

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5365847号
(P5365847)

(45) 発行日 平成25年12月11日(2013.12.11)

(24) 登録日 平成25年9月20日(2013.9.20)

(51) Int. Cl. F 1
G 0 6 F 13/10 (2006.01) G O 6 F 13/10 3 3 O C
G 0 6 F 13/14 (2006.01) G O 6 F 13/14 3 2 O B

請求項の数 4 (全 15 頁)

<p>(21) 出願番号 特願2009-51861 (P2009-51861) (22) 出願日 平成21年3月5日(2009.3.5) (65) 公開番号 特開2010-205124 (P2010-205124A) (43) 公開日 平成22年9月16日(2010.9.16) 審査請求日 平成23年3月10日(2011.3.10)</p>	<p>(73) 特許権者 000004237 日本電気株式会社 東京都港区芝五丁目7番1号 (74) 代理人 100130029 弁理士 永井 道雄 (74) 代理人 100166338 弁理士 関口 正夫 (74) 代理人 100152054 弁理士 仲野 孝雅 (72) 発明者 小池 毅 東京都港区芝五丁目7番1号 日本電気株式会社社内 審査官 坂東 博司</p>
---	--

最終頁に続く

(54) 【発明の名称】 仮想化装置における物理デバイスのコンフィグレーション処理方法及びコンピュータシステム

(57) 【特許請求の範囲】

【請求項1】

オペレーティングシステムと、

前記オペレーティングシステムから、物理デバイスへのコンフィグレーションアクセスを捕捉し、前記オペレーティングシステムがアクセスした前記物理デバイスのデバイスアドレスから、前記物理デバイスが所属するパーティションを特定し、

前記物理デバイスが自パーティションに所属するデバイスであった場合には、前記オペレーティングシステムが要求した前記物理デバイスへのコンフィグレーションレジスタリードまたはライトをそのまま実行し、

前記物理デバイスが自パーティションに所属していないデバイスであった場合には、前記物理デバイスを特定するコンフィグレーションレジスタの内容を、実在しないダミーデバイスの情報に置き換えて前記オペレーティングシステムに通知する仮想化手段と、

を備え、

前記オペレーティングシステムは前記ダミーデバイスに対応したダミーデバイスドライバを有し、該ダミーデバイスドライバによって、前記自パーティションに所属しない物理デバイスを、名称の付与された有効なデバイスとして認識し、

前記ダミーデバイスドライバは、前記オペレーティングシステムから要求される全物理デバイスに共通な処理のうち、他パーティションに所属している物理デバイスの動作に擾乱を与えるものについては、その実行を拒絶すること、を特徴とするコンピュータシステム。

10

20

【請求項 2】

請求項 1 に記載のコンピュータシステムにおいて、前記仮想化手段は、前記物理デバイスが自パーティションに所属していないデバイスであった場合には、前記物理デバイスの製造元および種類を特定するコンフィグレーションレジスタの内容を、実在しないダミーデバイスの製造元または種類に置き換えて前記オペレーティングシステムに通知するコンピュータシステム。

【請求項 3】

仮想化装置における物理デバイスのコンフィグレーション処理方法において、

仮想化ソフトウェア上で動作するオペレーティングシステムから、物理デバイスへのコンフィグレーションアクセスを捕捉し、

前記オペレーティングシステムがアクセスした前記物理デバイスのデバイスアドレスから、前記物理デバイスが所属するパーティションを特定し、

前記物理デバイスが自パーティションに所属するデバイスであった場合には、前記オペレーティングシステムが要求した前記物理デバイスへのコンフィグレーションレジスタリードまたはライトを実行し、前記物理デバイスが自パーティションに所属していないデバイスであった場合には、前記物理デバイスを特定するコンフィグレーションレジスタの内容を、実在しないダミーデバイスの情報に置き換えて前記オペレーティングシステムに通知し、

前記ダミーデバイスに対応したダミーデバイスドライバによって、前記自パーティションに所属しない物理デバイスを、名称の付与された有効なデバイスとして前記オペレーティングシステムに認識させ、

前記オペレーティングシステムから要求される全物理デバイスに共通な処理のうち、他パーティションに所属している物理デバイスの動作に擾乱を与えるものについては、その実行を前記ダミーデバイスドライバにて拒絶するコンフィグレーション処理方法。

【請求項 4】

請求項 3 に記載のコンフィグレーション処理方法において、前記物理デバイスが自パーティションに所属していないデバイスであった場合には、前記物理デバイスの製造元および種類を特定するコンフィグレーションレジスタの内容を、実在しないダミーデバイスの製造元または種類に置き換えて前記オペレーティングシステムに通知するコンフィグレーション処理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、仮想化装置における物理デバイスのコンフィグレーション処理方法及びコンピュータシステムに係わる。

【背景技術】

【0002】

仮想化技術は、ハードウェアとソフトウェアの結びつきを抽象化し、OS（オペレーティングシステム）などのソフトウェアがハードウェアを直接制御するのではなく、抽象化レイヤを介して動作できるようにする技術である。

【0003】

仮想化の技術自体はメインフレームが全盛の時代から存在し、1つのハードウェア上で複数の基幹システムを動作させることで、高価なハードウェアを、より効率的に使用することを目的としていた。

【0004】

しかし現在は、テクノロジーの進化とオープンソースソフトウェアの浸透により、安価なサーバ上でも容易に仮想化環境が構築できるようになった。

【0005】

本願発明に関連する技術として、特許文献 1 には、PCI デバイスにおいて、メモリマップド I/O レジスタのアドレス空間内にコンフィグレーションレジスタのアドレス空間を

10

20

30

40

50

規定し、メモリマップド I/O レジスタのアドレス空間内に規定されたコンフィギュレーションレジスタのアドレスを指定して、コンフィギュレーションレジスタに対しアクセスすることの記載がある。

【0006】

また、特許文献 2 には、PCI デバイスと非 PCI デバイスとを実装するコンピュータシステムにおいて、非 PCI デバイスに対応させた仮想的な PCI デバイスの所定の識別情報を少なくとも記憶した記憶領域をシステム BIOS 側に用意しておき、システム BIOS が参照可能な媒体上に、仮想的な PCI デバイスのコンフィギュレーション空間のヘッダ領域に仮想的な PCI デバイスの所定の識別情報が格納されることの記載がある。

