼動していたサーバを、一つの高性能なブレードサーバシステムに集約し、複数のOSとアプリケーションを稼動させることになる。ここで、複数のサーバを一つのブレードサーバシステムに統合する際には、独立して稼動していたサーバのI/Oカード(またはI/Oデバイス)も統合する必要がある。各サーバが統合前に使用していたI/Oカードをそのまま利用しようとする、ブレードサーバシステムのバックプレーン(またはI/Oスロット)が不足する。

【0005】

ブレードサーバシステムの筐体にI/O拡張用の筐体を付加して、必要な数のI/Oカードを装着することも可能ではあるが、筐体のサイズが大形化するためより大きな設置スペースが必要になってしまい、サーバ統合の効果が薄れてしまう。

【0006】

このため、ブレードサーバシステムを用いたサーバ統合では、サーバブレード間でI/Oカードを共有し、I/Oカードの数を削減する必要が生じる。

【0007】

I/Oカードを複数のサーバで共有する技術としては、一つのコンピュータを任意論理区画に分割する仮想計算機が知られている。これは、ホストOS上で複数のゲストOSを稼動させ、各ゲストOSを論理区画として提供し、ゲストOSホストOSのアプリケーションとして稼動させ、ゲストOSからのI/O要求はホストOSが一元的にI/O要求を処理することで、I/Oデバイスの共有を行っている(例えば、特許文献1)。

【0008】

あるいは、AS(Advanced Switching)等に代表されるI/Oスイッチを用いて複数のサーバからI/Oカードを共有する技術も知られている。これは、PCI-EXPRESS規格のI/Oカードとサーバ間を専用のプロトコルで通信を行うスイッチであり、一つのI/Oカードを複数のサーバで切り換えて使用し、共有することができる。

【特許文献1】米国特許第6,496,847号

【発明の開示】

【発明が解決しようとする課題】

【0009】

しかしながら、上記特許文献1のような従来例では、ゲストOSとI/Oカードとの間のI/Oアクセスは、常にホストOSが中継することになる。このため、ゲストOSとI/Oカードの間でDMA転送を行う場合などでは、ホストOSのメモリ領域からゲストOSのメモリ領域へ転送する処理が必要となり、このホストOSの処理がオーバーヘッドとなってI/Oアクセスの性能(転送速度やレスポンス)が低下する、という問題があった。

【0010】

また、上記従来例のASでは、PCI-EXPRESSを拡張した規格(専用のプロトコル)に対応する専用のI/Oカードが必要となり、従来のPCIやPCI-x、PCI-EXPRESSといった汎用のインターフェースで構成されたI/Oカードをそのまま利用することはできない。したがって、ASを利用してI/Oカードを共有するには、従来のI/Oカードを置き換えるために多額の費用が必要になるという問題がある。さらに、ASに対応したI/Oカードでは、専用のプロトコルを使用するため、デバイスドライバなども新規に導入する必要が生じる。このため、サーバ統合の際には、従来のOSに組み込んでいたデバイスドライバを変更する作業が必要となり、上記費用の問題に加えてサーバ統合に要する労力が増大する、という問題があった。

【0011】

そこで本発明は、上記問題点を鑑みてなされたもので、I/Oアクセスの性能低下を防いで複数のサーバ間でI/Oカードの共有を実現することを目的とし、さらに、共有する

10

20

30

40

50

I/Oカードを汎用のインターフェースで構成可能にすることで、サーバ統合を行う際のコストを抑制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、複数のサーバと、該複数のサーバとI/Oカードを接続するスイッチと、を備え、前記サーバが、前記I/Oカードから書き込み可能なメモリ空間を備えた計算機システムにおいて、前記スイッチとI/Oカードの間に配置されて、前記複数のサーバと前記I/Oカードとの間のアクセス要求信号及び応答信号を操作するI/Oカード共有部と、前記I/OカードとI/Oカード共有部にアクセス可能であって、前記I/Oカード共有部から書き込み可能なメモリと、前記I/Oカード共有部からの割り込みを受け付けるプロセッサとを備えて、前記複数のサーバに対する前記I/Oカードの割り当てを管理するI/Oプロセッサと、を備え、前記スイッチは、前記サーバからI/Oカードへの指令と前記メモリ空間を指定するベースアドレスとを含んだアクセス要求信号に、当該アクセス要求信号の宛先と要求元の経路情報を前記アクセス要求信号のヘッダ情報に付加して前記I/Oカード共有部へ送信し、前記I/Oカードからサーバへの応答と前記メモリ空間を指定するベースアドレスとを含んだ応答信号を、当該応答信号のヘッダ情報に含まれる宛先のサーバに転送するヘッダ処理部を有し、前記I/Oカード共有部は、前記スイッチから受信したアクセス要求信号に含まれる指令が、予め設定した第1の指令のときには前記I/Oカード共有部のメモリに当該アクセス要求信号を書き込む要求信号書込部と、前記スイッチから受信した前記アクセス要求信号に含まれる指令が、予め設定した第2の指令のときには前記I/Oプロセッサへ割り込みをかける割り込み発生部と、前記I/Oカードからサーバへの応答信号を受信したときには、当該応答信号に含まれるベースアドレスから要求元の経路情報を抽出して、当該抽出した要求元の経路情報を応答信号のヘッダ情報の宛先に設定するヘッダ修正部と、前記応答信号のベースアドレスに埋め込まれた要求元の経路情報を削除するベースアドレス修正部と、当該応答信号を前記スイッチへ送信する送信部と、を有し、前記I/Oプロセッサは、前記割り込みがあったときに、前記メモリに書き込まれたアクセス要求信号に含まれるベースアドレスの所定の上位ビットに、前記アクセス要求信号のヘッダ情報から前記要求元の経路情報を設定するアドレス変換部と、前記割り込みに基づいて前記要求元のヘッダ情報を上位ビットに設定した前記ベースアドレスを前記I/Oカードの応答先として設定する応答アドレス設定部と、前記割り込みに基づいて前記アクセス要求信号をI/Oカードへ送信し、前記I/Oカードの動作を起動するI/Oカード起動部と、を有し、前記I/Oカードは、前記アクセス要求信号に応じた処理の応答信号を、前記I/Oカード共有部に返信する

【発明の効果】

【0013】

したがって、本発明は、I/Oカードからサーバへ向かう応答信号を中継するI/Oカード共有部は、ベースアドレスに埋め込まれた要求元のサーバの識別子をヘッダ情報の宛先に付け替えることで、汎用バスのI/Oカードを共有することができる。そして、I/Oカードの起動後は、I/Oプロセッサは応答信号の転送へ介入せず、I/Oカード共有部のハードウェアがアドレス付け替えを行うため、前記従来例のようなソフトウェア処理によるオーバーヘッドを防いで、I/Oアクセスの性能低下を防ぎながら複数のサーバ間でI/Oカードの共有を実現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】

【0014】

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。

【0015】

図1は、第1の実施形態を示すブレードサーバシステムのブロック図である。

【0016】

ブレードサーバシステムは、複数のサーバブレード10-1~nと、各種I/Oインターフェースを備えたI/Oカード501、502と、サーバブレード10-1~nとI/O

10

20

30

40

50

カードを接続するスイッチ200と、I/Oカード501、502を複数のサーバブレード10-1~nで共有するためのI/Oカード共有機構400と、I/Oカードの共有を管理するI/Oプロセッサブレード600とを備える。そして、これらの各ブレードとスイッチ200及びI/Oカード共有機構400は一つの筐体(図示省略)内に収納される。

【0017】

サーバブレード10-1~nは、それぞれCPU101とメモリ102がチップセット(あるいはI/Oブリッジ)103を介して接続され、また、チップセット103は汎用バス11-1~nを介してスイッチ200に接続される。ここで、汎用バス11-1~nとしては、例えば、PCI-EXPRESS(図中PCI-ex)を適用した場合を示す。

10

【0018】

CPU101はメモリ102にロードしたOSやアプリケーションを実行することでサーバ#1~nを提供する。そして、CPU101は、チップセット103からスイッチ200及びI/Oカード共有機構400を介してI/Oカード501、502にアクセスを行う。

【0019】

スイッチ200は、サーバブレード10-1~nとI/Oカード501、502との間で送受信されるパケットのヘッダ情報を付加し、このヘッダ情報に基づいてパケットを転送するヘッダ処理部210を有する。

20

【0020】

ヘッダ処理部210は、サーバブレード10-1~nからI/Oカード501、502へ向けて送信されたパケット(アクセス要求信号)にヘッダ情報を付加して、このヘッダ情報に含まれる宛先のノード(I/Oカード)に転送する。このヘッダ情報は、サーバブレード10-1~nのアドレス情報(識別子)を要求元に設定し、I/Oカードのアドレス情報を宛先に設定する。また、スイッチ200のヘッダ処理部210は、I/Oカードからサーバブレード10-1~nへ向けたパケット(応答信号)を、このパケットに含まれる宛先(サーバ識別子)のサーバブレード10-1~nに転送する。ここで、本実施形態のパケットは、汎用バスとしてPCI-EXPRESSを採用しているため、PCIトランザクション(PCI-Tx)である。

30

【0021】

スイッチ200とI/Oカード501、502の間には、汎用バス301、311、312を介してI/Oカード501、502を複数のサーバブレード10-1~nで共有するためのI/Oカード共有機構400が接続される。また、I/Oカード共有機構400は汎用バス401を介してI/Oプロセッサブレード600に接続されており、I/Oプロセッサブレード600は後述するように、I/Oカードの共有に関するアドレス変換や共有状態などの管理を実行する。なお、I/Oプロセッサブレード600には、コンソール5が接続されており、管理者等の入力によりI/Oカード501、502の共有状態を設定する。

