

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5608243号
(P5608243)

(45) 発行日 平成26年10月15日(2014.10.15)

(24) 登録日 平成26年9月5日(2014.9.5)

(51) Int.Cl. F I
G 0 6 F 1 3 / 1 0 (2 0 0 6 . 0 1) G 0 6 F 1 3 / 1 0 3 3 0 C

請求項の数 22 (全 16 頁)

(21) 出願番号	特願2012-545042 (P2012-545042)	(73) 特許権者	591003943 インテル・コーポレーション
(86) (22) 出願日	平成21年12月24日 (2009.12.24)		アメリカ合衆国 95054 カリフォル ニア州・サンタクララ・ミッション カレ ッジ ブレーバード・2200
(65) 公表番号	特表2013-515983 (P2013-515983A)	(74) 代理人	110000877 龍華国際特許業務法人
(43) 公表日	平成25年5月9日 (2013.5.9)	(72) 発明者	ドン、ヤオズ アメリカ合衆国 95052 カリフォル ニア州・サンタクララ・ミッション カレ ッジ ブレーバード・2200 インテル ・コーポレーション内
(86) 国際出願番号	PCT/CN2009/001543		
(87) 国際公開番号	W02011/075870	審査官	木村 貴俊
(87) 国際公開日	平成23年6月30日 (2011.6.30)		
審査請求日	平成24年6月25日 (2012.6.25)		

最終頁に続く

(54) 【発明の名称】 仮想化環境において I/O 処理を行う方法および装置

(57) 【特許請求の範囲】

【請求項 1】

サービス仮想マシンが実行する方法であって、

前記サービス仮想マシンのデバイスモデルによって、前記サービス仮想マシンのデバイスドライバを呼び出して、入出力 (I/O) デバイスの仮想機能インタフェースを制御して、I/O 処理を、前記 I/O 処理に関連し、1 以上の I/O 記述子を有する I/O 情報であって、ゲスト仮想マシンによって書き込まれた I/O 情報を利用して実行する段階を備え、

前記 1 以上の I/O 記述子は、I/O 処理の種類と前記仮想機能インタフェースによって読出または書き込みを行うゲストメモリアドレスとを示すデータを含み、

前記デバイスモデルは、前記サービス仮想マシンのホスト OS 上で動作し、仮想 I/O デバイスまたは物理 I/O デバイスをエミュレートし、

前記デバイスモデル、前記デバイスドライバ、および、前記 I/O デバイスの前記仮想機能インタフェースは、前記ゲスト仮想マシンに割り当てられ、

前記 I/O デバイスの前記仮想機能インタフェースが前記ゲスト仮想マシンのアーキテクチャに準拠して動作できない場合、

前記デバイスドライバによって、前記ゲスト仮想マシンのアーキテクチャに準拠している前記 I/O 情報を、前記 I/O デバイスの前記仮想機能インタフェースのアーキテクチャに準拠しているシャドウ I/O 情報に変換する段階と、

前記デバイスドライバによって、前記 I/O デバイスの前記仮想機能インタフェースの

アーキテクチャに準拠している更新済みのシャドー I / O 情報を、前記ゲスト仮想マシンのアーキテクチャに準拠している更新済みの I / O 情報に変換する段階と

をさらに備え、

前記更新済みのシャドー I / O 情報は、前記 I / O 処理の実行に応じて、前記 I / O デバイスの前記仮想機能インタフェースによって更新される方法。

【請求項 2】

前記 I / O 処理が実行された後に、前記デバイスドライバによって、前記 I / O デバイスの前記仮想機能インタフェースのステータスを維持する段階をさらに備える請求項 1 に記載の方法。

【請求項 3】

前記 I / O 処理が実行された旨を、前記デバイスモデルが前記ゲスト仮想マシンに通知する段階をさらに備える請求項 1 または請求項 2 に記載の方法。

【請求項 4】

前記 I / O 情報は、前記 I / O デバイスの前記仮想機能インタフェースによって制御可能なヘッドポインタを先頭とするバッファに書き込まれる請求項 1 から請求項 3 のいずれか 1 つに記載の方法。

【請求項 5】

I / O 情報の終端を示すテイルポインタを、前記ゲスト仮想マシンで更新し、制御は、前記テイルポインタが更新されることに対応して、前記ゲスト仮想マシンから前記サービス仮想マシンに移行し、

制御が前記ゲスト仮想マシンから前記サービス仮想マシンに移行されることに対応して、前記サービス仮想マシンの前記デバイスモデルによって、前記サービス仮想マシンの前記デバイスドライバを呼び出す、請求項 1 から請求項 4 のいずれか 1 つに記載の方法。

【請求項 6】

デバイスモデルと、
デバイスドライバと
を備え、

前記デバイスモデルは、前記デバイスドライバを呼び出して、入出力 (I / O) デバイスの仮想機能インタフェースを制御させて、I / O 処理を、前記 I / O 処理に関連し、1 以上の I / O 記述子を含む I / O 情報であってゲスト仮想マシンによって書き込まれる I / O 情報を利用して実行させ、

前記 1 以上の I / O 記述子は、I / O 処理の種類と前記仮想機能インタフェースによって読みまたは書き込みを行うゲストメモリアドレスとを示すデータを含み、

前記デバイスモデルは、サービス仮想マシンのホスト OS 上で動作し、仮想 I / O デバイスまたは物理 I / O デバイスをエミュレートし、

前記デバイスモデル、前記デバイスドライバ、および、前記 I / O デバイスの前記仮想機能インタフェースは、前記ゲスト仮想マシンに割り当てられ、

前記 I / O デバイスの前記仮想機能インタフェースが前記ゲスト仮想マシンのアーキテクチャに準拠して動作できない場合、

前記デバイスドライバは、

前記ゲスト仮想マシンのアーキテクチャに準拠している前記 I / O 情報を、前記 I / O デバイスの前記仮想機能インタフェースのアーキテクチャに準拠しているシャドー I / O 情報に変換し、

前記デバイスドライバによって、前記 I / O デバイスの前記仮想機能インタフェースのアーキテクチャに準拠している更新済みのシャドー I / O 情報を、前記ゲスト仮想マシンのアーキテクチャに準拠している更新済みの I / O 情報に変換し、

前記更新済みのシャドー I / O 情報は、前記 I / O 処理の実行に応じて、前記 I / O デバイスの前記仮想機能インタフェースによって更新される装置。

【請求項 7】

前記デバイスドライバはさらに、前記 I / O 処理が実行された後に、前記 I / O デバイ

10

20

30

40

50

スの前記仮想機能インタフェースのステータスを維持する請求項 6 に記載の装置。

【請求項 8】

前記デバイスモデルはさらに、前記 I / O 処理が実施された旨を、前記ゲスト仮想マシンに通知する請求項 6 または請求項 7 に記載の装置。

【請求項 9】

前記 I / O 情報は、前記 I / O デバイスの前記仮想機能インタフェースによって制御可能なヘッドポインタを先頭とするバッファに書き込まれる請求項 6 から請求項 8 のいずれか 1 つに記載の装置。