【先行技術文献】

【特許文献】

【0007】

【特許文献 1】特開 2002 - 318779 号公報

【特許文献 2】特開 2005 - 250975 号公報

【発明の概要】

【発明が解決しようとする課題】

【0008】

近年、仮想化の目的も、ハードウェアの使用効率の追求から、ある程度の使用効率の低下は許容しながらも、より機構が単純で、一度開発したソフトウェア資産を長期に渡って使用し続けられる仮想化のインフラストラクチャが、重要視されるようになってきた。

【0009】

本願発明は、前述の仮想化技術の変遷を鑑み、I/O デバイスの仮想化の分割粒度を PCI デバイス単位に制限することで、仮想化ソフトウェアの構造の単純化と、パフォーマンスの向上を目的とするものである。

【課題を解決するための手段】

【0010】

本発明に係わるコンピュータシステムは、オペレーティングシステムと、前記オペレーティングシステムから、物理デバイスへのコンフィギュレーションアクセスを捕捉し、前記オペレーティングシステムがアクセスした物理デバイスのデバイスアドレスから、前記物理デバイスが所属するパーティションを特定し、

前記物理デバイスが自パーティションに所属するデバイスであった場合には、前記オペレーティングシステムが要求した前記物理デバイスへのコンフィギュレーションレジスタリードまたはライトをそのまま実行し、

前記物理デバイスが自パーティションに所属していないデバイスであった場合には、前記物理デバイスを特定するコンフィギュレーションレジスタの内容を、実在しないダミーデバイスの情報に置き換えて前記オペレーティングシステムに通知する仮想化手段と、

を備え、

前記オペレーティングシステムは前記ダミーデバイスに対応したダミーデバイスドライバを有し、該ダミーデバイスドライバによって、前記自パーティションに所属しない物理デバイスを、名称の付与された有効なデバイスとして認識し、

前記ダミーデバイスドライバは、前記オペレーティングシステムから要求される全物理デバイスに共通な処理のうち、他パーティションに所属している物理デバイスの動作に擾乱を与えるものについては、その実行を拒絶すること、を特徴とするコンピュータシステムである。

【0011】

本発明に係わる物理デバイスのコンフィギュレーション処理方法は、仮想化装置における物理デバイスのコンフィギュレーション処理方法において、

仮想化ソフトウェア上で動作するオペレーティングシステムから、物理デバイスへのコンフィギュレーションアクセスを捕捉し、

前記オペレーティングシステムがアクセスした物理デバイスのデバイスアドレスから、

10

20

30

40

50

前記物理デバイスが所属するパーティションを特定し、

前記物理デバイスが自パーティションに所属するデバイスであった場合には、前記オペレーティングシステムが要求した前記物理デバイスへのコンフィグレーションレジスタリードまたはライトを実行し、前記物理デバイスが自パーティションに所属していないデバイスであった場合には、前記物理デバイスを特定するコンフィグレーションレジスタの内容を、実在しないダミーデバイスの情報に置き換えて前記オペレーティングシステムに通知し、

前記ダミーデバイスに対応したダミーデバイスドライバによって、前記自パーティションに所属しない物理デバイスを、名称の付与された有効なデバイスとして前記オペレーティングシステムに認識させ、

前記オペレーティングシステムから要求される全物理デバイスに共通な処理のうち、他パーティションに所属している物理デバイスの動作に擾乱を与えるものについては、その実行を前記ダミーデバイスドライバにて拒絶するコンフィグレーション処理方法である。

【発明の効果】

【0012】

本発明によれば、仮想化ソフトウェア上で動作するOSからは、仮想メモリマップドI/O空間と実メモリマップドI/O空間のPCIデバイスのレジスタ割り付けが同一に見えるため、OSが標準で装備するPCIデバイスドライバを、そのまま利用することができ、仮想化ソフトウェアのプログラム規模の縮小と、再利用性を向上させる効果がある。

【0013】

また、OSがPCIデバイスのI/Oレジスタにアクセスする際に、仮想メモリマップドI/O空間から実メモリマップドI/O空間へのアドレス変換を行う必要がないため、仮想化によるI/Oスループットの低下を防ぐ効果がある。

【図面の簡単な説明】

【0014】

【図1】図1は本発明の一実施形態に係わるコンピュータシステムの構成図である。

【図2】準仮想化方式の構成を示すブロック図である。

【図3】完全仮想化方式の構成を示すブロック図である。

【図4】ホストバスに複数のCPUが接続されているマルチプロセッサシステムを示す図である。

【図5】PCIバスモデルのシステムを、仮想化によって2つのパーティションに分割した場合を示す図である。

【図6】仮想化を行っていないシステムでコンフィグレーション処理を行った場合の、実アドレス空間のメモリマップドI/Oのレジスタ割り付けと、仮想化されたパーティション1およびパーティション2のメモリマップドI/Oのレジスタ割り付けを示す図である。

【図7】デバイス管理テーブルを示す図である。

【図8】PCI仕様2.0で規定されたコンフィグレーション・サイクルの発生メカニズムを表す図である。

【図9】割り付けが行われた、実アドレスとパーティション1および2のメモリマップドI/O空間のコンフィグレーション結果を表す図である。

【図10】本実施形態の仮想化ソフトウェアの処理フローを示すフローチャートである。

【発明を実施するための形態】

【0015】

以下、本発明の実施の形態について図面を用いて詳細に説明する。

【0016】

まず、本発明の実施形態の理解を容易にするために、I/Oを仮想化する技術について説明する。I/Oを仮想化する技術には、準仮想化方式と完全仮想化方式がある。

【0017】

図2は、準仮想化方式の構成を示すブロック図である。ここでは、2つのゲストOS(

10

20

30

40

50

オペレーティングシステム) 11-1、11-2とホスト12を有する場合について説明する。準仮想化方式では、ゲストOS 11-1、11-2に仮想化ソフトウェア専用のデバイスドライバを組み込む。このドライバをフロントエンドドライバ43-1、43-2という。ゲストOS 11-1、11-2が実行したI/O処理は、それぞれフロントエンドドライバ43-1、43-2から仮想化ソフトウェア20の内部通信機構60を經由して、ホストOS 12(サービス・パーティション、あるいは管理ドメインなどとも言う)のバックエンドドライバ44に通知される。ホストOS 12からは、全ての物理ハードウェア30が参照できるため、仮想化ソフトウェア20は、ゲストOS 11-1、11-2から依頼されたI/O処理を、ホストOS 12上の標準デバイスドライバ41で実行可能な形式に変換して、物理デバイス51-1、51-2にアクセスする。

