

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4864817号
(P4864817)

(45) 発行日 平成24年2月1日(2012.2.1)

(24) 登録日 平成23年11月18日(2011.11.18)

(51) Int.Cl. F I
G O 6 F 9/46 (2006.01) G O 6 F 9/46 3 5 0

請求項の数 16 (全 38 頁)

(21) 出願番号	特願2007-164892 (P2007-164892)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成19年6月22日(2007.6.22)	(74) 代理人	100075513 弁理士 後藤 政喜
(65) 公開番号	特開2009-3749 (P2009-3749A)	(74) 代理人	100114236 弁理士 藤井 正弘
(43) 公開日	平成21年1月8日(2009.1.8)	(74) 代理人	100120260 弁理士 飯田 雅昭
審査請求日	平成22年2月15日(2010.2.15)	(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
		(72) 発明者	服部 直也 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】 仮想化プログラム及び仮想計算機システム

(57) 【特許請求の範囲】

【請求項1】

物理プロセッサとメモリを備えた物理計算機で実行されて、複数の仮想プロセッサを提供する仮想化プログラムにおいて、

第1の仮想マシンマネージャが第1の仮想プロセッサを生成する手順と、

前記第1の仮想プロセッサ上で実行される第2の仮想マシンマネージャが第2の仮想プロセッサを生成する手順と、

前記第2の仮想プロセッサがユーザプログラムを実行する手順と、

前記第1の仮想マシンマネージャが、前記物理プロセッサから当該第1の仮想マシンマネージャの呼び出しを受信する手順と、

前記第1の仮想マシンマネージャが、前記第1の仮想マシンマネージャ呼び出しの要因を解析する手順と、

前記第1の仮想マシンマネージャが、前記解析の結果に基づいて前記第2の仮想マシンマネージャと、前記ユーザプログラムのいずれを実行するかを判定する手順と、

前記第1の仮想マシンマネージャが、前記判定の結果に基づいて前記第1の仮想プロセッサに前記第2の仮想マシンマネージャまたは前記ユーザプログラムの一方を実行させる手順と、

を前記物理プロセッサに実行させることを特徴とする仮想化プログラム。

【請求項2】

前記第1の仮想マシンマネージャが、前記ユーザプログラムを実行する際の第1の仮想

プロセッサの状態を規定する第1の制御情報を前記メモリに設定する手順と、

前記第1の仮想マシンマネージャが、前記第2の仮想マシンマネージャを実行する際の前記第1の仮想プロセッサの状態を規定する第2の制御情報を前記メモリに設定する手順と、

前記第1の仮想マシンマネージャが、前記第2の仮想マシンマネージャまたは前記ユーザプログラムを実行する際の前記物理プロセッサの状態を規定する第3の制御情報を前記メモリに設定する手順と、

をさらに含み、

前記第1の仮想マシンマネージャが、判定の結果に基づいて前記第1の仮想プロセッサに前記第2の仮想マシンマネージャまたは前記ユーザプログラム的一方を実行させる手順は、

10

前記第1の仮想マシンマネージャが、前記第1の仮想プロセッサに前記ユーザプログラムの実行を開始させるときには、前記第1の制御情報を参照する手順と、

前記第1の仮想マシンマネージャが、前記第1の仮想プロセッサに前記第2の仮想マシンマネージャの実行を開始させるときには、前記第2の制御情報を参照する手順と、

前記第1の仮想マシンマネージャが、前記参照した前記第1の制御情報または第2の制御情報の一方で前記第3の制御情報を更新する手順と、

前記第1の仮想マシンマネージャが、前記物理プロセッサに対して前記第3の制御情報に従って命令の実行を開始するよう指示する第1の制御命令を発行する手順と、

を含むことを特徴とする請求項1に記載の仮想化プログラム。

20

【請求項3】

前記物理プロセッサは、インテル・アーキテクチャ32の命令セットに準拠するプロセッサであって、

前記第1の仮想マシンマネージャが、前記解析の結果に基づいて前記第2の仮想マシンマネージャと、前記ユーザプログラムのいずれを実行するかを判定する手順は、

前記呼び出し要因の解析の結果が、前記第2の仮想マシンマネージャからのVM-entry命令の発行である場合に、前記ユーザプログラムの実行開始と判定し、

前記呼び出し要因の解析の結果が、前記第1の仮想プロセッサから第2の仮想マシンマネージャに対してVM-exitを通知する場合に、前記第2の仮想マシンマネージャの実行開始と判定することを特徴とする請求項2に記載の仮想化プログラム。

30

【請求項4】

前記物理プロセッサは、AMD64の命令セットに準拠するプロセッサであって、

前記第1の仮想マシンマネージャが、前記解析の結果に基づいて前記第2の仮想マシンマネージャと、前記ユーザプログラムのいずれを実行するかを判定する手順は、

前記呼び出し要因の解析の結果が、前記第2の仮想マシンマネージャからのVMRUN命令の発行である場合に、前記ユーザプログラムの実行開始と判定し、

前記呼び出し要因の解析の結果が、前記第1の仮想プロセッサから第2の仮想マシンマネージャに対して#VMEEXITを通知する場合に、前記第2の仮想マシンマネージャの実行開始と判定することを特徴とする請求項2に記載の仮想化プログラム。

【請求項5】

40

前記物理プロセッサは、アイテニウム・プロセッサ・ファミリの命令セットに準拠するプロセッサであって、

前記第1の仮想マシンマネージャが、前記解析の結果に基づいて前記第2の仮想マシンマネージャと、前記ユーザプログラムのいずれを実行するかを判定する手順は、

前記呼び出し要因の解析の結果が、前記第2の仮想マシンマネージャからのPAL_VPS_RESUME_HANDLERプロシジャ呼び出しの発行である場合に、前記ユーザプログラムの実行開始と判定し、

前記呼び出し要因の解析の結果が、前記第1の仮想プロセッサから第2の仮想マシンマネージャに対して仮想化例外を通知する場合に、前記第2の仮想マシンマネージャの実行開始と判定することを特徴とする請求項2に記載の仮想化プログラム。

50

【請求項 6】

前記第 1 の仮想マシンマネージャが第 1 の仮想プロセッサを生成する手順は、
前記第 1 の仮想プロセッサを複数生成する手順を含み、
前記物理プロセッサの状態を規定する第 3 の制御情報を前記メモリに設定する手順は、
前記第 1 の仮想プロセッサ毎に、前記第 3 の制御情報をそれぞれ前記メモリに設定する
手順を含み、

前記第 1 の仮想マシンマネージャが、前記物理プロセッサに対して前記第 3 の制御情報
に従って命令の実行を開始するよう指示する第 1 の制御命令を発行する手順は、

前記複数の第 3 の制御情報のうちの一つを選択する第 2 の制御命令を発行する手順を含
むことを特徴とする請求項 2 に記載の仮想化プログラム。

10

【請求項 7】

前記第 2 の仮想マシンマネージャが、前記第 1 の仮想プロセッサで前記第 1 の制御命令
を受付可能か否かを確認する第 3 の制御命令を発行する手順をさらに含み、

前記第 1 の仮想マシンマネージャが、前記第 1 の仮想マシンマネージャ呼び出しの要因
を解析する手順は、

前記第 1 の仮想マシンマネージャ呼び出しの要因が、前記第 2 の仮想マシンマネージャ
からの前記第 3 の制御命令であるときに、予め設定された応答可否情報を参照する手順と

、
前記応答可否情報に応じて、前記第 3 の制御命令に対する応答を決定する手順と、

を含むことを特徴とする請求項 2 に記載の仮想化プログラム。

20

【請求項 8】

前記第 1 の仮想マシンマネージャが、前記第 1 の仮想プロセッサに対する設定要求を受
信する手順と、

前記第 1 の仮想マシンマネージャは、前記設定要求に基づいて前記応答可否情報を設定
する手順と、

をさらに含むことを特徴とする請求項 7 に記載の仮想化プログラム。

【請求項 9】

前記物理プロセッサは、インテル・アーキテクチャー 32 の命令セットに準拠するプロ
セッサであって、

前記第 3 の制御命令は、CPUID 命令であることを特徴とする請求項 7 に記載の仮想
化プログラム。

30

【請求項 10】

前記第 2 の仮想マシンマネージャが、前記第 1 の仮想マシンマネージャから前記第 3 の
制御命令に対する応答を受信する手順と、

前記第 2 の仮想マシンマネージャが、前記第 1 の制御命令が前記第 1 の仮想プロセッサ
にて受付可能である場合に、前記第 1 の仮想プロセッサに対して前記第 1 の制御命令を受
付可能に設定する命令を発行する手順と、

を含むことを特徴とする請求項 9 に記載の仮想化プログラム。

【請求項 11】

前記第 1 の仮想マシンマネージャが、前記ユーザプログラムを実行する際の第 1 の仮想
プロセッサの状態を規定する第 1 の制御情報を前記メモリに設定する手順と、

前記第 1 の仮想マシンマネージャが、前記第 2 の仮想マシンマネージャを実行する際の
前記第 1 の仮想プロセッサの状態を規定する第 2 の制御情報を前記メモリに設定する手順
と、

前記第 1 の仮想マシンマネージャが、前記第 2 の仮想マシンマネージャまたは前記ユー
ザプログラムを実行する際の前記物理プロセッサの状態を規定する第 3 の制御情報を前記
メモリに設定する手順と、

をさらに含み、

前記第 1 の仮想マシンマネージャが、前記解析の結果に基づいて前記第 2 の仮想マシン
マネージャと、前記ユーザプログラムのいずれを実行するかを判定する手順は、

40

50

前記解析の結果が、前記第2の仮想マシンマネージャから前記第1の制御情報に対して前記第1の仮想プロセッサの動作状態を反映させる第4の制御命令を第1の仮想マシンマネージャ呼び出しの要因のときに、

前記第1の仮想マシンマネージャが前記物理プロセッサに対して、前記第3の制御情報に前記物理プロセッサの動作状態を反映させる第4の制御命令を発行し、

前記第1の仮想マシンマネージャが前記物理プロセッサの動作状態を反映した前記第3の制御情報を参照し、

前記第1の仮想マシンマネージャが、前記参照した第3の制御情報に基づいて、前記第1の制御情報を更新した後に、前記第2の仮想マシンマネージャの実行を判定することを特徴とする請求項1に記載の仮想化プログラム。

10

【請求項12】

前記物理プロセッサは、インテル・アーキテクチャー32の命令セットに準拠するプロセッサであって、

前記第4の制御命令がVMCLEAR命令であることを特徴とする請求項11に記載の仮想化プログラム。

【請求項13】

物理プロセッサとメモリを備えた物理計算機上で複数の仮想プロセッサを提供する仮想計算機システムにおいて、

前記物理プロセッサ上で動作して、第1の仮想プロセッサを提供する第1の仮想マシンマネージャと、

20

前記第1の仮想プロセッサ上で動作して、ユーザプログラムを実行する第2の仮想プロセッサを提供する第2の仮想マシンマネージャと、を備え、

前記第2の仮想マシンマネージャは、

前記ユーザプログラムを実行する際の第1の仮想プロセッサの状態を規定する第1の制御情報と、前記第2の仮想マシンマネージャを実行する際の前記第1の仮想プロセッサの状態を規定する第2の制御情報とを格納する仮想プロセッサ制御データを有し、

前記第1の仮想マシンマネージャは、

前記第2の仮想マシンマネージャまたは前記ユーザプログラムを実行する際の前記物理プロセッサの状態を規定する第3の制御情報を前記メモリに格納する物理プロセッサ制御データと、

30

前記物理プロセッサから当該第1の仮想マシンマネージャの呼び出しを受信する受信部と、

前記第1の仮想マシンマネージャ呼び出しの要因を解析し、前記第2の仮想マシンマネージャと、前記ユーザプログラムのいずれを実行するかを判定する判定部と、

前記判定の結果に基づいて、前記仮想プロセッサ制御データの第1の制御情報と第2の制御情報の何れかを選択する選択部と、

前記選択された第1の制御情報または第2の制御情報で、前記物理プロセッサ制御データの第3の制御情報を更新する更新部と、

前記第3の制御情報に基づいて、前記第1の仮想プロセッサに前記第2の仮想マシンマネージャまたは前記ユーザプログラムの一方を実行させる命令を発行する命令発行部と、を有することを特徴とする仮想計算機システム。

40

【請求項14】

前記物理プロセッサは、インテル・アーキテクチャー32の命令セットに準拠するプロセッサであって、

前記判定部は、

前記第1の仮想マシンマネージャの呼び出し要因の解析の結果が、前記第2の仮想マシンマネージャからのVM-entry命令の発行である場合に、前記ユーザプログラムの実行開始と判定し、

前記第1の仮想マシンマネージャの呼び出し要因の解析の結果が、前記第1の仮想プロセッサから第2の仮想マシンマネージャに対してVM-exitを通知する場合に、前記

50

第2の仮想マシンマネージャの実行開始と判定することを特徴とする請求項13に記載の仮想計算機システム。

【請求項15】

前記物理プロセッサは、AMD64の命令セットに準拠するプロセッサであって、
前記判定部は、

前記第1の仮想マシンマネージャの呼び出し要因の解析の結果が、前記第2の仮想マシンマネージャからのVMRUN命令の発行である場合に、前記ユーザプログラムの実行開始と判定し、

前記第1の仮想マシンマネージャの呼び出し要因の解析の結果が、前記第1の仮想プロセッサから第2の仮想マシンマネージャに対して#VMEXITを通知する場合に、前記第2の仮想マシンマネージャの実行開始と判定することを特徴とする請求項13に記載の仮想計算機システム。

10

【請求項16】

前記物理プロセッサは、アイテニウム・プロセッサ・ファミリの命令セットに準拠するプロセッサであって、

前記判定部は、

前記第1の仮想マシンマネージャの呼び出し要因の解析の結果が、前記第2の仮想マシンマネージャからのPAL_VPS_RESUME_HANDLERプロシジャ呼び出しの発行である場合に、前記ユーザプログラムの実行開始と判定し、

前記第1の仮想マシンマネージャの呼び出し要因の解析の結果が、前記第1の仮想プロセッサから第2の仮想マシンマネージャに対して仮想化例外を通知する場合に、前記第2の仮想マシンマネージャの実行開始と判定することを特徴とする請求項13に記載の仮想計算機システム。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、仮想計算機システムに関し、特に、仮想化支援機能を備えたプロセッサを用いた仮想計算機システムに関する。