【請求項 10】

I / O 情報の終端を示すテイルポインタを、前記ゲスト仮想マシンで更新し、
制御は、前記テイルポインタが更新されることに対応して、前記ゲスト仮想マシンから前記サービス仮想マシンに移行し、

制御が前記ゲスト仮想マシンから前記サービス仮想マシンに移行されることに対応して、前記サービス仮想マシンの前記デバイスモデルによって、前記サービス仮想マシンの前記デバイスドライバを呼び出す、請求項 6 から請求項 9 のいずれか 1 つに記載の装置。

【請求項 11】

システムに、

サービス仮想マシンのデバイスモデルによって、前記サービス仮想マシンのデバイスドライバを呼び出して、入出力 (I / O) デバイスの仮想機能インタフェースを制御して、I / O 処理を、前記 I / O 処理に関連し、1 以上の I / O 記述子を含む I / O 情報であって、ゲスト仮想マシンによって書き込まれた I / O 情報を利用して実行する段階

を実行させるためのプログラムであり、

前記 1 以上の I / O 記述子は、I / O 処理の種類と前記仮想機能インタフェースによって読みまたは書き込みを行うゲストメモリアドレスとを示すデータを含み、

前記デバイスモデルは、前記サービス仮想マシンのホスト OS 上で動作し、仮想 I / O デバイスまたは物理 I / O デバイスをエミュレートし、

前記デバイスモデル、前記デバイスドライバ、および、前記 I / O デバイスの前記仮想機能インタフェースは、前記ゲスト仮想マシンに割り当てられ、

前記 I / O デバイスの前記仮想機能インタフェースが前記ゲスト仮想マシンのアーキテクチャに準拠して動作できない場合、

前記デバイスドライバによって、前記ゲスト仮想マシンのアーキテクチャに準拠している前記 I / O 情報を、前記 I / O デバイスの前記仮想機能インタフェースのアーキテクチャに準拠しているシャドー I / O 情報に変換する段階と、

前記デバイスドライバによって、前記 I / O デバイスの前記仮想機能インタフェースのアーキテクチャに準拠している更新済みのシャドー I / O 情報を、前記ゲスト仮想マシンのアーキテクチャに準拠している更新済みの I / O 情報に変換する段階と

をさらに前記システムに実行させ、

前記更新済みのシャドー I / O 情報は、前記 I / O 処理の実行に応じて、前記 I / O デバイスの前記仮想機能インタフェースによって更新される

プログラム。

【請求項 12】

前記システムに、

前記 I / O 処理が実行された後に、前記デバイスドライバによって、前記 I / O デバイスの前記仮想機能インタフェースのステータスを維持する段階をさらに実行させる請求項 11 に記載のプログラム。

【請求項 13】

前記システムに、前記 I / O 処理が実施された旨を、前記デバイスモデルから前記ゲスト仮想マシンに通知する段階をさらに実行させる請求項 11 または請求項 12 に記載のプログラム。

【請求項 14】

10

20

30

40

50

前記 I / O 情報は、前記 I / O デバイスの前記仮想機能インタフェースによって制御可能なヘッドポインタを先頭とするバッファに書き込まれる請求項 11 から請求項 13 のいずれか 1 つに記載のプログラム。

【請求項 15】

I / O 情報の終端を示すテイルポインタを、前記ゲスト仮想マシンで更新し、
制御は、前記テイルポインタが更新されることに対応して、前記ゲスト仮想マシンから前記サービス仮想マシンに移行し、

制御が前記ゲスト仮想マシンから前記サービス仮想マシンに移行されることに対応して、前記サービス仮想マシンの前記デバイスモデルによって、前記サービス仮想マシンの前記デバイスドライバを呼び出す、請求項 11 から請求項 14 のいずれか 1 つに記載のプログラム。

10

【請求項 16】

入出力デバイス (I / O デバイス) を有するハードウェアマシンと、
前記ハードウェアマシンと複数の仮想マシンとの間でやり取りを行うための仮想マシンモニタと

を備え、

前記複数の仮想マシンは、

入出力処理 (I / O 処理) に関連し、1 以上の I / O 記述子を含む I / O 情報を書き込むゲスト仮想マシンと、

デバイスモデルおよびデバイスドライバを含むサービス仮想マシンと、

20

を有し、

前記デバイスモデルは、前記デバイスドライバを呼び出して、前記 I / O デバイスの仮想機能インタフェースを制御させて、前記 I / O 処理を前記 I / O 情報を利用して実行させ、

前記 1 以上の I / O 記述子は、I / O 処理の種類と前記仮想機能インタフェースによって読出したまたは書込みを行うゲストメモリアドレスとを示すデータを含み、

前記デバイスモデルは、前記サービス仮想マシンのホスト OS 上で動作し、仮想 I / O デバイスまたは物理 I / O デバイスをエミュレートし、

前記デバイスモデル、前記デバイスドライバ、および、前記 I / O デバイスの前記仮想機能インタフェースは、前記ゲスト仮想マシンに割り当てられ、

30

前記 I / O デバイスの前記仮想機能インタフェースが前記ゲスト仮想マシンのアーキテクチャに準拠して動作できない場合、

前記サービス仮想マシンの前記デバイスドライバはさらに、

前記ゲスト仮想マシンのアーキテクチャに準拠している前記 I / O 情報を、前記 I / O デバイスの前記仮想機能インタフェースのアーキテクチャに準拠しているシャドー I / O 情報に変換し、

前記 I / O デバイスの少なくとも前記仮想機能インタフェースのアーキテクチャに準拠している更新済みのシャドー I / O 情報を、前記ゲスト仮想マシンのアーキテクチャに準拠している更新済みの I / O 情報に変換し、

前記更新済みのシャドー I / O 情報は、前記 I / O 処理の実行に応じて、前記 I / O デバイスの前記仮想機能インタフェースによって更新されるシステム。

40

【請求項 17】

前記ゲスト仮想マシンは、前記 I / O 情報を、前記 I / O デバイスの前記仮想機能インタフェースによって更新されるヘッドポインタを先頭とするバッファに書き込む請求項 16 に記載のシステム。

【請求項 18】

前記ゲスト仮想マシンは、前記 I / O 情報の終端を示すテイルポインタを更新する請求項 16 または請求項 17 に記載のシステム。

【請求項 19】

前記仮想マシンモニタは、前記テイルポインタが更新されたことを検出すると、前記シ

50

システムの制御を、前記ゲスト仮想マシンから前記サービス仮想マシンへと移行し、

前記デバイスモデルは、制御が前記ゲスト仮想マシンから前記サービス仮想マシンに移行されることに対応して、前記サービス仮想マシンの前記デバイスドライバを呼び出す、請求項 18 に記載のシステム。

【請求項 20】

前記 I/O デバイスの前記仮想機能インタフェースは、前記 I/O 処理が実行されたことに応じて前記 I/O 情報を更新する請求項 16 から請求項 19 のいずれか 1 つに記載のシステム。

【請求項 21】

前記デバイスドライバは、前記 I/O 処理が実行された後、前記 I/O デバイスの前記仮想機能インタフェースのステータスを維持する請求項 16 から請求項 20 のいずれか 1 つに記載のシステム。