10

【0018】

準仮想化方式は、物理デバイス51-1、51-2にアクセスする際に、ホストOS 12に付属の標準デバイスドライバ41を使用するため、システムの置き換え等によりハードウェアが更新された場合でも、フロントエンドドライバ43-1、43-2やバックエンドドライバ44を修正しなくて済むという利点がある。その反面、ゲストOS 11-1、11-2が実行したI/O処理は、必ず仮想化ソフトウェア20の内部通信機構60を經由してホストOS 12に集約されるため、処理のオーバーヘッドが発生するという問題がある。

【0019】

図3に示す完全仮想化方式は、準仮想化による処理オーバーヘッドとホストOS(オペレーティングシステム)12のようなサービス・パーティションを無くすよう改善した方式である。完全仮想化方式は、準仮想化方式のフロントエンドドライバ43-1、43-2の代わりに、仮想化ソフトウェア専用の仮想デバイスドライバ61-1、61-2を、OSの標準デバイスドライバ41-1、41-2の配下に組み込む方式である。

20

【0020】

仮想デバイスドライバ61-1、61-2は、その組み込み位置によって、仮想デバイス方式と直接割り当て方式に大別される。仮想デバイス方式は、仮想化ソフトウェア20が用意するデバイスドライバをOS 10-1、10-2内に組み込む方式である。このデバイスドライバを通じて、OS 10-1、OS 10-2が実行したI/O処理が仮想化ソフトウェア20に通知され、仮想化ソフトウェア20が適切な形でI/Oを再現する。もう一つの直接割り当て方式は、仮想デバイス方式のデバイスドライバ部分を仮想化ソフトウェア20内に取り込み、OS 10-1、10-2には直接仮想化されたハードウェアを見せる方式である。OS 10-1、10-2は標準デバイスドライバ41-1、41-2を使用して、仮想化ソフトウェア20内の仮想デバイスドライバ61-1、61-2にアクセスする。図3は直接割り当て方式を例示したものである。

30

【0021】

完全仮想化方式は、準仮想化方式と比べてゲストOS 11-1、11-2からホストOS 12への内部通信機構60を使用しない分、仮想化による処理オーバーヘッドは軽減される。但し、仮想デバイス方式・直接割り当て方式を問わず、完全仮想化方式は、OS 10-1、10-2に対応した仮想デバイスドライバ61-1、61-2を、ハードウェア

40

【0022】

仮想化技術において、物理デバイスと仮想デバイスの対応付けのために、専用の仮想デバイスドライバを必要とする理由に、オープンOSが採用しているPCIデバイス仕様には則ったコンフィグレーション処理がある。以下、物理ハードウェアとして図4に示すPCIバスモデルを使用して、オープンOSが採用しているコンフィグレーション処理について説明する。

【0023】

図4のシステムは、ホストバス5にCPU 1a~1dが接続されているマルチプロセッサシステムである。ホストバス5にはI/Oデバイスにアクセスするためのホストブリッ

50

ジ 2 が接続されている。本モデルは 2 つの P C I バス 6 a と 6 b を有する。P C I バス 6 a にはバス番号 0 (B = 0)、P C I バス 6 b にはバス番号 1 (B = 1) が付与されている。P C I バス 6 a と P C I バス 6 b は P C I - P C I ブリッジ 4 によって接続されている。P C I バス 6 a には 1 つの P C I デバイス 3 a が接続されており、P C I バス 6 b には 3 つの P C I デバイス 3 b ~ 3 d が接続されている。

【 0 0 2 4 】

P C I デバイスは、システムの起動時やホットプラグ時に、O S によるコンフィグレーション処理によってバス番号とデバイス番号の割り付けが決定される。

【 0 0 2 5 】

コンフィグレーション処理におけるバス番号の割り付け手順は以下の通りである。

10

(1) 処理 P 1

ホストブリッジの二次側バスにバス番号 0 を付与し、このバス番号を記憶する。ホストブリッジの二次側バスを検査し P C I - P C I ブリッジの存在を確認する。

(2) 処理 P 2

見つかった P C I - P C I ブリッジの二次側バスに、記憶してあったバス番号に 1 を加算したバス番号を付与し、このバス番号を記憶する。

(3) 処理 P 3

上記 (2) の処理 P 2 で見つけた P C I - P C I ブリッジの二次側バスを検査し、P C I - P C I ブリッジの配下に、更に P C I - P C I ブリッジが存在しないか確認する。新たな P C I - P C I ブリッジが見つかったら上記処理 P 2 の処理を行う。

20

(4) 処理 P 4

P C I - P C I ブリッジの配下に、新たな P C I - P C I ブリッジが見つからなくなったら、ホストブリッジに接続されている次の P C I - P C I ブリッジについて上記処理 P 2 を行う。

【 0 0 2 6 】

コンフィグレーション処理におけるデバイス番号の割り付け手順は以下の通りである。

(1) 処理 P 1 1

ホストブリッジまたは P C I - P C I ブリッジの二次側バスを検査する前に、デバイス番号を 0 に初期化し、このデバイス番号を記憶する。

(2) 処理 P 1 2

バスに接続されている P C I デバイスまたは P C I - P C I ブリッジを検査する。見つかった P C I デバイスまたは P C I - P C I ブリッジに、記憶してあったデバイス番号を割り当てて、デバイス番号に 1 を加算して記憶する。

30

(3) 処理 P 1 3

バスに接続されている P C I デバイスまたは P C I - P C I ブリッジを見つけるたびに、上記処理 P 1 2 を繰り返す。

(4) 処理 P 1 4

バスの配下に、新たな P C I デバイスまたは P C I - P C I ブリッジが見つからなくなったら、次の P C I - P C I ブリッジの二次側バスについて上記処理 P 1 1 を行う。

