

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5116497号
(P5116497)

(45) 発行日 平成25年1月9日(2013.1.9)

(24) 登録日 平成24年10月26日(2012.10.26)

(51) Int. Cl. F I
G06F 13/14 (2006.01) G06F 13/14 310D
G06F 13/10 (2006.01) G06F 13/10 330C

請求項の数 12 (全 26 頁)

(21) 出願番号	特願2008-20923 (P2008-20923)	(73) 特許権者	000005108
(22) 出願日	平成20年1月31日 (2008.1.31)		株式会社日立製作所
(65) 公開番号	特開2009-181418 (P2009-181418A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成21年8月13日 (2009.8.13)	(74) 代理人	100114236
審査請求日	平成22年10月13日 (2010.10.13)		弁理士 藤井 正弘
		(74) 代理人	100075513
			弁理士 後藤 政喜
		(74) 代理人	100120260
			弁理士 飯田 雅昭
		(72) 発明者	沖津 潤
			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】 情報処理システム、I/Oスイッチ及びI/Oパスの交替処理方法

(57) 【特許請求の範囲】

【請求項1】

プロセッサとメモリを備えた物理サーバと、
 前記物理サーバの計算機資源を仮想化して仮想サーバを実行する仮想化部と、
 前記仮想サーバを実行する物理サーバを複数備え、これら物理サーバと1つ以上のI/Oデバイスを接続するI/Oスイッチと、を備えて、前記仮想化部が前記仮想サーバのマイグレーションを行う情報処理システムであって、
 前記I/Oスイッチは、
 前記I/Oデバイスから仮想サーバへのトランザクション発行の抑止を指示する抑止指示情報を格納するレジスタと、
 前記レジスタに前記抑止指示情報が格納されたときに、前記I/Oデバイスからマイグレーションを行う仮想サーバへのトランザクションの発行を抑止し、前記トランザクションの抑止前に前記I/Oデバイスから発行されたトランザクションの完了を保証するトランザクション抑止制御部と、
 前記仮想サーバのメモリアドレスと、当該仮想サーバを実行する物理サーバのメモリ上のアドレスとの対応関係を保持するアドレス変換部を備えて、前記仮想サーバのメモリアドレスを前記物理サーバのメモリ上のアドレスに変換する仮想化アシスト部と、
 当該I/Oスイッチの構成を管理するスイッチ管理部と、
 前記I/Oデバイスから該物理サーバへのトランザクションを保持する第1のバッファと、

10

20

前記物理サーバから前記 I / O デバイスへのトランザクションを保持する第 2 のバッファと、
を備え、

前記仮想化部は、

前記物理サーバ上の仮想サーバと他の仮想サーバに割り付けられた I / O デバイスの対応を管理する I / O デバイス管理部と、

前記 I / O デバイス管理部からマイグレーション対象の仮想サーバに割り付けられた I / O デバイスを特定し、当該 I / O デバイスに対応する前記 I / O スイッチの前記レジスタに対して、前記抑止指示情報を設定および解除するトランザクション指示部と、

前記スイッチ管理部に対して当該 I / O スイッチの構成変更を指示し、前記仮想化アシスト部で管理されるアドレス変換部に対して物理サーバの前記アドレスの変更を指令する構成変更指示部と、を備え、

前記トランザクション指示部は、マイグレーションの開始時に前記抑止指示情報を前記レジスタに設定して、前記 I / O デバイスからマイグレーション対象の仮想サーバへのトランザクション発行を抑止し、前記マイグレーションが完了すると前記レジスタに設定した抑止指示情報を解除して前記レジスタに設定して前記 I / O デバイスからマイグレーションが完了した前記仮想サーバへのトランザクション発行を許可し、

前記トランザクション抑止制御部は、

前記第 1 のバッファからのトランザクションの発行を抑止するトランザクション抑止部と、

レスポンス付きトランザクションを生成して物理サーバに発行するレスポンス付きトランザクション発行部と、

前記第 2 のバッファを監視し、前記レスポンス付きトランザクションの完了を確認する応答確認部と、

前記応答確認部が前記レスポンス付きトランザクションの完了を確認したときに、トランザクション抑止の完了を前記仮想化部に通知する完了通知部と、
を含むことを特徴とする情報処理システム。

【請求項 2】

前記情報処理システムは、

前記スイッチ管理部を制御する I / O マネージャをさらに有し、

前記 I / O マネージャは、前記仮想化部と通信を行う設定インタフェースを備え、

前記仮想化部が当該設定インタフェースを介して前記 I / O マネージャに対して構成変更を指示すると、前記 I / O マネージャが前記スイッチ管理部を制御して I / O スイッチの構成を変更することを特徴とする請求項 1 に記載の情報処理システム。

【請求項 3】

前記レジスタは、

前記抑止指示情報を格納するフィールドと、アドレス情報を格納するフィールドと、を有することを特徴とする請求項 1 に記載の情報処理システム。

【請求項 4】

前記アドレス情報を格納するフィールドは、前記マイグレーション対象の仮想サーバのメモリのアドレスを格納することを特徴とする請求項 3 に記載の情報処理システム。

【請求項 5】

前記情報処理システムは、

前記物理サーバを管理するサーバマネージャをさらに有し、

前記サーバマネージャは、

前記仮想化部に対して、前記仮想サーバのマイグレーションの開始を指示する開始指示部を備えたことを特徴とする請求項 1 に記載の情報処理システム。

【請求項 6】

前記レスポンス付きトランザクションは、先行するトランザクションを追い抜かないことを特徴とする請求項 1 に記載の情報処理システム。

10

20

30

40

50

【請求項 7】

1つ以上の物理サーバと1枚以上のI/Oデバイスを接続するI/Oスイッチにおいて、

前記I/Oスイッチは、
前記物理サーバを接続する1つ以上の第1のポートと、前記I/Oデバイスを接続する1つ以上の第2のポートを備え、
前記第2のポートは、
前記I/Oデバイスからのトランザクション発行の抑止を指示する抑止指示情報を保持するレジスタと、
当該第2のポートに接続する前記I/Oデバイスから物理サーバへのトランザクションを保持する第1のバッファと、
前記物理サーバから該第2のポートに接続するI/Oデバイスへのトランザクションを保持する第2のバッファと、
前記レジスタに前記抑止指示情報が設定されたことを契機に前記第1のバッファからのトランザクション発行を抑止するトランザクション抑止部と、
該トランザクション抑止部によるトランザクション発行の抑止後に、前記抑止指示情報に基づきレスポンス付きトランザクション要求を発行し、前記レスポンス付きトランザクションの完了を確認する滞留トランザクション完了確認部と、を備え、
前記滞留トランザクション完了確認部は、
前記レスポンス付きトランザクション要求を生成して前記物理サーバに発行するレスポンス付きトランザクション発行部と、
前記第2のバッファを監視し、前記レスポンス付きトランザクションの完了を確認する応答確認部と、
を含んで、前記I/Oデバイスから前記物理サーバへ発行されたトランザクションの完了を保証することを特徴とするI/Oスイッチ。

【請求項 8】

前記レジスタは、
前記抑止指示情報を格納するフィールドと、
アドレス情報を保持するフィールドと、
を有することを特徴とする請求項7に記載のI/Oスイッチ。

【請求項 9】

前記物理サーバは、
該物理サーバ上で稼動する1つ以上の仮想サーバと、
該仮想サーバに割当てられたI/Oデバイスの対応を管理する仮想化部と、を備え、
該仮想化部が、前記レジスタに対して前記抑止指示情報を設定することを特徴とする請求項7に記載のI/Oスイッチ。

【請求項 10】

前記第2のポートは、
前記レジスタの前記抑止指示情報が解除されたのを契機に、前記I/Oデバイスから該物理サーバへのトランザクションを保持する第1のバッファからのトランザクション発行を再開するトランザクション再開部を含むことを特徴とする請求項7に記載のI/Oスイッチ。

【請求項 11】

前記レスポンス付きトランザクションは、先行するトランザクションを追い抜かないことを特徴とする請求項7に記載のI/Oスイッチ。

【請求項 12】

1つ以上のサーバと1枚以上のI/Oデバイスを1つ以上のI/Oスイッチで接続する情報処理システムにおいて、前記I/Oスイッチに設定したI/Oパスを切り替えるI/Oパスの交替処理方法であって、
前記I/Oスイッチは、

10

20

30

40

50

前記 I / O デバイスからのトランザクション発行の抑止を指示する抑止指示情報を保持するレジスタと、

前記 I / O デバイスから前記サーバへのトランザクションを抑止し、前記トランザクションの抑止前に前記 I / O デバイスから発行されたトランザクションの完了を通知するトランザクション抑止制御部と、

前記 I / O スイッチの構成制御を管理するスイッチ管理部と、
を備え、

前記サーバを停止するステップと、

前記レジスタに対して前記抑止指示情報を設定するステップと、

前記 I / O デバイスから該物理サーバへのトランザクションを保持する第 1 のバッファからのトランザクションの発行を抑止するステップと、

レスポンス付きトランザクションを生成して物理サーバに発行するステップと、

前記物理サーバから前記 I / O デバイスへのトランザクションを保持する第 2 のバッファを監視して、前記レスポンス付きトランザクションの完了を確認するステップと、

前記レスポンス付きトランザクションの完了を確認したときに、トランザクション抑止の完了を通知するステップと、

前記完了通知を受け、前記サーバに割り付けられた I / O デバイスの構成変更を行うステップと、

前記レジスタに設定された該抑止指示情報を解除するステップと、

前記サーバの停止を解除するステップと、

を含むことを特徴とする I / O パスの交替処理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、サーバ計算機と I / O デバイスを P C I スイッチで接続する情報処理システムに関し、特に、サーバ計算機を仮想化して他の物理計算機へ移動させる技術に関する。

【背景技術】

【0002】

近年、IT システムを構成するサーバ台数の増加による運用管理コスト増大と、CPU マルチコア化等による物理サーバ（物理計算機）の高性能化を背景として、1 台の物理サーバを論理的に分割して仮想サーバとして運用するサーバ仮想技術を用いて、物理サーバ台数を削減するサーバ統合が注目を集めるようになってきている。