【請求項 22】

前記デバイスモデルは、前記 I/O 処理が実行された旨を前記ゲスト仮想マシンに通知する請求項 16 から請求項 21 のいずれか 1 つに記載のシステム。

【発明の詳細な説明】

【背景技術】

【0001】

仮想マシンアーキテクチャは、物理マシンを論理的に分割して、マシンの下層ハードウェアを共有し、1 以上の独立して動作する仮想マシンに見えるようにするアーキテクチャである。入出力 (I/O) 仮想化 (IOV) によって、複数の仮想マシンが利用する一の I/O デバイスの機能を実現するとしてよい。

【0002】

ソフトウェア完全デバイスエミュレーションは、I/O 仮想化の一例であるとしてよい。I/O デバイスの完全エミュレーションによって、仮想マシンは既存のデバイスドライバを再利用することができるようになるとしてよい。シングル・ルート・I/O 仮想化 (SR-IOV) または任意のその他のリソースパーティション化方法は、I/O 仮想化の別の一例であるとしてよい。I/O デバイス機能 (例えば、データ移動に関する I/O デバイス機能) を複数の仮想インターフェース (VI) にパーティション化して、各仮想インターフェースを 1 つの仮想マシンを割り当てる場合、ソフトウェアエミュレーションレイヤの I/O オーバーヘッドが少なくなる場合がある。

【図面の簡単な説明】

【0003】

添付図面では、本明細書で説明する発明を図示しているが、一例に過ぎず限定するものではない。図示を簡略化して分かりやすくするべく、図面に図示する構成要素は必ずしも実寸に即したものではない。例えば、一部の構成要素の寸法は、分かりやすいように、他の構成要素と比較して強調しているとしてもよい。また、適切であると考えられる場合には、対応する構成要素または同様の構成要素を示す際に同じ参照符号を複数の図面にわたって繰り返し用いる。

【図 1】ゲスト仮想マシンを発行元とする I/O 処理を制御するサービス仮想マシンを含むコンピューティングプラットフォームの実施形態を示す図である。

【図 2 A】I/O 処理のための I/O 記述子を格納する記述子リング構造の実施形態を示す図である。

【図 2 B】I/O 処理のための I/O 記述子を格納する記述子リング構造およびシャドー記述子リング構造の実施形態を示す図である。

【図 3】I/O デバイスによるダイレクトメモリアクセス (DMA) のための入出力メモリ管理ユニット (IOMMU) テーブルの実施形態を示す図である。

【図 4】ゲスト仮想マシンによる I/O 処理に関連する I/O 情報を書き込む方法の実施形態を示す図である。

【図 5】サービス仮想マシンによって I/O 情報に基づき I/O 処理を実行する方法の実

10

20

30

40

50

施形態を示す図である。

【図 6 A】サービス仮想マシンによって I / O 情報に基づき I / O 処理を実行する方法の別の実施形態を示す図である。

【図 6 B】サービス仮想マシンによって I / O 情報に基づき I / O 処理を実行する方法の別の実施形態を示す図である。

【発明を実施するための形態】

【 0 0 0 4 】

以下に記載する説明では、仮想化環境において I / O 処理を実行する方法を説明する。以下に記載する説明では、論理実施例、疑似コード、オペランド指定手段、リソースのパーティション化 / 共有 / 複製の実施例、システム構成要素の種類および相関関係、および論理のパーティション化 / 統合化に関する選択肢等、具体的且つ詳細な内容を数多く記載する。これによって、本発明を深く理解していただきたい。しかし、以下に記載する具体的且つ詳細な内容を採用することなく本発明を実施し得る。また、制御構造、ゲートレベル回路および完全なソフトウェア命令列は、本発明をあいまいにすることを避けるべく、詳細な説明を省略している。当業者であれば、本明細書に記載した内容に基づき、過度な実験を必要とすることなく適切な機能を実現可能であろう。

【 0 0 0 5 】

本明細書において「一実施形態」、「ある実施形態」、「実施形態例」という場合、当該実施形態は特定の特徴、構造または特性を含むが、どの実施形態でもその特定の特徴、構造または特性を必ずしも含むものではない。さらに、上記の表現は、必ずしも同じ実施形態を意味しているものではない。さらに、特定の特徴、構造または特性がある実施形態に関連付けて説明されている場合、明示的な言及の有無に関わらず、この特徴、構造または特性を他の実施形態に関連付けて実施することは当業者の想到する範囲内であると考えられる。

【 0 0 0 6 】

本発明の実施形態は、ハードウェア、ファームウェア、ソフトウェアまたはこれらの任意の組み合わせで実施され得る。本発明の実施形態はまた、1 以上のプロセッサで読み出されて実行される機械可読媒体に格納されている命令として実現されるときにもよい。機械可読媒体は、機械（例えば、コンピューティングデバイス）が読み出し可能な形式で情報を格納または送信する任意のメカニズムを含むとしてよい。例えば、機械可読媒体は、リードオンリーメモリ（ROM）、ランダムアクセスメモリ（RAM）、磁気ディスク格納媒体、光格納媒体、フラッシュメモリデバイス、電氣的、光学的、音響的またはその他の方法で伝搬する信号（例えば、搬送波、赤外信号、デジタル信号等）等を含むとしてよい。

【 0 0 0 7 】

仮想化環境において I / O 処理を行うコンピューティングプラットフォーム 1 0 0 の実施形態を図 1 に示す。コンピューティングシステム 1 0 0 の例を列挙すると、これらで全てを網羅しているわけではないが、分散型コンピューティングシステム、スーパーコンピュータ、コンピューティングクラスタ、メインフレームコンピュータ、ミニコンピュータ、パーソナルコンピュータ、ワークステーション、サーバ、ポータブルコンピュータ、ラップトップコンピュータ、および、データを送受信して処理するその他のデバイスがある。

【 0 0 0 8 】

当該実施形態によると、コンピューティングプラットフォーム 1 0 0 は、1 以上のプロセッサ 1 1 1、メモリシステム 1 2 1、チップセット 1 3 1、I / O デバイス 1 4 1 およびその他の構成要素を有する下層ハードウェアマシン 1 0 1 を備えるとしてよい。1 以上のプロセッサ 1 1 1 は、プロセッサバス（図 1 には不図示）等の 1 以上のバスを介して、さまざまな構成要素（例えば、チップセット 1 3 1）に通信可能に結合されているとしてよい。プロセッサ 1 1 1 は、適切なアーキテクチャでコードを実行する 1 以上のプロセッシングコアを含む集積回路（IC）として実現されるときにもよい。

【 0 0 0 9 】

メモリシステム 1 2 1 は、プロセッサ 1 1 1 が実行すべき命令およびデータを格納するとしてよい。メモリ 1 2 1 の例としては、シンクロナスダイナミックランダムアクセスメモリ (S D R A M) デバイス、 R A M B U S ダイナミックランダムアクセスメモリ (R D R A M) デバイス、ダブルデータレート (D D R) メモリデバイス、スタティックランダムアクセスメモリ (S R A M) およびフラッシュメモリデバイス等の半導体デバイスのうち、1つまたは任意の組み合わせを含むとしてよい。