【 0 0 2 7 】

図 4 の括弧内の数字は、本手順に従って割り付けられた P C I デバイスまたは P C I - P C I ブリッジの、一次側のバス番号 (B) とデバイス番号 (D) を表す。

40

【 0 0 2 8 】

ここで、本 P C I バスモデルのシステムを、仮想化によって 2 つのパーティションに分割した場合を考える。図 5 は、P C I デバイス 3 a と 3 b をパーティション 1 に、P C I デバイス 3 c と 3 d をパーティション 2 に割り当てた構成である。

【 0 0 2 9 】

仮想化ソフトウェアは、O S からのコンフィグレーション処理に対して、自パーティションに所属しないデバイスは未実装状態に見せて O S に割り付けを行わせる。

【 0 0 3 0 】

50

図5のPCIバスモデルの括弧内の数字は、パーティショニングされた2つのOS 10 a、10 bのそれぞれから、前記コンフィグレーション処理に従ってバス番号とデバイス番号を付与した状態を表している。

【0031】

図5において、パーティション1に所属するPCIデバイス3 bと、パーティション2に所属するPCIデバイス3 cは、共にバス番号1 / デバイス番号0 (B = 1 / D = 0) を有する。その結果物理デバイスにこのバス番号 / デバイス番号を付与するとバス競合が発生する。

【0032】

この場合、バス番号 / デバイス番号のバス競合を回避するためには、仮想化ソフトウェアにて仮想デバイスと物理デバイスのバス番号 / デバイス番号の変換を行う必要がある。

【0033】

次に、図6を使用して、仮想化を行っていないシステムでコンフィグレーション処理を行った場合の、実アドレス空間のメモリマップドI/Oのレジスタ割り付けと、仮想化されたパーティション1およびパーティション2のメモリマップドI/Oのレジスタ割り付けを示す。

【0034】

オープンOSにおいては、前記コンフィグレーション処理で検出したPCIデバイスがレジスタ領域を必要とする場合、メモリマップドI/O空間の最上位アドレスから順番に、検出したPCIデバイスのレジスタを並べて配置するように動作する。

【0035】

同様の手順で、仮想化されたパーティションのOSがコンフィグレーション処理を行うと、自パーティションに所属していないデバイスは、仮想化ソフトウェアによって隠蔽されて未検出状態となるため、自パーティションの仮想メモリマップドI/O上にはマッピングされず、仮想メモリマップドI/Oと実メモリマップドI/Oのアドレス割り付けに差分が生じる。その結果、仮想メモリマップドI/Oと実メモリマップドI/Oを同一視してPCIデバイスにアクセスすると、I/Oレジスタのアドレス競合が発生する。

【0036】

仮想化によるPCIデバイスのI/Oレジスタのアドレス競合を回避するためには、仮想化ソフトウェアにて仮想メモリマップドI/Oと実メモリマップドI/Oのアドレス変換を行う必要があった。

【0037】

また、ゲストOSやPCIデバイスに付属するデバイスドライバの中には、仮想メモリマップドI/Oに対応していないデバイスドライバもあるため、これらについては、仮想化ソフトウェア専用のデバイスドライバをOSの配下に組み込む必要があった。

【0038】

次に、本実施形態について説明する。

【0039】

図1は本発明の一実施形態に係わるコンピュータシステムの構成図である。

【0040】

図1の実施形態のシステムの構成は、物理ハードウェア30を物理デバイス51 - 1、51 - 2の単位で、パーティション1とパーティション2の2つの仮想システムに分割したものである。

【0041】

パーティション1の物理デバイス51 - 1は、パーティション1に所属する物理CPU 50 - 1上で動作するOS (オペレーティングシステム) の制御下に置かれ、パーティション1のOS 10 - 1に付属する標準デバイスドライバ41 - 1を使用してアクセスされる。

【0042】

同様に、パーティション2の物理デバイス51 - 2は、パーティション2に所属する物

10

20

30

40

50

理CPU50-2上で動作するOSの制御下に置かれ、パーティション2のOS10-2に付属する標準デバイスドライバ41-2を使用してアクセスされる。

【0043】

仮想化ソフトウェア20は、物理CPU50-1、50-2に対応する仮想CPU40-1、40-2からのコンフィグレーションアクセスを捕捉し、例えば、自パーティションに所属していない物理デバイス51-2へ仮想CPU40-1がアクセスした場合に、物理デバイス51-2のベンダIDおよびデバイスIDを、実在しないダミーデバイスのIDに置き換えて仮想CPU40-1に通知する機能を有する。仮想化ソフトウェアは仮想化手段となる。

【0044】

ダミーデバイスドライバ42-1、42-2は、ダミーデバイスに対応したデバイスドライバであって、自パーティションに所属しない物理デバイスを、名称の付与された有効なデバイスとしてOSに認識させる機能を有する。例えば、ダミーデバイスドライバ42-1は、パーティション1に所属しない物理デバイス51-2を、名称の付与された有効なデバイスとしてOS10-1に認識させる。また、ダミーデバイスドライバ42-1、42-2は、OS10-1、10-2から要求される全PCIデバイスに共通な処理のうち、他パーティションに所属しているPCIデバイスの動作に擾乱を与えるものについては、その実行を拒絶する機能を有する。

【0045】

なお、本実施形態の構成では、物理CPU50-1、50-2を仮想化する手段として仮想CPU40-1、40-2を定義したが、CPUに関する仮想化技術については、ここでは説明を省略する。

【0046】

物理デバイス51-1、51-2と各パーティションの対応付けは、図7に示すデバイス管理テーブルにて行う。デバイス管理テーブルは、仮想化システムを構築する際に予め決定しておくべきテーブルであって、仮想化ソフトウェア20内に保持される。

【0047】

デバイス管理テーブルは、物理デバイス51-1、51-2のバス番号、デバイス番号、機能番号を入力に、当該物理デバイス51-1、51-2がどのパーティションに所属するPCIデバイスであるのかを規定する。