【背景技術】

【0002】

近年、オープンサーバの普及に伴い多数のサーバが企業の情報システムに導入された。特に価格性能比が高いiA(Intel Architectures)-32サーバの乱立は、電力やハードウェアの保守費用などを含むサーバの運用管理コストを増大させており、サーバを運用する企業において問題となっている。

30

【0003】

サーバの運用管理コストを低減するため、複数台のサーバを1台の物理的なサーバ(物理サーバ)に集約するサーバ統合が有望である。サーバ統合の実現方法としては、計算機資源を仮想化する機能を提供する仮想化ソフトウェアが注目されている。仮想化ソフトウェアは1台の物理サーバのCPU(プロセッサ)やI/O等の計算機資源を分割し、複数の仮想的なサーバ(仮想サーバ)に分割した計算機資源を割付ける制御ソフトウェアである。各仮想サーバ上では各々1つのOS(ゲストOS)が稼働できる。仮想化ソフトウェアを用いると、従来、複数の物理的なサーバで稼働していたOSやアプリケーションを各仮想サーバに割当てて、一つの物理計算機で複数のサーバを提供するサーバ統合を実現することができる。

40

【0004】

仮想化ソフトウェアにおける仮想サーバへの計算機資源の割付方針について説明する。計算機資源のうちCPUの割付について、iA-32向けの仮想化ソフトウェアではVT-x(Virtualization Technology for Xeon)機能(または、AMD-V)等の仮想化支援機能を備えたプロセッサを用いる方法が主流である。VT-xは仮想化ソフトウェアとゲストOSに異なる動作権限を割当てる機能であり、プロセッサのハードウェアに実装さ

50

れている（例えば、特許文献1または非特許文献1、2）。VT-xに対応したCPUは、ゲストOS間と仮想化ソフトウェアの動作権限の移行を検出してCPUのレジスタ状態の退避や回復を行い、仮想サーバ毎に独立した動作環境を提供する。

【0005】

一方、I/Oの割付方針は仮想化ソフトウェアによって異なる。仮想化ソフトウェアのI/Oの割り付けは、

- (1) 物理サーバのI/Oデバイスを直接使用させるI/O直接割付型と、
 - (2) 物理サーバのI/Oデバイスの種類やリビジョンを隠蔽する仮想I/O割付型と、
- に大別される。

【0006】

上記(1)のI/O直接割付型は、物理サーバで動作済みのI/Oについて、ファイルシステムの再構築等を行うことなく容易にサーバ統合を実現できる長所を有する。一方、上記(2)の仮想I/O割付型では、物理サーバのI/O種類に依らず一定のI/O構成をゲストOSに提供できる長所を有する。

【0007】

ここで、上述のような仮想化ソフトウェアとしては、次のようなものが知られている。まず、iA-32サーバ向けの仮想化ソフトウェアの一つとして、VMware（登録商標）のESX Serverが知られている。これは、iA-32のCPUを含む物理サーバで、上述したVT-x機能を利用して、既存のOSを複数動作させることが可能である。

【0008】

汎用機の技術に基づく仮想化ソフトウェアとしては、論理分割運転機能が知られている（IBM System370等）。これは、単一の物理計算機を複数の論理区画（LPAR）に分割し、各LPAR上で既存のOS及び仮想化管理部（VMM=Virtual Machine Manager）を実行する。論理分割運転機能では、汎用機（物理計算機）が上記iA-32のVT-xに相当する機能（LPARモード）を用いて、仮想サーバ（LPAR）上での従来のOSおよびVMMを稼働させている。

【0009】

この他、シミュレータによってiA-32のVT-x機能を提供するものとしてSimOS（<http://simos.stanford.edu/introduction.html>）等が知られている。この種のシミュレータは、仮想サーバ上のゲストOSの命令列を解釈して任意のサーバおよびCPU機能を提供するソフトウェアである。

【0010】

さらに、iA-64（IPF=Itanium Processor Family）のプロセッサに実装された仮想化支援機能としては、非特許文献3に記載されるVT-i機能が知られている。

【特許文献1】特表2005-529401号公報

【非特許文献1】「Intel 64 and IA-32 Architectures Software Developer's Manual VOL 2B」、[online]、Intel corp. 発行、[平成19年5月1日検索]、インターネット <URL: <http://www.intel.com/design/processor/manuals/253667.pdf> >

【非特許文献2】「AMD-Virtualization (AMD-V)」、[online]、Advanced Micro Devices, Inc. 発行、[平成19年5月1日検索]、インターネット <URL: http://www.amd.com/jp-ja/assets/content_type/DownloadableAssets/jp-ja_AMD-V_general_200702.pdf >

【非特許文献3】「Intel Itanium Architecture Software Developer's Manual」、[online]、Intel corp. 発行、[平成19年5月1日検索]、インターネット <URL: <ftp://download.intel.com/design/Itanium/manuals/24531705.pdf> >

【発明の開示】

【発明が解決しようとする課題】

【0011】

10

20

30

40

50

ところで、近年のサーバ統合の要求から、次世代のサーバOS (Windows (登録商標) Server Longhorn またはWindows (登録商標) Server 2008等) では上述の(2)のような仮想I/O割付型の仮想化ソフトウェアの機能を搭載することが検討されている。

【0012】

しかしながら、上記従来の仮想化ソフトウェアでは、次のような問題がある。

【0013】

上記ESXサーバ等の従来の仮想化ソフトウェアでは、物理計算機を構成するCPUの仮想化支援(VT-x)機能を用いて、複数の仮想サーバ上で既存のOSを稼働させることは可能であるが、仮想サーバに対してVT-x機能を提供することはできない。このため、ESXサーバなどでは次世代のサーバOSのように仮想化ソフトウェア(仮想化機能)を統合したOSを実行することが難しい、という問題がある。

10

【0014】

また、上記論理分割運転機能では、LPAR上のOSやVMMがLPARモードを利用することはできない。このため、仮想化ソフトウェアを統合したOSをLPAR上で稼働させるのは難しい、という問題があった。

【0015】

さらに、上記シミュレータの場合は、VT-x機能を仮想サーバに対して提供することができ、仮想化ソフトウェアを統合した次世代のOSを稼働させることは可能である。しかしながら、上記シミュレータは、仮想サーバ(ゲストOS)の命令列を解釈し、物理計算機のCPUが実行可能な命令列に変換(バイナリトランスレーション)するためのオーバーヘッドが発生するため、仮想サーバの性能(処理能力)が低下する、という問題があり、複数の仮想サーバを稼働させるサーバ統合を行うことは現実的ではない。

20

【0016】

また、仮想I/O割り付け型の仮想化機能を統合した次世代OSでは、開発環境の提供など同一のI/O構成の仮想サーバを多数生成する用途には向いているが、既存のサーバで動作済みの従来のOS(VT-x機能を利用しないOS=NTサーバ等)を統合する用途には不向きである。特に、現在、企業が運用するソフトウェア資産は、ほとんどが仮想化機能を含まないOSで稼働するソフトウェアである。このため、今後、サーバ統合を進めていく過程では、仮想化機能を含まない既存のOSと、新たに導入する仮想化機能を含んだ次世代OSが混在するサーバ統合環境が必要となる。しかし、上述のように、従来例では仮想化機能を統合した次世代OSを仮想計算機上に統合することは難しく、また、次世代OSで既存のOSを統合することは難しい、という問題があった。

30

【0017】

そこで本発明は、上記問題点に鑑みてなされたもので、仮想計算機の性能を低下させることなく、仮想化機能を統合したOSを仮想サーバで実行可能な仮想計算機を提供することを目的とし、さらに、仮想化機能を統合した次世代のOSと、プロセッサの仮想化支援機能を利用しない従来のOSを一つの物理計算機に統合することを目的とする。

【課題を解決するための手段】

【0018】

本発明は、物理プロセッサとメモリを備えた物理計算機で実行されて、複数の仮想プロセッサを提供する仮想化プログラムにおいて、第1の仮想マシンマネージャが第1の仮想プロセッサを生成する手順と、前記第1の仮想プロセッサ上で実行される第2の仮想マシンマネージャが第2の仮想プロセッサを生成する手順と、前記第2の仮想プロセッサがユーザプログラムを実行する手順と、前記第1の仮想マシンマネージャが、前記物理プロセッサから当該第1の仮想マシンマネージャの呼び出しを受信する手順と、前記第1の仮想マシンマネージャが、前記第1の仮想マシンマネージャ呼び出しの要因を解析する手順と、前記第1の仮想マシンマネージャが、前記解析の結果に基づいて前記第2の仮想マシンマネージャと、前記ユーザプログラムのいずれを実行するかを判定する手順と、前記第1の仮想マシンマネージャが、前記判定の結果に基づいて前記第1の仮想プロセッサに前記第2の仮想マシンマネージャまたは前記ユーザプログラムの一方を実行させる手順と、を

40

50

前記物理プロセッサに実行させる。

【 0 0 1 9 】

また、前記第 1 の仮想マシンマネージャが、前記ユーザプログラムを実行する際の第 1 の仮想プロセッサの状態を規定する第 1 の制御情報を前記メモリに設定する手順と、前記第 1 の仮想マシンマネージャが、前記第 2 の仮想マシンマネージャを実行する際の前記第 1 の仮想プロセッサの状態を規定する第 2 の制御情報を前記メモリに設定する手順と、前記第 1 の仮想マシンマネージャが、前記第 2 の仮想マシンマネージャまたは前記ユーザプログラムを実行する際の前記物理プロセッサの状態を規定する第 3 の制御情報を前記メモリに設定する手順と、をさらに含み、前記第 1 の仮想マシンマネージャが、判定結果に基づいて前記第 1 の仮想プロセッサに前記第 2 の仮想マシンマネージャまたは前記ユーザプログラム
10
の一方を実行させる手順は、前記第 1 の仮想マシンマネージャが、前記第 1 の仮想プロセッサに前記ユーザプログラムの実行を開始させるときには、前記第 1 の制御情報を参照する手順と、前記第 1 の仮想マシンマネージャが、前記第 1 の仮想プロセッサに前記第 2 の仮想マシンマネージャの実行を開始させるときには、前記第 2 の制御情報を参照する手順と、前記第 1 の仮想マシンマネージャが、前記参照した前記第 1 の制御情報または第 2 の制御情報の一方で前記第 3 の制御情報を更新する手順と、前記第 1 の仮想マシンマネージャが、前記物理プロセッサに対して前記第 3 の制御情報に従って命令の実行を開始するよう指示する第 1 の制御命令を発行する手順と、を含む。

【 発明の効果 】

【 0 0 2 0 】

したがって、本発明は、第 1 の仮想マシンマネージャの呼び出しの要因に応じて、第 2 の仮想プロセッサ上のユーザプログラムを実行するための第 1 の制御情報と、第 2 の仮想マシンマネージャを実行するための第 2 の制御情報の何れかを選択し、物理プロセッサを制御するための第 3 の制御情報を更新することで、第 1 の仮想プロセッサ上で仮想化機能を備えた次世代 OS をユーザプログラムとして実行することが可能となる。また、ユーザプログラムと第 2 の仮想マシンマネージャの切り換えは、第 1 または第 2 の制御情報で第 3 の制御情報を更新すればよいので、仮想計算機の性能を低下させることなく次世代 OS の稼働を実現できる。

【 0 0 2 1 】

そして、第 2 の仮想マシンマネージャを持たない既存の OS と、第 2 の仮想マシンマネージャを含む次世代 OS を、ひとつの計算機に統合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 】

【 0 0 2 2 】

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。

【 0 0 2 3 】

図 1 は、第 1 の実施形態を示し、本発明を適用する仮想計算機システムのブロック図である。物理サーバ（物理計算機）101 は、仮想化支援機能を備えて演算処理を実行する物理 CPU（プロセッサ）104 と、データやプログラムを格納するメモリ 105 と、物理サーバ 101 の外部の装置とデータの送受信を行うための I/O 装置 106 と、を含んで構成される。なお、I/O 装置 106 は、ネットワークインターフェースやホストバスアダプタ等で構成される。また、物理 CPU 104 は、複数の CPU で構成しても良いし、複数の演算コアを備えた CPU で構成しても良い。また、物理 CPU 104 の仮想化支援機能（VMX：Virtual Machine Extensions）は、上述の VT-x（Virtualization Technology for IA-32 Processors）または AMD-V（AMD Virtualization technology）あるいは VT-i（Virtualization Technology for Itanium architecture processors）を含むものである。なお、仮想化支援機能を利用する動作モードを VMX モードとし、仮想化支援機能を利用しない通常の特権レベルによる動作モードは通常動作モードとする。

【 0 0 2 4 】

物理サーバ 101 では、複数の仮想サーバ 102 a ~ 102 n を稼働させるため、物理
50

サーバ101の物理的な計算機資源を仮想化した計算機資源に変換し、各仮想サーバ102a~102nへ仮想計算機資源を割り当てるホストVMM(仮想マシンマネージャ:Virtual Machine Manager)10が実行される。このホストVMMは、メモリ105に読み込まれて物理CPU104で実行されるプログラムとして提供される。ホストVMM10は、各仮想サーバ102a~102nに対して、仮想CPU108a~108nを提供し、また、メモリ105とI/O装置106を各仮想サーバ102a~102nに割り当てる。なお、各仮想サーバ102a~102nに、物理サーバ101の計算機資源を割り当てる手法については、周知または公知のものを適宜用いればよいので、ここでは詳述しない。

【0025】

物理サーバ101は、ユーザーインターフェース301を提供する管理コンソール300に接続され、管理者などが計算機資源の割り当てなどの設定をユーザーインターフェース301を介してホストVMM10へ入力する。また、ユーザーインターフェース301は、ホストVMM10から受信した設定状態等を管理コンソール300の表示装置へ出力する。

【0026】

物理サーバ101のホストVMM10上で稼動する仮想サーバ102a~102nは、ユーザプログラム110a~110nとして、ゲストOS111a~111nがそれぞれ稼動し、各ゲストOS111a~111n上ではアプリケーション112a~112nが実行されている。各ゲストOS111a~111nは、ホストVMM10が提供する仮想CPU108a~108nで実行される。なお、仮想CPU108a~108nは、一つの仮想サーバに対して複数の仮想CPUを割り当てるようにすることができる。