【0022】

I/Oカード501、502は、SCSI(またはSAS)やFC(Fibre Channel)あるいはEthernet(登録商標)等のインターフェースを備えている。I/Oカード501、502には、サーバブレード10-1~nのメモリ102へ直接アクセスするDMA(Direct Memory Access)コントローラ513をそれぞれ備えている。そして、I/Oカード501、502は、DMAコントローラ513でDMAを行うサーバブレード10-1~n上のメモリ102のMMIO(Memory Mapped I/O)のベースアドレスを指定するベースアドレスレジスタ511と、I/Oカード501、502に対する指令を指定するコマンドレジスタ512を備えている。DMAコントローラ513は、ベースアドレスレジスタ511に書き込まれたメモリ102のアドレスに対して、コマンドレジスタ512に書き込まれたコマンドに対応する動作を実行する。なお、I/Oカード501、502は

40

50

、P C I規格に準拠した図示しないレジスタ（コンフィグレーションレジスタ、レイテンシタイマレジスタ等）を有するものである。

【0023】

次に、I/Oプロセッサブレード600は、CPU602とメモリ603がチップセット（あるいはI/Oブリッジ）601を介して接続され、また、チップセット601は汎用バス401を介してI/Oカード共有機構400に接続される。そして、I/Oプロセッサブレード600では、後述するように、所定の制御プログラムが実行され、サーバブレード10-1～nからのI/Oアクセスに対してアドレス変換などの処理を実行する。

【0024】

< I/Oカード共有機構 >

次に、本発明のI/Oカード共有機構400の詳細について、図2のブロック図を参照しながら以下に詳述する。

【0025】

I/Oカード共有機構400は、サーバブレード10-1～nとI/Oカード501、502の間に位置して、サーバブレード10-1～nからのI/Oアクセスパケットのアドレスの変換を行うことで、ひとつのI/Oカードを複数のサーバブレード10-1～nで共有することを実現する。ここで、スイッチ200と汎用バス及びI/Oカード501、502はP C I - E X P R E S Sに準拠しており、以下、I/OアクセスパケットをP C Iトランザクションという。

【0026】

I/Oカード共有機構400の主な機能は、

- 1) サーバブレード10-1～nからI/Oカード501、502へのP C Iトランザクションを、I/Oプロセッサブレード600のメモリ603へ書き込む機能
 - 2) I/Oカード501、502のコマンドレジスタ512への書き込み要求に基づいて、I/Oプロセッサブレード600のCPU602へ割り込みを要求する機能
 - 3) I/Oカード501、502からサーバブレード10-1～nへのD M AによるP C Iトランザクションの宛先を変換する機能
- の3つである。

【0027】

図2において、I/Oカード共有機構400は、後述するアドレス情報テーブル411を格納する連想メモリ410と、サーバブレード10-1～nからのP C Iトランザクションからヘッダ情報を分離するヘッダ情報抽出部406と、ヘッダ情報を除くP C Iトランザクションの本体を解析してI/Oプロセッサブレード600に指令を送る第1のトランザクションデコーダ402（図中T xデコーダ1）と、I/Oプロセッサブレード600からの信号を解析してI/Oカード501、502に指令を送る第2のトランザクションデコーダ403（図中T xデコーダ2）と、I/Oカード501、502からのP C Iトランザクションを解析して、D M A転送であれば宛先アドレスを修正（変換）する第3のトランザクションデコーダ404（図中T xデコーダ3）と、トランザクションデコーダ402からの指令によりI/Oプロセッサブレード600のCPU602へ割り込みをかける割り込み発生部407と、トランザクションデコーダ402からの指令によりI/Oプロセッサブレード600のメモリ603へ書き込みを行うメモリ書込部408と、第1のトランザクションデコーダ402の指令に基づいてI/Oプロセッサブレード600へ出力する信号を選択する信号選択部412と、第3のトランザクションデコーダ404の指令に基づいてサーバ#1～#n側（スイッチ200）へ出力する信号を選択する信号選択部413とを主体にして構成されている。

【0028】

ここで、P C Iトランザクションは、図3で示すように、コマンドやオーダーまたはサーバアドレス（M M I Oベースアドレス）462などのデータを格納するP C Iトランザクション本体461と、経路情報を格納するヘッダ情報451から構成される。そして、ヘッダ情報451は、宛先452を先頭にして、P C Iトランザクションを発行した要求

10

20

30

40

50

元453が格納される。例えば、サーバブレード10-1からI/Oカード501へのPCIトランザクションは、宛先452にI/Oカード501への経路情報(アドレスなど)が設定され、要求元453にサーバブレード10-1の経路情報が設定され、PCIトランザクション本体461には、コマンドやオーダーに加えてサーバブレード10-1のI/Oレジスタのアドレス情報がMMIOベースアドレス462に設定される。

【0029】

図2において、301-1はスイッチ200(サーバブレード10-1~n側)とI/Oカード共有機構400を接続する汎用バス301上で、スイッチ200からI/Oカード共有機構400に向かう下りのPCIトランザクションを示し、301-2は汎用バス301上でI/Oカード共有機構400からスイッチ200(サーバブレード10-1~n側)に向かう上りのPCIトランザクションを示す。同様に、401-2はI/Oカード共有機構400とI/Oプロセッサブレード600を接続する汎用バス401上で、I/Oカード共有機構400からI/Oプロセッサブレード600へ向かう下りの指令信号を示し、401-1は、汎用バス401上でI/Oプロセッサブレード600からI/Oカード共有機構400へ向かう上りの指令信号(またはPCIトランザクション)を示す。さらに、311-1、312-1はI/Oカード共有機構400とI/Oカード501、502を接続する汎用バス311、312上でI/Oカード501、502からI/Oカード共有機構400へ向かう上りのPCIトランザクションを示し、311-2、312-2は汎用バス311、312上でI/Oカード共有機構400からI/Oカード501、502へ向かう下りのPCIトランザクションを示す。

【0030】

I/Oカード共有機構400は、スイッチ200(サーバ側)から下りのPCIトランザクション301-1を受信すると、ヘッダ情報抽出部406がPCIトランザクションを図3に示したヘッダ情報451とPCIトランザクション本体461に分離する。また、ヘッダ情報抽出部406はPCIトランザクションに設定されたMMIOベースアドレスからオフセットを抽出する。そして、ヘッダ情報抽出部406は、ヘッダ情報451とオフセットを連想メモリ410へ入力し、PCIトランザクション本体461を第1のトランザクションデコーダ402へ入力する。

【0031】

連想メモリ410は、CAM(Contents Addressable Memory)等で構成され、I/Oプロセッサブレード600により設定されたアドレス情報テーブル411を保持する。このアドレス情報テーブル411には、後述するように、I/Oカード共有機構400に接続されたI/Oカード501、502に対する各サーバ#1~#nのアクセス許可情報(割り当て情報)が格納されている。

【0032】

そして、連想メモリ410は、検索を行うアドレス(ヘッダ情報451)を検索キー(ビット列)として入力し、入力した検索キーに対応するアドレスを予め設定したテーブル(アドレス情報テーブル411)から出力するものである。後述するように、連想メモリ410は、入力されたヘッダ情報451から該当するMMIOのベースアドレスとI/Oプロセッサブレード600上のメモリ603のアドレスを出力する。

【0033】

第1のトランザクションデコーダ402は、ヘッダ情報抽出部406から受信したPCIトランザクション本体461と、連想メモリ410から受信したMMIOベースアドレスとから、PCIトランザクション本体461の指令を解析してI/Oプロセッサブレード600へ出力する下りの指令信号401-2を選択する。また、トランザクションデコーダ402は、PCIトランザクション本体461の指令が所定の指令でない場合には、I/Oカード共有機構400が受信したPCIトランザクションをそのままI/Oカード501または502へ下りのPCIトランザクションとして転送する。

【0034】

ここで、複数のサーバブレード10-1~n(サーバ#1~#n)間でひとつのI/O

10

20

30

40

50

カードを共有する際の、メモリ空間について説明する。以下の例では、3つのサーバ#1～#3がひとつのI/Oカード501を共有する場合を示す。

【0035】

MMIOのアドレス空間は、ひとつのI/Oカードについて各サーバ#1～#3が自身のメモリ102にそれぞれ設定したI/O領域である。例えば、図4で示すように、サーバ#1～#3(サーバブレード10-1～3)には、使用(または共有)するI/Oカード(ここでは、I/Oカード501)毎にそれぞれMMIOベースアドレスが0xA、0xB、0xC、メモリ空間のサイズを示すオフセットが0>X、0>Y、0>ZのMMIO領域が設定されている。なお、これらのMMIOベースアドレスとオフセットは、各サーバブレード10-1～nで起動したBIOSまたはOSが決定するものである。

10

【0036】

これら各サーバ#1～#3のMMIOに対応して、プロセッサブレード600のメモリ603には、図5で示すように、各サーバ#1～#3で共有するI/Oカード501のメモリ空間6031～6032が後述するように設定される。図5では、サーバ#1のMMIOベースアドレス0xAに対応してプロセッサブレード600のメモリ603にメモリ空間6031(0>P)が設定される。なお、I/Oプロセッサブレード600は、メモリ603に設定したI/Oカード共有設定テーブル(図7参照)に基づいて、共有対象のI/Oカードを共有するサーバのMMIOについてのみ、メモリ603上にメモリ空間を設定する。同様にI/Oカード501を共有するサーバ#2のMMIOベースアドレス0xBに対応してメモリ空間6032(0>Q)が設定され、サーバ#3のMMIOベースアドレス0xCに対応してメモリ空間6032(0>R)が設定される。

20

【0037】

そして、上記サーバ#1～#3のMMIOベースアドレス=0xA、0xB、0xCのI/O領域がI/Oカード501を共有するため、連想メモリ410のアドレス情報テーブル411は、図6で示すようにI/Oプロセッサブレード600により設定される。