【0048】

本実施形態では、説明を容易化するため、物理デバイス51-1、51-2に対応した実デバイスのベンダID/デバイスIDと、これに対応するダミーデバイスのベンダID/デバイスIDを、デバイス管理テーブル内に保持する。但し、実デバイスのベンダID/デバイスIDは、物理デバイス51-1、51-2のコンフィグレーションレジスタをリードすることで、またダミーデバイスのベンダID/デバイスIDは、固定値を使用することで代用可能なため、これらの情報は必須の構成要素ではない。

【0049】

本実施形態の動作について、図1と図4および図7から図9を使用して説明する。

【0050】

図8はPCI仕様2.0で規定されたコンフィグレーション・サイクルの発生メカニズムを表す図である。AT互換機系のシステムでは、PCIデバイス3a~3dおよびPCI-PCIブリッジ4のコンフィグレーションレジスタにアクセスするためのI/Oポートを、ホストブリッジ2に内蔵する。ホストブリッジ2のコンフィグレーション・サイクル発生回路内には、CONFIG_ADDRレジスタと呼ばれる4バイト幅のレジスタが存在し、CPU1a~1dがI/OアドレスCF8hに対してI/Oリードまたはライトを実行すると、これをCONFIG_ADDRレジスタに対するリード/ライトとして扱う。CONFIG_ADDRレジスタは、ホストブリッジ2がアクセスすべきPCIバスのバス番号、デバイス番号、機能番号、およびPCIデバイスのレジスタアドレスを指し示すためのものであって、その内容はCONFIG_ADDRレジスタにライトした最後

10

20

30

40

50

の値を保持する。

【0051】

ホストブリッジ2は、もう1つCONFIG__DATAレジスタと呼ばれる4バイト幅のレジスタを内蔵する。CONFIG__DATAレジスタは、CPU1a~1dがI/OアドレスCFChからCFhまでの4バイト以内のI/Oリードまたはライトを実行すると、これをCONFIG__DATAレジスタに対するリード/ライトとして扱い、CONFIG__ADDRレジスタで指定したPCIデバイスに対して、コンフィグレーション・サイクルを発生させる。

【0052】

本実施形態の仮想化ソフトウェア20は、図1における仮想CPU40-1、40-2からのCONFIG__DATAレジスタへのアクセスを捕捉し、このときのCONFIG__ADDRレジスタの値から、各パーティションのOS10-1、10-2が、どの物理デバイスのコンフィグレーションレジスタにアクセスしようとしたのかを判断する。

10

【0053】

次に、仮想化ソフトウェア20は、図7に示すデバイス管理テーブルを索引し、物理デバイスがどのパーティションに所属するデバイスなのかを判断する。

【0054】

仮想化ソフトウェア20は、仮想CPUが発行したコンフィグレーションアクセスが、自パーティションに所属する物理デバイスに対するものであった場合は、そのまま仮想CPUが発行したコンフィグレーションアクセスを実行する。例えば、仮想CPU40-1が発行したコンフィグレーションアクセスが、自パーティションに所属する物理デバイス50-1に対するものであった場合は、そのまま仮想CPU40-1が発行したコンフィグレーションアクセスを実行する。

20

【0055】

また、仮想CPUが発行したコンフィグレーションアクセスが、自パーティションに所属しない物理デバイスに対するものであった場合は、CONFIG__ADDRレジスタの内容から、アクセス先の物理デバイスのレジスタアドレスを判定し、ベンダIDあるいはデバイスIDへのリード要求の場合は、実在しないダミーデバイスのベンダIDおよびデバイスIDに差し替えて、リードしたデータを仮想CPUに通知する。例えば、仮想CPU40-1が発行したコンフィグレーションアクセスが、自パーティションに所属しない物理デバイス51-2に対するものであった場合は、CONFIG__ADDRレジスタの内容から、アクセス先の物理デバイス51-2のレジスタアドレスを判定し、ベンダIDあるいはデバイスIDへのリード要求の場合は、実在しないダミーデバイスのベンダIDおよびデバイスIDに差し替えて、リードしたデータを仮想CPU40-1に通知する。

30

【0056】

本実施形態では、ダミーデバイスが使用するベンダIDとして5853hを、またデバイスIDとして9999h(PCIデバイス用)と8888h(PCI-PCIブリッジ用)を定義する。

【0057】

図7に、図4で示したPCIバスモデルの、全PCIデバイス3a~3dとPCI-PCIブリッジ4のベンダID/デバイスIDについて、実デバイスとダミーデバイスのそれぞれのIDを記述する。

40

【0058】

本実施形態の仮想化ソフトウェア20上で動作するOS10-1、10-2は、前述の機構によって、自パーティションに所属しない物理デバイス51-2、51-1を、物理的/論理的に実装されているものと認識しながらも、使用することのない未知のデバイス(ダミーデバイス)として扱う。

【0059】

その結果、OS10-1、10-2の起動時に実施されるコンフィグレーション処理においても、何れのパーティションのOSからも、同じPCIバス構成が参照できるため、

50

仮想化を行わなかったときと同じ図 4 に示すバス番号 / デバイス番号が、各物理デバイス 5 1 - 1、5 1 - 2 に付与されることとなる。

【 0 0 6 0 】

図 9 は、割り付けが行われた、実アドレスとパーティション 1 および 2 のメモリマップド I / O 空間のコンフィグレーション結果を表す。

【 0 0 6 1 】

図 7 および図 9 において、P C I デバイス # 1 はベンダ I D : 1 0 1 4 h / デバイス I D : 1 0 8 0 h を有する。本デバイスはパーティション 1 に所属しているため、パーティション 1 には実デバイスと同じベンダ I D / デバイス I D を持つ P C I デバイスが、実アドレス空間と同じメモリマップド I / O の位置に割り付けられる。

10

【 0 0 6 2 】

一方、パーティション 2 からは、P C I デバイス # 1 はベンダ I D : 5 8 5 3 h / デバイス I D : 9 9 9 9 h を有する、実在しない未知のデバイス (ダミーデバイス) として認識されるが、本ダミーデバイスは P C I デバイス # 1 と同じベースアドレスレジスタのサイズを有しているため、実アドレス空間と同じメモリマップド I / O の位置に、P C I デバイス # 1 と同じ I / O レジスタサイズを持つダミーデバイスとして割り付けられる。