【 0 0 1 0 】

チップセット 1 3 1 は、1以上のプロセッサ 1 1 1、メモリ 1 2 1 およびその他の構成要素、例えば、I/Oデバイス 1 4 1の間をつなぐ1以上の通信経路を提供するとしてよい。I/Oデバイス 1 4 1 は、これらに限定されないが、ペリフェラル・コンポーネント・インターコネクタ (P C I) バスまたは P C I E x p r e s s (P C I e) バスを介してホストマザーボードと接続されている P C I デバイスおよび/または P C I e デバイスを含むとしてよい。I/Oデバイス 1 4 1 の例は、ユニバーサルシリアルバス (U S B) コントローラ、グラフィクスアダプタ、オーディオコントローラ、ネットワークインターフェースコントローラ (N I C)、ストレージデバイス等を含むとしてよい。

10

【 0 0 1 1 】

コンピューティングプラットフォーム 1 0 0 はさらに、仮想マシンモニタ (V M M) 1 0 2 を備えるとしてよい。VMM 1 0 2 は、下層のハードウェアおよび上層の下層マシン (例えば、サービス仮想マシン 1 0 3、ゲスト仮想マシン 1 0 3₁ - 1 0 3_n) のインターフェースとして機能して、仮想マシン (例えば、サービス仮想マシン 1 0 3 のホストオペレーティングシステム 1 1 3、ゲスト仮想マシン 1 0 3₁ - 1 0 3_n のゲストオペレーティングシステム 1 1 3₁ - 1 1 3_n) の複数のオペレーティングシステム (O S) の動作を円滑化および管理して、下層物理リソースを共有する。仮想マシンモニタの例としては、Xen、ESXサーバ、仮想PC、Virtual Server、Hyper-V、Parallel、OpenVZ、Qemu等が挙げられるとしてよい。

20

【 0 0 1 2 】

ある実施形態によると、I/Oデバイス 1 4 1 (例えば、ネットワークカード) は、複数の機能部にパーティション化されるとしてよい。これらの機能部は、入出力仮想化 (I O V) アーキテクチャ (例えば、シングル・ルート I O V) をサポートするコントロールエンティティ (C E) 1 4 1₀ と、専用アクセス用のランタイムリソース (例えば、ネットワークデバイスにおけるキュー対) を持つ複数の仮想機能インターフェース (V I) 1 4 1₁ - 1 4 1_n とを含む。CE および V I の例は、シングル・ルート I/O 仮想化アーキテクチャまたはマルチ・ルート I/O 仮想化アーキテクチャにおける物理機能および仮想機能を含むとしてよい。CE はさらに、V I 機能を設定および管理するとしてよい。ある実施形態によると、複数のゲスト仮想マシン 1 0 3₁ - 1 0 3_n が、C E 1 4 1₀ が制御する複数の物理リソースを共有する一方、ゲスト仮想マシン 1 0 3₁ - 1 0 3_n のそれぞれは、V I 1 4 1₁ - 1 4 1_n のうち1以上に割り当てられるとしてよい。例えば、ゲスト仮想マシン 1 0 3₁ は V I 1 4 1₁ に割り当てられるとしてよい。

30

【 0 0 1 3 】

他の実施形態は I/O デバイス 1 4 1 の構造について他の技術を採用し得るものと考えられたい。ある実施形態によると、I/O デバイス 1 4 1 は、1以上の V I を含むが C E を含まないとしてもよい。例えば、パーティション化機能を持たないレガシー N I C は、N U L L C E 条件において機能する1つの V I を含むとしてよい。

40

【 0 0 1 4 】

サービス仮想マシン 1 0 3 には、デバイスモデル 1 1 4、C E ドライバ 1 1 5 および V I ドライバ 1 1 6 のコードがロードされているとしてよい。デバイスモデル 1 1 4 は、実際の I/O デバイス 1 4 1 のソフトウェアエミュレーションであってもなくてもよい。C E ドライバ 1 1 5 は、コンピューティングプラットフォーム 1 0 0 の初期化およびランタイムにおいて、I/O デバイスの初期化および設定に関連する C E 1 4 1₀ を管理すると

50

してよい。V I ドライバ 1 1 6 は、管理ポリシーに応じて、V I 1 4 1₁ - V I 1 4 1_nのうち1以上を管理するデバイスドライバであってよい。ある実施形態によると、管理ポリシーに基づいて、V I ドライバは、V I ドライバがサポートしているゲストVMに割り当てられているリソースを管理するとしてよく、C E ドライバは、グローバル動作を管理するとしてよい。

【0015】

ゲスト仮想マシン 1 0 3₁ - 1 0 3_n はそれぞれ、V M M 1 0 2 が提示する仮想デバイスを管理するゲストデバイスドライバ、例えば、ゲスト仮想マシン 1 0 3₁ のゲストデバイスドライバ 1 1 6₁ またはゲスト仮想マシン 1 0 3_n のゲストデバイスドライバ 1 1 6_n のコードがロードされているとしてよい。ゲストデバイスドライバは、V I 1 4 1 およびそれらのドライバ 1 1 6 に適応しているモードで動作することが可能または不可能であるとしてよい。ある実施形態によると、ゲストデバイスドライバは、レガシードライバであってよい。

【0016】

ある実施形態によると、ゲスト仮想マシンのゲストオペレーティングシステム（例えば、ゲストVM 1 0 3₁ のゲストOS 1 1 3₁）がゲストデバイスドライバ（例えば、ゲストデバイスドライバ 1 1 6₁）をロードすることに応じて、サービスVM 1 0 3 がV I ドライバ 1 1 6 およびデバイスモデル 1 1 4 のインスタンスを実行するとしてよい。例えば、デバイスモデル 1 1 4 のインスタンスは、ゲストデバイスドライバ 1 1 6₁ にサービスを提供し、一方で、V I ドライバ 1 1 6 のインスタンスは、ゲストVM 1 0 3₁ に割り当てられたV I 1 4 1₁ を制御するとしてよい。例えば、ゲストデバイスドライバ 1 1 6₁ が、8 2 5 7 1 E B ベースのN I C (Intel Corporation (米国カリフォルニア州サンタクララ) 社製のネットワークコントローラ) のレガシードライバであり、ゲストVM 1 0 3₁ に割り当てられているV I 1 4 1₁ が、8 2 5 7 1 E B ベースのN I C であるか、または、8 2 5 7 1 E B ベースのN I C に準拠した、または、していないその他の種類のN I C である場合、サービスVM 1 0 3 は、仮想8 2 5 7 1 E B ベースのN I C を表すデバイスモデル 1 1 4 のインスタンス、および、V I 1 4 1₁ を制御するV I ドライバ 1 1 6 のインスタンス、つまり、8 2 5 7 1 E B ベースのN I C、または、8 2 5 7 1 E B ベースのN I C に準拠した、または、していないその他の種類のN I C を実行するとしてよい。