【0003】

上記のサーバ仮想化技術には幾つかの実現方法が存在するが、中でもハイパバイザを用いた仮想化方式は、仮想化に伴うオーバヘッドが小さいことが特徴として知られている（例えば、非特許文献 1）。この非特許文献 1 に示されるサーバ仮想化方法は、物理サーバのファームウェアとして実装されるハイパバイザと、ハードウェアとして実装される仮想化アシスト機能から構成される。ハイパバイザは、物理サーバ上の各仮想サーバ（仮想計算機）を制御し、仮想化アシスト機能は I / O デバイスからのメモリアクセスのアドレスを変換する。非特許文献 1 に示されるサーバ仮想化方法では、ハイパバイザが仮想サーバに I / O デバイスを直接割当てることから、物理サーバと仮想サーバ間で同一の OS 環境を利用できる。また、ハードウェアベースの仮想化アシスト機能が I / O デバイスのアドレス変換処理を行うため、I / O 処理に伴うアドレス変換処理オーバヘッドを削減して、I / O スループットを向上できる。

【0004】

一方で、CPU マルチコア化等による物理サーバの高性能化が進むと、物理サーバの性能に応じた I / O スループットが求められる。物理サーバの性能に応じた I / O スループットを向上する技術としては、複数の物理サーバに対して複数の I / O デバイスを接続可能な P C I スイッチが有望である。

【0005】

10

20

30

40

50

さらに、上記サーバ仮想化技術とP C Iスイッチを組合せ、仮想サーバとI / Oデバイスの対応付けを柔軟化するI / O仮想化技術の標準化も進められてきている(S R / M R - I o V (非特許文献2))。

【0006】

このような背景のもと、今後、サーバ仮想化技術とP C Iスイッチを組み合わせた情報処理システムが主流になると考えられる。

【0007】

ところで、サーバ仮想化技術を適用した計算機システムの柔軟性や可用性を高める機能の一つにライブマイグレーションがある(非特許文献3)。

【0008】

ライブマイグレーションは、稼働中の仮想サーバを他の物理サーバに移動させる機能である。ライブマイグレーションによって、物理サーバの負荷に応じて稼働中の仮想サーバの配置を変更したり、メンテナンス対象の物理サーバから稼働中の仮想サーバを退避できる。その結果、システムの柔軟性や可用性を高めることができる。

【0009】

上記に述べたようなサーバ仮想化技術とP C Iスイッチを組み合わせた情報処理システムにおいてライブマイグレーションを実現する場合、稼働中の仮想サーバの3つの状態を保持して引き継ぐ必要がある。3つの状態は、(a)CPU稼働状態、(b)メモリ内容、(c)I / Oデバイス状態である。(a)のCPU稼働状態は、ハイパバイザが仮想サーバの動作を停止させることで保持することができる。

【0010】

しかしながら、通常のハイパバイザは、仮想サーバに割当てられたI / OデバイスからのDMA等のメモリアクセスや外部I / Oデバイス発のトランザクション処理を止められないため、(b)メモリ内容や(c)I / Oデバイス状態を保持することができない。そのため、上記に述べたようなサーバ仮想化技術においてライブマイグレーションを実現する場合は、稼働するI / Oデバイスの影響を排除し(b)メモリ内容と、(c)I / Oデバイス状態を保持する必要がある。

【0011】

これらの状態を保持する技術としては幾つか存在する。

【0012】

例えば、特許文献1には、物理サーバとI / OデバイスがP C Iスイッチで接続されP C Iマネージャによって管理されたシステムにおいて、仮想I / Oデバイスのマイグレーションを実現するために、マイグレーション中にI / Oデバイスからのトランザクション処理を停止させる方法が開示されている。

【0013】

この特許文献1では、I / Oデバイスのコンフィギュレーション空間上にマイグレーションビットを持ち、I / Oデバイスが該マイグレーションビットを参照して仮想I / Oデバイスがマイグレーション中かどうかを判断する。

【0014】

I / Oデバイスは、マイグレーション中であれば自身の処理を停止し、I / Oデバイス発メモリ行きのトランザクションの発行を抑止する。また、P C Iマネージャは、マイグレーション中であればP C Iスイッチ内のI / Oデバイス発メモリ行きのトランザクションを特定のストレージ領域に退避し、マイグレーション完了後にストレージ領域に退避したトランザクションをリストアする。

【0015】

これにより、マイグレーション中のI / Oデバイス発トランザクションがメモリ内容を書き替えたり、I / Oデバイス状態が変更されることを防止する。

【0016】

また、特許文献2では、ホストOSがI / Oデバイスをエミュレーションしたデバイスエミュレータを仮想サーバに提供し、仮想サーバのゲストOSが該デバイスエミュレータ

10

20

30

40

50

を利用して I / O デバイスに間接的にアクセスする方法が開示されている。

【 0 0 1 7 】

ホスト O S が、マイグレーション中であることを認識し、マイグレーション中であれば、該エミュレーションデバイスの処理を停止させることで、マイグレーション中の仮想サーバのメモリ内容と I / O デバイス状態を保持できる。

【特許文献 1】米国特許出願公開第 2 0 0 7 / 0 1 8 6 0 2 5 号明細書

【特許文献 2】米国特許第 6 4 9 6 8 4 7 号明細書

【非特許文献 1】上野仁 他共著、「情報システムの運用効率を向上する「BladeSymphony」のサーバ仮想化機構「Virtage」」、2007年7月、インターネット <<http://www.hitachiyoron.com/2007/07/pdf/07a10.pdf>>

10

【非特許文献 2】Michael Krause 他 共著、「I/O Virtualization and Sharing」、2006年11月、http://www.pcisig.com/developers/main/training_materials/

【非特許文献 3】Christopher Clark 他 共著、「Live Migration of Virtual Machines」、2005年5月、NSDI (Networked Systems Design and Implementation) '05

【発明の開示】

【発明が解決しようとする課題】

【 0 0 1 8 】

しかしながら、上記従来例のような 前述の特許文献 1 を用いた方法では、I / O デバイスが仮想サーバのマイグレーション中か否かを判定する機能を有している場合にのみ適用可能である。したがって P C 用途などに広く流通している汎用 I / O カードを対象に出来ないという問題がある。

20

【 0 0 1 9 】

また、前述の特許文献 2 を用いた方法では、仮想サーバが I / O デバイスにアクセスを行うために、ゲスト O S とホスト O S の切り替えが必要となるため、この切り換えがオーバヘッドとなって処理性能が低下するという問題がある。さらには、仮想サーバに I / O デバイスを直接割り付ける場合には、デバイスエミュレータの方式は適用できないという問題もある。

【 0 0 2 0 】

上記問題を解決するため、物理サーバと I / O デバイスが P C I スイッチを介して接続された情報システムにおいて、汎用的な I / O デバイスが仮想サーバに直接割当てられている場合であっても、処理のオーバヘッドを削減しつつ、マイグレーション中仮想サーバのメモリ内容と I / O デバイス状態を保持することが必要である。

30

【 0 0 2 1 】

本発明の課題は、汎用的な I / O デバイスが仮想サーバに直接割当てられている場合であっても、処理オーバヘッドを削減しつつ、マイグレーション中の仮想サーバのメモリ内容と I / O デバイス状態を保持する機構を提供する事である。

【課題を解決するための手段】

【 0 0 2 2 】

本発明の一実施形態によれば、複数の物理サーバと 1 枚以上の I / O デバイスを接続する P C I スイッチを有する情報処理システムにおいて、物理サーバは、仮想サーバと該仮想サーバに割り付けられた I / O デバイスの対応を管理する仮想化部を備える。

40

【 0 0 2 3 】

I / O スイッチは、I / O デバイスから仮想サーバへのトランザクション発行の抑止要求を指示するレジスタと、I / O デバイスから仮想サーバへのトランザクションを抑止し、抑止前に発行された I / O デバイスからのトランザクションの完了を保證するトランザクション抑止制御部と、仮想サーバのアドレスを物理サーバのメモリ上のアドレスに変換する仮想化アシスト部と、I / O スイッチの構成を管理するスイッチ管理部を備える。

【 0 0 2 4 】

仮想化部は、I / O スイッチのレジスタに対して、トランザクション抑止要求 (抑止指示情報) と仮想サーバのメモリアドレスを設定するトランザクション指示部と、I / O ス

50

イッチからの完了通知を受け、スイッチ管理部に対する構成変更の指示と、仮想化アシスト部にアドレス変換部の変更を指示する構成変更指示部を備える。

【0025】

トランザクション指示部によって、I/Oデバイスから仮想サーバへのトランザクションの抑止と抑止前に発行されたI/Oデバイスからトランザクションの完了保証処理を行い、I/Oデバイスからのトランザクションが仮想サーバのメモリ状態を書き換えないようにする。また、構成変更指示部によって、I/Oデバイスが保持するDMA等のメモリアドレスが仮想サーバの移動後も有効になるよう仮想化アシスト部のアドレス変換部を更新し、I/Oデバイスの状態を維持する。

【発明の効果】

10

【0026】

したがって、本発明は、汎用的なI/Oデバイスが仮想サーバに直接割当てられている場合であっても、処理のオーバーヘッドを削減しつつマイグレーション中の仮想サーバのメモリ内容とI/Oデバイス状態を保持できる。これにより、PCIスイッチ等のI/Oスイッチを介して汎用的なI/Oデバイスを仮想計算機に割り当てた状態で、仮想計算機を他の物理計算機へ円滑に移動させることが可能となる。