【 0 0 6 3 】

P C I デバイス # 2 から # 4 についても、同様の手順でメモリマップド I / O への割り付けが行われる。

【 0 0 6 4 】

20

また、本実施形態では、P C I - P C I ブリッジを何れのパーティションにも所属していない共通のデバイスとして定義している。このようなデバイスは、図 7 に示すデバイス管理テーブル上はパーティション番号 0 を付与され、ベンダ I D : 5 8 5 3 h / デバイス I D : 8 8 8 8 h を有するダミーブリッジとして、パーティション 1 および 2 の O S に認識される。

【 0 0 6 5 】

図 9 のメモリマップド I / O 空間のデバイス割り付け状態を参照すると、パーティション 1 および 2 からは、自パーティションに所属する P C I デバイスは、実アドレス空間のメモリマップド I / O と全く同じアドレスに割り付けられていることが分かる。

【 0 0 6 6 】

30

また、自パーティションに所属しない P C I デバイスあるいは P C I - P C I ブリッジは、ダミーデバイスまたはダミーブリッジとして割り付けられていることが分かる。

【 0 0 6 7 】

本機構によって、O S は標準のデバイスドライバを使用して、自パーティションに所属する P C I デバイスを制御することが可能となる。

【 0 0 6 8 】

なお、P C I 仕様では、図 8 に示した P C I デバイス用のコンフィグレーションレジスタ配置 (タイプ 0 0 h) 以外に、P C I - P C I ブリッジ用のレジスタ配置 (タイプ 0 1 h) や、更に後位の P C I 仕様に準拠したレジスタ配置などがある。

【 0 0 6 9 】

40

これらの中には、ベンダ I D やデバイス I D に類似したサブシステム・ベンダ I D / サブシステム I D、あるいはベースアドレスレジスタに類似した拡張 R O M ベースアドレスレジスタなどがあるが、これらについても本発明は同様の効果を奏するため、本実施形態では説明を省略する。

【 0 0 7 0 】

次に、図 1 に示すダミーデバイスドライバ (またはダミーブリッジドライバ) の役割について説明する。

【 0 0 7 1 】

O S が仮想化ソフトウェアを使用して、自パーティションに所属しない P C I デバイスを参照した場合、その P C I デバイスは未知のデバイスとして O S に認識されるため、自

50

パーティションに所属しないP C Iデバイスの機能をO Sが使用することはない。

【0072】

しかし、全P C Iデバイスに共通のO S操作については、O Sはこれを実行する可能性がある。特に、このような処理が、他パーティションに所属しているP C Iデバイスに対して実行された場合には、他パーティションのシステム運用に影響を与える場合もある。例えばP C Iデバイスのホットプラグに先立って実施される、デバイスの切り離し可否の確認がこれに該当する。

【0073】

本実施形態のデバイスドライバは、全P C Iデバイスに共通のO S操作を防止するため、自パーティションに所属しないP C Iデバイスを、名称の付与された有効なデバイスとしてO Sに認識させ、他パーティションのシステム運用に影響を与えるようなO S操作が、ダミーデバイスドライバに要求された場合には、その実行を拒絶する役割を果たす。

【0074】

図10に、本実施形態の仮想化ソフトウェアの処理フローを示す。

【0075】

仮想化ソフトウェアは、ステップS1で、仮想CPUからのCONF I G _ D A T Aレジスタへのアクセスを捕捉すると、ステップS2で、CONF I G _ A D D Rレジスタの設定値から、アクセスすべきP C Iデバイスのバス番号、デバイス番号、機能番号を取得し、ステップS3で本情報を基にデバイス管理テーブルを索引し、P C Iデバイスの所属するパーティション番号を特定する。

【0076】

ステップS4で、もしP C Iデバイスが自パーティションに所属するデバイスであった場合は、仮想CPUからのP C Iデバイスへのレジスタリードまたはレジスタライトをそのまま実行する。ステップS4にて、P C Iデバイスが自パーティションに所属しないデバイスであった場合は、ステップS6でCONF I G _ A D D Rレジスタの設定値から、アクセスされたコンフィグレーションレジスタのアドレスを特定し、ベンダIDまたはデバイスIDへのリードかどうかを判断する。もしこれに該当する場合は、仮想CPUに、デバイス管理テーブルで設定されているダミーデバイスのベンダID / デバイスIDを返却する。

【0077】

ステップS6において、それ以外のコンフィグレーションレジスタへのアクセスの場合は、ステップS8においてP C Iデバイスがシステムの起動時やホットプラグ時に既に初期化されていたかどうかを判断し、初期化済みの場合は、ステップS9においてP C Iデバイスの動作に擾乱を与えないよう、仮想CPUからのコンフィグレーションレジスタのリードのみを実行し、そうでない場合は、ステップS10においてP C Iデバイスのコンフィグレーションレジスタへのリードおよびライトを実行する。

【0078】

なお、ステップS8、S9、S10については、コンフィグレーションレジスタの種類によって、個々に異なるレジスタアクセスの制御を必要とするが、本発明の目的は仮想マップドI / O空間と実マップドI / O空間の、P C Iデバイスのレジスタ割り付けを一致させることにあるため、本実施形態の説明では簡略化したフローを示す。

【0079】

以上説明した、本実施形態では、サーバ仮想化(バーチャル・マシン)技術における、P C Iデバイスのコンフィグレーション処理において、自パーティションに所属しないP C Iデバイスに対して、O Sがコンフィグレーションアクセスした場合に、P C IデバイスのベンダIDおよびデバイスIDを、実在しないダミーデバイスのIDに置き換えてO Sに通知することによって、O Sからは自パーティションに所属していない物理デバイスの種類を隠蔽しながらも、あたかもダミーデバイスが実装されているかのようにコンフィグレーションさせることで、仮想マップドI / O空間と実マップドI / O空間のP C Iデバイスのレジスタ割り付けを一致させる機能を提供する。