【0038】

図6のアドレス情報テーブル411には、I/Oカード共有機構400がサーバ#1～#3から受信したPCIトランザクションのヘッダ情報451と比較を行うヘッダ4111と、ヘッダ情報451とアドレス情報テーブル411のヘッダ4111が一致したときに出力する各サーバのMMIOベースアドレス4112と、ヘッダ情報451とアドレス情報テーブル411のヘッダ4111が一致したときに出力するI/Oプロセッサブレード600のメモリ空間のアドレス(図中IoP ADDR)4113と、I/Oプロセッサブレード600のCPU602への割り込みを行う際に、連想メモリ410へ入力されたオフセットと比較を行うためのオフセット4114が予め設定される。

30

【0039】

上述のようにサーバ#1～#3でI/Oカード501を共有する場合、図6のアドレス情報テーブル411のヘッダ4111には、I/Oカード501とサーバ#1～#3のアドレス情報が設定され、MMIOベースアドレス4112には、図4に示した各サーバ#1～#3のMMIOベースアドレスが設定され、メモリ空間アドレス4113には図5で示したように各サーバ#1～#3のMMIOベースアドレスに対応したメモリ空間6031～6033のアドレス情報が設定される。また、オフセット4114は、サーバ#1～#3のメモリ空間のサイズを求めるためにMMIOベースアドレスからの差分が設定される。

40

【0040】

ここで、ヘッダ4111には、ヘッダ情報451の宛先452がI/Oカード501のアドレス情報を示す「Io1」と、I/Oアクセスを要求した要求元453のサーバ#1～#3のアドレス情報を示す「SV1」～「SV3」が一对の比較用アドレスとしてヘッダ4111に設定される。

【0041】

例えば、サーバ#1(サーバブレード10-1)からI/Oカード501へのPCIト

50

ランザクションには、ヘッダ情報 4 5 1 の宛先 4 5 2 に「I o 1」が設定され、要求元 4 5 3 にサーバ # 1 のアドレス情報である「S V 1」が設定されている。ヘッダ情報抽出部 4 0 6 で抽出されたヘッダ情報 4 5 1 を、連想メモリ 4 1 0 のアドレス情報テーブル 4 1 1 へ入力すると、サーバ # 1 の M M I O ベースアドレス $0 \times A$ と、サーバ # 1 が I / O カード 5 0 1 を共有するための I / O プロセッサブレード 6 0 0 のメモリ空間 6 0 3 1 のアドレス = $0 \times P$ が連想メモリ 4 1 0 から出力される。

【 0 0 4 2 】

連想メモリ 4 1 0 から出力された M M I O ベースアドレスは第 1 のランザクションデコード 4 0 2 へ入力され、メモリ空間 6 0 3 1 のアドレスは信号選択部 4 1 2 へ入力される。

10

【 0 0 4 3 】

以上のようなメモリ空間により、各ランザクションデコード 4 0 2 ~ 4 0 4 の動作について以下に説明する。

【 0 0 4 4 】

第 1 のランザクションデコード 4 0 2 は、ヘッダ情報抽出部 4 0 6 から受信した P C I トランザクション本体 4 6 1 より I / O カード 5 0 1、5 0 2 に対する指令を抽出し、この指令の内容に応じて、I / O カード共有機構 4 0 0 が I / O プロセッサブレード 6 0 0 または I / O カード 5 0 1、5 0 2 へ出力する信号を以下のように決定する。

【 0 0 4 5 】

A) P C I トランザクション本体 4 6 1 から抽出した指令が、I / O カード 5 0 1、5 0 2 のコマンドレジスタ 5 1 2 への書き込み指令 (I / O カードの動作を起動させる指令、例えば、D M A 転送開始コマンド) であれば、ランザクションデコード 4 0 2 は割り込み発生部 4 0 7 に割り込みの出力を指令し、信号選択部 4 1 2 へ割り込み発生部 4 0 7 の出力を選択して汎用バス 4 0 1 の下りの指令信号 4 0 1 - 2 として I / O プロセッサブレード 6 0 0 へ出力するよう指令する。この割り込み信号には、連想メモリ 4 1 0 が出力した I / O プロセッサブレード 6 0 0 のアドレス 4 1 1 3 を含める。また、割り込みを行う場合には、連想メモリ 4 1 0 へ入力したヘッダ情報 4 5 1 とオフセットが、アドレス情報テーブル 4 1 1 のヘッダ 4 1 1 1 とオフセット 4 1 1 4 に一致している必要がある。

20

【 0 0 4 6 】

B) P C I トランザクション本体 4 6 1 から抽出した指令が、I / O カード 5 0 1、5 0 2 のコマンドレジスタ 5 1 2 以外のレジスタ (例えば、ベースアドレスレジスタ 5 1 1 等) への書き込み指令 (例えば、D M A の初期化要求等で I / O カードの動作を起動しない指令) であれば、ランザクションデコード 4 0 2 はメモリ書込部 4 0 8 に対して I / O プロセッサブレード 6 0 0 の所定のメモリ空間へ、P C I トランザクション (ヘッダ情報 4 5 1 と P C I トランザクション本体 4 6 1) を書き込むように指令し、信号選択部 4 1 2 へメモリ書込部 4 0 8 の出力を選択して汎用バス 4 0 1 の下りの指令信号 4 0 1 - 2 として I / O プロセッサブレード 6 0 0 へ出力するよう指令する。なおメモリ書込部 4 0 8 は、連想メモリ 4 1 0 から出力された I / O プロセッサブレード 6 0 0 のアドレス 4 1 1 3 に対して、ヘッダ情報 4 5 1 と P C I トランザクション本体 4 6 1 を書き込む。

30

【 0 0 4 7 】

C) P C I トランザクション本体 4 6 1 から抽出した指令が、I / O カード 5 0 1、5 0 2 のレジスタへの書き込み要求以外の指令であれば、ランザクションデコード 4 0 2 は信号線 4 2 0 から受信した P C I トランザクションを、I / O カード共有機構 4 0 0 と I / O カード 5 0 1、5 0 2 を接続する汎用バス 3 1 1、3 1 2 から下りの P C I トランザクション 3 1 1 - 2、3 1 2 - 2 としてそのまま出力する。この場合、ランザクションデコード 4 0 2 は、P C I トランザクションのヘッダ情報 4 5 1 の宛先 4 5 2 に基づいて、該当する I / O カード 5 0 1 または 5 0 2 が接続された汎用バス 3 1 1 または 3 1 2 を選択して出力する。

40

【 0 0 4 8 】

なお、上記 A) ~ C) において、ランザクションデコード 4 0 2 は後述する I / O カ

50

ード共有設定405を参照して、要求元のサーバに要求先のI/Oカードが割り当てられていない場合にはアクセスを禁止する。また、トランザクションデコーダ402は後述するI/Oカード共有設定405を参照して、割り当て状態(属性)が占有を示す「Dedicate」の場合には、連想メモリ410によるアドレスのデコードを行わず、上記C)の機能によりPCIトランザクションをそのまま宛先のI/Oカードに転送し、サーバ#1~#nとI/Oカード501、502が直接I/O処理を行う通常のアクセスを行う。

【0049】

以上のように、I/Oカード共有機構400のトランザクションデコーダ402は、サーバ#1~#nとI/Oカード501、502の間に介在して、I/Oカード501、502のレジスタへの書き込みをI/Oプロセッサブレード600への操作(メモリ空間への書き込み処理、割り込み処理)に変換し、共有機能のないI/Oカード501、503を共有可能にするのである。例えば、サーバ#1~#nがI/Oカード501、502へDMA転送要求(DMA初期化要求)を行うと、上記B)の機能により、I/Oカード共有機構400がDMA転送を行うMMIOベースアドレスをI/Oプロセッサブレード600のメモリ空間に書き込む。次に、サーバ#1~#nがDMA転送の開始を指令すると、上記A)の機能によりI/Oプロセッサブレード600に割り込みをかけ、後述するように、I/Oプロセッサブレード600のCPU602がサーバ#1~#nに代わってDMAを実行させるI/Oカード501または502のコマンドレジスタ512、ベースアドレスレジスタ511に書き込みを行う。そして、I/Oカード501または502は、サーバ#1~#nの代理でI/Oアクセス要求を行ったI/Oプロセッサブレード600の指令により、要求元のサーバ#1~#nにDMA転送を行う。なお、I/Oプロセッサブレード600の詳細な動作については後述する。また、サーバブレード10-1~nから起動時の初期化のためにI/Oカードのコンフィグレーションレジスタ(図示省略)参照要求があったときには、トランザクションデコーダ402は、I/Oプロセッサブレード600のCPU602へ割り込みをかけ、さらに、PCIトランザクションをメモリ603へ書き込む。

【0050】

次に、I/Oカード共有機構400の第2のトランザクションデコーダ403の主な機能は、汎用バス401を介してI/Oプロセッサブレード600から受信した上りの指令信号401-1を、共有されているI/Oカード501または502に対してのみ出力するようにフィルタリングを行うことである。

【0051】

このため、I/Oカード共有機構400は、I/Oプロセッサブレード600のメモリ603上に設定したI/Oカード共有設定テーブル610(図7参照)を参照する領域としてI/Oカード共有設定405を格納するレジスタ430を備える。

【0052】

ここで、I/Oカード共有設定テーブル610は、図7で示すように、I/Oカード毎に、どのサーバが割り当てられているか(利用可能か)を示す属性のテーブルであり、コンソール5などから管理者が設定するものである。このテーブルは、I/Oカードのアドレス情報(デバイス番号など)などで構成される識別子611と、I/Oカードの機能などを示す種別612と、各サーバ#1~#3の割り当て状態613~615から構成される。

【0053】

図7はサーバ#1~#3とI/Oカード501、502の関係(属性)を示し、識別子611のI/Oカード1はI/Oカード501を示し、種別612がSCSIカードであり、属性は「Share」でサーバ#1、#2、#3がI/Oカード1を共有していることを示している。なお、稼動するサーバの数に応じて割り当て状態613~615の数も増減する。