【0017】

図1に図示されている実施形態は説明のためのものであり、他の技術を採用すればコンピューティングシステム 1 0 0 の他の実施形態が実施され得るものと考えられたい。例えば、デバイスモデル 1 1 4 は、V I ドライバ 1 1 6 またはC E ドライバに組み込まれているとしてよく、または、全て1つにまとめられているとしてもよい。OSカーネル等の特権モードで動作するとしてもよいし、OSユーザランド等の非特権モードで動作するとしてもよい。サービスVMは、さらに複数のVMに分割され、1つのVMがC E を実行する一方で、別のVMが、複数のVM間での通信を十分に実行しつつ、デバイスモデルおよびV I ドライバまたは任意のその他の組み合わせを実行するとしてよい。

【0018】

ある実施形態によると、ゲストVM 1 0 3₁ で実行されているアプリケーション（例えば、アプリケーション 1 1 7₁）によってI/O処理が指示された場合、ゲストデバイスドライバ 1 1 6₁ は、当該I/O処理に関連するI/O情報を、ゲストVM 1 0 3₁ に割り当てられているバッファ（図1には不図示）に書き込むとしてよい。例えば、ゲストデバイスドライバ 1 1 6₁ は、I/O記述子を図2Aに示すリング構造に書き込むとしてよい。この場合、リンク構造の1つのエントリを、1つのI/O記述子に使う。ある実施形態によると、I/O記述子は、データパケットに関連するI/O処理を示すとしてよい。例えば、ゲストアプリケーション 1 1 7₁ がゲストメモリアドレス x x x - y y y との間で100個のパケットの読出または書込を指示すると、ゲストデバイスドライバ 1 1 6₁ は、100個のI/O記述子を図2Aの記述子リングに書き込むとしてよい。ゲストデバ

10

20

30

40

50

イスドライバ116₁は、ヘッドポインタ201を先頭として、記述子を記述子リングに書き込むとしてよい。ゲストデバイスドライバ116₁は、I/O処理に関連する記述子の書き込みを完了した後、テイルポインタ202を更新するとしてよい。ある実施形態によると、ヘッドポインタ201およびテイルポインタ202は、ヘッドレジスタおよびテイルレジスタ(不図示)に格納しているとしてよい。

【0019】

ある実施形態によると、記述子は、データ、I/O処理の種類(読出か書込か)、VI141₁がデータの読出または書込を行うゲストメモリアドレス、I/O処理のステータス、I/O処理に必要な他の情報を含むとしてよい。

【0020】

ある実施形態によると、ゲストデバイスドライバ116₁がゲストVM103₁に割り当てられているVI141₁に適應したモードで動作できない場合、例えば、VI141₁およびゲストデバイスドライバ116₁がサポートしているビット形式および/または意味論が異なるために、ゲストデバイスドライバ116₁が書き込んだ記述子に基づいてVI141₁がI/O処理を実施できない場合、VIドライバ116は、シャドーリング(図2Bに示す)を生成して、ゲストVM103₁のアーキテクチャに準拠している記述子、ヘッドポインタおよびテイルポインタを、VI141₁のアーキテクチャに準拠しているシャドー記述子(S記述子)、シャドーヘッドポインタ(Sヘッドポインタ)およびシャドーテイルポインタ(Sテイルポインタ)に変換するとしてよい。この結果、VI141₁がシャドー記述子に基づいてI/O処理を実施できるようになる。

【0021】

図2Aおよび図2Bに図示されている実施形態は説明のためのものであり、他の技術を採用すればI/O情報の他の実施形態を実施し得ると考えられたい。例えば、I/O情報は、図2Aおよび図2Bのリング構造以外のデータ構造、例えば、ハッシュテーブル、リンクテーブル等で書き込まれるとしてもよい。別の例を挙げると、1つのリングを受信および送信の両方に用いるとしてよく、複数の異なるリングを受信または送信に用いるとしてよい。

【0022】

IOMMUまたは同様の技術によって、I/Oデバイス141は、記述子リングまたはシャドー記述子リングの記述子から取得されたゲストアドレスをホストアドレスに再マッピングすることによって、メモリシステム121に直接アクセスすることができるとしてよい。図3は、IOMMUテーブルの実施形態を示す図である。ゲストVM103₁等のゲスト仮想マシンは、ゲストVMのアーキテクチャに準拠したゲストメモリアドレスと、ホストコンピューティングシステムのアーキテクチャに準拠しているホストメモリアドレスとの間の対応関係を示すIOMMUテーブルを少なくとも1つ持つとしてよい。VMM102およびサービスVM103は、全てのゲスト仮想マシンのIOMMUテーブルを管理するとしてよい。更に、IOMMUページテーブルは、さまざまな方法でインデックスが付与されるとしてよく、例えば、デバイス識別子(例えば、PCIeシステムにおけるバス:デバイス:機能番号)でインデックスが付与されたり、ゲストVM番号、または、IOMMU実施例で特定されるその他の任意の方法でインデックスが付与されるとしてよい。

【0023】

他の実施形態では他の技術を用いてメモリアクセスを実現し得るものと考えられたい。ある実施形態によると、IOMMUは、例えば、ソフトウェアソリューションによってゲストアドレスがホストアドレスと等しくなる場合には、利用されないとしてよい。別の実施形態によると、ゲストデバイスドライバは、VMM102と協働して、IOMMUテーブルと同様のマッピングテーブルを利用して、ゲストアドレスをホストアドレスに変換するとしてよい。

【0024】

図4は、ゲスト仮想マシンによるI/O処理に関連するI/O情報を書き込む方法の実

10

20

30

40

50

施形態を示す図である。以下の説明は、ゲストVM103₁を一例に挙げて行う。他のゲストVMについても同一または同様の技術が適用可能であると理解されたい。

【0025】

ブロック401では、ゲストVM103₁で実行されているアプリケーション117₁が、I/O処理、例えば、ゲストメモリアドレスxxx-yyyへの100個のパケットの書込を指示するとしてよい。ブロック402では、ゲストデバイスドライバ116₁は、I/O処理に関するI/O記述子を生成して、ゲストVM103₁の記述子リング（例えば、図2Aまたは図2Bに示す記述子リング）に書き込むとしてよい。ブロック403では、I/O処理に関連する全ての記述子が記述子リングに書き込まれる。ある実施形態によると、ゲストデバイスドライバ116₁は、ヘッドポインタ（例えば、図2Aに示すヘッドポインタ201または図2Bに示すヘッドポインタ2201）を先頭として、I/O記述子を書き込むとしてよい。ブロック404において、ゲストデバイスドライバ116₁は、I/O処理に関する記述子が全てバッファに書き込まれた後、テイルポインタ（例えば、図2Aのテイルポインタ202または図2Bに示すテイルポインタ2202）を更新するとしてよい。