【0027】

そして、仮想サーバ102aでは、ゲストOS111aとして仮想化機能(ゲストVMM20)を統合した次世代OSが実行され、仮想サーバ102nでは、ゲストOS111nとして仮想化機能を利用しない既存のOS(例えば、NTサーバ)が実行される。

【0028】

ホストVMM10は、既存のOSを実行する仮想サーバ102nに対しては、仮想CPU108nと、管理コンソール300から設定された計算機資源を割り当て、既存のゲストOS111nと、アプリケーション112nを実行する。

【0029】

一方、ホストVMM10は、次世代OSを実行する仮想サーバ102aへ割り当てる仮想CPU108aに、仮想化支援機能を提供する。仮想CPU108a上ではゲストVMM(第2仮想マシンマネージャ)20が稼働し、このゲストVMM20は仮想CPU208a~208iを提供する。次世代OSが稼働する仮想サーバ102aでは、第1の仮想CPU108a上で第2の仮想CPU208a~208iが複数提供され、各仮想CPU208a~208iでは、複数のユーザプログラム110a(ゲストOS111a及びアプリケーション112a)~110i(ゲストOS111i及びアプリケーション112i)が実行される。

【0030】

以下、本第1実施形態では、物理CPU104にVT-x機能を備え、仮想サーバ102aのゲストOS111aが仮想化機能を統合した次世代OSの例について説明する。

【0031】

VT-x機能を利用するホストVMM10は、メモリ105の所定の領域に仮想サーバ102a~102nの状態と物理CPU104を制御するための制御情報を格納するホストVMM保持データ11を格納する。そして、ホストVMM保持データ11には、物理CPU104を制御するための物理CPU制御データ13が格納される。物理CPU制御データ13は、仮想化支援機能を利用する仮想CPU108a~108nのステータスを示すデータ構造で、VMCB(Virtual Machine Control Block)またはVMCS(Virtual Machine Control Struc

10

20

30

40

50

ture)という。なお、本実施形態では、ホストVMM保持データ11内の物理CPU制御データ13を、図2で示すようにシャドウVMCB_0~_nとし、ゲストVMM20が操作する仮想CPU制御データ21をゲストVMCBとして区別する。また、物理CPU104は、物理CPU制御データ13のシャドウVMCB_0~_nを参照するためのポインタ115を備える。

【0032】

ホストVMM10は、物理CPU制御データ13内のシャドウVMCB_0~_nを書き換えることで、物理CPU104の動作モードを、ユーザプログラム110aまたはゲストVMM20を実行する動作モード(VMXノンルートモード)と、ホストVMM10を実行する動作モード(VMXルートモード)の何れかに設定する。

10

【0033】

仮想サーバ102aでは、ホストVMM10が提供する仮想CPU108a上で、次世代OSであるゲストOS111aに統合された仮想化機能(仮想化ソフトウェア)が、ゲストVMM20として稼動する。ゲストVMM20には、仮想CPU108aの仮想化支援機能を制御するためのVMCB(ゲストVMCB22)を含む仮想CPU制御データ21が格納される。ゲストVMCB22は、ホストVMM10によって割り当てられたメモリ105上の所定の領域に格納される。

【0034】

また、仮想CPU108aは、ゲストVMM20の仮想CPU制御データ21(ゲストVMCB22)を参照するためのポインタ116を備える。このポインタ116は、VT-x機能に対応するゲストOS111aが保持する仮想CPU108aの制御データ構造(ゲストVMCB22)へのポインタであり、ゲストOS111aからのVMPTRLD命令発行時に初期化される。なお、ポインタ116の初期化は、ゲストVMM20が実行する。

20

【0035】

図1において、次世代OSが稼動する仮想サーバ102aの仮想CPU208a~208iは、次世代OSに統合されたゲストVMM20が提供するものであり、各仮想CPU208a~208i上でユーザプログラム110a~110iを実行することができる。なお、仮想サーバ102aのゲストVMM20は、ゲストOS111aのアドインソフトウェアとして機能するものであっても良い。

30

【0036】

<本発明の概要>

仮想化を支援するVT-x機能では、ホストVMM10が物理サーバ101のメモリ105上に確保するシャドウVMCBを用いて、物理CPU104の動作モードを制御する。仮想化支援機能としてVT-x機能を備えた物理CPU104は、通常の動作モードと仮想化支援機能を提供するVMX(Virtual Machine Extensions)モードを備え、VMXモードでは、ホストVMM10が動作するホストモード(以下、VMXルートモード)と、ゲストVMM20またはユーザプログラム(ゲストOS111aまたはアプリケーション112a)が動作するゲストモード(以下、VMXノンルートモード)の何れかに切り替える。

40

【0037】

物理CPU制御データ13のシャドウVMCB中には仮想サーバ102a上のユーザプログラム110aの動作状態を規定するフィールド(ゲスト状態エリア131)が1種類しかなく、単純には次世代OSの仮想化機能であるゲストVMM20と、ユーザプログラム110a(ゲストOS111aまたはアプリケーション112a)を区分できない。

【0038】

そこで、本発明では、1つの仮想CPUが、ゲストVMM20とユーザプログラム110aを同時に実行することがないことに着目し、ホストVMM10が仮想化機能を統合したゲストOS111aを実行する仮想サーバ102aのゲストVMM20と、ユーザプログラム110aの実行の切替を監視して、動作モードの切替時にホストVMM保持データ

50

11のシャドウVMMCB_0~_n内のゲスト状態エリア131を書き換えることにより、仮想サーバ102a上で仮想化機能を稼動させるものである。

【0039】

このため、ホストVMM10は、仮想化機能を統合した次世代OSであるゲストOS111aを監視して、物理CPU104の動作モードを、ホストVMM10が動作するVMXルートモードと、ゲストVMM20またはユーザプログラム110aが動作するVMXノンルートモードに切り替え、所定の条件(VM-exit)でVMXルートモードではホストVMM10がゲストVMM20やゲストOSの命令をエミュレートする。これにより、ゲストOS111aに対して、仮想CPU108aが仮想化支援機能を提供するように見せかける。

10

【0040】

<ホストVMMの構成>

ホストVMM10は、上述のホストVMM保持データ11の他に、仮想サーバ102a~102nを監視して、物理CPU104の動作モードをVMXルートモードとVMXノンルートモードの何れかに切り替えるCPU制御部12を備える。

【0041】

さらに、ホストVMM10は、CPU制御部12が仮想サーバ102a~102nの状態を取得し、各仮想サーバ102a~102nへ指令を送信するためのインターフェースである制御・通信インターフェース14と、CPU制御部12が物理CPU104へ命令を発行するための命令発行インターフェース15と、物理CPU104が物理CPU制御データ13を参照または更新するための参照・更新インターフェース16と、I/O装置106からの割り込み要求などを受け付けて、この要求に応答するためのI/O要求・応答インターフェース17を備える。

20

【0042】

図2は、ホストVMM10とゲストVMM20の機能を示すブロック図である。ホストVMM10は、物理CPU104の仮想化支援機能(VT-x機能)を利用するため、メモリ105の所定の領域にホストVMM保持データ11を設定する。

【0043】

ホストVMM保持データ11は、ゲストOS111a~111nの仮想化支援機能の利用の有無や、仮想CPU108a~108nの状態を示すフラグ類を格納する領域と、仮想CPU108a~108n毎のステータス等を格納するシャドウVMMCB_0~_nを保持する物理CPU制御データ13の領域を含んで構成される。

30

【0044】

ホストVMM保持データ11のゲストOSや仮想CPUの状態を示すフラグとしては、例えば、ゲストOS111a~111n毎に物理CPU104の仮想化支援機能を利用可能か否かを識別する仮想化機能有効フラグ141と、仮想CPU108a~108n毎に仮想化支援機能を利用中か否かを設定するVMXONフラグ142と、仮想CPU108a~108n仮想化支援機能がVMXルートモードまたはVMXノンルートモードの何れで動作しているのかを示す動作モードフラグ143が挙げられる。

【0045】

仮想化機能有効フラグ141は、ゲストOS111a~111n毎に設定されて「1」であれば該当するゲストOSが仮想化支援機能を利用可能であることを示し、「0」であればゲストOSが仮想化支援機能を利用しないことを示す。この仮想化機能有効フラグ141は、ゲストOS111a~111n毎に管理コンソール300から設定したり、あるいは、予め設定したファイルなどから設定する。

40

【0046】

VMXONフラグ142は、各仮想CPU108a~108n毎にVMXモードの動作状態を示すフラグで、「1」のときには該当する仮想CPU108a~108nの動作モードがVMXモードであることを示し、「0」のときには該当する仮想CPU108a~108nの動作モードが仮想化支援機能を利用しない通常動作モードであることを示す。

50

VMXONフラグ142は、ゲストOS111a~111nがVMXON命令を発行したときにホストVMM10が「1」をセットし、ゲストOS111a~111nがVMXOFF命令を発行したときにホストVMM10が「0」にリセットする。

【0047】

動作モードフラグ143は、仮想CPU108a~108n上で動作するプログラムの動作モードを追跡するフラグである。この動作モードフラグ143は、ゲストOS111a~111nのVM-entry時にホストVMM10が「1」にセットし、ゲストOSのVM-exit時にホストVMM10が「0」にリセットする。つまり、この動作モードフラグ143は、VMXONフラグ142が「1」のときに各仮想CPU108a~108n毎のVMXモードの種別を示すもので、動作モードフラグ143が「0」のときには仮想CPU108a~108nがゲストVMM20を実行する状態（仮想CPUのVMXルートモード）を示し、動作モードフラグ143が「1」のときには仮想CPUがユーザプログラム110a（ゲストOS111aまたはアプリケーション112a）を実行する状態（仮想CPUのVMXノンルートモード）を示す。

10

【0048】

ここで、VMXモードのVMXルートモードとVMXノンルートモードの遷移については、上記非特許文献1に記載されるとおりである。ここでは、概要のみについて説明する。通常動作モードからVMXモードへ移行する際には、ホストVMM10がVMXON命令を発行し、物理CPU104の動作モードをVMXモードに移行させる。そして、VMXモードへ移行したホストVMM10は、該当する仮想CPU108a~108nのシャドウVMCB_0~_n-1にユーザプログラム110aを実行させるための情報を書き込んでからVM-entry命令（VMLAUNCH命令またはVMRESUME命令）を発行し、VMXルートモードからVMXノンルートモードへ移行する。このVMXルートモードからVMXノンルートモードへの遷移をVM-entryという。

20

【0049】

逆に、VMXノンルートモードからVMXルートモードへの遷移は、VM-exitという。このVM-exitは、ゲストOS111a~111nが特権命令を発行したときなどの所定の要因で物理CPU104がホストVMM10へVM-exitを通知する。ホストVMM10のCPU制御部12がVM-exitを検知すると、所定のエミュレーションを行ってゲストVMM20またはゲストOSの処理を完了した後、シャドウVMCBを必要に応じて書き換えてからVM-entry命令（第1制御命令）を発行してVMXノンルートモードからVMXルートモードへ切り替える。

30

【0050】

本発明では、次世代OSであるゲストOS111aのVM-entryの際に、ホストVMM10がゲストVMCB22のゲスト状態エリア221とホスト状態エリア222を読み込んで、ゲストOS111aの動作に応じた一方のエリアの内容をシャドウVMCB_0のゲスト状態エリア131へ設定することで、仮想サーバ102aでゲストOS111aの仮想化機能を実現するのである。

【0051】

VT-x機能では以上のような、VM-entryとVM-exitの遷移により、ホストVMM10とゲストVMM20またはユーザプログラム110aを切り替える。このため、VM-entryとVM-exitの前後の物理CPU104のステータス等を保持するために、物理CPU制御データ13のデータ構造であるシャドウVMCB_0~_nを使用する。

40

【0052】

物理CPU制御データ13は、仮想CPU108a~108n毎にシャドウVMCB_0~_nが設定され、各シャドウVMCBには次のようなデータが格納される。

【0053】

ゲスト状態エリア131には、図3で示すように、仮想CPU108a~108nのレジスタ状態などのステータスが格納される。すなわち、後述するように、ゲストVMM2

50

0のステータスまたはユーザプログラム110aのステータスが選択的に格納される。

【0054】

ホスト状態エリア132には、図3で示すように、ホストVMM10の物理CPU104のレジスタ状態などのステータスが格納される。VM実行制御領域133には、仮想サーバ102a~102nの設定情報が格納され、例えば、例外ビットマップやI/Oビットマップなどの情報が含まれる。VM-exit制御領域134には、VM-exitの要因などの情報が格納される。VM-entry制御領域135には、VM-entryの動作を制御するための情報が格納される。そして、VM-exit情報領域136には、VM-exitを発生する要因(命令またはイベント)を格納する。VM-exitを発生させる要因としては、例えば、図4に示す「description」に示すものがVM-exit情報領域136に設定される。

10

【0055】

以上のようなシャドウVMCB_0~_nにより、ホストVMM10は各仮想サーバ102a~102nの制御を行う。

【0056】

一方、仮想サーバ102aのゲストVMM20が管理する仮想CPU制御データ21には、上記物理CPU制御データ13のシャドウVMCB_0と同様のデータ構造であるゲストVMCB22が格納される。

【0057】

ゲストVMM20は、ユーザプログラム110a(ゲストOS111aまたはアプリケーション112a)を実行する仮想CPU108aのレジスタ状態などのステータス等を格納するゲスト状態エリア221と、ゲストVMM20を実行する仮想CPU108aのレジスタ状態などのステータスを格納するホスト状態エリア222と、仮想サーバ102aの設定情報を格納するVM実行制御領域223と、仮想サーバ102aにおけるVM-exitの要因などの情報を格納するVM-exit制御領域224と、仮想サーバ102aにおけるVM-entryの動作を制御するための情報を格納するVM-entry制御領域225と、仮想サーバ102aにおけるVM-exitの原因を特定する情報を格納するVM-exit情報領域226とを備える。

20

【0058】

ゲスト状態エリア221には、上記図3と同様に、仮想CPU108aのレジスタ状態などのステータスが格納される。すなわち、後述するように、ゲストOS111aのステータスまたはアプリケーション112aのステータスが選択的に格納される。

30

【0059】

仮想サーバ102aの状態を制御するホストVMM10のシャドウVMCB_0では、ゲスト状態エリア131にゲストVMM20のステータスまたはユーザプログラム110a(ゲストOS111aまたはアプリケーション112a)のステータスが格納され、ホスト状態エリア132にはホストVMM10のステータスが格納される。