【0054】

10

20

30

40

50

また、識別子 6 1 1 の I/O カード 2 は I/O カード 5 0 2 を示し、種別 6 1 2 が NIC カードであり、サーバ # 2 のみに割り当てられて占有（非共有）されている属性「Dedicate」を示し、この I/O カード 2 が他のサーバ # 1、# 3 に共有されていないことを示している。なお、稼動するサーバの数に応じて割り当て状態 6 1 3 ~ 6 1 5 の数も増減する。また、I/O カード 2 は、サーバ # 1、# 3 について割り当てられていない（利用可能ではない）ため、I/O カード 2 に対してこれらのサーバからのアクセスは禁止される。

【 0 0 5 5 】

トランザクションデコード 4 0 3 は、I/O プロセッサブレード 6 0 0 から PCI トランザクションを受信すると、PCI トランザクション本体 4 6 1 からサーバアドレス（MMIO ベースアドレス）を抽出し、アドレス情報テーブル 4 1 1 の MMIO ベースアドレス 4 1 1 2 と比較し、一致する MMIO ベースアドレスがあればそのエントリのヘッダ 4 1 1 1 から宛先と要求元を取得する。次に、I/O カード共有設定 4 0 5 の識別子と取得した宛先を比較し、一致するエントリで取得した要求元と同一のサーバを検索する。該当するサーバに宛先の I/O カードが割り当てられていれば、トランザクションデコード 4 0 3 が受信した PCI トランザクションは正当なものであるため、宛先の I/O カードに出力する。一方、該当するサーバに宛先の I/O カードが割り当てられていなければ、不平等な I/O アクセス要求であるため、PCI トランザクションを破棄する。なお、I/O カード共有機構 4 0 0 は、PCI トランザクションを破棄してから要求元のサーバにエラーの通知を行っても良い。

【 0 0 5 6 】

次に、I/O カード共有機構 4 0 0 の第 3 のトランザクションデコード 4 0 4 の主な機能は、汎用バス 3 1 1、3 1 2 を介して I/O カード 5 0 1、5 0 2 から受信した上りの PCI トランザクション 3 1 1 - 1、3 1 2 - 1 に対して、I/O アクセスの要求元のサーバ # 1 ~ # n へ返信するためにヘッダ情報 4 5 1 と PCI トランザクション本体 4 6 1 のサーバアドレス 4 6 2 を変換することである。

【 0 0 5 7 】

また、トランザクションデコード 4 0 4 は、I/O カード側から受信した PCI トランザクション 3 1 1 - 1、3 1 2 - 1 が、アドレス変換が必要なトランザクション（DMA 等）であるか、アドレス変換が不要なトランザクション（例えば、割り込みなどのイベント）であるかを判定し、信号選択部 4 1 3 でトランザクションデコード 4 0 4 の出力または上りの PCI トランザクション 3 1 1 - 1、3 1 2 - 1 の何れかを選択する。

【 0 0 5 8 】

トランザクションデコード 4 0 4 は、I/O カード側から受信した PCI トランザクション 3 1 1 - 1、3 1 2 - 1 が DMA 転送などであれば、アドレス変換が必要であると判定してトランザクションデコード 4 0 4 の出力を選択するよう信号選択部 4 1 3 に指令する。一方、アドレス変換が不要な PCI トランザクションであれば、受信した PCI トランザクション 3 1 1 - 1、3 1 2 - 1 をそのまま出力するよう信号選択部 4 1 3 へ指令する。

【 0 0 5 9 】

ここで、アドレス変換の有無の判定は、後述するように、図 3 で示した PCI トランザクション本体 4 6 1 の MMIO ベースアドレス 4 6 2 の未使用領域に設定された所定の上位ビットにサーバ # 1 ~ # n の識別子（アドレス情報など）が含まれているか否かに応じて、トランザクションデコード 4 0 4 が決定する。つまり、トランザクションデコード 4 0 4 は、PCI トランザクション本体 4 6 1 の MMIO ベースアドレス 4 6 2 の上位ビットにサーバ # 1 ~ # n の識別子（以下、サーバ識別子）が含まれていれば、アドレス変換が必要な PCI トランザクションであると判定し、サーバ識別子が含まれていなければアドレス変換が不要な PCI トランザクションであると判定する。

【 0 0 6 0 】

< I/O プロセッサブレード >

次に、I/Oプロセッサブレード600の機能について以下に説明する。図8は、I/Oプロセッサブレード600を主体とする機能ブロック図である。

【0061】

図8において、I/Oプロセッサブレード600のメモリ603には、上記図7で示したI/Oカード共有設定テーブル610と、複数のサーバが共有するI/Oカードのメモリ空間6031~6032(図中603x)が格納され、さらに、I/Oカード共有機構400からの割り込み(図中INT)により起動する割り込み処理部620が図示しないROMなどからロードされている。また、サーバブレード10-1~nの起動時には、I/Oカード共有機構400からの割り込みによって、初期化処理部630が図示しないROMなどからメモリ603へロードされる。

10

【0062】

I/Oカード共有設定テーブル610は、上記したようにI/Oプロセッサブレード600に接続されたコンソール5から、管理者などにより適宜設定されるもので、I/Oカードとサーバ#1~#nの割り当てを規定する。メモリ603のメモリ空間603xは、後述するように、サーバ#1~#nの起動時にCPU602により設定される。そして、I/Oカード共有機構400は、受信したスイッチ200からのPCIトランザクションが、上記B)に該当するとき、例えば、I/Oカードのベースアドレスレジスタ511への書き込みコマンド(DMA初期化要求)を含むとき、このPCIトランザクション本体461とヘッダ情報451を、アクセスを行うI/Oカードと要求元のサーバに対応するメモリ空間603xに書き込む。

20

【0063】

その後、I/Oカード共有機構400は、受信したスイッチ200からのPCIトランザクションが、上記A)に該当するとき(コマンドレジスタ512への書き込み)、I/Oプロセッサブレード600のCPU602に割り込みをかけて、割り込み処理部620を起動する。

【0064】

割り込み処理部620は、I/Oカード共有機構400からの割り込み指令に含まれるメモリ空間603xのアドレスに基づいて、メモリ空間603xに予め書き込まれたヘッダ情報451とPCIトランザクション本体461を読み込む。PCIトランザクション本体461に含まれるコマンドがDMA転送の場合、ヘッダ情報451とPCIトランザクション本体461に含まれるMMIOベースアドレス462を後述するように一時的に変換し、起動するI/Oカードのアドレスレジスタ511に変換したMMIOベースアドレスを書き込む。

30

【0065】

次に、割り込み処理部620は、ヘッダ情報451の宛先452となっているI/Oカードのコマンドレジスタ512へ、割り込みを発生させたPCIトランザクション本体461に含まれる指令(例えば、DMA転送開始)を書き込んでI/Oカードの動作を起動する。

【0066】

次に、割り込み処理部620がDMA転送の場合に行う上記アドレス変換について以下に説明する。

40

【0067】

図3で示すように、PCI-EXPRESSまたはPCIにおいてはPCIトランザクションのMMIOのアドレス空間として64ビット(図中0~63bit)が定義されている。また、サーバ#1~#n(サーバブレード10-1~n)のCPU101も64ビットのアドレッシングが可能なものが普及しつつある。しかし、64ビットのアドレス空間を使い切るようなメモリ102をサーバブレード10-1~nに実装するのは現実的ではなく、現状では数十GBを搭載可能なメモリ空間の上限とするのが一般的である。このため、CPU101等のアドレスバスも、図3の使用領域のように、64ビット未満の所定値、例えば、52ビット(0~51ビット)等に設定されている。

50

【 0 0 6 8 】

したがって、MMIOのアドレス空間としては64ビットが規定されてはいるものの、サーバブレード10-1~nの実装上は、アドレス空間の上位ビットは未使用となっている。

【 0 0 6 9 】

上述のように、下位52ビットをアクセス可能なアドレス空間とした場合、PCIトランザクション本体461に格納されるMMIOベースアドレス462の上位ビットのうち、図3の52~63bitは未使用領域463となる。ブレードサーバシステムでは、実装可能なサーバブレードの数を数十程度とすると、図3の未使用領域は12ビットあるので、この未使用領域の12ビットのうち少なくとも6ビットなどを使用すれば、筐体内の全てのサーバを識別することができる。

10

【 0 0 7 0 】

一方、PCI-EXPRESSに準拠したI/Oカードは、I/Oカード側で複数のサーバ#1~#nを識別することはできず、一旦DMA転送が開始されると、初期化時のMMIOベースアドレス462しか認識できないので、複数のサーバ#1~#nにサーバ#1~#nにDMA転送を行うことはできない。

【 0 0 7 1 】

そこで、本発明では、DMA転送の際に、PCIトランザクションのMMIOベースアドレス462の未使用の上位ビットを、サーバ識別子(アドレス情報)を格納する領域として利用し、I/Oプロセッサブレード600の割り込み処理部620がMMIOベースアドレス462の上位ビットにサーバ識別子となる要求元453を埋め込んでアドレス変換を行う。

20

【 0 0 7 2 】

そして、割り込み処理部620がI/Oカードのベースアドレスレジスタ511にアドレス変換を行ったMMIOベースアドレス462を書き込み、コマンドレジスタ512へDMA転送開始を書き込んでI/OカードのDMAを起動する。

【 0 0 7 3 】

I/OカードのDMA転送開始後には、I/Oカード共有機構400のトランザクションデコーダ404が、I/OカードからのPCIトランザクションを受信すると、MMIOベースアドレス462の上位ビットの未使用領域463に埋め込まれたサーバ識別子を抽出し、PCIトランザクションのヘッダ情報451の宛先452に書き込む。そして、トランザクションデコーダ404は、MMIOベースアドレス462のサーバ識別子の領域に「0」を書き込んで埋め込まれた要求元の経路情報を削除してから、このPCIトランザクションをスイッチ200へ送信する。スイッチ200は、PCIトランザクションのヘッダ情報に基づいて宛先に指定されたサーバ、つまりDMA転送を要求したサーバ#1~#nへPCIトランザクションを転送する。