【発明を実施するための最良の形態】

【0027】

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。

【0028】

20

以下に、本発明の実施の形態を添付の図面に基づいて説明する。

【0029】

図1は、本発明の第1の実施の形態である仮想サーバの状態を保持する機構を備えた情報処理システム100の構成の一例を示すブロック図である。

【0030】

情報処理システム100は、1つ以上の物理サーバ110a~110bと、1つ以上のI/Oデバイス120a~120bと、物理サーバ上の仮想サーバを制御するサーバマネージャ140と、I/Oデバイス120a~120bと物理サーバ110a~110bを接続するPCIスイッチ150と、PCIスイッチ150を管理するPCIMマネージャ130とを含む。物理サーバ110a~110bと、PCIMマネージャ130と、サーバマネージャ140と、I/Oデバイス120a~120bは、PCIスイッチ150を介して接続されている。また、物理サーバ110a~110bと、PCIMマネージャ130と、サーバマネージャ140はEthernet(登録商標)あるいはI2C(Integrated Circuit)などの管理用ネットワーク102により接続されている。あるいは、PCIスイッチ150を介してインバウンドでアクセスするための仕組みを備えても良い。物理サーバ110a~110bと、PCIMマネージャ130と、サーバマネージャ140の間で情報の受け渡しが出来ればどのような接続方法でも良い。また、本発明の第1の実施の形態では、2台の物理サーバと、2枚のI/OデバイスがPCIスイッチに接続されている例を示しているが、この数には限定されない。

30

【0031】

40

物理サーバ110aはハードウェア116aとハイパバイザ111aを備え、物理サーバ110a上では仮想サーバ115aが動作する。

【0032】

ハードウェア116aは、物理サーバ上のハードウェアリソースであるCPU(プロセッサ)、チップセット、メモリ等を備える。ハードウェア116aの物理的な接続構成は後述する(図3参照)。なお、PCIMマネージャ130とサーバマネージャ140も物理サーバ110aのハードウェア116aと同様にCPUとチップセット及びメモリ等から構成された計算機である。

【0033】

ハイパバイザ111aは、物理サーバ110a上に実装されたファームウェアもしくはは

50

アプリケーションであり、物理サーバ110a上の仮想サーバ115aを管理する。仮想サーバ115aには、ハイパバイザ111aにより管理されるハードウェア116aのCPUやメモリ等のリソースが割当てられる。仮想サーバ115aへのCPUリソースの割当を行うために、ハイパバイザ111aはCPUリソースの仮想サーバ115aへの割当を管理するCPUスケジューラ(図示省略)を保持する。CPUスケジューラは周知または公知の技術を用いれば良く、本発明の実施の形態では詳細説明を省略する。

【0034】

本発明の第1の実施の形態では、仮想サーバ115aのみを示すが、仮想サーバ数は一つに制限されない。ハイパバイザ111aが必要に応じて仮想サーバを生成する。

【0035】

ハイパバイザ111aは、I/Oデバイス管理表117aと、I/O発T×抑止指示部112aと、I/O構成変更指示部113aと、I/O発T×再開指示部114aと、設定インタフェース101aを備える。

【0036】

I/Oデバイス管理表117aは、仮想サーバ115aと仮想サーバ115aに割り付けられたI/Oデバイス120a、120bの対応を管理する。I/Oデバイス管理表117aの構成は後述する(図2参照)。

【0037】

I/O発T×(トランザクション)抑止指示部112aは、PCIスイッチ150に対して、仮想サーバ115aに割当てられたI/Oデバイス120a、120bが発行するトランザクションを抑止するよう指示する。I/O発T×抑止指示部112aは、例えば、PCIスイッチ150のコンフィギュレーションレジスタ158に設けた設定レジスタ161に対して書込みトランザクションを発行する。

【0038】

I/O構成変更指示部113aは、PCIマネージャ130に対してPCIスイッチ150の構成を変更するよう指示する。また、I/O構成変更指示部113aは、PCIスイッチ150の仮想化アシスト部153に対してアドレス変換表152を変更するよう指示する。

【0039】

I/O発T×再開指示部114aは、PCIスイッチ150に対して、仮想サーバ115aに割当てられたI/Oデバイス120a、120bからのトランザクション(I/O発トランザクション)を発行するよう指示する。具体的には、PCIスイッチ150に設けた設定レジスタ161に対して書込みトランザクションを発行する。

【0040】

設定インタフェース101aは、サーバマネージャ140およびPCIマネージャ130とハイパバイザ111aとの間で設定情報をやり取りするためのインタフェースである。設定情報は、通常のネットワークを用いて情報のやり取りをしても良いし、情報共有用のレジスタをハイパバイザ111aに用意し、レジスタアクセスを介して情報をやり取りしても良い。やり取りするための情報とは、例えばI/O構成変更をするためのトリガとなる情報である。

【0041】

物理サーバ110aは、以上のように構成されており、物理サーバ110bも物理サーバ110aと同様に構成されたハイパバイザ111b、仮想サーバ115b、ハードウェア116bを備える。

【0042】

I/Oデバイス120a~120bは、NICやファイバチャネルのHBAなどの汎用的なI/Oインタフェースカードである。

【0043】

PCIマネージャ130は、PCIスイッチ150を管理する管理プログラムを含む計算機である。PCIマネージャ130は、CPUおよびメモリなどを備えたハードウェア

10

20

30

40

50

上に実装されている。PCIマネージャ130は、PCIスイッチ150の構成を変更するためのI/O構成変更部131と、ハイパバイザ111a~111bやサーバマネージャ140との間で設定情報のやり取りを行う設定インタフェース101dを備える。また、本発明の第1の実施の形態では、単一のPCIマネージャ130のみ示すが、信頼性向上のため複数のPCIマネージャを備えても良い。その場合、複数のPCIマネージャ間では整合性が取れるように情報の制御を行う。

【0044】

サーバマネージャ140は、情報処理システム100全体を管理するプログラムを含む計算機であり、CPUおよびメモリなどを備えたハードウェア上に実装されている。サーバマネージャ140は、仮想サーバ115a(115b)のマイグレーションの処理開始を指示する処理開始指示部141と、ハイパバイザ111a~111bやPCIマネージャ130との間で設定情報のやり取りを行う設定インタフェース101eとを備える。本発明の第1の実施の形態では、単一のサーバマネージャ140のみを示すが、信頼性向上のため複数のサーバマネージャを備えても良い。この場合、複数のサーバマネージャ間では整合性が取れるように情報を制御する。

10

【0045】

PCIスイッチ150は、複数の物理サーバ110a~110bとI/Oデバイス120a~120bとを接続するスイッチファブリックであり、物理サーバ110a、110bとI/Oデバイス120a、120bの間でトランザクションの送受信を行う。トランザクションの詳細については後述する(図10参照)。

20

【0046】

本発明の第1の実施の形態で示す情報処理システム100では、単一のPCIスイッチ150のみを示しているが、複数のPCIスイッチを備えていてもよい。それぞれのPCIスイッチが別々のPCIマネージャで管理されていても良いし、単一のPCIマネージャで管理されていても良い。また、本発明の第1の実施の形態では、PCIスイッチ150はPCI-Expressプロトコルを対象としたスイッチファブリックとする。ただし、PCIスイッチ150は、PCIプロトコルやPCI-Xプロトコルといった他のプロトコルを対象としたスイッチファブリック(I/Oスイッチ)であってもよい。

【0047】

PCIスイッチ150は、1つ以上のUpstreamポート151a~151cと、1つ以上のDownstreamポート160a~160bと、スイッチング部157と、PCIスイッチ管理部154と、設定レジスタ161及びルーティングテーブルなどを格納するコンフィギュレーションレジスタ158とを備え、スイッチング部157は仮想化アシスト部153を含む。本発明の第1の実施の形態では、3つのUpstreamポートと2つのDownstreamポートを示すが、ポート数はこの数には限定されない。

30

【0048】

Upstreamポート151a~151cは物理サーバ110a、110b、PCIマネージャ130を接続するポートである。

【0049】

Downstreamポート160a~160bはI/Oデバイス120a、120bを接続するポートである。Downstreamポート160aは、Tx抑止制御部162と、設定レジスタ161を備える。Tx抑止制御部162は、設定レジスタ161に設定された情報に従い、仮想サーバ115a、115bに割当てられたI/Oデバイス120a、120bからのトランザクションの抑止と再開を行う。設定レジスタ161は、ハイパバイザ111aのI/O発Tx抑止指示部112aから発行される抑止指示要求を保持する。Downstreamポート160aの詳細は後述する(図5参照)。

40

【0050】

PCIスイッチ管理部154は、PCIスイッチ150のスイッチング部157を制御する機能要素であり、PCIマネージャ130と連携して動作する。PCIスイッチ管理

50

部154はPCIスイッチ150を1つ又は複数のPCIツリーに分割して管理する。PCIツリーは、1つのUpstreamポートと1つ又は複数のDownstreamポートの組からなる。PCIスイッチ管理部154は、あるPCIツリーのポートに接続する物理サーバ110a、110bやI/Oデバイス120a、120bが該PCIツリーに接続していない物理サーバやI/Oデバイスにアクセスできないよう、UpstreamポートとDownstreamポート間のトランザクションのルーティングを行う。PCIスイッチ管理部154は、複数のPCIツリーが存在する場合は、PCIツリー識別子を用いてPCIツリーを一意に特定する。また、PCIスイッチ管理部154は、PCIマネージャ130のI/O構成変更部131からの指示により、PCIスイッチ150のコンフィギュレーションレジスタ158を書き換えて、PCIツリーに属するUpstreamポートと、Downstreamポートの設定変更を行う。

10

【0051】

スイッチング部157は、PCIスイッチ管理部154が管理するツリー構造に従って、UpstreamポートとDownstreamポートとの間でトランザクションを転送する。

【0052】

仮想化アシスト部153は、Upstreamポート151a~151cとDownstreamポート160a~160bを接続する経路155、156上(スイッチング部157)に位置する制御回路であり、アドレス変換表152を含む。アドレス変換表152は、仮想サーバ115a、115bのアドレス(仮想アドレス)と物理サーバ110a、110bのアドレス(物理アドレス)の対応関係を管理する表である。アドレス変換表152の詳細は後述する(図4参照)。仮想化アシスト部153は、アドレス変換表152を参照して、I/Oデバイス120a、120bから仮想サーバ110a、110bに発行されるトランザクションのアドレスを変換して、物理サーバ110a、110bのメモリにトランザクションを発行する。また、仮想化アシスト部153は、ハイパバイザ111a、111bと連携して動作し、アドレス変換表152に設定されている仮想アドレスと物理アドレスの対応を設定し、必要に応じて変更する。本発明の第1の実施の形態では、仮想化アシスト部153を単一の制御要素として示したが、UpstreamポートもしくはDownstreamポート毎に分散して設けても良い。