【0060】

一方、ゲストVMM20のゲストVMCB22では、ゲスト状態エリア221にゲストOS111aのステータスまたはアプリケーション112aのステータスが格納され、ホスト状態エリア222にはゲストVMM20のステータスが格納される点でシャドウVMCB_0と相違する。

40

【0061】

次に、ホストVMM10のCPU制御部12の構成について説明する。CPU制御部12は、管理コンソール300からの入力などに基づいて各仮想サーバ102a~102nへ物理サーバ101の計算機資源を割り当てるリソース管理部(図示省略)に加えて、ゲストOS111aの仮想化機能(ゲストVMM20)を稼働させるため、仮想CPUのステータスの読み込み先を選択する状態エリア選択部121と、エミュレータ122と、シャドウVMCB参照・更新部123と、VM-exitハンドラ124と、VM-entry命令発行部125と、ユーザインターフェース301とを備える。なお、リソース管

50

理部については、上述のように公知または周知の技術を適用すればよいので、本実施形態では詳述しない。

【0062】

VM-exitハンドラ124は、物理CPU104からVM-exitを受信して、エミュレータ122を起動する。

【0063】

エミュレータ122は、VM-exitハンドラ124から受信したVM-exitの発生要因となった命令またはイベントを特定し、後述するように特定した発生要因に対応するモジュールを起動してゲストVMM20やユーザプログラム110aに代わって処理を実行する。

10

【0064】

エミュレータ122が、仮想CPU108aの動作モードの変更（VMXルートとVMXノンルートの切り替え）を検出した場合は、特定した要因に応じたモジュール（VM-entry命令実行モジュール、VMCLEAR命令実行モジュール、CPUID命令実行モジュール、VM-exit条件検出モジュール）を起動して、状態エリア選択部121を起動する。

【0065】

ここで、VM-exitの発生要因は、図4に示すように、シャドウVMCB_0のVM-exit情報領域136の「Exit reason」に設定されたものである。図4に示すVM-exitの要因リストには、VM-entry命令の発行に起因する要因1361と、VMCLEAR命令（第4制御命令）の発行に起因する要因1362と、ゲストVMM20へのVM-exit通知の有無を設定した通知条件1363が予め設定される。

20

【0066】

例えば、仮想CPU108aの動作モードを切り替える際のVM-exitの要因として、VM-entry命令（VMLAUNCH命令またはVMRESUME命令）を検出した場合には、図4の要因1361に該当するので、ホストVMM10はVM-entry命令実行モジュールでエミュレーションを行い、ゲストVMM20またユーザプログラム110aに代わって処理を行う。

【0067】

あるいは、VM-exitの要因が、図4の通知条件1363に該当するゲストVMMへの通知条件の場合には、VM-exit条件検出モジュールを起動して同様にエミュレーションを行う。

30

【0068】

状態エリア選択部121は、ホストVMM保持データ11からVM-exitを発生した仮想CPU（この例では108a）の動作モードフラグ143を読み込んで、VMXルートモードとVMXノンルートモードの何れであるかを判定する。動作モードフラグ143が「0」のVMXルートモードであればゲストVMM20が動作していたので、状態エリア選択部121は、VM-exitを発生した仮想サーバ（この例では102a）のゲストVMCB22からホスト状態エリア222を読み込む。

40

【0069】

一方、動作モードフラグ143が「1」のVMXノンルートモードであればユーザプログラム110aが動作していたので、状態エリア選択部121は、該当する仮想サーバ102aのゲストVMCB22からゲスト状態エリア221を読み込む。

【0070】

ゲストVMCB22からの読み込みが完了すると、CPU制御部12ではシャドウVMCB参照・更新モジュール123を起動する。シャドウVMCB参照・更新モジュール123は、状態エリア選択部121が読み込んだゲストVMCB22の情報を、VM-exitの処理対象である仮想CPU108aに対応するシャドウVMCB_0のゲスト状態エリア131に書き込んでシャドウVMCBを更新する。

50

【 0 0 7 1 】

シャドウVMCB_0のゲスト状態エリア131の更新が完了すると、CPU制御部12は仮想CPU108aの動作モードをVMXルートモードからVMXノンルートモードへ切り替えるため、動作モードフラグ143を「1」に更新し、物理CPU104のポインタ115に操作対象の仮想CPU108aのシャドウVMCB_0のアドレスをセットする。

【 0 0 7 2 】

そして、CPU制御部12は、VM-entry命令発行モジュール125を起動して、物理CPU104に対してVM-entry命令（VMRESUME命令）を発行する。

10

【 0 0 7 3 】

物理CPU104はVM-entry命令を受け付けると、ポインタ115が指し示すシャドウVMCB_0のゲスト状態エリア131を読み込んで、状態エリア選択部121が選択した仮想サーバ102aのゲストVMM20またはユーザプログラム110aを実行する。

【 0 0 7 4 】

以上のように、仮想サーバ102aで仮想化機能を統合した次世代OSのゲストOS111aが稼働している場合、ホストVMM10のCPU制御部12は、物理CPU104のVM-exitを検出すると、シャドウVMCB_0の動作モードフラグ143を参照して、仮想CPU108aで実行中のプログラムがゲストVMM20とユーザプログラム110aの何れであるかを判定する。そして、CPU制御部12は、仮想CPU108aで実行していたプログラムに応じてゲスト状態エリア221またはホスト状態エリア222の内容をホストVMM10のシャドウVMCB_0のゲスト状態エリア131へ書き込んでからVM-entry命令を実行する。

20

【 0 0 7 5 】

こうして、仮想サーバ102aのゲストVMM20が保持するゲストVMCB22の情報で、ホストVMM10が保持するシャドウVMCB_0のゲスト状態エリア131を更新することにより、ゲストOS111aの仮想化機能を稼働させる仮想サーバ102aと、仮想化機能を用いない既存のOSを実行する仮想サーバ102nをひとつの物理サーバ101に統合することが可能となる。

30

【 0 0 7 6 】

<ホストVMMの処理の詳細>

次に、ホストVMM10で行われる処理について、図5を参照しながら以下に説明する。図5は、仮想サーバ102a～102nの稼働中に物理CPU104からのVM-exitを受信したときにホストVMM10のCPU制御部12が行う処理の一例を示すフローチャートである。なお、この例では、VM-exitの要因がVM-entry命令である場合を示す。

【 0 0 7 7 】

まず、S1では、ホストVMM10のCPU制御部12は、VM-exitの要因となった操作対象のゲストOS111a～111n（以下、ゲストOSとする）に対応するホストVMM保持データ11の仮想化機能有効フラグ141を参照し、VM-exitを発生したゲストOSがVT-x機能を利用可能であるか否かを判定する。操作対象のゲストOSがVT-x機能を利用可能である場合にはS2へ進む。一方、操作対象のゲストOSがVT-x機能を利用しない場合（NTサーバまたは2000サーバなど）には、S15へ進んで、ホストVMM10は前記従来例に示した特許文献1または非特許文献1に記載されるような仮想マシン処理を行う。

40

【 0 0 7 8 】

S2では、ホストVMM10がVM-exitの要因となった仮想CPU108a～108n（以下、仮想CPUとする）のVMXONフラグ142を参照し、この操作対象の仮想CPUがVMXモードであるか否かを判定する。VMXONフラグ142が「1」で

50

あれば仮想CPUはVMXモードであるのでS3へ進む。一方、VMXONフラグ142が「0」の場合には、仮想CPUが通常の動作モードであるので、上記と同様にS15へ進み、ホストVMM10は既存の仮想マシン処理を実行する。

【0079】

S3では、ホストVMM10のCPU制御部12が、操作対象の仮想CPUの動作モードフラグ143を参照し、仮想CPUの動作モードがVMXルートモードとVMXノンルートモードの何れであるかを判定する。動作モードフラグ143が「0」の場合には仮想CPUがVMXルートモードであると判定してS4へ進む。一方、動作モードフラグ143が「1」の場合には、仮想CPUがVMXノンルートモードであると判定してS11へ進む。

10

【0080】

S4では、CPU制御部12が物理CPU104から受信したVM-exitの発生要因を特定する。この例では、仮想CPUの動作モードがVMXルートモードであるので、CPU制御部12は、ゲストVMM20の状態(ステータス)を格納したゲスト状態エリア131を参照し、ゲストVMM20が発行したVM-entry命令をVM-exitの要因として特定する。

【0081】

次に、S5ではCPU制御部12は、エミュレータ122からVM-entry命令実行モジュールを実行し、操作対象の仮想CPUをVMXノンルートモードへ切り替えるのに必要な所定の処理(状態エリア選択部121の起動など)をゲストVMM20に代わってエミュレートする。

20

【0082】

次に、S6ではホストVMM10のCPU制御部12は、仮想CPUの動作モードをVMXルートモードからVMXノンルートモードへ移行させるので、操作対象の仮想CPUの動作モードフラグ143を「1」に更新する。そして、S7では、CPU制御部12が操作対象のゲストVMM20の仮想CPU制御データ21からゲスト状態エリア221に格納されているゲストOSまたはアプリケーションのステータスを読み込む。

【0083】

次に、S8において、CPU制御部12は、物理CPU104に対してVMPTLDR命令を発行して、操作対象の仮想CPUに対応するシャドウVMCBをアクティブに設定し、アクティブにしたシャドウVMCBのアドレスをポインタ115に設定する。このVMPTLDR命令(第2制御命令)により、ホストVMM10は、複数のシャドウVMCB_0~_n-1の中から操作対象の仮想CPU(仮想サーバ)のシャドウVMCBを選択する。

30

【0084】

S9では、CPU制御部12が、上記S7で読み込んだゲスト状態エリア221の情報で操作対象のシャドウVMCB_0のゲスト状態エリア131を更新する。そして、S10でCPU制御部12は、物理CPU104に対してVM-entry命令を発行する。

【0085】

VM-entry命令を受け付けた物理CPU104は、ポインタ115で指定されたシャドウVMCBのゲスト状態エリア131の内容に基づいて、操作対象の仮想サーバのユーザプログラム(ゲストOSまたはアプリケーション)を実行する。なお、物理CPU104は、VM-entry命令を受け付けると、実行していたホストVMM10のステータスをシャドウVMCBのホスト状態エリア132へ格納し、次の呼び出しに備える。

40

【0086】

一方、上記S3で動作モードフラグ143が「1」と判定された場合は、操作対象の仮想CPUがユーザプログラムを実行中のVMXノンルートモードであるので、S11の処理を行う。

【0087】

50

S 1 1では、CPU制御部12が図4のVM - e x i tの要因リストを参照し、通知条件1363からゲストVMMへのVM - e x i t通知条件を検索する。この例では、VM - e x i tの要因となったVM - e n t r y命令(VMLAUNCH、VMRESUME命令)であるので、VM - e x i t通知条件に一致する。

【0088】

そして、CPU制御部12は、S12で図2のVM - e x i t条件検出モジュールを実行し、操作対象の仮想CPUをVMXルートモードへ切り替えるのに必要な所定の処理(状態エリア選択部121の起動など)をユーザプログラムに代わってエミュレートする。

【0089】

次に、S13では、操作対象の仮想CPUの動作モードをVMXノンルートモードからVMXルートモードへ移行させるため、CPU制御部12が操作対象の仮想CPUの動作モードフラグ143を「0」にリセットする。そして、S14では、CPU制御部12が操作対象のゲストVMM20の仮想CPU制御データ21からホスト状態エリア221に格納されているゲストVMM20のステータスを読み込む。

【0090】

S14の処理が完了すると、上記S8～S10を実行し、CPU制御部12は、上記S14で読み込んだホスト状態エリア222の情報で操作対象のシャドウVMCBのゲスト状態エリア131を更新し、ゲストVMM20のステータスをセットしてから物理CPU104に対してVM - e n t r y命令を発行する。

【0091】

この結果、VM - e n t r y命令を受け付けた物理CPU104は、ポインタ115で指定されたシャドウVMCBのゲスト状態エリア131の内容に基づいて、操作対象の仮想サーバのゲストVMM20を実行する。

【0092】

このように、ゲストOSが仮想化機能を統合した次世代OSである場合、ホストVMM10は、仮想CPUの動作モードと、発生したVM - e x i tの要因に応じてシャドウVMCBのゲスト状態エリア131に書き込むステータスを、ゲストVMMまたはユーザプログラムの何れか一方から選択する。そして、ホストVMM10が物理CPU104(第1の仮想CPU)に対してVM - e n t r y命令を発行することで、仮想サーバ上でゲストVMMと、このゲストVMMが提供する第2の仮想CPU上で稼動するユーザプログラムを切り替えて実行することが可能となり、ゲストVMMは仮想サーバ上に複数の仮想化環境(ユーザプログラム)を提供することができる。

【0093】

なお、上記S1の判定でゲストOSが仮想化機能を利用しない場合、または上記S2の判定で、仮想CPUが物理CPU104のVT - x機能を利用しない場合は、S15において、前記従来例に示した特許文献1または非特許文献1の記載に準じた仮想マシン処理をホストVMM10で行えばよい。

【0094】

S15の仮想マシン処理は、例えば、図1に示した仮想サーバ102nのゲストOS111nは、既存のOSであり、このゲストOS111nまたはアプリケーション112n(ユーザプログラム110n)が特権命令などの所定の命令を実行すると、上述したように物理CPU104はVM - e x i tの発生をホストVMM10に通知する。

【0095】

ホストVMM10は、物理CPU104からのVM - e x i tの通知を受けると、ユーザプログラム110n(仮想CPU108n)のステータスをシャドウVMCB_n - 1のゲスト状態エリア131に格納する。そして、ホストVMM10は、ポインタ115のアドレスをホストVMM10のステータスを格納したホスト状態エリアに設定して、所定の処理を実行する。

【0096】

ホストVMM10は、特権命令など所定の処理が完了すると、ホストVMM10のステ

10

20

30

40

50

ータスをホスト状態エリア132へ格納し、ポインタ115のアドレスにゲスト状態エリア131を設定してから、VM-entry命令(VMRESUME命令)を発行して、制御を仮想CPU108nに移動してユーザプログラム110nの実行を再開する。

【0097】

このように、本発明によれば、仮想化機能を統合した次世代OSと、既存のOSをひとつの物理サーバ101へ統合することが可能となつて、物理サーバの数を削減してサーバの運用管理コストを削減することが可能となる。