10

【0026】

図5は、サービスVM103によるI/O処理を行う方法の実施形態を示す。当該実施形態は、ゲスト仮想マシンのゲストデバイスドライバが、ゲスト仮想マシンに割り当てられているVIおよび/またはVIのドライバに準拠したモードで動作可能であるという条件の場合に適用されるとしてよい。例えば、ゲストデバイスドライバが82571EBベースのNICのレガシードライバである一方、VIは、82571EBベースのNICまたは82571EBベースのNICに準拠した他の種類のNIC、例えば、82576EBベースのNICの仮想機能である。以下の説明は、ゲストVM103₁を一例に挙げて行う。他のゲストVMにも同一または同様の技術を適応可能であると理解されたい。

20

【0027】

ブロック501において、ゲストVM103₁がテイルポインタ（例えば、図2Aのテイルポインタ202）を更新することによって、仮想マシンイグジット（例えば、VMExit）がトリガされるとしてよい。仮想マシンイグジットがVMM102によって取得されると、VMM102は、システムの制御をゲストVM103₁のゲストOS113₁からサービスVM103のデバイスモデル114に移行するとしてよい。

30

【0028】

ブロック502において、デバイスモデル114は、テイル更新に応じて、VIドライバ116を呼び出すとしてよい。ブロック503-506では、VIドライバ116が、ゲストVM103₁に割り当てられたVI1141を制御して、ゲストVM103₁が書き込んだI/O記述子（例えば、図2AのI/O記述子）に基づいてI/O処理を実施するとしてよい。具体的には、ブロック503において、I/O記述子を準備するべく、VIドライバ116がVI1141を呼び出すとしてよい。ある実施形態によると、VIドライバ116は、テイルレジスタ（不図示）を更新することによって、VI1141を呼び出すとしてよい。ブロック504において、VI1141は、ゲストVM103₁（例えば、図2Aに図示した記述子リング）の記述子リングから記述子を読み出して、I/O記述子の記述内容にしたがってI/O処理を実行するとしてよい。例えば、パケットを受信して、当該パケットをゲストメモリアドレスxxxに書き込むとしてよい。ある実施形態によると、VI1141は、記述子リングのヘッドポインタ（例えば、図2Aのヘッドポインタ201）が指し示しているI/O記述子を読み出すとしてよい。

40

【0029】

ある実施形態によると、VI1141は、IOMMUまたは同様の技術を利用して、I/O処理のためのダイレクトメモリアクセス（DMA）を実行するとしてよい。例えば、VI1141は、ゲストVM103₁のために生成されたIOMMUテーブルから、ゲストメモリアドレスに対応するホストメモリアドレスを取得して、メモリシステム121に対してパケットの読出または書込を直接実行するとしてよい。別の実施形態によると、V

50

I 1 1 4 1 は、ゲストアドレスとホストアドレスとの間の固定マッピングにおいてゲストアドレスがホストアドレスに等しい場合には、I O M M U テーブルを利用することなくダイレクトメモリアクセスを実行するとしてよい。ブロック 5 0 5 では、V I 1 1 4 1 がさらに、I / O 記述子を更新するとしてよい。例えば、I / O 記述子に含まれている I / O 処理のステータスを更新して、I / O 記述子が実施された旨を示すとしてよい。ある実施形態によると、V I 1 1 4 1 は、I / O 記述子の更新において、I O M M U テーブルを利用してもしなくてもよい。V I 1 1 4 1 はさらに、ヘッドポインタを更新して、ヘッドポインタを進めて、記述子リングの次の I / O 記述子を指し示すようにするとしてよい。

【 0 0 3 0 】

ブロック 5 0 6 において、V I 1 1 4 1 は、テイルが指し示す I / O 記述子に到達したか否かを判断するとしてよい。到達していない場合、ブロック 5 0 4 および 5 0 5 において、V I 1 1 4 1 は、記述子リングからの I / O 記述子の読出を継続して行い、I / O 記述子によって指示された I / O 処理を実行するとしてよい。到達している場合、ブロック 5 0 7 において、V I 1 1 4 1 は、I / O 処理が完了した旨を、例えば、V M M 1 0 2 への割り込みを通知することによって、V M M 1 0 2 に通知するとしてよい。ブロック 5 0 8 において、V M M 1 0 2 は、I / O 処理が完了した旨を、例えば、サービス V M 1 0 3 への割り込みを挿入することによって、V I ドライバ 1 0 6 に通知するとしてよい。

【 0 0 3 1 】

ブロック 5 0 9 において、V I ドライバ 1 1 6 は、V I 1 1 4₁ のステータスを維持して、I / O 処理が完了した旨をデバイスモデル 1 1 4 に通知するとしてよい。ブロック 5 1 0 において、デバイスモデル 1 4 がゲスト V M 1 1 3₁ に仮想割り込みを通知して、ゲストデバイスドライバ 1 1 6₁ が、イベントを処理して I / O 処理が実行された旨をアプリケーション 1 1 7₁ に通知するとしてよい。例えば、ゲストデバイスドライバ 1 1 6₁ は、データが受信され利用する準備が整っている旨をアプリケーション 1 1 7₁ に通知するとしてよい。ある実施形態によると、デバイスモデル 1 1 4 はさらに、ヘッドレジスタ（不図示）を更新して、記述子リングの制御がゲストデバイスドライバ 1 1 6₁ に戻された旨を示すとしてよい。ゲストデバイスドライバ 1 1 6₁ への通知は他の方法で実行され、その方法は、デバイス/ドライバポリシー、例えば、ゲストデバイスドライバがデバイス割り込みをディセーブルした場合に作成したデバイス/ドライバポリシーによって決まり得ると考えられたい。

【 0 0 3 2 】

上述した実施形態は説明のためのものであり、他の技術を採用すれば他の実施形態が実現され得るものと考えられたい。例えば、V M M メカニズムによっては、V I 1 1 4 1 は、I / O 処理が完了した旨を上層のマシンに通知する方法を異ならせるとしてよい。ある実施形態によると、V I 1 1 4 1₁ は、V M M 1 0 2 を介するのではなく、直接サービス V M 1 0 3 に通知するとしてよい。別の実施形態によると、V I 1 1 4 1 は、記述子リングに列挙されている I / O 処理の全てではなく 1 以上が完了した時点で上層のマシンに通知するとしてよい。この構成では、ゲストアプリケーションに、I / O 処理の一部が完了した旨が遅滞なく通知されるとしてよい。

【 0 0 3 3 】

図 6 A および図 6 B は、サービス V M 1 0 3 によって I / O 処理を行う方法の別の実施形態を示す図である。当該実施形態は、ゲスト仮想マシンのゲストデバイスドライバが、V I および / またはゲスト仮想マシンに割り当てられている V I のドライバに準拠したモードで動作できないという条件において適用されるとしてよい。以下の説明は、ゲスト V M 1 0 3₁ を一例に挙げて行う。同一または同様の技術を他のゲスト V M に適用可能であると理解されたい。

【 0 0 3 4 】

ブロック 6 0 1 では、V M M は、例えば、ゲストデバイスドライバ 1 1 6 が仮想デバイス（例えば、デバイスモデル 1 1 4）にアクセスするとゲスト V M 1 0 3₁ が発生させる仮想マシンイグジット（例えば、V M E x i t）を取得するとしてよい。ブロック 6 0 2