20

【0053】

図3は、物理サーバ110aのハードウェア116aの構成を説明した図である。ハードウェア116aは、CPU301、チップセット302、メモリ303、Rootポート304を含む。CPU301は、プログラムを実行するプロセッサである。CPU301は、メモリアドレス領域の一部にPCIスイッチ150のレジスタをマッピングしたMMIO(メモリマップドI/O)を介して、PCIスイッチ150のレジスタにアクセスする。また、CPU301およびチップセット302は、CPU発Tx抑止要求(Quiescence要求)とCPU発Tx抑止解除要求(Dequiesce要求)の処理をサポートし、CPU301発トランザクションの抑止と抑止解除が可能である。また、チップセット302は、CPU301、メモリ303、Rootポート304を接続する。メモリ303は、物理サーバ110aの主記憶領域であり、一部の領域が仮想サーバ115aのメモリ領域として利用される。Rootポート304は、PCIスイッチ150のUpstreamポート151bと接続し、接続するPCIツリーのルートとなる。

30

40

【0054】

図10は、物理サーバ110a、110bとI/Oデバイス120a、120b間で送受信されるトランザクションの構成である。トランザクションは、ヘッダ1201とペイロード1202で構成されるパケットである。

【0055】

ヘッダ1201は、PCIスイッチ150がトランザクションをルーティングするために必要な情報であり、トランザクション送信元識別子1203(Requester ID)と、送信先アドレス1204と、トラフィッククラス1205を含む。また、複数の

50

PCIツリーがある場合には、ヘッダ1201にPCIツリー識別子が追加される。ヘッダ情報は、標準化仕様で決まるため、本発明の実施の形態では、詳細説明を省略する。ここでは本発明に関連するヘッダ情報のみ説明する。

【0056】

トランザクション送信元識別子1203は、送信元のI/OデバイスやRootポートのPCIツリーにおけるバス番号、デバイス番号、ファンクション番号から構成される識別子であり、送信元I/Oデバイスに接続するDownstreamポートやRootポートに接続するUpstreamポートを一意に識別できる。

【0057】

送信先アドレス1204は、トランザクション送信先のメモリのアドレスである。トランザクション送信先アドレス1204は、トランザクションが通過するPCIスイッチ150内のVirtual Channelを一意に特定するための情報である。本発明の実施の形態では、トランザクション送信先アドレスは0から7までの値を設定可能である。

【0058】

ペイロード1202は、トランザクションが保持するデータを格納する。例えば、メモリライトトランザクションの場合は、メモリに書き込むデータを格納する。

【0059】

図2は、I/Oデバイス管理表117aの構成を示す図である。なお、物理サーバ110bのI/Oデバイス管理表117bも同様に構成される。

【0060】

I/Oデバイス管理表117aは、ハイパバイザ111aにより管理され、仮想サーバ識別子201と、I/Oデバイス識別子202を含む。仮想サーバ識別子201はI/Oデバイス管理表117aを保持する物理サーバ110a内の仮想サーバ115aを一意に識別する番号である。I/Oデバイス識別子202は、物理サーバ110aに割り当てられたI/Oデバイスを一意に識別する識別子である。本発明の第1の実施の形態では、I/Oデバイス識別子202として図10に示したトランザクション送信元識別子1203を用いる。トランザクション送信元識別子は、トランザクションのヘッダ1201に含まれており、I/Oデバイスが接続するDownstreamポートを一意に特定可能である。ただし、トランザクションに含まれ、I/Oデバイスの接続するDownstreamポートを一意に識別できるならば他の識別子を用いても良い。

【0061】

図4は、アドレス変換表152の詳細を示した図である。アドレス変換表152は、PCIツリー毎に管理され、I/Oデバイス識別子401、変換オフセット402から成る組を保持する。I/Oデバイス識別子401は、I/O発トランザクションのヘッダに含まれているトランザクション送信元識別子1203を用い、PCIツリー内のI/Oデバイスを一意に識別する。変換オフセット402は、該I/Oデバイスが割り当てられた仮想サーバ115aのメモリ領域の、物理サーバ110aのメモリ領域における開始アドレスを示す。

【0062】

仮想化アシスト部153は、アドレス変換表152を参照してI/O発トランザクションに含まれるトランザクション送信元識別子1203に対応する変換オフセットを求め、該変換オフセットを用いてI/O発トランザクションの送信先アドレス1204を仮想アドレスから物理アドレスに変換する。本発明の第1の実施の形態では、アドレス変換表152はPCIツリー毎に分割して管理されているが、複数のPCIツリーをまとめて一つのアドレス変換表で管理しても良い。

【0063】

図5に、Downstreamポート160aのブロック構成を示す。Downstreamポート160aは、設定レジスタ161と、Tx抑止制御部162と、接続するI/OデバイスからのトランザクションをUpstreamポートへ送信する受信バッファ507と、Upstreamポートからのトランザクションを、このDownstream

10

20

30

40

50

mポートに接続されたI/Oデバイスへ送信する送信バッファ508を備え、I/Oデバイス120a、120bと物理サーバ110a間を流れるトランザクションの送受信を行う。

【0064】

設定レジスタ161は、図11で示すようにPCIスイッチ150のコンフィギュレーションレジスタ158に設けたもので、Downstreamポート160a、160b毎に設けたものであり、設定レジスタ161がDownstreamポート160aに対応し、設定レジスタ161bがDownstreamポート160bに対応する。各設定レジスタ161、161bは、物理サーバ110a、110bやPCIマネージャ130からMMIO経由でアクセス可能なレジスタであり、Downstreamポート毎に抑止ビット509とアドレス510を格納するフィールドを含む。抑止ビット509は、ハイパバイザ111aのI/O発T×抑止指示部112aおよびI/O発T×再開指示部114aからの指示に従い、設定、解除される。本発明の第1の実施の形態では、PCIのコンフィギュレーション空間への書込みを行うトランザクションを用いる。例えば、I/O発T×抑止指示部112は、コンフィギュレーション空間への書込みトランザクションにより、抑止ビット501に1をセットし、I/O発T×再開指示部114aは、コンフィギュレーション空間への書込みトランザクション用いて、抑止ビット509の値を0クリアする。設定レジスタ161のアドレス510には、I/O発T×抑止指示部112aからの指示に従い、マイグレーションの対象となる仮想サーバの仮想アドレスが設定される。また、I/O発T×再開指示部112aは、コンフィギュレーション空間への書込みトランザクション用いて、アドレス510の値をクリアする。

【0065】

T×抑止制御部162は、T×抑止部501と、滞留T×完了確認部502と、T×再開部503を備える。

【0066】

T×抑止部501は、受信バッファ507内のI/O発トランザクションをUpstreamポート507へ発行されないように制御する。例えば、PCIスイッチ150内の経路156（スイッチング部157）と受信バッファ507間で行われるフロー制御の仕組みを利用し、T×抑止部501は受信バッファ507が発行したトランザクションに対するACK（応答）を返送しないように制御することで、受信バッファ507のトランザクションの発行を抑止する。

【0067】

滞留T×完了確認部502は、レスポンス付きT×発行部504と、T×応答確認部505と、T×完了通知部506を備え、抑止前に発行されたI/O発トランザクションの完了を保証する。

【0068】

レスポンス付きT×発行部504は、設定レジスタ161のアドレス510を送信先アドレスとするメモリリードトランザクションを生成し、生成したメモリリードトランザクションを発行する。メモリリードトランザクションは、レスポンス付きT×の一種である。

【0069】

T×応答確認部505は、レスポンス付きT×発行部504が発行したメモリリードトランザクションの応答の受信を確認する。本発明の第1の実施の形態では、T×応答確認部505が送信バッファ508を監視し、レスポンス付きT×発行部504が発行したメモリリードトランザクションの応答を送信バッファ508が受信したことを確認する。

【0070】

T×完了通知部506は、T×応答確認部505が応答を確認した事をうけて、Downstreamポート160aと接続したI/Oデバイスに割当てられた物理サーバのハイパバイザ111aに、T×抑止制御の完了を通知する。

【0071】

Tx再開部503は、設定レジスタ161の抑止ビット509がクリアされた事を検出すると受信バッファ507内のI/O発トランザクションの発行を再開する。例えば、Tx再開部503が、受信バッファ507が発行したトランザクションに対するACKの返送の抑止を解除することで、受信バッファ507のトランザクションの発行は再開される。

【0072】

なお、Downstreamポート160bも上記Downstreamポート160aと同様に構成され、設定レジスタ161bとTx抑止制御部162bを備える。また、計算機と接続されるUpstreamポート151a~151cは、少なくとも送信バッファと受信バッファを備えていれば良く、設定レジスタやTx抑止制御部が設定されていなくても良い。

10

【0073】

図6は、図5で説明したTx抑止制御部162（または162b）で行われる処理の一例を示すフローチャートである。また、図7は、Tx抑止制御部162の処理に関するトランザクションの流れを示した図である。以下、図6および図7を参照してTx抑止制御部の処理フローを説明する。なお、以下の説明では、Downstreamポート160aに関する処理の一例を示すが、Downstreamポート160bも同様の処理を行うことができる。

【0074】

図6、図7の本処理は、I/Oデバイス120aがDownstreamポート160aに接続されている場合に開始される(S600)。

20

【0075】

まず、Tx抑止制御部162のTx抑止部401は、設定レジスタ161の抑止ビット509を監視し、抑止ビットが設定されたかどうかをチェックする(S601)。具体的には、Tx抑止制御部162は、設定レジスタ161の抑止ビット509の0から1への遷移を検出するまで、設定レジスタ161から抑止ビット509の値を繰り返して読み出す(図7の701)。抑止ビット509に1が設定された場合、次のステップS602を行う。抑止ビット509が0である場合には、ステップS601を繰り返す。

【0076】

次に、受信バッファ507内のトランザクションの発行を抑止する(S602)。具体的には、Tx抑止部501が、受信バッファ507がUpstreamポート151a~151cへ向けて発行したトランザクションに対するACK(応答)の返送を抑止する(図7の702)。