【0098】

さらに、上述のように、ホストVMM10は次世代OSに対して仮想CPUがVT-x機能を提供するように見せかけることができ、仮想化ソフトウェアを統合したOSを確実に稼働させることが可能となる。また、本発明のホストVMM10によれば、前記従来例のシミュレータのように、命令列の変換などのオーバーヘッドがなく、仮想計算機の性能を低下させることなく、仮想化機能を統合したOSを仮想サーバで実行することが可能となるのである。

【0099】

また、本発明では、複数の仮想CPU108a~108nに対応して、複数のシャドウVMCB_0~_n-1を設けたので、物理サーバ101で複数のゲストVMM20を実行する場合でも、シャドウVMCB_0~_n-1を切り替えるだけで迅速に処理を切り替えることが可能となる。これにより、物理サーバ101に複数の次世代OSを統合した場合でも、仮想サーバの性能を維持することが可能となる。

【0100】

図6は、仮想サーバの稼働中にVMCLEAR命令の実行に起因して、物理CPU104からのVM-exitを受信したときにホストVMM10のCPU制御部12が行う処理の一例を示すフローチャートである。

【0101】

S21~S23は、上記図5のS1~S3と同様であり、ホストVMM10のCPU制御部12が、VM-exitの要因となった操作対象のゲストOSがVT-x機能を利用可能であるか否かを判定し、操作対象の仮想CPUのVMXONフラグ142を参照して、操作対象の仮想CPUがVT-x機能を使用中か否かを判定する。操作対象のゲストOSがVT-x機能を利用しない場合や、VMXONフラグ142が「0」である仮想CPUが通常の動作モードの場合は、S35へ進み、ホストVMM10は既存の仮想マシン処理を実行する。

【0102】

さらに、CPU制御部12は、操作対象の仮想CPUの動作モードフラグ143を参照し、仮想CPUの動作モードがVMXルートモードとVMXノンルートモードの何れであるかを判定し、VMXルートモードであればS24へ進み、VMXノンルートモードであればS30へ進む。

【0103】

S24では、CPU制御部12は、物理CPU104から受信したVM-exitから要因を特定する。この例では、CPU制御部12がゲスト状態エリア131を参照して、VM-exitの要因としてゲストVMM20のVMCLEAR命令を特定する。

【0104】

次に、S25ではCPU制御部12のエミュレータ122からVMCLEAR命令実行モジュールを起動してエミュレーションを行う。S26では、エミュレーションにより、物理CPU104の動作状態(ステータス)を読み込んで、操作対象の仮想CPUに対応するシャドウVMCBを更新する。これにより、物理CPU104のステータスをシャドウVMCBに反映させる。

【0105】

S27では、CPU制御部12が操作対象の仮想CPUに対応するシャドウVMCBのゲスト状態エリア131に格納されているステータスを読み込む。

10

20

30

40

50

【0106】

次に、S28において、CPU制御部12は、シャドウVMCBのゲスト状態エリア131から読み込んだステータスをゲストVMM20のゲスト状態エリア221へ書き込んで更新する。これにより、ゲストVMM20のゲスト状態エリア221は、物理CPU104のステータスと同期する。

【0107】

S29では、CPU制御部12が、ポインタ115のアドレスを操作対象の仮想CPUに対応するシャドウVMCBのゲスト状態エリア131に設定してから、物理CPU104に対してVM-entry命令を発行する。

【0108】

VM-entry命令を受け付けた物理CPU104は、ポインタ115で指定されたシャドウVMCBのゲスト状態エリア131の内容に基づいて、操作対象の仮想サーバのユーザプログラム（ゲストOSまたはアプリケーション）を実行する。

【0109】

一方、上記S23で動作モードフラグ143が「1」と判定された場合は、操作対象の仮想CPUがVMXノンルートモードであるので、S30の処理を行う。S30では、図4のVM-exitの要因リストを参照し、VM-exitの要因を解析する。CPU制御部12は、通知条件1363からゲストVMMへのVM-exit通知条件を検索する。この例では、VM-exitの要因となったVMCLEAR命令は、VM-exit通知条件に一致する。

【0110】

そして、S31ではCPU制御部12がゲストVMM20に対して規定外命令エラーを通知する。すなわち、ユーザプログラムの権限ではVMCLEAR命令を実行できなかったことを、ゲストVMM20へ通知する。

【0111】

次に、S32ではCPU制御部12が、操作対象の仮想CPUの動作モードをVMXノンルートモードからVMXルートモードへ移行させるので、操作対象の仮想CPUの動作モードフラグ143を「0」にリセットする。

【0112】

S33では、CPU制御部12が操作対象のゲストVMM20の仮想CPU制御データ21から、ホスト状態エリア221に格納されているゲストVMM20のステータスを読み込む。

【0113】

S34では、CPU制御部12が、上記S33で読み込んだホスト状態エリア222のステータスを、操作対象のシャドウVMCBのゲスト状態エリア131に書き込んでシャドウVMCBを更新する。その後、上記S29では、CPU制御部12がポインタ115のアドレスをゲスト状態エリア131に設定してから、物理CPU104に対してVM-entry命令を発行する。

【0114】

この結果、VM-entry命令を受け付けた物理CPU104は、ポインタ115で指定されたシャドウVMCBのゲスト状態エリア131の内容に基づいて、操作対象の仮想サーバのゲストVMM20を実行する。

【0115】

なお、上記S21の判定でゲストOSが仮想化機能を利用しない場合、または上記S22の判定で、仮想CPUが物理CPU104のVT-x機能を利用しない場合は、S35において、前記図5のS15と同様にして、従来例に示した特許文献1または非特許文献1の記載に準じた仮想マシン処理をホストVMM10で行えばよい。

【0116】

このように、ゲストVMM20がVMCLEAR命令を発行した場合には、ホストVMM10は、物理CPU104の動作状態（ステータス）を取得してシャドウVMCBのゲ

10

20

30

40

50

スト状態エリア131を、取得した動作状態で更新する。そして、ホストVMM10は、シャドウVMCBのゲスト状態エリア131のステータスを、ゲストVMCBのゲスト状態エリア221に反映させることで、物理CPU104のステータスとホストVMM10のゲスト状態エリア131及びゲストVMM20上のゲスト状態エリア221のステータスを同期させることができる。

【0117】

これにより、仮想化機能を統合した次世代OSは、仮想サーバ上で仮想環境を円滑に提供することが可能となる。

【0118】

図7は、仮想サーバの稼働中にCPUID命令(第3制御命令)の実行に起因して、物理CPU104からのVM-exitを受信したときにホストVMM10のCPU制御部12が行う処理の一例を示すフローチャートである。

10

【0119】

まず、S41では、図5のS3と同様に、ホストVMM10のCPU制御部12は、VM-exitの要因となった操作対象の仮想CPUの動作モードフラグ143を参照し、仮想CPUの動作モードがVMXルートモードとVMXノンルートモードの何れであるかを判定する。動作モードフラグ143が「0」の場合には仮想CPUがVMXルートモードであると判定してS42へ進む。一方、動作モードフラグ143が「1」の場合には、仮想CPUがVMXノンルートモードであると判定してS41へ進む。

【0120】

20

S42では、図5のS4と同様にして、CPU制御部12が物理CPU104から受信したVM-exitの発生要因をゲスト状態エリア131の値から特定する。この例では、VMXルートモードであるので、CPU制御部12は、ゲストVMM20のCPUID命令をVM-exitの発生要因として特定する。

【0121】

S43では、図5のS1と同様に、CPU制御部12は、VM-exitの要因となった操作対象のゲストOSに対応するホストVMM保持データ11の仮想化機能有効フラグ141を参照し、VM-exitを発生したゲストOSがVT-x機能を利用可能であるか否かを判定する。操作対象のゲストOSがVT-x機能を利用可能である場合にはS44へ進む。一方、操作対象のゲストOSがVT-x機能を利用しない場合には、S45へ進む。

30

【0122】

S44では、CPUID命令の戻り値を示すリターンレジスタ(IA-32の場合ECX)の所定のビット(CPUID.1.ECX[5]=1)を1に設定して、VT-x機能が有効であることを設定する。

【0123】

一方、S45では、CPUID命令の戻り値を示すリターンレジスタ(IA-32の場合ECX)の所定のビットを0に設定して、VT-x機能が無効であることを設定する。

【0124】

次に、S46においてCPU制御部12は、エミュレータ122からCPUID命令実行モジュールを起動して、ゲストVMM20もしくはユーザプログラムへ通知すべき例外のチェック等を行い、チェック結果をゲスト状態エリア131に格納する。

40

【0125】

S47では、CPU制御部12が物理CPU104に対してVM-entry命令を発行して、仮想サーバに制御を移す。

【0126】

一方、上記S41の判定で、動作モードフラグ143が「1」となるユーザプログラムを実行していた場合(VMXノンルートモード)にはS48へ進み、ユーザプログラムを実行している期間中のVM-exitの要因を解析する。S48では、CPU制御部12が物理CPU104から受信したVM-exitの発生要因をゲスト状態エリア131の

50

値から特定する。この例では、VMXノンルートモードであるので、CPU制御部12は、ゲストOSまたはアプリケーションのCPUID命令をVM-exitの発生要因として特定する。

【0127】

S49では、CPU制御部12が、ユーザプログラム実行中のCPUID命令に起因するVM-exitについて、ゲストVMM20への通知が必要であるか否かを判定する。この判定は、CPU制御部12が、仮想CPU制御データ21のゲストVMCB22のVM実行制御領域223に予め設定された情報を参照して、判定を行う。CPUID命令に伴うVM-exitをゲストVMM20へ通知する必要がある場合はS50へ進む。一方、ゲストVMM20へのVM-exitの通知が不要であればS43へ進み、上述のS43移行の処理を実行する。

10

【0128】

S50では、ゲストVMM20へVM-exitの通知を行うため、CPU制御部12はVM-exitエミュレータを実行し、S51で操作対象の仮想CPUの動作モードフラグ143を「0」にリセットしてVMXルートモードへ切り替える。

【0129】

次に、S52でCPU制御部12は、操作対象のゲストVMCB22のホスト状態エリア222のステータスを読み込む。S53ではCPU制御部12が、上記読み込んだステータスを、操作対象の仮想CPUのシャドウVMCBのゲスト状態エリア131に書き込んで更新を行う。その後、S47で、CPU制御部12が物理CPU104に対してVM

20

【0130】

以上のように、ゲストVMM20がCPUID命令を発行した場合と、CPUID命令に起因するVM-exitをゲストVMM20へ通知しない場合では、仮想化機能有効フラグ141の設定に応じてリターンレジスタの値を設定した後、CPU制御部12は、VM-entry命令を発行して、ゲストVMM20またはユーザプログラムへ制御を戻すことができる。また、ユーザプログラムがCPUID命令を発行し、ゲストVMM20へVM-exitの通知が必要な場合では、シャドウVMCBのゲスト状態エリア131を、ゲストVMCB22のホスト状態エリア222の内容で更新して、制御をゲストVMM20へ移すことができる。

30

【0131】

以上のように、本第1実施形態によれば、ホストVMM10がゲストOSやアプリケーションの状態と仮想CPUの状態を監視してシャドウVMCBのゲスト状態エリア131を書き換えることで、仮想サーバの性能を低下させることなく、仮想化ソフトウェアを統合したOSを確実に稼働させることが可能となる。そして、ひとつの物理サーバ101で仮想化機能を統合した次世代OSと、既存のOSを共存させることが可能となって、サーバ統合を効率よく行うことで、サーバの運用コストを削減することが可能となる。

【0132】

なお、ホストVMM10が管理コンソール300に提供するユーザーインターフェースとしては、例えば、図8に示すようなGUIで構成される。図8は、各仮想サーバ毎に仮想化機能を有効にするか無効にするかを設定するユーザーインターフェースを示す。図中「VM_」は、仮想サーバの識別子を示し、「現在設定」は「ON」であれば仮想化機能が有効であることを示し、「OFF」であれば仮想化機能が無効であることを示す。そして、「新規設定」に「ON」または「OFF」を設定してから、画面左下の「OK」をクリックすることにより、所望の仮想サーバについて仮想化機能の有効または無効を切り替えることができる。すなわち、ホストVMM10は、ユーザーインターフェース301から入力された「新規設定」の値に基づいて、ホストVMM保持データ11のうち、該当する仮想サーバの仮想化機能有効フラグ141に「1」または「0」を設定する。これにより、ユーザーインターフェース301から入力された設定値を任意の仮想サーバへ反映させることができる。

40

50

【 0 1 3 3 】

< 第 2 実施形態 >

図 9 は、第 2 の実施形態を示し、前記第 1 実施形態に示した物理 CPU 1 0 4 の VT - x 機能を、前記従来例（非特許文献 2）に示した AMD - V（Virtualization）に置き換えたものである。

【 0 1 3 4 】

図 9 において、AMD - V を実装した物理 CPU 1 0 4 a は、前記第 1 実施形態に示したホスト VMM 1 0 のシャドウ VMCB のアドレスを指し示すポインタ 1 1 5 に加えて、ホスト状態エリア 1 3 2 のアドレスを指し示すポインタ 1 1 8 を備える。

【 0 1 3 5 】

AMD - V を備える物理 CPU 1 0 4 a は、VM - e n t r y 命令に代わって、VMRUN 命令を実行するとホスト VMM 1 0 を実行するホストモードからゲスト VMM 2 0 またはユーザプログラムを実行するゲストモードへの切り換えを行う。そして、ホスト VMM 1 0 は、AMD - V の仕様に依りて、次の点が前記第 1 実施形態の VT - x 機能と相違する。

【 0 1 3 6 】

ホスト VMM 保持データ 1 1 のうち、前記第 1 実施形態の VMXON フラグ 1 4 2 を削除し、仮想サーバ毎の仮想化機能有効フラグ 1 4 1 と、仮想 CPU 1 0 8 a ~ n 毎の動作モードフラグ 1 4 3 を備える。そして、ホスト VMM 保持データ 1 1 の物理 CPU 制御データは、図 1 0 で示すように、仮想 CPU 毎のホスト状態エリア 1 3 2 と、仮想 CPU 毎のシャドウ VMCB_0 ~ _n - 1 に、ゲスト状態エリア 1 3 1 と制御エリア 1 3 7 を備える。なお、シャドウ VMCB_0 ~ _n - 1 の制御エリア 1 3 7 は、前記第 1 実施形態の VM - e x i t 情報領域 1 3 6 に相当する情報を格納する。