10

20

30

40

50

において、VMM102は、システムの制御を、ゲストVM103₁のゲストOS113₁からサービスVM103のデバイスモデル114へと移行させるとしてよい。ブロック603において、デバイスモデル114は、仮想マシンイグジットが、ゲストデバイスドライバ116₁がI/O処理に関連するI/O記述子の記述子リング(例えば、図2Bの記述子リング)への書き込みを完了したことに応じてトリガされているか否かを判断するとしてよい。ある実施形態によると、ゲストVM113₁は、テイルポインタ(例えば、図2Bのテイルポインタ2202)を更新して、I/O記述子の終端を示すとしてよい。この場合、デバイスモデル114は、テイルポインタの更新によって仮想マシンイグジットがトリガされたか否かを判断するとしてよい。

【0035】

10

I/O記述子の書込をゲストデバイスドライバ116₁が完了したことによって仮想マシンイグジットがトリガされていない場合、図6Aおよび図6Bの方法は、ブロック601に戻るとしてよい。つまり、VMMは次のVMイグジットを取得するとしてよい。I/O記述子の書込をゲストデバイスドライバ116₁が完了したことによって仮想マシンイグジットがトリガされている場合、ブロック604において、デバイスモデル114は、VIDライバ116を呼び出して、ゲストVM103₁のアーキテクチャに準拠しているI/O記述子を、ゲストVM103₁に割り当てられているVI141₁のアーキテクチャに準拠しているシャドーI/O記述子へと変換して、シャドーI/O記述子をシャドー記述子リング(例えば、図2Bに示すシャドー記述子リング)に格納する。

【0036】

20

ブロック605において、VIDライバ116は、ゲストVM103₁のアーキテクチャに準拠したテイルポインタを、VI141₁のアーキテクチャに準拠したシャドーテイルポインタに変換するとしてよい。

【0037】

ブロック606-610では、VIDライバ116は、VI1141を制御して、ゲストVM103₁が書き込んだI/O記述子に基づいてI/O処理を実施するとしてよい。具体的には、ブロック606において、VIDライバ116は、シャドー記述子の準備を整えるべくVI1141を呼び出すとしてよい。ある実施形態によると、VIDライバ116は、シャドーテイルポインタ(不図示)を更新することによってVI1141を呼び出すとしてよい。ブロック607において、VI1141は、シャドー記述子リングからシャドーI/O記述子を読み出して、シャドーI/O記述子の記述内容に従ってI/O処理を実施するとしてよい。例えば、パケットを受信して、当該パケットをゲストメモリアドレスxxxに書き込むか、または、ゲストメモリアドレスxxxからパケットを読み出して当該パケットを送信するとしてよい。ある実施形態によると、VI1141は、シャドー記述子リングのシャドーヘッドポインタ(例えば、図2Bのシャドーヘッドポインタ2201)が指し示すI/O記述子を読み出すとしてよい。

30

【0038】

ある実施形態によると、VI1141は、IOMMUまたは同様の技術を用いて、I/O処理のためのダイレクトメモリアクセスを実行するとしてよい。例えば、VI1141は、ゲストVM103₁のために生成されたIOMMUテーブルから、ゲストメモリアドレスに対応するホストメモリアドレスを取得して、受信したパケットをメモリシステム121に直接書き込むとしてよい。別の実施形態によると、VI1141は、ゲストアドレスとホストアドレスとの間の固定マッピングにおいてゲストアドレスがホストアドレスに等しい場合には、IOMMUテーブルを利用することなくダイレクトメモリアクセスを実行するとしてよい。ブロック608では、VI1141はさらに、シャドーI/O記述子を更新するとしてよい。例えば、シャドーI/O記述子に含まれているI/O処理のステータスを更新して、I/O記述子が実施された旨を示すとしてよい。ある実施形態によると、VI1141は、I/O記述子の更新のためにIOMMUテーブルを利用するとしてよい。VI1141はさらに、シャドーヘッドポインタを更新して、シャドーヘッドポインタを進めて、シャドー記述子リングの次のシャドーI/O記述子を指し示すようにする

40

50

としてよい。

【0039】

ブロック609において、V I ドライバ116は、更新済みのシャドーI/O記述子およびシャドーヘッドポインタをI/O記述子およびヘッドポインタに変換し直して、この新しいI/O記述子およびヘッドポインタで記述子リングを更新するとしてよい。ブロック610において、V I 1141は、シャドーテイルポインタが指し示すシャドーI/O記述子に到達するか否かを判断するとしてよい。到達していない場合、ブロック607 - 609において、V I 1141は、シャドー記述子リングからのシャドーI/O記述子の読出を継続して行い、シャドーI/O記述子に記述されているI/O処理を実行するとしてよい。到達している場合、ブロック611において、V I 1141は、I/O処理が完了した旨を、例えば、VMM102への割り込みを通知することによって、VMM102に通知するとしてよい。VMM102はこの後、I/O処理が完了した旨を、例えば、サービスVM103に割り込みを挿入することによって、V I ドライバ106に通知するとしてよい。

10

【0040】

ブロック612において、V I ドライバ116は、V I 1141のステータスを維持して、I/O処理が完了した旨をデバイスモデル114に通知するとしてよい。ブロック613において、デバイスモデル114は、ゲストデバイスドライバ116₁に仮想割り込みを通知して、ゲストデバイスドライバ116₁がイベントを処理してI/O処理が実施された旨をアプリケーション117₁に通知するとしてよい。例えば、ゲストデバイスドライバ116₁は、データを受信して利用する準備が整っている旨をアプリケーション117₁に通知するとしてよい。ある実施形態によると、デバイスモデル114はさらに、ヘッドレジスタ(不図示)を更新して、記述子リングの制御がゲストデバイスドライバ116₁に再び移行される旨を示すとしてよい。ゲストデバイスドライバ116₁への通知は他の方法で実行され、その方法は、デバイス/ドライバポリシー、例えば、ゲストデバイスドライバがデバイス割り込みをディセーブルした場合に作成したデバイス/ドライバポリシーによって決まり得ると考えられたい。

20

【0041】

上述した実施形態は説明のためのものであり、他の技術を採用した場合は他の実施形態が実施され得るものと考えられたい。例えば、VMMメカニズムによっては、V I 1141は、I/O処理が完了した旨を上層のマシンに通知する方法を異ならせるとしてよい。ある実施形態によると、V I 1141₁は、VMM102を介するのではなく、直接サービスVM103に通知するとしてよい。別の実施形態によると、V I 1141は、記述子リングに列挙されているI/O処理の全てではなく1以上が完了した時点で上層のマシンに通知するとしてよい。この構成では、ゲストアプリケーションに、I/O処理の一部が完了した旨が遅滞なく通知されるとしてよい。