30

【0077】

次に、S602によるトランザクションの発行抑止前に発行されたI/O発トランザクションの完了を保証するために、設定レジスタ161が保持するアドレス510に対してレスポンス付Txを生成して発行する(S603)。I/Oデバイスと物理サーバ間のパス(I/Oパス)が複数のパスに分かれている場合は、全てのI/Oパスについてレスポンス付きトランザクションを発行する。具体的には、滞留Tx完了確認部502のレスポンス付きTx発行部504が、設定レジスタ161からアドレス510を取得し(図7の703)、該アドレスに対してレスポンス付きTxの1つであるメモリリードトランザクションを生成し(図7の704)、Upstreamポートへ発行する(図7の705)。また、PCIスイッチ150が複数のVirtual Channelを備える場合には、PCIスイッチ150の全ての利用可能なVirtual ChannelについてI/O発トランザクションの完了を保証するために、トラフィッククラスが0から7のヘッダを付加した8種類のメモリリードトランザクションを発行する。全てのトラフィッククラスに対してメモリリードトランザクションを発行することで、利用可能な全てのVirtual ChannelについてI/O発トランザクションの保証が行える。

40

【0078】

次に、S603で発行したレスポンス付きトランザクションの完了を確認する(S60

50

4)。PCI-Expressのオーダリングルールでは、レスポンス付きトランザクションは、先行するメモリライトトランザクションを追い抜かないため、発行したレスポンス付きトランザクションが完了すれば、抑止前に発行されたI/O発メモリライトトランザクションが完了していることが保証される。すなわち、S603で発行したレスポンス付きトランザクションの完了を確認した後は、仮想サーバ115aのメモリ内容が保持される。滞留Tx完了確認部502のTx応答確認部505が、メモリリードトランザクションの応答(図7の706)を全て確認するまで、ステップS604を繰り返す(図7の707)。

【0079】

次に、Tx抑止制御部162がハイパバイザ111aにトランザクション抑止完了通知を送信する(S605)。具体的には、滞留Tx完了確認部502のTx完了通知部506が、ハイパバイザ111aに処理完了通知708を発行する。処理完了通知708の発行手段は、物理サーバ110aに対する割込みでも、物理サーバ110aのメモリ303に対する書込みでもよい。

【0080】

次に、I/O発トランザクションの抑止指示情報が解除されるまで待つ(S606)。具体的には、Tx再開部503は、設定レジスタ161の抑止ビット509が0にクリアされたかどうかをチェックし(図7の701)、抑止ビット509が0に遷移したことを受けて次の処理に進む。

【0081】

最後に、受信バッファ507内のI/O発トランザクションの抑止を解除する(S607)。具体的には、Tx再開部503が、受信バッファ507が発行したトランザクションに対するACKの返送の抑止を解除することで、受信バッファ507からUpstreamポート151a~151cへ向けたトランザクションの発行を再開させる(図7の709)。

【0082】

以上の一連の処理を持って、Tx抑止制御部162の処理を完了する(S608)。以上のように、Tx抑止制御部162がI/O発トランザクションの抑止と完了保証を行うことで、物理サーバ110aのメモリ内容をI/O発トランザクションから保護できる。すなわち、I/O発トランザクションの抑止を開始した後に、設定レジスタ161のアドレス510にセットした移動元の仮想サーバのアドレスに対してレスポンス付きトランザクションを発行し、発行したレスポンス付きトランザクションの完了を確認することで、抑止開始以前に発行したトランザクションが完了していることを保証することができるのである。

【0083】

図8は、情報処理システム100で仮想サーバのライブマイグレーションを実現する処理のフローチャートである。また、図9は、ライブマイグレーションの処理に関するトランザクションの流れを示した情報処理システム100の要部を示すブロック図である。

【0084】

図8のフローチャートは、I/Oデバイス120aが割当てられている仮想サーバ115aを、物理サーバ110aから物理サーバ110bにマイグレーションするという想定で説明する。

【0085】

マイグレーションの対象となる仮想サーバ115aをマイグレーション対象仮想サーバと呼ぶ。マイグレーション対象仮想サーバが存在する物理サーバ110aをマイグレーション元物理サーバと呼び、マイグレーション元物理サーバ110a上のハイパバイザ111aをマイグレーション元ハイパバイザと呼ぶ。また、マイグレーション先の物理サーバ110bをマイグレーション先物理サーバと呼び、マイグレーション先物理サーバ上110bのハイパバイザ111bをマイグレーション先ハイパバイザと呼ぶ。

【0086】

10

20

30

40

50

図8のフローチャートは、サーバマネージャ140の処理開始指示部141がマイグレーション元ハイパバイザ111aに対してマイグレーション対象仮想サーバ115aのマイグレーション開始要求を発行してから、マイグレーション先ハイパバイザ111bがサーバマネージャ140にマイグレーション完了を通知するまでの処理フローを示す。

【0087】

ライブマイグレーションの処理は、マイグレーション元ハイパバイザ111aがサーバマネージャ140の処理開始指示部141によるマイグレーション開始要求901を受付けたことを契機に開始される(S801)。

【0088】

図9のマイグレーション開始要求901には、マイグレーション元仮想サーバ115aの仮想サーバ識別子および、マイグレーション先物理サーバ110bの識別子を含む。マイグレーション元ハイパバイザ111aは、マイグレーション開始要求901を受けて以下の処理を行う。

【0089】

まず、マイグレーション元ハイパバイザ111aは、マイグレーション対象仮想サーバ115aを停止する(S802)。実現にあたっては既知の手段を用いる。本発明の第1の実施の形態では、まず、マイグレーション元ハイパバイザ111aは、ハードウェア116aのCPU301およびチップセット302に対してCPU発T_x抑止要求を。次に、ハイパバイザ111aが、仮想サーバ115aにCPUリソースを割当てないようにCPUスケジューラの設定を変更し、仮想サーバ115aの動作を停止させる。最後にCPU発T_x抑止解除要求を行う。以上の処理により、仮想サーバ115aの停止と仮想サーバ115a発のトランザクションの完了が保証される。

【0090】

次に、マイグレーション元ハイパバイザ111aは、PCIスイッチ150に対して、マイグレーション対象仮想サーバ115aに接続するI/Oデバイス120aからトランザクションの抑止を指示する(S803)。具体的には、マイグレーション元ハイパバイザ111aのI/O発T_x抑止指示部112aが、I/Oデバイス管理表117aを参照して、仮想サーバ115aに割当てられているI/Oデバイス120aを抽出する。次に、I/O発T_x抑止指示部112aは、I/Oデバイス120aに接続するDownstreamポート160aの設定レジスタ161に対して、PCIスイッチ150のMMIOを経由してコンフィギュレーションレジスタ158への書込みを行う(図9の902)。具体的には、設定レジスタ161の抑止ビット509に1をセットし、アドレス510に仮想サーバ115aが用いるメモリアドレスの一部をセットする。その結果、Downstreamポート160aのT_x抑止制御部162により、仮想サーバ115aに割当てられたI/Oデバイス120aからのトランザクションの抑止処理が開始される。このI/Oデバイスからのトランザクションの抑止処理フローは図6で説明した通りである。

【0091】

次に、マイグレーション元ハイパバイザ111aは、PCIスイッチ150のT_x抑止制御部162からのI/Oデバイス120a発のトランザクションの抑止完了通知を待つ(S804)。マイグレーション元ハイパバイザ111aが、抑止完了通知708を受け取ることで、マイグレーション対象の仮想サーバ115a宛の全てのトランザクションが完了したことを確認でき、I/O発トランザクションによって仮想サーバ115aのメモリ内容が書き換わらないことを保証できる。

【0092】

次に、マイグレーション元ハイパバイザ111aは、マイグレーション対象仮想サーバ115aをマイグレーション元物理サーバ110aからマイグレーション先物理サーバ110bに移動させる(S805)。実現にあたっては既知の手段を用いる。具体的には、ハイパバイザ111aが、仮想サーバ115aのOS及びアプリケーションのイメージをマイグレーション先物理サーバ110bにコピーする。本発明の第1の実施の形態では、

ハイパバイザ 1 1 1 a が、管理用ネットワーク 1 0 2 を利用したアウトバウンド通信を利用して、仮想サーバ 1 1 5 a のメモリ内容とハイパバイザ 1 1 1 a が保持する仮想サーバ 1 1 5 a の構成情報を、マイグレーション先物理サーバ 1 1 0 b にコピーする。ただし、P C I スイッチ 1 5 0 を介したインバウンド通信による仮想サーバの移動を行っても良い。

【 0 0 9 3 】

次に、マイグレーション元ハイパバイザ 1 1 1 a は、マイグレーション対象仮想サーバ 1 1 5 a に割当てられた I / O デバイス 1 2 0 a の引継ぎを指示する (S 8 0 6)。具体的には、ハイパバイザ 1 1 1 a の I / O 構成変更指示部 1 1 3 a は、P C I スイッチ管理部 1 5 4 の設定を変更し、I / O デバイス 1 2 0 a が接続する P C I ツリーを変更する。すなわち、I / O デバイス 1 2 0 a の割り付けを物理サーバ 1 1 0 a から物理サーバ 1 1 0 b の仮想サーバ 1 1 5 a へ切り替える。また、仮想化アシスト部 1 5 3 にアクセスし、アドレス変換表 1 5 2 における物理アドレスと仮想アドレスの対応を変更する。この変更は、アドレス変換表 1 5 2 の I / O デバイス識別子 4 0 1 が I / O デバイス 1 2 0 a に一致するエントリの変換オフセット 4 0 2 を、仮想サーバ 1 1 5 a の移動先である物理サーバ 1 1 0 b のメモリの物理アドレスに対応するオフセット値に更新する。つまり、ハイパバイザ 1 1 1 a の構成変更指示部 1 1 3 a が I / O デバイス識別子 4 0 1 と新たな物理サーバ 1 1 0 b の変換オフセット 4 0 2 を P C I スイッチ管理部 1 5 4 に指令する。その後、ハイパバイザ 1 1 1 a は仮想サーバ 1 1 5 a が物理サーバ 1 1 0 a から削除されたので、I / O デバイス管理表 1 1 7 a から I / O デバイス 1 2 0 a に関する情報を削除する。