30

【 0 0 7 4 】

すなわち、I/Oプロセッサブレード600の割り込み処理部620が、MMIOの未使用領域463に要求元453のアドレスをサーバ識別子として埋め込んだものをI/Oカードのアドレスレジスタ511に書き込んで、I/OカードのDMA転送を起動し、I/Oカードから出力されるDMA転送はI/Oカード共有機構400によって、PCIトランザクション内のMMIOベースアドレス462の未使用領域463からサーバ識別子を抽出して、ヘッダ情報451の宛先452に設定する。これにより、I/Oカード自体には複数のサーバ#1~#nを識別する機能がなくても、I/Oカード共有機構400とI/Oプロセッサブレード600によりI/Oカードの共有が可能になるのである。

40

【 0 0 7 5 】

< I/Oカード共有処理 >

次に、PCIトランザクションを主体にしたI/Oカード共有機構400とI/Oプロセッサブレード600によるI/Oカード共有の処理を図9に示す。

【 0 0 7 6 】

50

図9において、S1ではサーバ#1～#nが、I/Oアクセスを行うI/Oカードに対して、PCIトランザクションにDMA初期化コマンドなどを設定して送信する。サーバ#1～#nは、PCIトランザクションのヘッダ情報451の要求元453に自分のアドレス情報(サーバ識別子)を設定し、MMIOベースアドレス462にこのI/Oカードにサーバが割り当てたMMIOベースアドレスを設定する。

【0077】

S2では、I/Oカードとサーバ間に介在するI/Oカード共有機構400が、このPCIトランザクションを受信すると、I/Oカードのアドレスレジスタ511への書き込み指令を含むDMA初期化コマンドがあるので、このPCIトランザクションをI/Oプロセッサブレード600のメモリ空間603xに書き込む。この時点では、I/Oカードへのアクセスは行われない。

10

【0078】

S3では、サーバ#1～#nがI/Oカードに対してDMA転送開始のPCIトランザクションを送信すると、I/Oカード共有機構400はI/Oカードの動作を起動する指令を含むためI/Oプロセッサブレード600に割り込みをかけて、割り込み処理部620を起動する。

【0079】

割り込み処理部620は、メモリ空間603xに書き込まれたPCIトランザクションからMMIOベースアドレス462を読み込んで、I/Oカードのアドレスレジスタ511に書き込む。このとき、割り込み処理部620は、PCIトランザクション内のMMIOベースアドレス462の未使用領域463に、要求元のサーバを示すヘッダ情報451の要求元453を埋め込んでおく。そして、割り込み処理部620がI/Oカードのコマンドレジスタ512へDMA転送開始を書き込んで、I/Oカードの動作を起動する。

20

【0080】

S4では、I/Oカードがアドレスレジスタ511に設定されたMMIOベースアドレス462に対してDMA転送(書き込みまたは読み込み)を行う。

【0081】

I/OカードからのDMAによるPCIトランザクションには、MMIOベースアドレス462の上位ビットに設定された未使用領域463にサーバ識別子が埋め込まれている。

30

【0082】

S5では、I/Oカードとサーバ#1～#nの間に介在するI/Oカード共有機構400が、PCIトランザクションを受信すると上記図2のトランザクションデコーダ404でDMAによるPCIトランザクションであるか否かを判定する。

【0083】

このトランザクションデコーダ404で行われるI/OカードからのPCIトランザクションがDMAであるか否かの判定は、上記MMIOベースアドレス462の未使用領域463の全ビットが0でなければ、サーバ識別子が埋め込まれていると判定して、DMAによるPCIトランザクションであると判定することができる。

【0084】

そして、DMAのPCIトランザクションの場合、トランザクションデコーダ404は、MMIOベースアドレス462の未使用領域463の内容をヘッダ情報451の宛先452に設定し、スイッチ200で識別可能なサーバ#1～#nのアドレス情報に変換する。この後、トランザクションデコーダ404は未使用領域463の内容を消去するため、未使用領域463の全ビットに0をセットしてからこのPCIトランザクションを送信する。

40

【0085】

この宛先452に基づいてスイッチ200は、DMAのPCIトランザクションを、宛先452に設定されたDMAの要求元のサーバに転送し、サーバ#1～#nに設定したMMIOに対して所定のアクセスが行われる。

50

【 0 0 8 6 】

このように、I/Oカードのアドレスレジスタ511に設定されたMMIOベースアドレスの所定の上位ビットに、DMAを要求したサーバ識別子(要求元453)が設定されているので、DMAを繰り返して行ってもI/Oカード共有機構400が各PCIトランザクションの宛先452をDMAの要求元であるサーバのアドレス情報に付け替えるので、汎用のI/Oカードを用いながらも複数のサーバブレード10-1~nで共有することが可能となるのである。

【 0 0 8 7 】

上記処理を時系列的に示したものが、図10のタイムチャートである。まず、S11ではサーバ#1~#nが、I/Oアクセスを行うI/Oカードに対して、PCIトランザクションにDMA初期化コマンドなどを設定して送信する。

10

【 0 0 8 8 】

S12では、I/Oカード共有機構400が、I/Oカードのコマンドレジスタ512以外のレジスタに対する書き込み要求であるので、I/Oプロセッサブレード600のメモリ空間にPCIトランザクションの内容を書き込む。

【 0 0 8 9 】

次に、S13では、サーバ#1~#nがI/Oアクセスを行うI/Oカードに対して、PCIトランザクションにDMA転送開始コマンドなどコマンドレジスタ512への書き込み指令を設定して送信する。

【 0 0 9 0 】

S14では、I/Oカード共有機構400が、I/Oカードのコマンドレジスタ512に対する書き込み指令であるので、I/Oプロセッサブレード600のCPU602へ割り込みを要求する。

20

【 0 0 9 1 】

S15では、I/Oプロセッサブレード600のCPU602が割り込み処理部620を起動して、アドレスレジスタ511等コマンドレジスタ512以外の内容を読み出す。

【 0 0 9 2 】

S16では、割り込み処理部620が読み込んだPCIトランザクションの内容がDMAであれば、ヘッダ情報451の要求元453をMMIOベースアドレス462の未使用領域463にセットする。そして、このアドレス変換を行ったMMIOベースアドレス462をI/Oカードのアドレスレジスタ511に書き込み、コマンドレジスタ512にDMA転送開始コマンドを書き込んでI/Oカードの動作を起動する。

30

【 0 0 9 3 】

S17では、I/OカードのDMAコントローラ513が、アドレスレジスタ511のメモリ空間にDMAアクセスを行う。

【 0 0 9 4 】

S18では、I/OカードからのPCIトランザクションを受信し、DMAであるかをトランザクションデコーダ404が上述のように判定する。そしてトランザクションデコーダ404は、DMAのPCIトランザクションであれば、PCIトランザクション本体461のMMIOベースアドレス462の未使用領域463のサーバ識別子を、ヘッダ情報451の宛先452に設定してアドレス情報の変換(復元処理)を実施する。そして、トランザクションデコーダ404が、スイッチ200を介してDMAを要求したサーバにPCIトランザクションを送信する。

40

【 0 0 9 5 】

以上の手順により、MMIOベースアドレス462の未使用領域463にDMAを要求したサーバの識別子を埋め込んでI/OカードのDMAコントローラ513に処理を実行させることで、一つのI/Oカードを複数のサーバ#1~#nで共有することが可能となる。

【 0 0 9 6 】

50

< アドレス情報の設定処理 >

次に、図 11 はサーバブレード 10 - 1 ~ n を起動したときに実行される、アドレス情報テーブル 411 の設定処理の一例を示すタイムチャートである。

【 0097 】

連想メモリ 410 に格納されるアドレス情報テーブル 411 は、サーバブレード 10 - 1 ~ n が起動するたびに図 11 の処理により更新される。なお、以下では、サーバブレード 10 - 1 を起動した場合について説明する。

【 0098 】

まず、S20 でサーバブレード 10 - 1 の電源を ON にする。S21 では、電源の投入により CPU 101 が図示しない BIOS を起動して、I/O カード (デバイス) の初期化を行うために各 I/O カードのコンフィグレーションレジスタに対して読み込みを要求する。

【 0099 】

S22 では、I/O カード共有機構 400 のトランザクションデコーダ 402 は、I/O カードのコンフィグレーションレジスタの読み込み要求を受けたので、上述したように、I/O プロセッサブレード 600 の CPU 602 へ割り込みをかけ、さらに、このコンフィグレーションレジスタの読み込み要求の PCI トランザクションをメモリ 603 へ書き込む。このとき、コンフィグレーションレジスタの読み込み要求を出したサーバブレード 10 - 1 では MMIO が未設定であるため、トランザクションデコーダ 402 は予め設定したアドレスに書き込みを行う。

【 0100 】

S23 では、I/O プロセッサブレード 600 の CPU 602 が、トランザクションデコーダ 402 の割り込みによって図 8 に示した初期化処理部 630 を起動する。初期化処理部 630 は、所定のアドレスに書き込まれたコンフィグレーションレジスタの読み込み要求から、I/O カード共有設定テーブル 610 を読み込んで当該サーバブレードに割り当てられた I/O カードを確認する。

【 0101 】

S24 では、I/O カードが当該サーバブレード 10 - 1 に割り当てられていなければ (アクセス禁止)、初期化処理部 630 は I/O カード共有機構 400 を介して割り当てがないことを通知する (マスターアポート扱い)。一方、I/O カードが当該サーバブレードに割り当てられていれば、初期化処理部 630 はこの I/O カードのコンフィグレーションレジスタの内容を読み込んで、サーバブレード 10 - 1 に応答する。なお、S24 の処理は、I/O カード共有 610 に設定されている全ての I/O カードについて順次実行する。