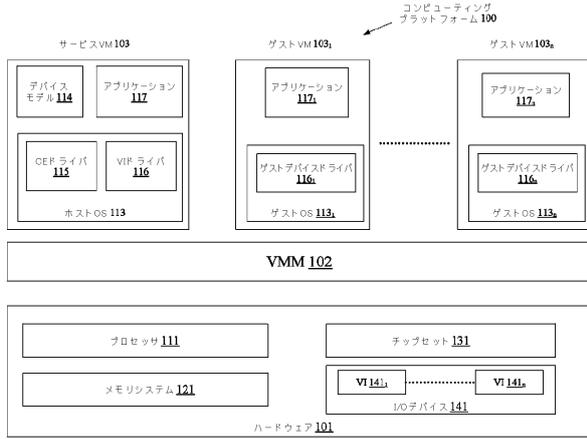
30

【0042】

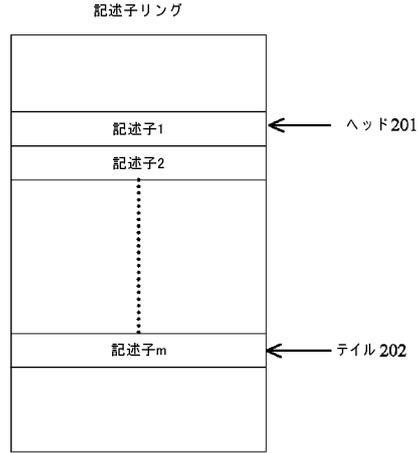
実施形態例に基づき本発明の具体的な特徴を説明してきたが、上記の説明は本発明を限定するものと解釈されるべきものではない。実施形態例のさまざまな変形例、および、本発明の他の実施形態は、本発明の関連技術分野の当業者に明らかであり、本発明の意図および範囲に含まれるものとする。

40

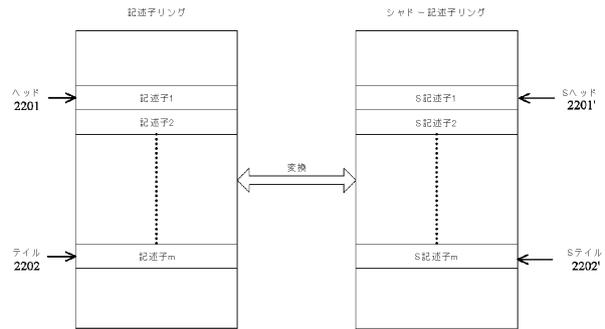
【図1】



【図2A】



【図2B】

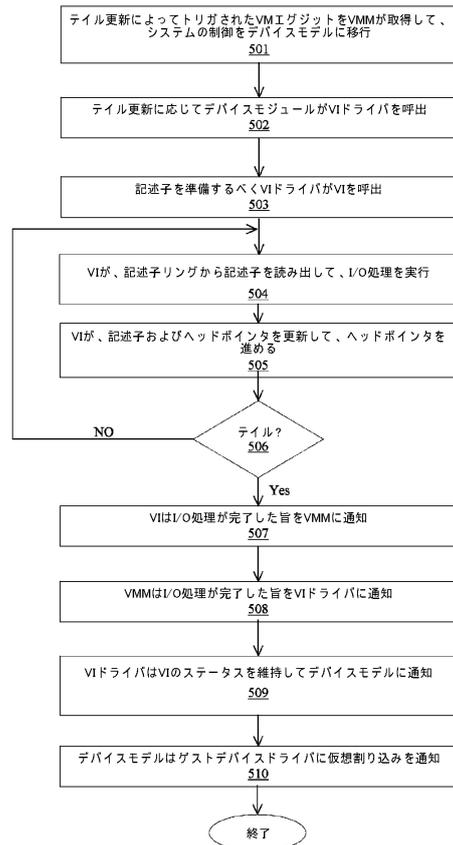


【図3】

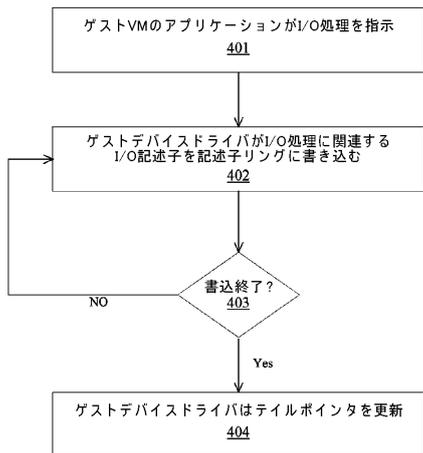
IOMMUテーブル

ゲストアドレス	ホストアドレス

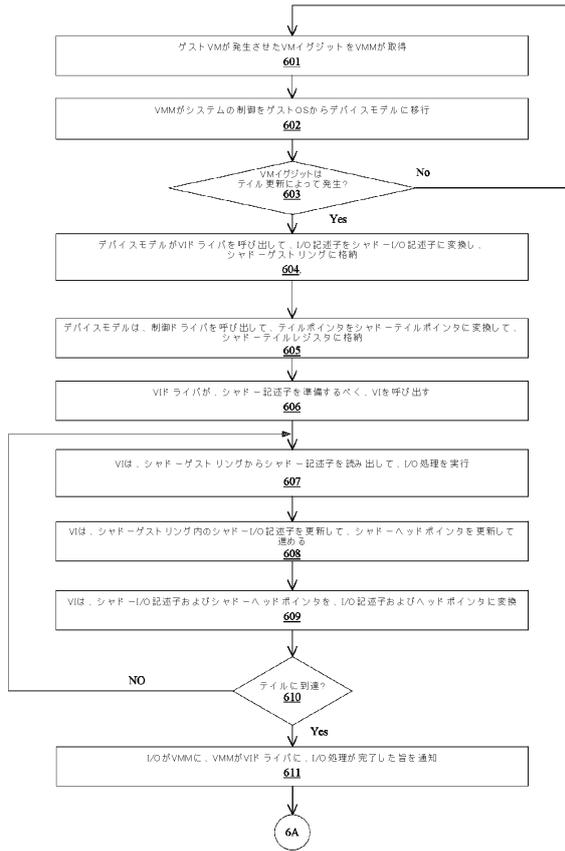
【図5】



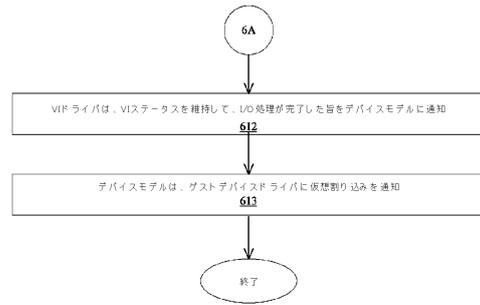
【図4】



【図 6 A】



【図 6 B】



フロントページの続き

(56)参考文献 国際公開第2007/115425(WO, A1)

米国特許出願公開第2007/0168641(US, A1)

Paul Barhamら, Xen and the Art of Virtualization, SOSP'03 Proceedings of the nineteenth ACM symposium on Operating systems principles, 米国, ACM, 2003年10月19日, Volume 37 Issue 5, pp.164-177, [online]インターネット<検索日:2013/09/05>, URL, http://dl.acm.org/ft_gateway.cfm?id=945462&ftid=231540&dwn=1&CFID=358346258&CFTOKEN=59321358

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08

G06F 12/00 - 12/16

G06F 13/00 - 13/42