【 0 1 3 7 】

すなわち、制御エリア 1 3 7 には、図 1 1 で示すように、予め設定された _VMEXIT（前記第 1 実施形態の VM - e x i t に相当）の発生要因のコードを示す EXITCODE 1 3 7 1 と、コードの名称 1 3 7 2 と、このコードが VMRUN 命令に起因することを示す判定結果 1 3 7 3 と、_VMEXIT をゲスト VMM 2 0 へ通知することを示す通知条件 1 3 7 4 が含まれている。

【 0 1 3 8 】

CPU 制御部 1 2 のエミュレータ 1 2 2 は、VM - e n t r y 命令に代わって、VMRUN 命令実行モジュールを備え、VMCLEAR 命令実行モジュールを削除した。また、CPU 制御部 1 2 は、前記第 1 実施形態の VM - e n t r y 命令発行部に代わって、物理 CPU 1 0 4 に VMRUN 命令を発行する VMRUN 命令発行部 1 2 5 A を備える。

【 0 1 3 9 】

次に、ホスト VMM 1 0 が仮想サーバ 1 0 2 a ~ 1 0 2 n へ提供する仮想 CPU 1 0 8 a ~ 1 0 8 n も、物理 CPU 1 0 4 a の構成の変更に応じて、ゲスト VMM 2 0 内のゲスト状態エリア 2 2 1 のアドレスを指し示すポインタ 1 1 6 に加えて、ホスト状態エリア 2 2 2 のアドレスを指し示すポインタ 1 1 7 を備える。なお、ゲスト VMM 2 0 は、前記第 1 実施形態の仮想サーバ 1 0 2 a と同じく、仮想化機能を統合した次世代 OS をゲスト OS 1 1 1 a に採用した場合に仮想サーバ 1 0 2 a で稼働するものである。

【 0 1 4 0 】

ゲスト VMM 2 0 の仮想 CPU 制御データ 2 1 は、ホスト VMM 1 0 のホスト VMM 保持データ 1 1 と同様に、AMD - V の仕様に依りてゲスト VMM 2 0 のステータスを保持するホスト状態エリア 2 2 2 と、ユーザプログラムのステータスを格納するゲスト状態エリア 2 2 1 及び制御エリア 2 2 6 を備えたゲスト VMCB 2 2 から構成される。

【 0 1 4 1 】

以上のような構成により、ホスト VMM 1 0 は、物理 CPU 1 0 4 a から _VMEXIT を受信すると、シャドウ VMCB_0 ~ _n - 1 のゲスト状態エリア 1 3 1 に、ゲスト VMM 2 0 のゲスト状態エリア 2 2 1 またはホスト状態エリア 2 2 2 の内容を書き込むこと

10

20

30

40

50

で、仮想サーバ102aで稼働する次世代OSに統合されたゲストVMM20と、このゲストVMM20上で稼働するユーザプログラムを切り替えることができる。

【0142】

次に、ホストVMM10で行われる処理について、図12を参照しながら以下に説明する。図12は、仮想サーバ102aの稼働中に物理CPU104aからの_VMEXITを受信したときにホストVMM10のCPU制御部12が行う処理の一例を示すフローチャートである。なお、この例では、_VMEXITの要因がVMRUN命令である場合を示す。

【0143】

まず、S61では、ホストVMM10のCPU制御部12は、_VMEXITの要因となった操作対象のゲストOS111a（以下、ゲストOSとする）に対応するホストVMM保持データ11の仮想化機能有効フラグ141を参照し、_VMEXITを発生したゲストOSがAMD-V機能を利用可能であるか否かを判定する。仮想化機能有効フラグ141が「1」であれば操作対象のゲストOSがAMD-V機能を利用可能であると判定してS62へ進む。一方、仮想化機能有効フラグ141が「0」であれば操作対象のゲストOSがAMD-V機能を利用しないと判定してS73へ進む。S73では前記第1実施形態の図5に示したS15と同様に、ホストVMM10が前記従来例に示した仮想マシン処理を行う。

10

【0144】

S62では、ホストVMM10のCPU制御部12が、操作対象の仮想CPUの動作モードフラグ143を参照し、仮想CPUの動作モードがホストモード（前記第1実施形態のVMXルートモード）とゲストモード（前記第1実施形態のVMXノンルートモード）の何れであるかを判定する。動作モードフラグ143が「0」の場合には仮想CPUがホストモードであると判定してS63へ進む。一方、動作モードフラグ143が「1」の場合には、仮想CPUがゲストモードであると判定してS69へ進む。

20

【0145】

S63では、CPU制御部12が物理CPU104aから受信した_VMEXITの発生要因を特定する。この例では、ホストモードであるので、CPU制御部12は、ゲストVMM20の状態（ステータス）を格納したゲスト状態エリア131を参照し、ゲストVMM20が発行したVMRUN命令を_VMEXITの要因として特定する。

30

【0146】

次に、S64ではCPU制御部12のエミュレータ122からVMRUN命令実行モジュールを実行し、操作対象の仮想CPUをゲストモードへ切り替えるのに必要な所定の処理（状態エリア選択部121の起動など）をゲストVMM20に代わってエミュレートする。

【0147】

次に、S65ではホストVMM10のCPU制御部12は、仮想CPUの動作モードをホストモードからゲストモードへ移行させるので、操作対象の仮想CPUの動作モードフラグ143を「1」に更新する。

【0148】

そして、S66では、CPU制御部12が操作対象のゲストVMM20のゲストVMCB22からゲスト状態エリア221に格納されているユーザプログラム（ゲストOSまたはアプリケーション）のステータスを読み込む。

40

【0149】

S67では、CPU制御部12が、上記S66で読み込んだゲスト状態エリア221の情報で操作対象のシャドウVMCB_0のゲスト状態エリア131を更新する。そして、S68でCPU制御部12は、物理CPU104a（操作対象の仮想CPU108a）に対してVMRUN命令を発行する。

【0150】

VMRUN命令を受け付けた物理CPU104aは、ポインタ115で指定されたシャ

50

ドウVMCBのゲスト状態エリア131の内容に基づいて、操作対象の仮想サーバのユーザプログラム(ゲストOSまたはアプリケーション)を実行する。なお、物理CPU104aは、VMRUN命令を受け付けると、ホストVMM10のステータスをホストVMM保持データ11のホスト状態エリア132へ格納し、次の呼び出しに備える。

【0151】

一方、上記S62で動作モードフラグ143が「1」と判定されたゲストモードの場合は、操作対象の仮想CPUがユーザプログラムを実行中であるので、S69の処理を行う。S69では、CPU制御部12がEXITCODEを取得して、図11に示した_VMEXITの要因リストを参照し、操作対象の仮想サーバでユーザプログラムを実行中であるので、通知条件1374からゲストVMMへの_VMEXIT通知条件を検索する。この例では、_VMEXITの要因がVMRUN命令が、_VMEXIT通知条件に一致する。そして、S70ではCPU制御部12がエミュレータ122から図9の_VMEXIT条件検出モジュールを実行し、操作対象の仮想CPUをホストモードへ切り替えるのに必要な所定の処理(状態エリア選択部121の起動など)をユーザプログラムに代わってエミュレートする。

10

【0152】

次に、S71ではCPU制御部12が、操作対象の仮想CPUの動作モードをゲストモードからホストモードへ移行させるので、操作対象の仮想CPUの動作モードフラグ143を「0」にリセットする。

【0153】

そして、S72では、CPU制御部12が操作対象のゲストVMM20から仮想CPU制御データ21のホスト状態エリア221に格納されているゲストVMM20のステータスを読み込む。

20

【0154】

S72の処理が完了すると、上記S67~S68を実行し、CPU制御部12は、上記S72で読み込んだホスト状態エリア222の情報で操作対象のシャドウVMCBのゲスト状態エリア131を更新し、物理CPU104a(操作対象の仮想CPU108a)に対してVMRUN命令を発行する。この結果、VMRUN命令を受け付けた物理CPU104aは、ポインタ115で指定されたシャドウVMCBのゲスト状態エリア131の内容に基づいて、操作対象の仮想サーバのゲストVMM20を実行する。

30

【0155】

このように、仮想サーバで稼働するゲストOSが仮想化機能を統合した次世代OSである場合、ホストVMM10は、仮想CPUの動作モードと、発生した_VMEXITの要因に応じてシャドウVMCBのゲスト状態エリア131に書き込むステータスを、ゲストVMMまたはユーザプログラムの何れか一方から選択する。そして、ホストVMM10が物理CPU104aに対してVMRUN命令を発行することで、仮想サーバ上で第1の仮想CPU上で稼働するゲストVMMと、このゲストVMMが提供する第2の仮想CPU上で稼働するユーザプログラムを切り替えて実行することが可能となり、ゲストVMMは仮想サーバ上で仮想化環境を提供することができる。

【0156】

なお、上記S61の判定でゲストOSが仮想化機能を利用しない場合は、S15において、前記従来例に示した特許文献1または非特許文献1の記載に準じた仮想マシン処理をホストVMM10で行えばよい。

40

【0157】

S73の仮想マシン処理は、例えば、図1に示した仮想サーバ102nのゲストOS111nは既存のOSであり、このゲストOS111nまたはアプリケーション112n(ユーザプログラム110n)が特権命令などの所定の命令を実行すると、上述したように物理CPU104aは_VMEXITの発生をホストVMM10に通知する。

【0158】

ホストVMM10は、物理CPU104aからの_VMEXITの通知を受けると、ユ

50

ーザプログラム 110n (仮想CPU 108n) のステータスをシャドウVMCB_n - 1 のゲスト状態エリア 131 に格納する。

【0159】

ホストVMM 10 は、特権命令など所定の処理が完了すると、ホストVMM 10 のステータスをホスト状態エリア 132 へ格納し、ポインタ 115 のアドレスにゲスト状態エリア 131 を設定してから、VMRUN 命令を発行して、制御を仮想CPU 108n に移動してユーザプログラム 110n の実行を再開する。

【0160】

このように、第2の実施形態によれば、前記第1実施形態と同様に、仮想化機能を統合した次世代OSをホストVMM 10 上で稼働させることが可能となる。そして、仮想化機能を統合した次世代OSと、既存のOSをひとつの物理サーバ 101 へ統合することが可能となって、物理サーバの数を削減してサーバの運用管理コストを削減することが可能となる。

10

【0161】

さらに、上述のように、ホストVMM 10 は第1の仮想CPU上で稼働するゲストVMMに対して第1の仮想CPUがAMD-V機能を提供するように見せかけることができ、仮想化ソフトウェアを統合したOSを確実に稼働させることが可能となる。また、本発明のホストVMM 10 によれば、前記従来例のシミュレータのように、命令列の変換などのオーバーヘッドがなく、仮想計算機の性能を低下させることなく、仮想化機能を統合したOSを仮想サーバで実行することが可能となるのである。

20

【0162】

<第3実施形態>

図13は、第3の実施形態を示し、前記第1実施形態に示した物理CPU 104の仮想化支援機能を、VT-x機能から前記従来例(非特許文献3)に示したIPF(Intanium Processor Family)のVT-i機能に置き換えたものである。

【0163】

図13において、VT-i機能を実装した物理CPU 104bでは、ゲスト(ゲストVMM 20またはユーザプログラム 110a)を実行中に予め設定した命令やイベントがあると仮想化例外(Virtualization Fault)を発生し、ホストVMM 10に制御を移す。ホストVMM 10では、仮想化例外を検知すると、予め設定された仮想化プロシジャ(PAL_VPS_RESUMEプロシジャなど)を実行(エミュレーション)し、仮想CPU 108aから08nを制御するためのシャドウVPD(Virtual Processor Descriptor)を更新してから物理CPU 104bへ仮想化プロシジャ(PAL_VPS_RESUMEプロシジャ)を発行し、ゲストに制御を移す。なお、PAL_VPS_RESUMEプロシジャは、PAL_VPS_RESUME_NORMALプロシジャ呼び出し、またはPAL_VPS_RESUME_HANDLERプロシジャ呼び出しの総称とする。

30

【0164】

すなわち、物理CPU 104bは、前記第1実施形態のVM-exitに代わって仮想化例外を発生し、VM-entry命令に代わって仮想化プロシジャを実行することで、ホストVMM 10とゲストの制御を切り替えるアーキテクチャである。詳細については前記非特許文献3に詳しいので、本アーキテクチャの詳細な説明は省略する。

40

【0165】

ホストVMM 10のホストVMM保持データ 11には、前記第1実施形態と同様の仮想化機能有効フラグ 141と動作モードフラグ 143に加え、仮想CPU 108a~108nのステータスを格納するアーキテクチャ状態 151を含む仮想CPU制御データ B150を設ける。

【0166】

そして、物理CPU制御データ 13は、VT-i機能の仕様に応じて以下の情報を格納する。物理CPU制御データ 13には、仮想CPU 108a~108n毎のステータスを

50

格納するアーキテクチャ状態 1310 を含むシャドウ VPD_{0 ~ n - 1} と、仮想 CPU 108a ~ 108n 毎に仮想化例外の発生を抑止を設定する VDC (Virtualization Disable Control) 領域_{0 ~ n - 1} (1320) と、仮想 CPU 108a ~ 108n が実行する命令またはイベントのうち、物理 CPU 104b で仮想化例外を発生する命令またはイベントを予め定義した仮想化例外要因 1330 と、を備える。

【0167】

ホスト VMM 保持データ 11 に格納されるアーキテクチャ状態は、図 14 で示すようになる。すなわち、仮想 CPU 制御データ B150 のアーキテクチャ状態 151 には、仮想 CPU 108a ~ 108n が実行するゲスト (ゲスト VMM 20 またはユーザプログラム) のステータスが格納され、シャドウ VPD のアーキテクチャ状態 1310 には、物理 CPU 104b で実行されるゲストまたはホスト VMM 10 のステータスが格納される。