【 0102 】

S25 では、I/O プロセッサブレード 600 から I/O カードのコンフィグレーションレジスタの内容を受信すると、サーバブレード 10 - 1 は取得した I/O カードの情報から MMIO 空間や IO 空間を設定し、MMIO ベースアドレスの設定などを行う。そして、サーバブレード 10 - 1 は、MMIO ベースアドレスを I/O カードに通知する。なお、この通知は I/O カード毎に実行される。

【 0103 】

S26 では、I/O カード共有機構 400 がサーバブレード 10 - 1 からの MMIO ベースアドレスを通知する PCI トランザクションを受信する。そして、I/O カード共有機構 400 は、MMIO ベースアドレスの設定通知であるので、I/O プロセッサブレード 600 の CPU 602 へ割り込みをかけ、さらに、MMIO ベースアドレスを通知する PCI トランザクションをメモリ 603 へ書き込む。なお、この処理は上記 S22 と同様に行われる。

【 0104 】

S27 では、I/O プロセッサブレード 600 の CPU 602 が、トランザクションデコーダ 402 の割り込みによって初期化処理部 630 を起動する。初期化処理部 630 は

10

20

30

40

50

、所定のアドレスに書き込まれたサーバブレード10-1のMMIOベースアドレス設定通知から、メモリ603にサーバブレード10-1の当該I/Oカードに対応するメモリ空間6031を割り当てる。そして、初期化処理部630は、サーバブレード10-1のMMIOベースアドレス及びオフセットと、I/Oプロセッサブレード600のメモリ空間6031のアドレスと、割り当てたI/Oカードのアドレス情報と、サーバブレード10-1のアドレス情報をI/Oカード共有機構400に通知し、連想メモリ410のアドレス情報テーブル411にこれらのアドレスを反映させる。なお、S27の処理は、サーバブレード10-1が利用するI/Oカード毎に繰り返して実行すればよい。

【0105】

以上の処理により、複数のサーバブレード10-1~nで共有されるI/Oカード501、502へ起動したサーバブレード10-1~nが直接アクセスすることなく、I/Oプロセッサブレード600が代理でコンフィグレーションレジスタの内容を取得し、メモリ空間603xの設定などを行う。ブレードサーバシステムでは、他のサーバブレードが稼働しているときに、新たなサーバブレードを設置し起動することができる。このような場合に、新たなサーバブレードの起動時にI/Oプロセッサブレード600が、I/Oカードに代わって応答をすることで、他の稼働中のサーバブレードのI/Oアクセスに影響を当てることなく新規のサーバブレードを起動することが可能となる。

【0106】

なお、上記ではサーバブレードに搭載したBIOSにより起動する例を示したが、図示はしないが、EFI(Extensible Firmware Interface)により起動する場合でも上記と同様に処理することが可能である。

【0107】

<I/Oカード共有設定テーブル>

図7に示したI/Oカード共有設定テーブル610について、以下に説明する。

【0108】

上述したように、I/Oカード共有設定テーブル610は、I/Oカード毎に、どのサーバに割り当てられているかを示すテーブルであり、コンソール5などから管理者が適宜設定するものである。なお、各サーバ#1~#nとI/Oカードの共有、占有、アクセス禁止の属性を、図7に示すインターフェースで、コンソール5のディスプレイへ表示し、図示しないマウスやキーボードなどのインターフェースで設定することができる。

【0109】

I/Oカード共有設定テーブル610の識別子611は、I/Oカードの追加、変更の度に、管理者が適宜設定する。また、I/Oカードの機能などを示す種別612は、I/Oカードのコンフィグレーションレジスタのクラスコード及びサブクラスコードを読み込んで、本テーブルに設定することができる。

【0110】

各サーバ#1~#3の割り当て状態613~615は、管理者がサーバ#1~#3やI/Oカードの特性や性能等から共有の有無、割り当ての有無を適宜設定する。図7において、複数のサーバに割り当てて共有する場合には「Share」を設定し、割り当てない場合には「No Assigned」を設定し、一つのサーバで占有する場合には「Dedicate」を設定して他のサーバからこのI/Oカードへのアクセスを禁止する。なお、I/Oカードの割り当てが「No Assigned」に設定されたサーバは、当該I/Oカードへアクセスするとアクセスが拒否される。この拒否されたアクセスが読み込みの場合には、I/Oカード共有機構400は全ビットが1のデータを応答し、上記アクセスが書き込みの場合には、マスターアポートを通知する。

【0111】

このI/Oカード共有設定テーブル610は、I/Oカード共有機構400のI/Oカード共有設定405に反映される。そして、I/Oプロセッサブレード600からI/OカードへのPCITランザクションを管理する第2のランザクションデコード403は、I/Oカード共有設定405に基づいて、正当なランザクションのみを許可し、不正

10

20

30

40

50

なトランザクション、すなわち、I/Oカード共有設定テーブル610で割り当てられていないサーバによるアクセスを禁止する。

【0112】

このI/Oカード共有設定テーブル610により、全てのI/Oカードを共有することもできるが、スループットを確保するためなどで一つのI/Oカードあるサーバに占有させ、他のI/Oカードを共有するように設定することができる。これにより、I/Oカードの共有と占有を混在させることが可能となって、ブレードサーバシステムのI/Oデバイスの構成に柔軟性を持たせることが可能となって、I/Oデバイスのリソースを有効に利用できるのである。

【0113】

<第2実施形態>

図12は、第2の実施形態を示すブレードサーバシステムのブロック図である。このブレードサーバシステムでは、前記第1実施形態のI/Oカードにコマンドチェーン制御部514を付加し、サーバ#1～#nのメモリ02に設定されたデータ構造体を順次読み込んでI/O処理を行うI/Oカード1501、1502で構成したもので、その他の構成は前記第1実施形態と同様である。

【0114】

I/Oカード1501、1502は、サーバブレード10-1～nのメモリ102に設定されたデータ構造体を順次読み込んで、各データ構造体の記述に従ってI/O動作を行うもので、いわゆるコマンドチェーン処理を行うものである。

【0115】

サーバブレード10-1～nで稼動するサーバ#1～#nは、I/Oカードを使用する際に、図13で示すように、各サーバ#1～#nのメモリ空間の所定のアドレスにデータ構造体1020～1022を設定しておく。

【0116】

例えば、サーバ#1のメモリ空間には、3つのデータ構造体1020(CCW1～3)がアドレス0×Dに設定され、サーバ#2のメモリ空間には、2つのデータ構造体1021(CCW11、12)がアドレス0×Eに設定され、サーバ#3のメモリ空間には、4つのデータ構造体1022(CCW21～24)がアドレス0×Fに設定される。なお、これらのデータ構造体は、各サーバ#1～#3のOSまたはアプリケーションが適宜設定するものである。

【0117】

先頭及び中間のデータ構造体には、次のデータ構造体があることを示すフラグが設定されている。例えば、データ構造体1020のCCW1とCCW2には、次のデータ構造体があることを示すフラグが設定され、CCW3にはこのフラグが設定されず、最後のデータ構造体であることを示す。

【0118】

そして、各サーバ#1～#3は、使用するI/Oカード1501または1502に対して動作を起動させる指令を送信し、メモリ空間に設定したデータ構造体1020～1022のアドレスを通知する。

【0119】

I/Oカード1501、1502は、サーバ#1～#nから起動の指令とともにメモリ空間のアドレスを受信すると、指定されたメモリ空間のデータ構造体を読み込んで、データ構造体1020～1022に記述されたI/Oアクセスを実行する。

【0120】

このようなコマンドチェーン処理を行うI/Oカード1501、1502を備えたブレードサーバシステムに、前記第1実施形態と同様のI/Oカード共有機構400と、I/Oプロセッサブレード600を適用する例を以下に示す。

【0121】

I/Oカード共有機構400とI/Oプロセッサブレード600によるI/Oカードの

10

20

30

40

50

共有処理のタイムチャートを図14に示す。なお、以下では、サーバ#1(サーバブレード10-1)がI/Oカード1501を使用する場合について説明する。

【0122】

まず、サーバ#1が、I/Oアクセスを行うI/Oカード1501に対して動作の起動を指令するPCIトランザクションを送信する(S31)。このPCIトランザクションは、サーバ#1がI/Oカード1501に割り当てたMMIOベースアドレスをMMIOベースアドレス462に設定し、動作の起動を指令するコマンドとデータ構造体1020のアドレス0xDをPCIトランザクション本体461内に設定する。

【0123】

I/Oカード共有機構400は、サーバ#1からのPCIトランザクションを受信するとPCIトランザクションの指令を解析し、動作の起動の指令はI/Oカード1501のコマンドレジスタ512への書き込み指令であるので、I/Oプロセッサブレード600のCPU602へ割り込みをかける(S32)。同時に、I/Oカード共有機構400は、I/Oプロセッサブレード600の所定のメモリ空間6031(図5参照)にPCIトランザクションの内容を書き込む(S33)。

【0124】

割り込みにより起動したCPU602は、前記第1実施形態の割り込み処理部620を起動する。そして割り込み処理部620が、メモリ空間に書き込まれたPCIトランザクションのデータ構造体1020のアドレス0xDを読み込んで、MMIOベースアドレス462の要求元453を取得する。次に、割り込み処理部620は、要求元453のサーバ#1のメモリ102に対して取得したアドレスからデータ構造体1020を読み込み、I/Oプロセッサブレード600の所定のメモリ空間にデータ構造体1020をコピーする(S34)。なお、このメモリ空間は図15で示すように、サーバ#1~#n毎に予め設定したデータ構造体を格納する領域であり、この例では、サーバ#1~#3のデータ構造体用メモリ空間として0xS、0xT、0xUが設定され、サーバ#1のデータ構造体1020は、図示のようにアドレス0xSから格納される。