【 0 0 9 4 】

P C I スイッチ管理部 1 5 4 の設定を変更するためには、既知の手段を用いる。具体的には、ハイパバイザ 1 1 1 a の I / O 構成変更指示部 1 1 3 a が、P C I マネージャ 1 3 0 の I / O 構成変更部 1 3 1 に対して、I / O 構成変更要求 9 0 6 を発行する。I / O 構成変更要求 9 0 6 は、P C I スイッチ 1 5 0 もしくは設定インタフェース 1 0 1 a、1 0 1 d を介して P C I マネージャ 1 3 0 に送られる。I / O 構成変更要求 9 0 6 には、P C I スイッチの識別子、引継ぎ対象 I / O デバイスの I / O デバイス識別子が含まれる。また、P C I スイッチ 1 5 0 内に複数の P C I ツリーが含まれる場合は、I / O 構成変更要求 9 0 6 には、引継ぎ元 P C I ツリーの P C I ツリー識別子、引継ぎ先 P C I ツリーの P C I ツリー識別子が含まれる。I / O 構成変更部 1 3 1 は、I / O 構成変更要求 9 0 6 を受けて、設定変更要求 9 0 7 を P C I スイッチ管理部 1 5 4 に発行する。

【 0 0 9 5 】

また、アドレス変換表 1 5 2 の設定を変更するために、I / O 構成変更指示部 1 1 3 a は、仮想化アシスト部 1 5 3 に対して、アドレス変換表更新要求 9 0 4 を発行する。アドレス変換表更新要求 9 0 4 には、I / O デバイス識別子、変換オフセットの情報が含まれる。P C I スイッチ 1 5 0 内に複数の P C I ツリーが含まれる場合は、アドレス変換表更新要求 9 0 4 には、P C I ツリー識別子が含まれる。仮想化アシスト部 1 5 3 は、アドレス変換表更新要求 9 0 4 を受けて、アドレス変換表 1 5 2 を更新する。

【 0 0 9 6 】

P C I スイッチ管理部 1 5 4 の設定と仮想化アシスト部 1 5 3 の設定を変更することで、D o w n s t r e a m ポート 1 6 0 a の受信バッファ 5 0 7 に存在するトランザクションは、マイグレーション先物理サーバ 1 1 0 b 上の仮想サーバ 1 1 5 a に割当てられたメモリ領域に書き込まれる。

【 0 0 9 7 】

次に、マイグレーション先ハイパバイザ 1 1 1 b は、I / O 構成変更の完了通知を待つ (S 8 0 7)。完了通知は、ハイパバイザ 1 1 1 b が P C I マネージャ 1 3 0 の I / O 構成変更部 1 3 1 から明示的に受け取っても良いし、ハイパバイザ 1 1 1 b がマイグレーション先物理サーバ 1 1 0 b に I / O デバイス 1 2 0 a が追加されたことを共有レジスタなどを介して検知するのでもよい。ハイパバイザ 1 1 1 b が I / O 構成変更が完了した事を認識できれば、何れの方法であっても良い。ハイパバイザ 1 1 1 b は、I / O 構成変更の

10

20

30

40

50

完了通知を受け、I/Oデバイス管理表117bにI/Oデバイス120aに関する情報を追加する。

【0098】

次に、マイグレーション先ハイパバイザ111bは、マイグレーション対象仮想サーバ115aの処理を再開させる(S808)。実現にあたっては既知の手段を用いる。例えば、ハイパバイザ111bが、仮想サーバ115aにCPUリソースを割当てようハイパバイザ111bのCPUスケジューラの設定を変更し、仮想サーバ115aの動作を再開させる。

【0099】

次に、マイグレーション先ハイパバイザ111bは、PCIスイッチ150に対して、I/Oデバイス発トランザクションの再開を指示する(S809)。具体的には、ハイパバイザ111bのI/O発Tx再開指示部114bは、I/Oデバイス管理表117bを参照して、マイグレーション対象仮想サーバ115aに割当てられているI/Oデバイス120aが接続するDownstreamポート160aを抽出する。次に、Downstreamポート160aに対応する設定レジスタ161に対して、コンフィギュレーション空間への書き込みトランザクションを用いてレジスタアクセスを行う(図9の908)。レジスタアクセス908では、設定レジスタ161の抑止ビット509とアドレス510を0にクリアする。Downstreamポート160aのTx再開部503は、設定レジスタ161の各種設定情報が0にクリアされた事を検知し、I/Oデバイス120aからのトランザクションの送付を再開させる。

【0100】

最後に、ステップ810(S810)ではマイグレーション先ハイパバイザ111bは、サーバマネージャ140にマイグレーションが完了した事を通知する(図9の909)。

【0101】

以上の一連の処理を持って、ライブマイグレーションを完了する(S811)。以上のように、ライブマイグレーションの処理を行うことで、マイグレーション対象仮想サーバ115aに割当てられたI/Oデバイス120aからのトランザクションがマイグレーション中の仮想サーバのメモリ領域に書き込まれることを防止できる。そのため、ライブマイグレーション対象仮想サーバ115aのメモリ状態とI/Oデバイス状態の保持を実現できる。

【0102】

<変形例1>

本発明の第1の実施の形態で述べた仮想サーバの状態を保持する機構は、仮想サーバのライブマイグレーション用途の他に、I/Oパス交替機能にも適用可能である。I/Oパス交替機能は、物理サーバと物理サーバに割り付けられたI/Oデバイスの間の経路(I/Oパス)を現用系と待機系の複数を用意し、I/Oパス上のポート等に障害が発生した際にI/Oパスを現用系から待機系にフェールオーバーする機能である。I/Oパス交替機能によって、PCIスイッチのポート故障による情報処理システムの停止を回避できるため、情報処理システムの可用性が向上する。

【0103】

図12は、本発明の変形例1におけるI/Oパス交替機能を備えた情報処理システム1000の構成を示すブロック図である。

【0104】

情報処理システム1000は、1つ以上の物理サーバ1010と、1つ以上のI/Oデバイス120aと、PCIマネージャ1030と、サーバマネージャ1040と、1つ以上のPCIスイッチ150a~150bとを含む。物理サーバ1010と、PCIマネージャ1030と、サーバマネージャ1040と、I/Oデバイス120aは、PCIスイッチ150a~150bを介して接続されている。PCIスイッチ150aとPCIスイッチ150bは、Upstreamポート151aとDownstreamポート160

c、Upstreamポート151bとDownstreamポート160dという2つの経路で接続されている。また、物理サーバ1010と、PCIマネージャ1030と、サーバマネージャ1040は管理用ネットワーク102により接続されている。なお、PCISwitch150a、150bは、前記第1実施形態と同様に、コンフィギュレーションレジスタとスイッチング部とを備えるが図12においては図示を省略した。また、Downstreamポート160a～160dはTx抑制制御部とコンフィギュレーションレジスタに設定された設定レジスタ161を備えるが、Downstreamポート160b～160dについては設定レジスタ等を省略した。

【0105】

物理サーバ1010には、I/Oデバイス120aが割り当てられており、物理サーバ1010とI/Oデバイス120aの間で現用系のI/OパスとしてUpstreamポート151d、Downstreamポート160d、Upstreamポート151b、Downstreamポート160aを経由するI/Oパスが設定されている。また、待機系のI/OパスとしてUpstreamポート151d、Downstreamポート160c、Upstreamポート151a、Downstreamポート160aを経由するI/Oパスが設定されている。

10

【0106】

情報処理システム1000を構成するコンポーネントの構成は本発明の第1の実施の形態で示した情報処理システム100を構成するコンポーネントの構成と類似するため、以下、情報処理システム1000と情報処理システム100の差分について説明する。

20

【0107】

物理サーバ1010は、CPU、チップセット、メモリ等から成るハードウェア116を備える。物理サーバ1010では、OS1015が動作し、OS1015ではアプリケーション1016が動作する。

【0108】

OS1015は、ドライバモジュール1017と、I/O障害検出部1011と、サーバ発Tx抑制部1012と、I/Oパス交替指示部1013と、サーバ発Tx再開部1014を備える。ドライバモジュール1017はI/Oデバイス120aのドライバである。I/O障害検出部1011と、サーバ発Tx抑制部1012と、I/Oパス交替指示部1013と、サーバ発Tx再開部1014は、OS1015の一機能として実装されているため、アプリケーション1016が意識することなく、I/Oパス交替処理を行える。

30

【0109】

I/O障害検出部1011は、OS1015が利用するI/OデバイスやPCISwitchからの障害通知を検出し、障害がI/Oパスに関する場合、I/Oパス交替処理を開始する。具体的には、I/O障害検出部1011は、例えばPCI-Express Switchが備えるAdvanced Error Reporting機能により受け取ったI/O障害通知を解析し、障害がI/Oパスに関する場合は、物理サーバ1010をリセットせずに、I/Oパス交替処理を開始する。

【0110】

サーバ発Tx抑制部1012は、障害が発生したI/Oパスを利用するI/Oデバイスへのトランザクション(サーバ発トランザクション)の発行を抑制する。実現にあたっては既知の手段を用いる。例えば、物理サーバ1010上のOS1015は、I/OデバイスのHot Plug機能を備え、Hot Plugの仕組みを用いて物理サーバ1010に割当てられているI/Oデバイスを切り離す。

40

【0111】

I/Oパス交替指示部1013は、PCIマネージャ1030に対して、障害が発生したI/Oパスの交替を指示する。具体的には、I/Oパス交替指示部1013は、PCIマネージャ1030に対してI/Oパス交替要求を発行する。I/Oパス交替要求は、障害が発生したPCISwitchの識別子と、障害発生ポートの識別子が含まれる。I/Oパ

50

ス交替要求は、物理サーバ1010からPCIスイッチ150bを介してPCIマネージャ1030に通知されても良いし、物理サーバ1010から管理用ネットワーク102を介してサーバマネージャ1040を経由してPCIマネージャ1030に通知されても良い。