10

【0168】

また、ホスト VMM 保持データ 11 の物理 CPU 制御データ 13 に格納される仮想化例外要因 1330 は、例えば、図 15 の仮想化例外要因リストで示すように設定される。図 15 において、仮想化例外要因リストには、仮想化例外の要因を示すコード 1331 と、コードに対応する仮想化例外の内容 1332 と、制御をゲストへ移すための PAL_VPS_RESUME プロシジャの呼び出し設定 1333 と、ゲスト VMM 20 への仮想化例外の通知設定 1334 から構成される。設定 1333 では、図中「 」が設定された命令のときに、制御をゲストへ移す PAL_VPS_RESUME プロシジャを呼び出すことを意味する。また、通知設定 1334 は、図中「 」が設定された命令のときに、ゲスト VMM 20 へ通知を行うことを意味し、ゲスト VMM 20 の VDC 領域 2202 が仮想化例外の発生を抑止していないときにゲスト VMM 20 への通知を行うことを意味する。

20

【0169】

ゲスト VMM 20 が保持する仮想 CPU 108a の状態は、仮想 CPU 制御データ A210 に格納される。仮想 CPU 制御データ A210 は、仮想 CPU 108a のステータスを格納するアーキテクチャ状態 2210 を含むゲスト VPD 220 と、ゲスト OS またはアプリケーションが実行した命令またはイベントのうち、仮想 CPU 108a で仮想化例外を発生する命令またはイベントを定義した仮想化例外要因 2201 と、ゲスト VMM 20 を実行する仮想 CPU 108a が VT-i 機能を利用するか否かを示す VDC 領域 2202 とを含む。

30

【0170】

CPU 制御部 12 は、物理 CPU 104b からの仮想化例外を検知し、エミュレータ 122 へ通知する仮想化例外ハンドラ 1240 と、物理 CPU 104b が提供する PAL (Processor Abstraction Layer) で仮想化例外となった処理を実行させるエミュレータ 122 と、仮想 CPU 108a の稼働状態に基づいて、ゲスト VPD のアーキテクチャ状態 2210 と、ホスト VMM 保持データ 11 の仮想 CPU 制御データ B150 のアーキテクチャ状態 151 の何れかを選択する状態エリア選択部 121 と、状態エリア選択部 121 で選択された情報でシャドウ VPD のアーキテクチャ状態 1310 を更新するシャドウ VPD 更新部 1250 と、ホスト VMM 10 からゲスト VMM 20 またはゲスト VMM 20 上の第 2 の仮想 CPU 208a 上で稼働するユーザプログラムへ制御を切り替えるため、物理 CPU 104b へ所定のプロシジャの実行を指令する仮想化プロシジャ指令部 1250 と、を備える。

40

【0171】

次に、ホスト VMM 10 の CPU 制御部 12 で行われる処理の一例について、図 16 のフローチャートを参照しながら以下に説明する。なお、以下の説明では、前記第 1 実施形態の図 1 に示した仮想サーバ 102a で次世代 OS のゲスト OS 111a とアプリケーション 112a 及びゲスト VMM 20 が稼働している例について説明する。

【0172】

図 16 は、仮想サーバ 102a の稼働中に物理 CPU 104b からの仮想化例外を受信したときにホスト VMM 10 の CPU 制御部 12 が行う処理の一例を示すフローチャート

50

である。なお、この例では、仮想化例外の要因が `PAL__VPS__RESUME` プロシジャからの `vm sw` 命令である場合を示す。

【0173】

まず、S81では、ホストVMM10のCPU制御部12が、仮想化例外の要因となった操作対象のゲストOS111a（以下、ゲストOSとする）に対応するホストVMM保持データ11の仮想化機能有効フラグ141を参照し、仮想化例外が発生したゲストOSがVT-i機能を利用可能であるか否かを判定する。仮想化機能有効フラグ141が「1」であれば操作対象のゲストOSがVT-i機能を利用可能であると判定してS82へ進む。一方、仮想化機能有効フラグ141が「0」であれば操作対象のゲストOSがVT-i機能を利用しないと判定してS94へ進む。S94では、ホストVMM10が前記従来例の非特許文献3に示したVT-i機能を用いて仮想マシン処理を行う。

10

【0174】

S82では、ホストVMM10のCPU制御部12が、操作対象の仮想CPUの動作モードフラグ143を参照し、仮想CPUの動作モードがホストモード（前記第1実施形態のVMXルートモード）とゲストモード（前記第1実施形態のVMXノンルートモード）の何れであるかを判定する。動作モードフラグ143が「0」の場合には仮想CPUがゲストVMM20を実行するホストモードであると判定してS83へ進む。一方、動作モードフラグ143が「1」の場合には、仮想CPUがユーザプログラムを実行するゲストモードであると判定してS90へ進む。

【0175】

S83では、CPU制御部12が物理CPU104bから受信した仮想化例外の発生要因を特定する。この例では、ホストモードであるので、CPU制御部12は、ゲストVMM20の状態（ステータス）を格納したホストVMM保持データ11内の仮想CPU制御データB150のアーキテクチャ状態151を参照し、ゲストVMM20が`PAL__VPS__RESUME` プロシジャから発行した`vm sw`命令を仮想化例外の要因として特定する。

20

【0176】

次に、S84ではCPU制御部12が、仮想化例外の要因となった`PAL__VPS__RESUME` プロシジャモジュールを実行し、所定の処理をゲストVMM20に代わって実行し、操作対象の仮想CPUをゲストモードへ切り替えるのに必要な所定の処理（状態エリア選択部121の起動など）を実行する。

30

【0177】

次に、S85ではCPU制御部12が、仮想CPUの動作モードをホストモードからゲストモードへ移行させるので、操作対象の仮想CPUの動作モードフラグ143を「1」に更新する。

【0178】

そして、S86では、CPU制御部12が操作対象のゲストVMM20の仮想CPU制御データA210内のゲストVPD220からアーキテクチャ状態2210に格納されているユーザプログラム（ゲストOSまたはアプリケーション）のステータスを読み込む。

【0179】

S87では、CPU制御部12が物理CPU104bに対して`PAL__VPS__SAVE/RESTORE` プロシジャを発行し、仮想化例外の発行要因となったシャドウVPD_0を選択する。

40

【0180】

次に、S88では、CPU制御部12が、上記S86で読み込んだアーキテクチャ状態2210のステータスを、上記S87で選択した操作対象のシャドウVPD_0のアーキテクチャ状態1310に書き込んで更新する。そして、S88においてCPU制御部12は、物理CPU104bに対して`PAL__VPS__RESUME` プロシジャを発行する。

【0181】

`PAL__VPS__RESUME` プロシジャを受け付けた物理CPU104bは、S87

50

で選択されたシャドウV P Dのアーキテクチャ状態1 3 1 0の内容に基づいて、操作対象の仮想サーバのユーザプログラム(ゲストOSまたはアプリケーション)を実行する。

【0 1 8 2】

一方、上記S 8 2で動作モードフラグ1 4 3が「1」と判定されたゲストモードの場合は、操作対象の仮想C P Uがユーザプログラムを実行中であるので、S 9 0の処理を行う。S 9 0では、C P U制御部1 2が仮想化例外の要因を取得して、図1 5に示した仮想化例外の要因リストを参照し、通知条件1 3 3 4からゲストV M Mへの仮想化例外通知条件を検索する。この例では、仮想化例外の要因がv m s w命令であるので、仮想化例外通知条件はゲストV M M 2 0のV D C領域2 2 0 2の設定によるものとなる。ここでは、ゲストV M M 2 0へ仮想化例外を通知するものとして扱う。そして、S 9 1では図1 3の仮想化例外条件検出モジュールを実行し、操作対象の仮想C P Uをホストモードへ切り替えるのに必要な所定の処理(状態エリア選択部1 2 1の起動など)をゲストV M M 2 0に代わってエミュレートする。

10

【0 1 8 3】

次に、S 9 2ではC P U制御部1 2が、操作対象の仮想C P Uの動作モードをゲストモードからホストモードへ移行させるので、操作対象の仮想C P Uの動作モードフラグ1 4 3を「0」にリセットする。

【0 1 8 4】

そして、S 9 3では、C P U制御部1 2がホストV M M保持データ1 1内の操作対象の仮想C P Uに対応する仮想C P U制御データB 1 5 0のアーキテクチャ状態1 5 1に格納されているゲストV M M 2 0のステータスを読み込む。

20

【0 1 8 5】

S 9 3の処理が完了すると、上記S 8 7～S 8 9を実行し、C P U制御部1 2は、上記S 9 3で読み込んだホストV M M保持データ1 1の仮想C P U制御データB 1 5 0のアーキテクチャ状態1 5 1の情報で操作対象のシャドウV P Dのアーキテクチャ状態1 3 1 0を更新し、物理C P U 1 0 4 bに対してP A L _ V P S _ R E S U M Eプロシジャを発行する。この結果、P A L _ V P S _ R E S U M Eプロシジャを受け付けた物理C P U 1 0 4 bは、シャドウV P D_0のアーキテクチャ状態1 3 1 0の内容に基づいて、操作対象の仮想サーバのゲストV M M 2 0を実行する。

【0 1 8 6】

このように、仮想サーバで稼働するゲストOSが仮想化機能を統合した次世代OSである場合、ホストV M M 1 0は、仮想C P Uの動作モードと、発生した仮想化例外の要因に応じて、シャドウV P Dのアーキテクチャ状態1 3 1 0に書き込むステータスを、ゲストV M Mまたはユーザプログラムの何れか一方から選択する。そして、ホストV M M 1 0が物理C P U 1 0 4 bに対してP A L _ V P S _ R E S U M Eプロシジャを発行することで、仮想サーバ上でゲストV M Mと、このゲストV M Mが提供する第2の仮想C P U 2 0 8 a上で稼働するユーザプログラムを切り替えて実行することが可能となり、ゲストV M Mは仮想サーバ上で仮想化環境を提供することができる。

30

【0 1 8 7】

なお、上記S 8 1の判定でゲストOSが仮想化機能を利用しない場合は、S 9 4において、前記従来例に示した特許文献1または非特許文献3の記載に準じた仮想マシン処理をホストV M M 1 0で行えばよい。

40

【0 1 8 8】

S 9 4の仮想マシン処理は、例えば、図1に示した仮想サーバ1 0 2 nのゲストOS 1 1 1 nは既存のOSであり、このゲストOS 1 1 1 nまたはアプリケーション1 1 2 n(ユーザプログラム1 1 0 n)が仮想化例外要因の命令を実行すると、上述したように物理C P U 1 0 4 bは仮想化例外の発生をホストV M M 1 0に通知する。

【0 1 8 9】

ホストV M M 1 0は、物理C P U 1 0 4 bからの仮想化例外の通知を受けると、ユーザプログラム1 1 0 n(仮想C P U 1 0 8 n)のステータスを仮想C P U制御データB 1 5

50

0のアーキテクチャ状態151に格納する。

【0190】

ホストVMM10は、仮想化例外の要因となった命令をエミュレートした後、PAL_VPS_RESUMEプロシジャを発行して、制御を仮想CPU108nに移動してユーザプログラム110nの実行を再開する。

【0191】

このように、第3の実施形態によれば、前記第1実施形態と同様に、仮想化機能を統合した次世代OSと、既存のOSをひとつの物理サーバ101へ統合することが可能となって、物理サーバの数を削減してサーバの運用管理コストを削減することが可能となる。

【0192】

さらに、上述のように、ホストVMM10は次世代OSに対して仮想CPUがVT-i機能を提供するように見せかけることができ、仮想化ソフトウェアを統合したOSを確実に稼働させることが可能となる。また、本発明のホストVMM10によれば、前記従来例のシミュレータのように、命令列の変換などのオーバーヘッドがなく、仮想計算機の性能を低下させることなく、仮想化機能を統合したOSを仮想サーバで実行することが可能となるのである。

【0193】

なお、上記各実施形態において、物理CPU104、104a、104bとして記載したプロセッサは、マルチコアプロセッサの構成でも良く、また、ホモジニアスなプロセッサやヘテロジニアスのプロセッサを採用することができる。すなわち、物理CPU104、104a、104bとして複数の汎用プロセッサコア(CPU)と特定用途プロセッサコアを含むヘテロジニアス・マルチコア・プロセッサを用いる場合は、汎用プロセッサコアが仮想化支援機能を備えていれば本発明を適用することができる。

【産業上の利用可能性】

【0194】

以上のように、本発明は、複数の仮想サーバを提供する仮想計算機システムに適用することができる。また、本発明は、物理計算機上で複数の仮想サーバを提供する仮想化管理(VMM)ソフトウェアに適用することができる。

【図面の簡単な説明】

【0195】

【図1】第1の実施形態を示し、本発明を適用する仮想計算機システムのブロック図である。

【図2】第1の実施形態を示し、ホストVMMとゲストVMMの機能を示すブロック図である。

【図3】第1の実施形態を示し、ホストVMM保持データのシャドウVMCBの制御データの一例を示すブロック図である。

【図4】第1の実施形態を示し、VM-exitの要因リストの一例を示す説明図である。

【図5】第1の実施形態を示し、ホストVMMで行われる処理の一例を示すフローチャートで、VM-entry命令に起因するVM-exitの処理を示す。

【図6】第1の実施形態を示し、ホストVMMで行われる処理の一例を示すフローチャートで、VMCLEAR命令に起因するVM-exitの処理を示す。

【図7】第1の実施形態を示し、ホストVMMで行われる処理の一例を示すフローチャートで、CPUID命令に起因するVM-exitの処理を示す。

【図8】第1の実施形態を示し、ユーザーインターフェースの一例を示す画面イメージである。

【図9】第2の実施形態を示し、物理CPUがAMD-V機能をサポートする場合のホストVMMとゲストVMMの機能を示すブロック図である。

【図10】第2の実施形態を示し、ホストVMM保持データのシャドウVMCBの制御データの一例を示すブロック図である。

10

20

30

40

50

【図 1 1】第 2 の実施形態を示し、_V M E X I T の要因リストの一例を示す説明図である。

【図 1 2】第 2 の実施形態を示し、ホスト V M M で行われる処理の一例を示すフローチャートで、V M R U N 命令に起因する _V M E X I T の処理を示す。

【図 1 3】第 3 の実施形態を示し、物理 C P U が V T - i 機能をサポートする場合のホスト V M M とゲスト V M M の機能を示すブロック図である。

【図 1 4】第 3 の実施形態を示し、ホスト V M M 保持データのシャドウ V P D 及び仮想 C P U 制御データ B の制御データの一例を示すブロック図である。