【産業上の利用可能性】

【0080】

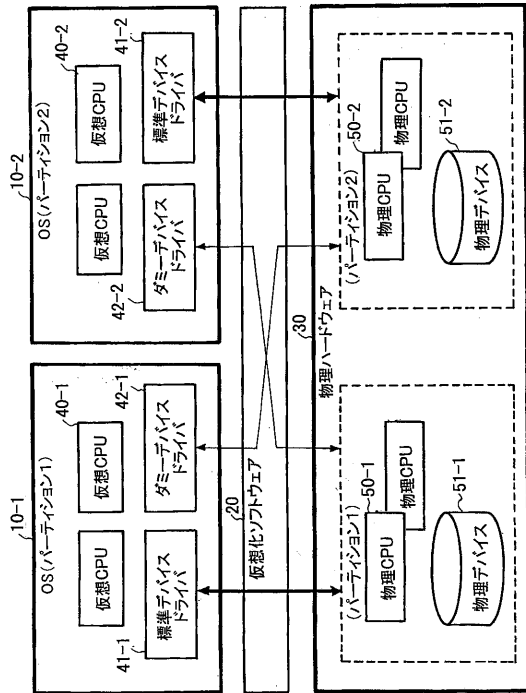
本発明は、P C I デバイスを含むコンピュータシステムに適用される。

【符号の説明】

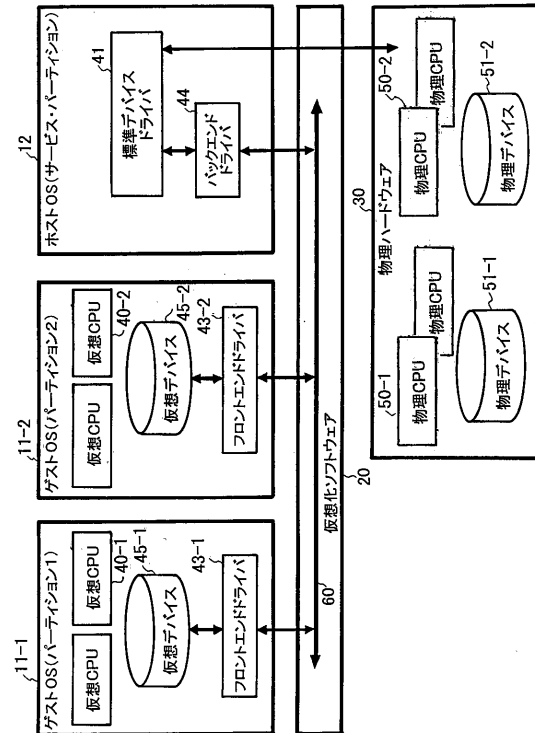
【0081】

- 10 - 1、10 - 2 OS
- 20 仮想化ソフトウェア
- 30 物理ハードウェア
- 41 - 1、41 - 2 標準デバイスドライバ
- 42 - 1、42 - 2 ダミーデバイスドライバ
- 50 - 1、50 - 2 物理CPU
- 51 - 1、51 - 2 物理デバイス