【0112】

サーバ発Tx再開部1014は、サーバ発Tx抑止部1012によって抑止されたI/Oデバイスへのトランザクションの発行を再開させる。実現にあたっては既知の手段を用いる。例えば、OS1015は、I/OデバイスのHot Plug機能を備え、Hot Plugの仕組みを用いて物理サーバ1010に割当てられているI/Oデバイスを接続する。

10

【0113】

サーバ発Tx抑止部1012およびサーバ発Tx再開部1014は、サーバ発トランザクションの制御を行うための機構であればHot Plug機能に限定されない。例えば、物理サーバ1010がハイパバイザを備える場合、S802やS808で述べた方法を用いて、ハイパバイザがサーバ発トランザクションの制御を行っても良い。また、物理サーバ1010に接続するPCIスイッチ150bが物理サーバ1010からのトランザクションを制御する機構を備えるならば、物理サーバ1010はサーバ発トランザクションの制御を行うための機構を備えなくても良い。

【0114】

PCIマネージャ1030は、I/O構成変更部131と、I/O発Tx抑止指示部112d、I/O構成変更指示部113d、I/O発Tx再開指示部114dと、I/Oパス交替完了通知部1031を備える。I/O発Tx抑止指示部112dと、I/O構成変更指示部113dと、I/O発Tx再開指示部114dは、本発明の第1の実施の形態で示した物理サーバ110aが備えるI/O発Tx抑止指示部112aと、I/O構成変更指示部113aと、I/O発Tx再開指示部114aと同一である。

20

【0115】

I/Oパス交替完了通知部1031は、I/Oパス交替を指示した物理サーバ1010に対して、I/Oパス交替完了を通知する。物理サーバ1010は、I/Oパス交替完了の通知を受けて、サーバ発トランザクションを再開させる。

【0116】

図13は、PCIスイッチ150a、150b上でI/Oパス交替を実現する処理のフローチャートである。以下、物理サーバ1010と物理サーバ1010に割当てられているI/Oデバイス120a間のI/OパスのUpstreamポート151bに障害が発生し、I/Oパスを現用系から待機系にフェールオーバーするという想定で、図12に示すI/Oパス交替処理を説明する。

30

【0117】

I/Oパス交替処理は、PCIスイッチ間のパスで障害が発生した際などに開始される(S1100)。

【0118】

まず、物理サーバ1010において、I/Oパスで障害が発生した事を検出する(S1101)。具体的には、I/O障害検出部1011は、例えばPCI-Expressスイッチが備えるAdvanced Error Reporting機能により、PCIスイッチ150aのUpstreamポート151bで障害が発生した事を検出する。

40

【0119】

次に、物理サーバ1010において、障害が発生したI/Oパスを利用するI/Oデバイス120aへのトランザクション(サーバ発トランザクション)の発行を抑止し、PCIマネージャ1030にI/Oパス交替を指示する(S1102)。具体的には、サーバ発Tx抑止部1012は、Hot Plugの仕組みを用いてI/Oデバイス120aを切り離す。その結果、物理サーバ1010から、I/Oパス交替対象のI/Oパスへのトランザクションが抑止される。また、I/Oパス交替指示部1013は、PCIマネージャ

50

ャ1030に対してI/Oパス交替要求を発行する。I/Oパス交替要求には、障害発生PCIスイッチ150aの識別子と、障害発生ポート151bの識別子が含まれる。

【0120】

PCIマネージャ1030は、I/Oパス交替要求を受け取ると、まず、障害発生I/Oパスに接続するI/Oデバイス120aからのトランザクションを抑止する(S1103)。具体的には、PCIマネージャ1030のI/O発T×抑止指示部112bは、I/Oデバイス120aが接続するDownstreamポート160aの設定レジスタ161に対して、PCIスイッチのMMIOを経由してコンフィギュレーションレジスタへの書込みを行う。その結果、上述と同様にDownstreamポート160aのT×抑止制御部162により、I/Oデバイス120aからのトランザクションの抑止処理が開始される。I/Oデバイスからのトランザクションの抑止処理フローは図6で説明した通りである。

10

【0121】

次に、PCIマネージャ1030は、T×抑止制御部162からの抑止完了通知を待つ(S1104)。PCIマネージャ1030が、抑止完了通知を受け取ることによってI/Oデバイス120aからのトランザクションが抑止されたことを保証できる。

【0122】

次に、PCIマネージャ1030は、PCIスイッチ150a~150bに対して、I/O構成変更を指示する(S1105)。具体的には、PCIマネージャ1030のI/O構成変更指示部113aは、障害発生ポート151bを避ける待機系I/Oパス情報を元に、PCIスイッチ150a~150bに関するI/Oパス交替情報を生成する。PCIスイッチ150aに関するI/Oパス交替情報は、交替前パス情報「Upstreamポート151b、Downstreamポート160a」、交替後パス情報「Upstreamポート151a、Downstreamポート160a」を含む。また、PCIスイッチ150bに関するI/Oパス交替情報は、交替前パス情報「Upstreamポート151b、Downstreamポート160a」、交替後パス情報「Upstreamポート151a、Downstreamポート160a」を含む。次に、I/O構成変更指示部113aは、I/O構成変更部131に対してI/O構成変更要求を発行する。I/O構成変更要求には、前述のPCIスイッチ150a~150bに関するI/Oパス交替情報が含まれる。

20

30

【0123】

I/O構成変更部131は、I/O構成変更要求を受けて、設定変更要求をPCIスイッチ管理部154a~154bに発行する。PCIスイッチ管理部154aに対する設定変更要求は、前述のPCIスイッチ150aに関するI/Oパス交替情報を含む。また、PCIスイッチ管理部154bに対する設定変更要求は、前述のPCIスイッチ150bに関するI/Oパス交替情報を含む。PCIスイッチ管理部154a~154bは、設定変更要求に従い、該当するPCIツリーに属するポートの構成を変更する。なお、I/Oパス交替処理では、仮想サーバの移動を伴わないため、アドレス変換表の設定の変更は発生しない。

【0124】

次に、PCIマネージャ1030は、PCIスイッチ150aに対して、I/Oパス交替対象I/Oデバイス120aからのトランザクションの再開を指示する(S1106)。具体的には、PCIマネージャ1030のI/O発T×再開指示部114bは、I/Oデバイス120aが接続するDownstreamポート160aの設定レジスタ160に対して、PCIスイッチ150aのMMIOを経由してコンフィギュレーションレジスタへの書込みを行う。その結果、Downstreamポート160aのT×抑止制御部162により、I/Oデバイス120aからのトランザクションが再開される。

40

【0125】

次に、PCIマネージャ1030は、物理サーバ1010に対し、I/Oパス交替完了を通知する(S1107)。具体的には、サーバマネージャ1040のI/Oパス交替完

50

了通知部 1031 は、管理用ネットワーク 102 を介して、物理サーバ 1010 に対して、I/Oパス交替完了を通知する。

【0126】

物理サーバ 1010 では、I/Oパス交替完了の通知を受けて、サーバ発トランザクションを再開させる (S1108)。具体的には、サーバ発Tx再開部 1014 は、I/Oパス交替完了の通知を受けて、Hot Plug の仕組みを用いて I/O デバイス 120a を接続する。その結果、物理サーバ 1010 から I/O デバイス 120a へのトランザクションの発行が再開される。

【0127】

以上の一連の処理を持って、I/Oパス交替処理を完了する (S1109)。以上のように、I/Oパス交替処理を行うことで、交替対象の I/Oパス上にトランザクションが存在しないことが保証された状態で I/Oパス交替を実現でき、I/Oパス交替処理に伴うトランザクションの喪失などを回避できる。

【0128】

<変形例 2>

変形例 2 は、上記変形例 1 と比べ、I/O デバイス 120a のドライバモジュール 1017 が、図 12 に示した変形例 1 の I/O 障害検出部 1011 と、サーバ発Tx抑止部 1012 と、I/Oパス交替指示部 1013 と、サーバ発Tx再開部 1014 を備え、さらに、I/Oパス交替を実現する機構を備える点が異なる。ドライバモジュール 1017 が備える I/Oパス交替を実現する機構は、例えば、I/O デバイス 120a と連携することで I/O デバイス 120a に関する I/Oパスを複数管理し、それら I/Oパスを任意に変更する機構である。

【0129】

ドライバモジュール 1017 の I/Oパス交替指示部 1013 は、PCI マネージャ 1030 に対して I/O 発トランザクションの抑止を指示すると共に、ドライバモジュール 1017 が備える I/Oパス交替を実現する機構を用いて、I/Oパスの交替を行う。

【0130】

ドライバモジュール 1017 が I/Oパス交替処理の機構を備えることで、OS 1015 やアプリケーション 1016 が I/Oパス交替処理の機構を備えなくても、I/Oパス交替を行える。一方で、I/O デバイス 120a およびドライバモジュール 1017 は、I/Oパス交替処理のための機構を備える必要があるため、汎用的な I/O デバイスに適用できない。

【0131】

<変形例 3>

変形例 3 は、上記変形例 1、変形例 2 と比べ、PCI マネージャ 1030 が I/O 障害検出部 1011 を備える点が異なる。

【0132】

本変形例 3 では、PCI マネージャ 1030 が図 12 に示した変形例 1 の I/O 障害検出部 1011 を有し、I/O 障害検出部 1011 は、I/Oパス障害を検出したら、物理サーバ 1010 に対して I/Oパス障害を通知する。物理サーバ 1010 のサーバ発Tx抑止部 1012 は、I/Oパス障害の通知を受けて、サーバ発トランザクションの発行を抑止する。

【産業上の利用可能性】

【0133】

以上のように、本発明では、I/Oスイッチを備えて計算機と I/O デバイスの接続を動的に変更する計算機システムや、I/Oスイッチ内の経路を動的に変更する計算機システムに適用することができる。