【0125】

次に割り込み処理部620は、DMA等で使用するMMIOの処理を実施する。メモリ空間6031に書き込んだPCIトランザクションから、ヘッダ情報451の要求元453をMMIOベースアドレス462の未使用領域463にセットしたものを、目的のI/Oカード1501のアドレスレジスタ511に書き込んで、アドレス変換を実施しておく(S35)。

【0126】

この後、割り込み処理部620は、受信したPCIトランザクションのI/Oカードの動作の起動指令に基づいて、I/Oカード1501のコマンドレジスタに動作の起動指令を書き込み、コマンドチェーン制御部514にデータ構造体1020のアドレス0xSを通知する(S36)。

【0127】

I/Oカード1501は、割り込み処理部620の起動指令に基づいて起動し、受信したI/Oプロセッサブレード600のメモリ603のアドレス0xSからひとつのデータ構造体(CCW1)1020を読み込む(S37)。なお、I/Oカード共有機構400は、I/Oカード1501からI/Oプロセッサブレード600への通信については、そのまま転送するものとする。

【0128】

I/Oカード1501はこのデータ構造体1020の記述に応じてI/O操作を行う。例えば、読み込んだデータ構造体1020がDMAの場合、I/Oカード1501はアドレスレジスタ511にセットされたMMIOベースアドレスに対してDMA転送を実施する(S38)。I/OカードからのDMAによるPCIトランザクションには、MMIOベースアドレス462の上位ビットに設定された未使用領域463にサーバ識別子が埋め込まれている。

10

20

30

40

50

【 0 1 2 9 】

I/Oカード1501とサーバ#1の間に介在するI/Oカード共有機構400が、PCIトランザクションを受信すると上記図2のトランザクションデコーダ404でDMAによるPCIトランザクションであるか否かを判定する。

【 0 1 3 0 】

このトランザクションデコーダ404で行われるI/OカードからのPCIトランザクションがDMAであるか否かの判定は、上記MMIOベースアドレス462の未使用領域463の全ビットが0でなければ、サーバ識別子が埋め込まれていると判定して、DMAによるPCIトランザクションであると判定することができる。

【 0 1 3 1 】

そして、DMAのPCIトランザクションの場合、トランザクションデコーダ404は、MMIOベースアドレス462の未使用領域463の内容をヘッダ情報451の宛先452に設定し、スイッチ200で識別可能なサーバ#1～#nのアドレス情報に変換する。この後、トランザクションデコーダ404は未使用領域463の内容を消去するため、未使用領域463の全ビットに0をセットしてからこのPCIトランザクションを送信する(S39)。

【 0 1 3 2 】

この宛先452に基づいてスイッチ200は、DMAのPCIトランザクションを、宛先452に設定されたDMAの要求元のサーバに転送し、サーバ#1～#nに設定したMMIOに対して所定のアクセスが行われる。

【 0 1 3 3 】

データ構造体(CCW1)1020で指定されたI/O操作が完了すると、I/Oカード1501は、次のデータ構造体(CCW2)を指定されたI/Oプロセッサブレード600のメモリ603のアドレス0x5から読み込んで、上記と同様に実行する。

【 0 1 3 4 】

このように、コマンドチェーン処理を行うI/Oカード1501、1502の場合では、I/Oプロセッサブレード600のメモリ603のメモリ空間に、各サーバ#1～#3のデータ構造体1020～1022をコピーしておき、サーバ#1～#3に代わってI/Oプロセッサブレード600がI/Oカード1501、1502からの読み込み要求に応答する。

【 0 1 3 5 】

したがって、I/Oカード1501、1502の起動時に、I/Oアクセスを要求したサーバのデータ構造体1020～1022を格納したメモリ603上のアドレスをI/Oカードへ通知することで、コマンドチェーン処理を行うI/Oカードを複数のサーバ#1～#3で共有することが可能となる。

【 0 1 3 6 】

< 第3実施形態 >

図16は、第3の実施形態を示すブレードサーバシステムのブロック図を示す。本実施形態は、前記第1または第2実施形態のスイッチ200に、I/Oカード共有機構400を組み込んで一体にしたものである。

【 0 1 3 7 】

スイッチ200Aは、I/Oカード共有機構400を含んで構成され、前記第1または第2実施形態で述べたように動作する。I/Oカード共有機構400をスイッチ200Aに包含することで、ブレードサーバシステムに搭載するスロットの数を低減でき、筐体をコンパクトに構成することが可能となる。

【 0 1 3 8 】

< まとめ >

以上のように、本発明によれば、複数のサーバブレード10-1～nで一つのI/Oカードを共有してDMA転送を行う場合は、DMA転送開始時にはI/Oプロセッサブレード600がI/Oカードのアドレスレジスタ511にDMAを要求したサーバの識別子を

10

20

30

40

50

MMIOベースアドレス内に埋め込んでおき、I/Oカードからサーバへ向かうPCIトランザクションを中継するI/Oカード共有機構400は、MMIOベースアドレスに埋め込まれたサーバ識別子をヘッダ情報451の宛先452に付け替えることで、汎用バスのI/Oカードを共有することができる。

【0139】

そして、DMA転送開始後は、I/Oプロセッサブレード600はPCIトランザクションの転送へ介入せず、I/Oカード共有機構400のハードウェアがアドレス変換を行うため、前記従来例のようなソフトウェア処理によるオーバーヘッドを防いで、I/Oアクセスの性能低下を防ぎながら複数のサーバ間でI/Oカードの共有を実現することが可能となる。

10

【0140】

さらに、複数のサーバブレード10-1~nで共有するI/Oカードは、汎用のバスインターフェースで構成することができるので、上述のようなサーバ統合を行う際には、従来使用していたI/Oカードをそのまま利用することができるので、サーバ統合の際のコストを抑制することが可能となる。

【0141】

また、サーバ統合の際には、複数のサーバブレード10-1~nで一つのI/Oカードを共有できるので、I/Oカードの数を低減することができ、前記従来例のようなI/Oカードの増大を防いでコンパクトな筐体を用いることができる。

【0142】

また、I/Oカード共有設定テーブル610により、共有するI/Oカードと、ひとつのサーバブレードのみに占有させるI/Oカードを混在させることができ、ブレードサーバシステムの柔軟な構成が可能となる。

20

【0143】

なお、上記各実施形態では、汎用バスとしてPCI-EXPRESSを適用した例を示したが、PCIやPCI-X等の汎用バスを採用してもよい。

【0144】

また、上記各実施形態では、サーバブレード10-1~nとサーバ#1~#nが1対1で対応した例を示したが、サーバ#1~#nが仮想計算機の論理区画で構成されるときは、サーバ識別子を論理区画番号とすればよい。

30

【0145】

また、上記各実施形態では、I/Oカード共有機構400に直接I/Oプロセッサブレード600を接続した例を示したが、I/Oプロセッサブレード600をスイッチ200に接続しても良い。さらに、I/Oプロセッサブレード600に代わって、サーバブレードがI/Oプロセッサブレード600の処理を実行するようにしても良い。

【産業上の利用可能性】

【0146】

以上のように、本発明は複数のサーバとI/OカードまたはI/Oデバイスを汎用バスで接続するコンピュータシステムに適用することができ、特に、サーバ統合を行うブレードサーバシステムに適用することで、コストの低減と性能の確保実現できる。

40

【図面の簡単な説明】

【0147】

【図1】第1の実施形態を示すブレードサーバシステムのブロック図。

【図2】I/Oカード共有機構のブロック図。

【図3】PCIトランザクションの説明図。

【図4】サーバのメモリ空間を示す説明図。

【図5】I/Oプロセッサブレードのメモリ空間を示す説明図。

【図6】アドレス情報テーブルの一例を示す説明図。

【図7】I/Oカード共有設定テーブルの一例を示す説明図。

【図8】I/Oカード共有機構とI/Oプロセッサブレードの関係を示すブロック図。

50

【図9】 I/Oカード共有処理の流れを示す説明図。

【図10】 I/Oカード共有処理の流れを示すタイムチャート。

【図11】 サーバブレードを起動したときに実行される、アドレス情報の設定処理の一例を示すタイムチャートである。

【図12】 第2の実施形態を示すブレードサーバシステムのブロック図。

【図13】 第2の実施形態を示し、サーバのメモリ空間を示す説明図。

【図14】 第2の実施形態を示し、I/Oカード共有処理の流れを示すタイムチャート。

【図15】 第2の実施形態を示し、I/Oプロセッサブレードのメモリ空間を示す説明図

【図16】 第3の実施形態を示すブレードサーバシステムのブロック図。

10

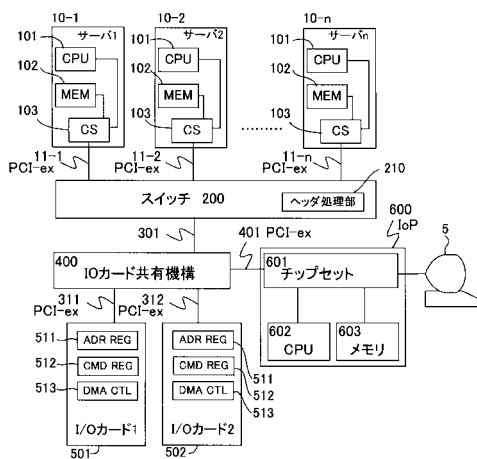
【符号の説明】

【0148】

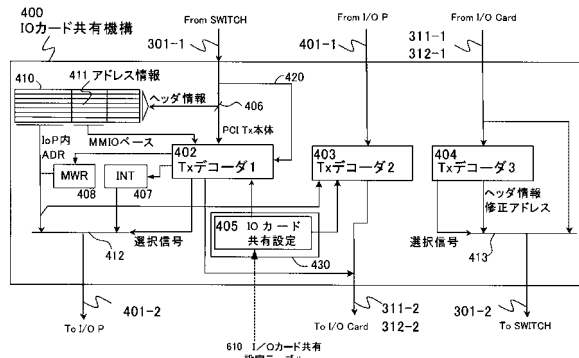
- 10 - 1 ~ 10 - n サーバブレード
- 200 スイッチ
- 400 I/Oカード共有機構
- 501、502 I/Oカード
- 600 I/Oプロセッサブレード
- 405 I/Oカード共有設定
- 411 アドレス情報テーブル
- 610 I/Oカード共有設定テーブル