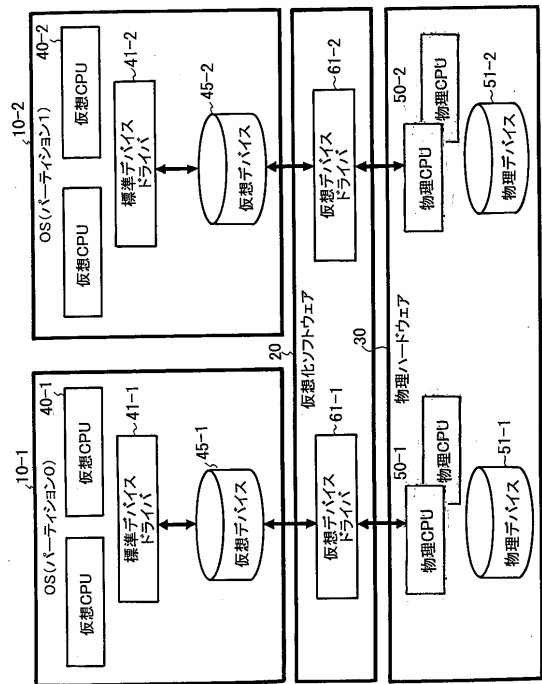
【図1】



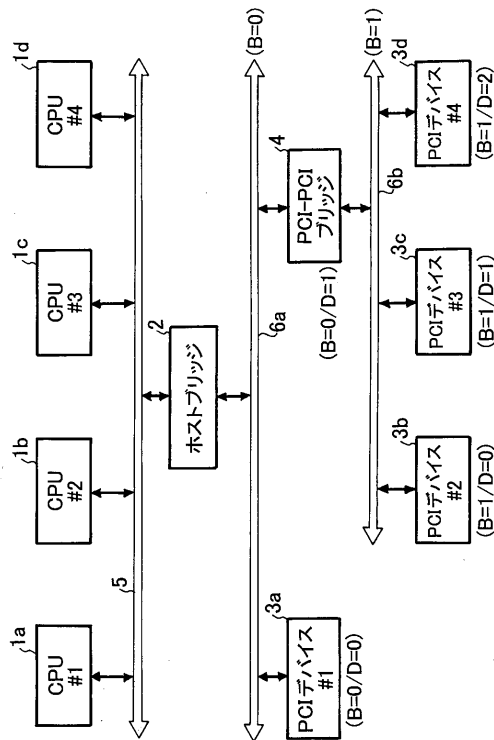
【図2】



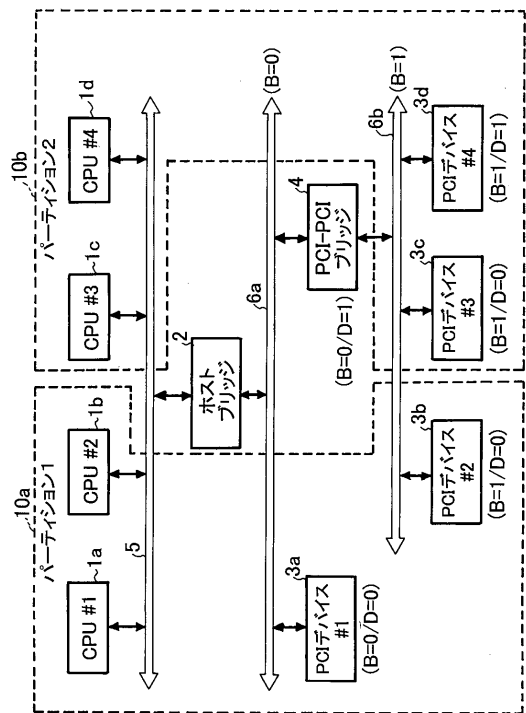
【 図 3 】



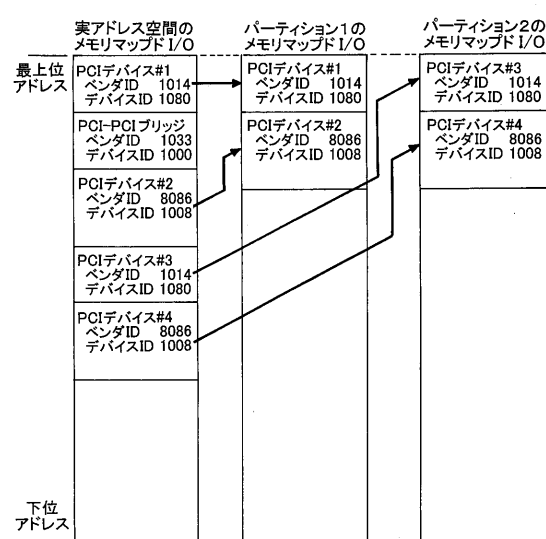
【 図 4 】



【 図 5 】



【 図 6 】

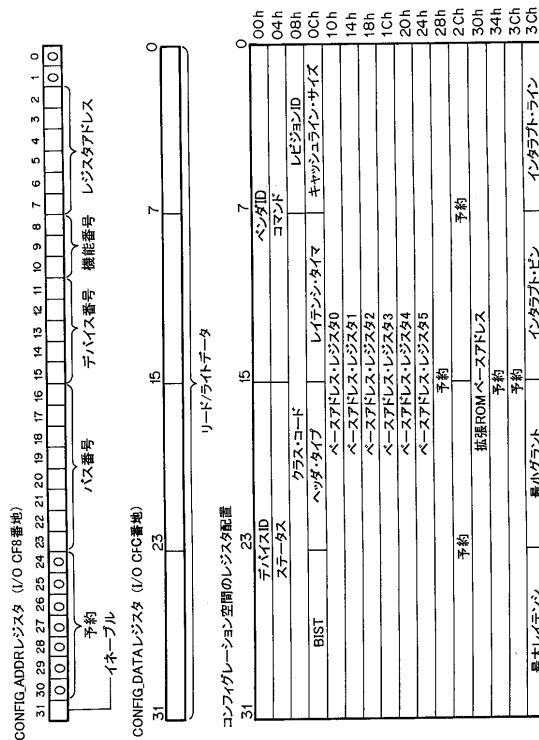


【 図 7 】

デバイス管理テーブル

CONFIG_ADDRの設定値		実デバイス		ダミーデバイス		
バス番号	デバイス番号	機能番号	ベンダID	デバイスID	ベンダID	デバイスID
00	00	00	1014	1080	5853	9999
00	01	00	1033	1000	5853	8888
01	00	00	8086	1008	5853	9999
01	01	00	1014	1080	5853	9999
01	02	00	8086	1008	5853	9999

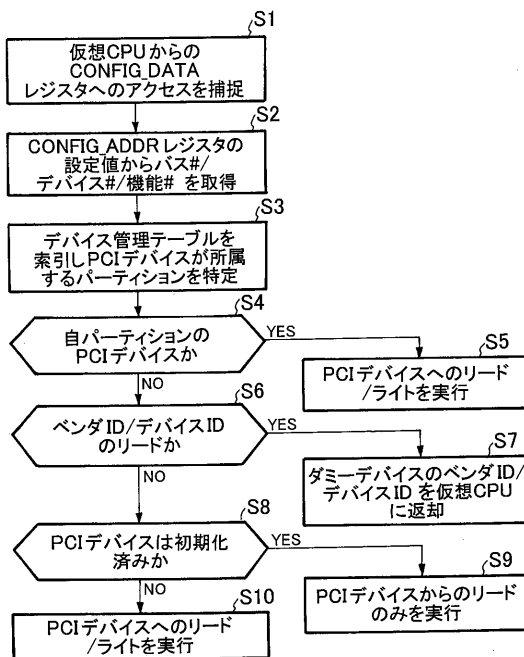
【 図 8 】



【 図 9 】

	実アドレス空間のメモリマップド I/O	パーティション1のメモリマップド I/O	パーティション2のメモリマップド I/O
最上位アドレス	PCIデバイス#1 ベンダID 1014 デバイスID 1080	PCIデバイス#1 ベンダID 1014 デバイスID 1080	ダミーデバイス ベンダID 5853 デバイスID 9999
	PCI-PCIブリッジ ベンダID 1033 デバイスID 1000	ダミーブリッジ ベンダID 5853 デバイスID 8888	ダミーブリッジ ベンダID 5853 デバイスID 8888
	PCIデバイス#2 ベンダID 8086 デバイスID 1008	PCIデバイス#2 ベンダID 8086 デバイスID 1008	ダミーデバイス ベンダID 5853 デバイスID 9999
	PCIデバイス#3 ベンダID 1014 デバイスID 1080	ダミーデバイス ベンダID 5853 デバイスID 9999	PCIデバイス#3 ベンダID 1014 デバイスID 1080
下位アドレス	PCIデバイス#4 ベンダID 8086 デバイスID 1008	ダミーデバイス ベンダID 5853 デバイスID 9999	PCIデバイス#4 ベンダID 8086 デバイスID 1008

【 図 10 】



フロントページの続き

- (56)参考文献 米国特許出願公開第2008/0162865(US,A1)
米国特許出願公開第2007/0044108(US,A1)
特開平10-275223(JP,A)
特開平11-161591(JP,A)
特開平11-143806(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/10

G06F 13/14