【図面の簡単な説明】

【0134】

【図 1】第 1 の実施の形態を示し仮想サーバを実行する情報処理システムの構成の一例を

10

20

30

40

50

示すブロック図である。

【図2】第1の実施の形態を示し、I/Oデバイス管理表の一例を示す説明図である。

【図3】第1の実施の形態を示し、物理サーバのハードウェア構成を示すブロック図である。

【図4】第1の実施の形態を示し、アドレス変換表の一例を示す説明図。

【図5】第1の実施の形態を示し、Downstreamポートの構成を示すブロック図である。

【図6】第1の実施の形態を示し、Tx抑止制御部で行われる処理の一例を示すフローチャートである。

【図7】第1の実施の形態を示し、Tx抑止制御部の処理に関するトランザクションの流れを示すブロック図である。 10

【図8】第1の実施の形態を示し、情報処理システムで行われる仮想サーバのライブマイグレーション処理の一例を示すフローチャートである。

【図9】第1の実施の形態を示し、ライブマイグレーションの処理に関するトランザクションの流れを示した情報処理システムの要部を示すブロック図である。

【図10】第1の実施の形態を示し、トランザクションの構成を示すブロック図である。

【図11】第1の実施の形態を示し、コンフィギュレーションレジスタに設定された設定された設定レジスタを示すブロック図。

【図12】変形例1を示し、I/Oパス交替機能を備えた情報処理システムの構成を示すブロック図である。 20

【図13】計算機システムで行われるPCIスイッチ150a、150bのI/Oパス交替処理のフローチャートである。

【符号の説明】

【0135】

100、1000 情報処理システム

110a、110b、110c 物理サーバ

111a、111b、111c ハイパバイザ

112a、112b、112d I/O発Tx抑止指示部

113a、113b、113d I/O構成変更指示部

114a、114b、114d I/O発Tx再開指示部 30

115a、115b 仮想サーバ

120a、120b I/Oデバイス

130 PCIマネージャ

140 サーバマネージャ

150 PCIスイッチ

153 仮想化アシスト部

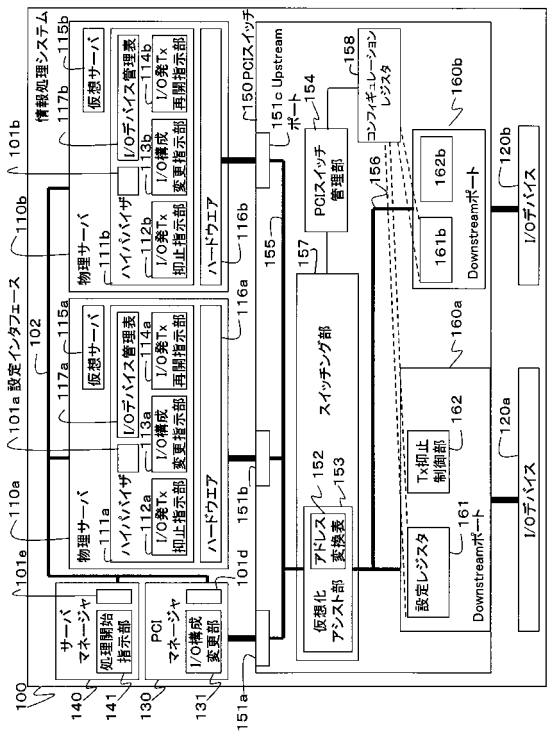
154 PCIスイッチ管理部

160a、160b Downstreamポート

161 設定レジスタ

152 Tx抑止制御部 40

【図1】

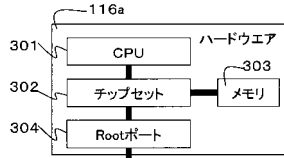


【図2】

I/Oデバイス管理表 117a

仮想サーバ識別子	I/Oデバイス識別子
1	10012
2	10011
...	11012
...	...

【図3】

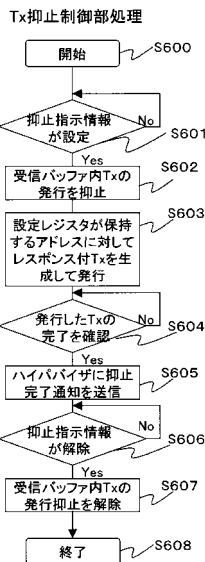


【図4】

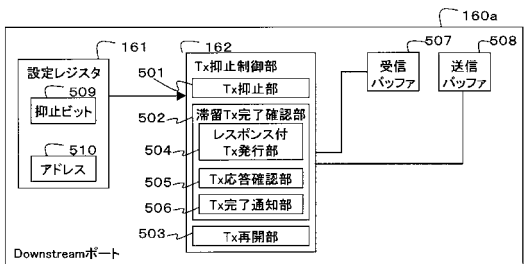
アドレス変換表 152

I/Oデバイス識別子	変換オフセット
10012	0x1000000
10011	0x2000000
...	...

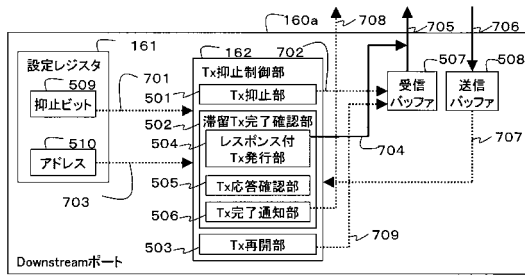
【図6】



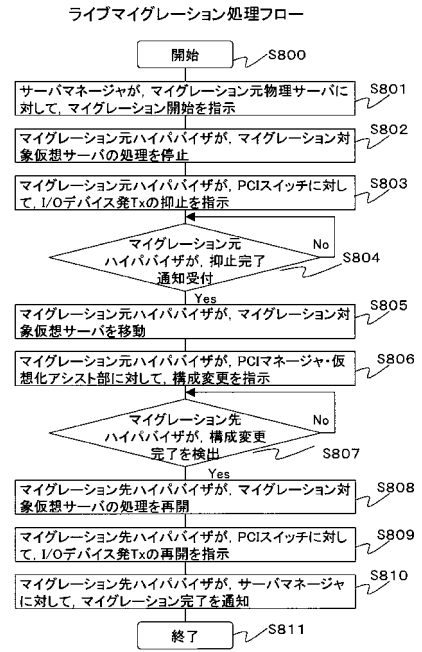
【図5】



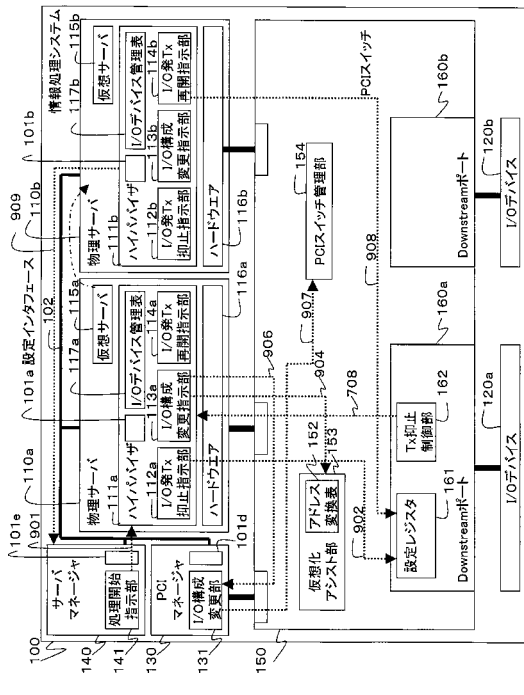
【図7】



【図8】



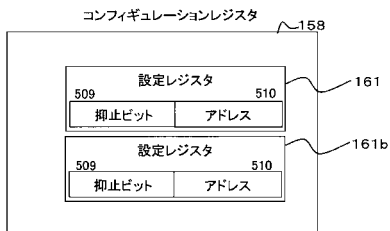
【図9】



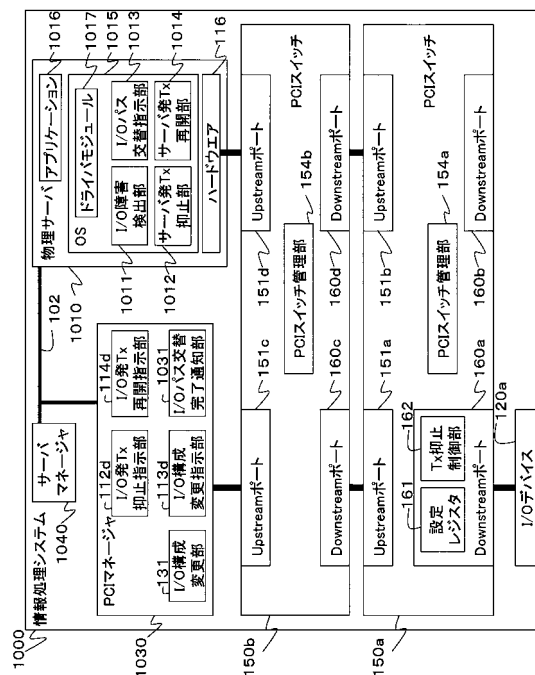
【図10】



【図 1 1】

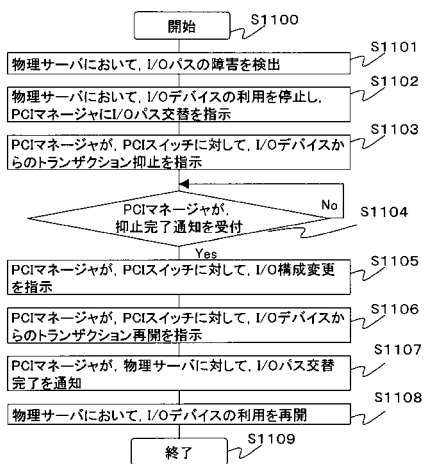


【図 1 2】



【図 1 3】

I/Oバス交替処理のフロー



フロントページの続き

- (72)発明者 馬場 貴成
東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所 中央研究所内
- (72)発明者 上原 敬太郎
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 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

審査官 木村 貴俊

- (56)参考文献 米国特許出願公開第2007/0186025 (US, A1)
特開平06-290067 (JP, A)
特開2004-032224 (JP, A)

- (58)調査した分野(Int.Cl., DB名)
G06F 3/06 - 3/08、11/20
G06F 13/00 - 13/42