20

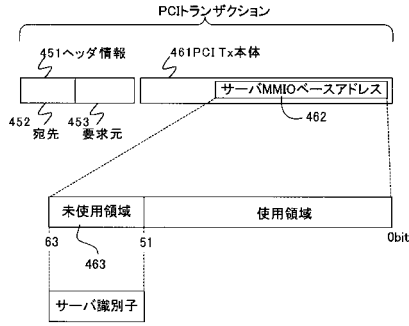
【図1】



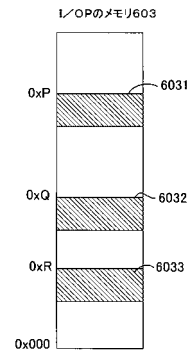
【図2】



【図3】



【図5】



【図6】

411 アドレス情報テーブル

HEADER	MMIO ADDR	IoP ADDR	OFFSET
Io1:SV1	0xA	0xP	0xX
Io1:SV2	0xB	0xQ	0xY
Io1:SV3	0xC	0xR	0xZ

4111

4112

4113

4114

【図7】

610 I/Oカード共有設定テーブル

	種別	サーバ1	サーバ2	サーバ3
IOカード1	SCSIカード	Share	Share	Share
IOカード2	NICカード	No Assigned	Dedicate	No Assigned

611

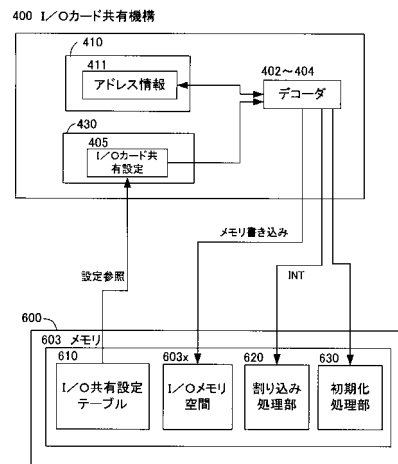
612

613

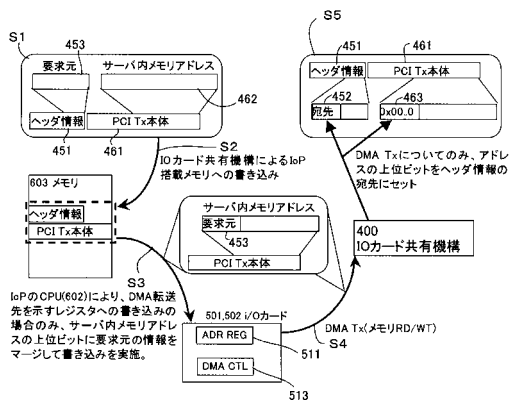
614

615

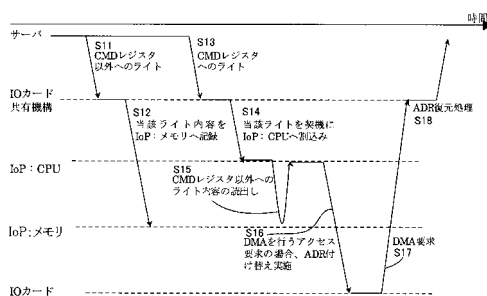
【図8】



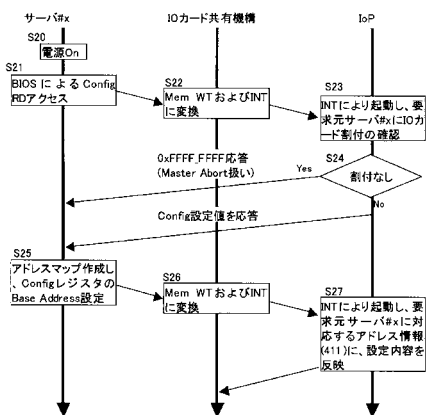
【図9】



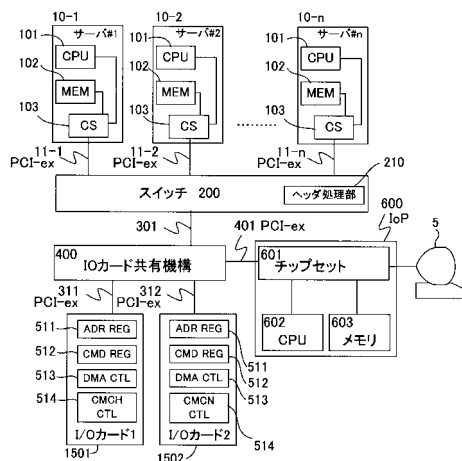
【図10】



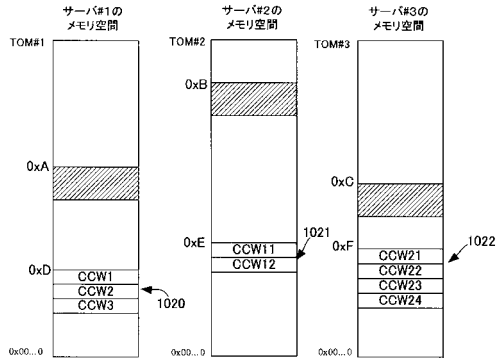
【図11】



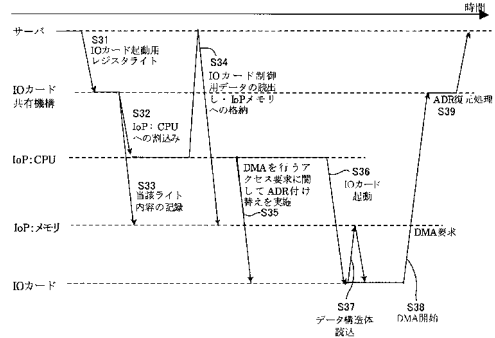
【図12】



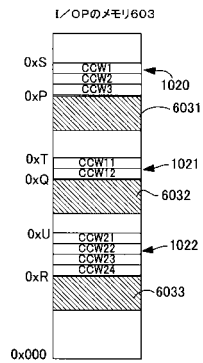
【図13】



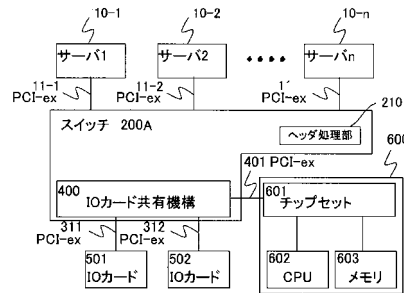
【図14】



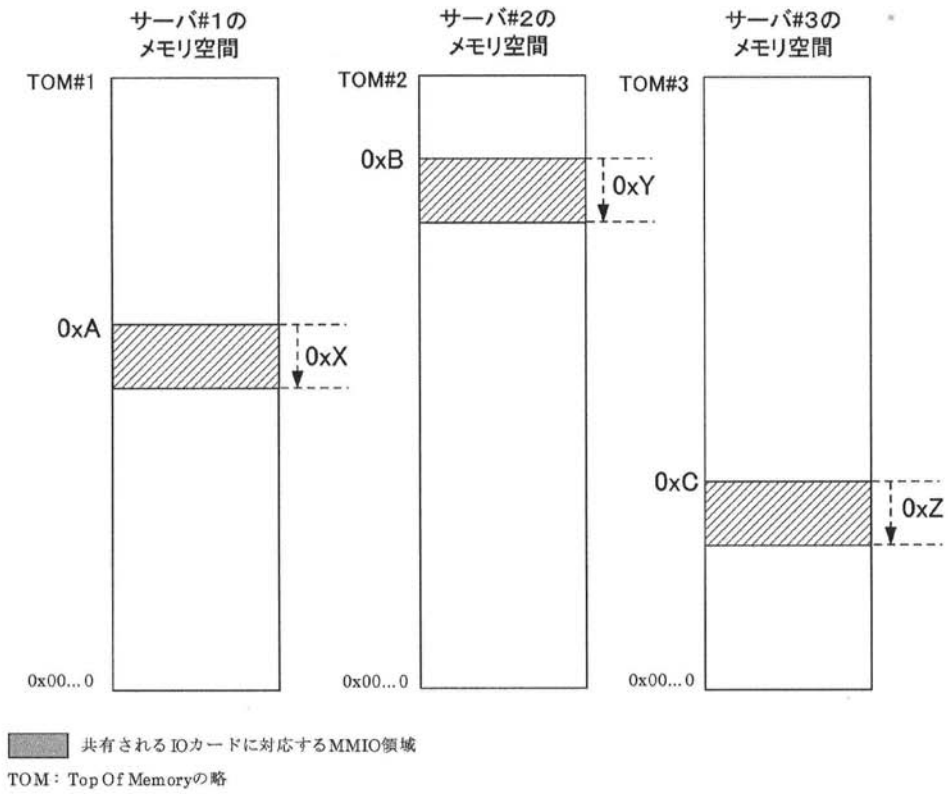
【図15】



【図16】



【 図 4 】



フロントページの続き

(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

審査官 坂東 博司

(56)参考文献 特表2007-502579(JP,A)

特開2004-54949(JP,A)

特開2005-209197(JP,A)

特表2007-507045(JP,A)

特表2006-506736(JP,A)

米国特許出願公開第2006/253619(US,A1)

(58)調査した分野(Int.Cl., DB名)

G06F 13/14