【図 1 5】第 3 の実施形態を示し、仮想化例外の要因リストの一例を示す説明図である。

【図 1 6】第 3 の実施形態を示し、ホスト V M M で行われる処理の一例を示すフローチャートで、P A L _ V P S _ R E S U M E プロシジャ実行に起因する仮想化例外処理を示す。

10

【符号の説明】

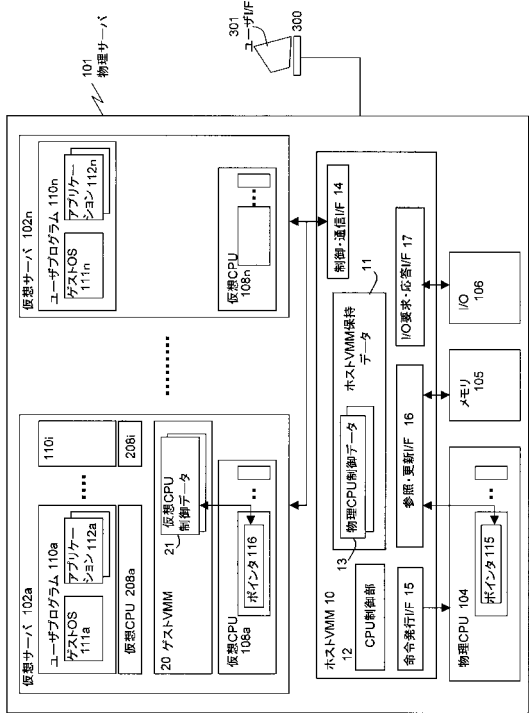
【 0 1 9 6 】

- 1 0 ホスト V M M
- 1 1 ホスト V M M 保持データ
- 1 2 C P U 制御部
- 1 3 物理 C P U 制御データ
- 2 0 ゲスト V M M
- 2 1 仮想 C P U 制御データ
- 1 0 1 物理サーバ (物理計算機)
- 1 0 2 a ~ 1 0 2 n 仮想サーバ
- 1 0 8 a ~ 1 0 8 n 仮想 C P U
- 1 0 4 C P U
- 1 0 5 メモリ
- 1 1 0 a ~ 1 1 0 n ユーザプログラム
- 1 1 1 a ~ 1 1 1 n ゲスト O S
- 1 1 2 a ~ 1 1 2 n アプリケーション
- 2 0 8 仮想 C P U
- 3 0 1 ユーザーインターフェース

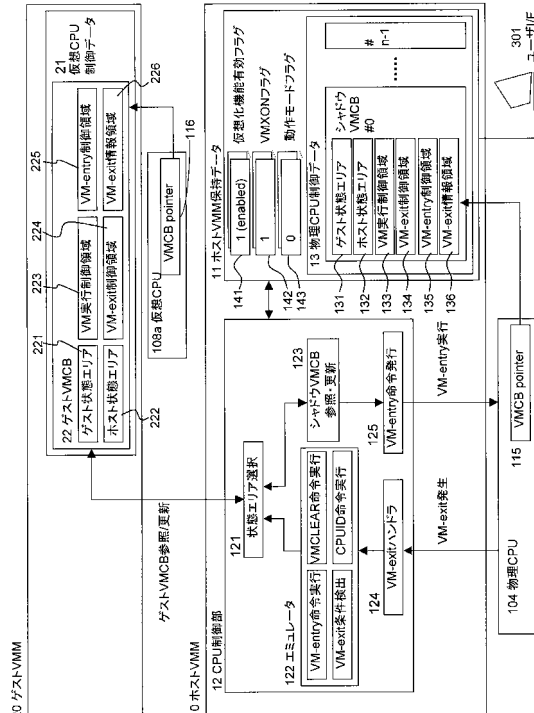
20

30

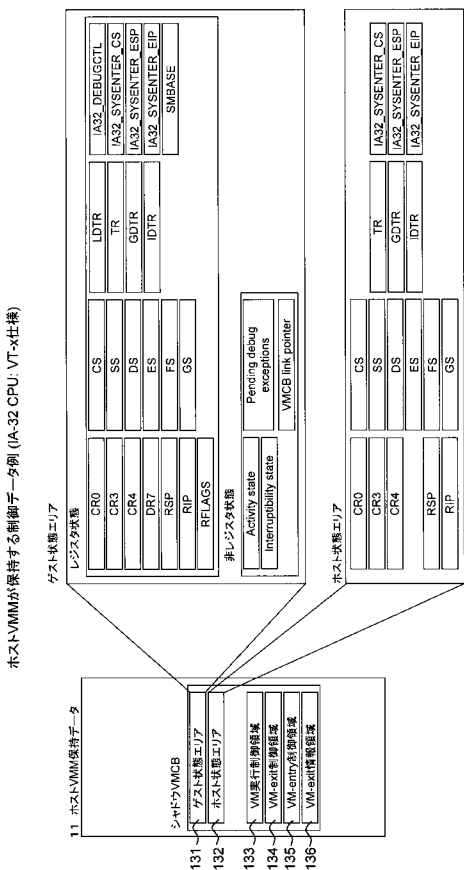
【図 1】



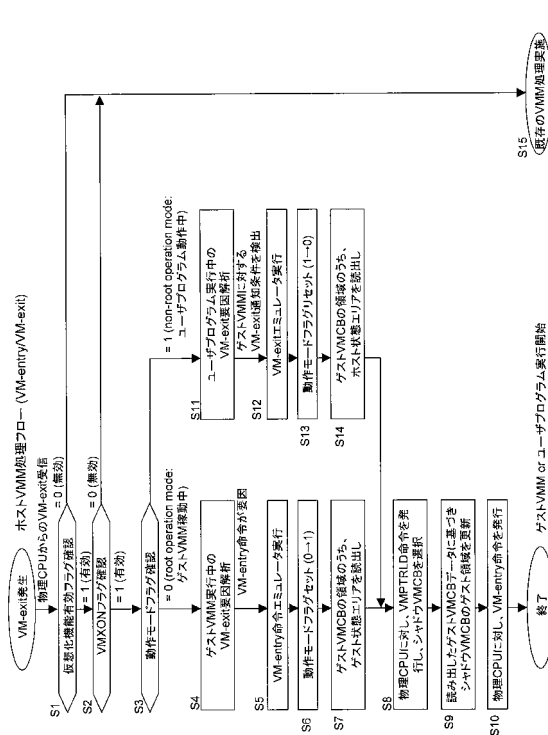
【図 2】



【図 3】

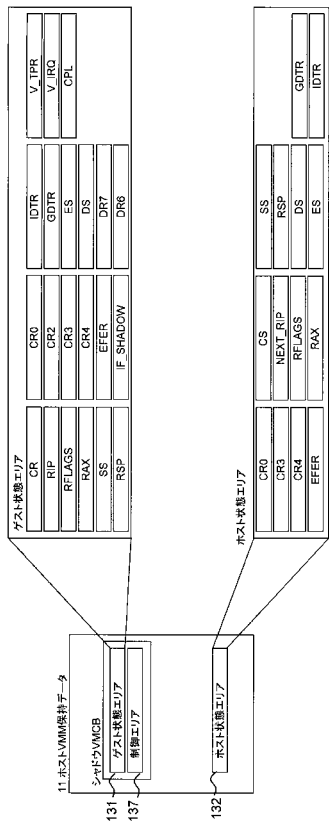


【図 5】

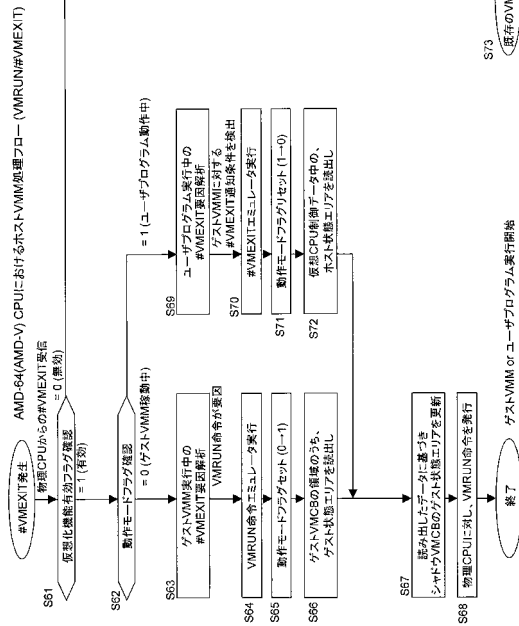


【図10】

ホストVMMが保持する制御データ例 (AMD-64 CPU: AMD-V仕様)

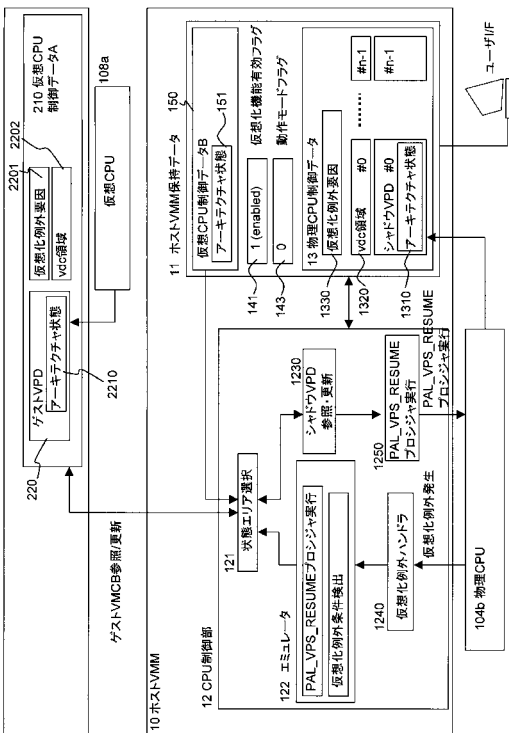


【図12】

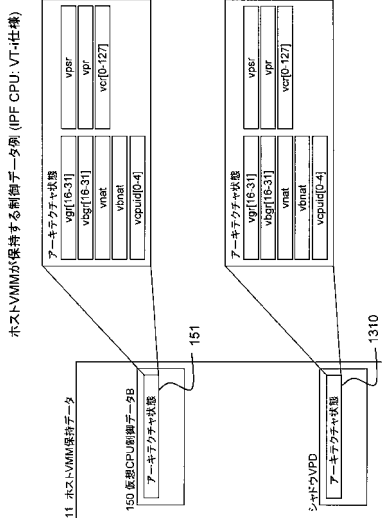


【図13】

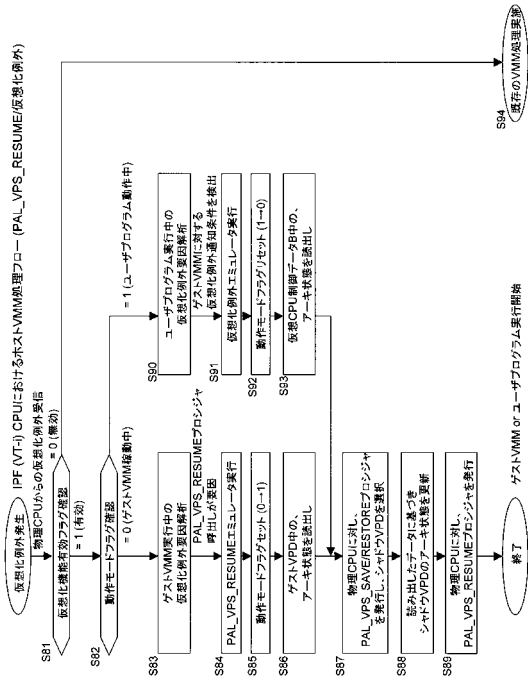
ホストVMMおよびゲストVMMの構成例 (for IPF CPU)



【図14】



【 図 16 】

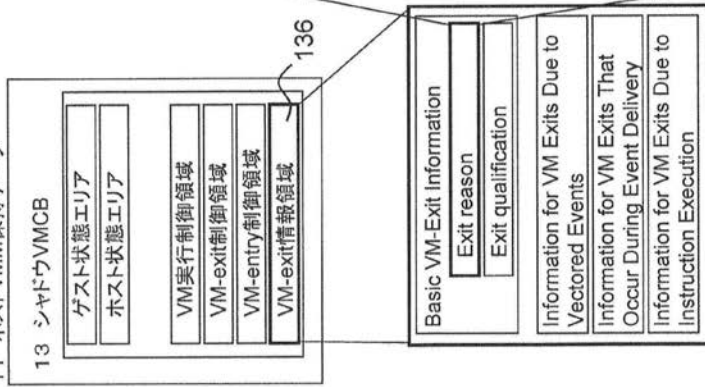


S84 (既存のVMM処理再開)

ゲストVMM or ユーザープログラム実行開始

【 図 4 】

11 ホストVMM保持データ



IA-32 (VT-x) CPUにおけるVM-exit要因一覧

VM-exit reason	description	判定		VM-exit通知条件
		VM-entry命令発行起因のVM-exit	VMCLEAR命令起因のVM-exit	
0	Exception or non-maskable interrupt(NMI)			ゲストVMCB設定による(*)
1	External Interrupt			ゲストVMCB設定による(*)
2	Triple Fault			
3	INIT signaled			
4	Start-up IPI (SIP)			
5	SMI			
6	Other SMI			
7	Interrupt Window			ゲストVMCB設定による(*)
9	Task switch			
10	CPUID			
12	HLT			ゲストVMCB設定による(*)
13	INVD			
14	INVLPG			ゲストVMCB設定による(*)
15	RDPMC			ゲストVMCB設定による(*)
16	RDTS			ゲストVMCB設定による(*)
17	RSM			
18	VMCALL			
19	VMCLEAR		○	
20	VMLAUNCH	○		
21	VMPTRLD			
22	VMPTRST			
23	VMREAD			
24	VMRESUME	○		
25	VMWRITE			
26	VMXOFF			
27	VMXON			
28	CR register access			ゲストVMCB設定による(*)
29	MOV DR			ゲストVMCB設定による(*)
30	I/O instruction			ゲストVMCB設定による(*)
31	RDMSR			ゲストVMCB設定による(*)
32	WRMSR			ゲストVMCB設定による(*)
33	VM-entry failure due to invalid guest state			
34	VM-entry failure due to MSR loading			
36	MWAIT			ゲストVMCB設定による(*)
39	MONITOR			ゲストVMCB設定による(*)
40	PAUSE			ゲストVMCB設定による(*)
41	VM-entry failure due to machine check			
43	TPR below threshold			ゲストVMCB設定による(*)

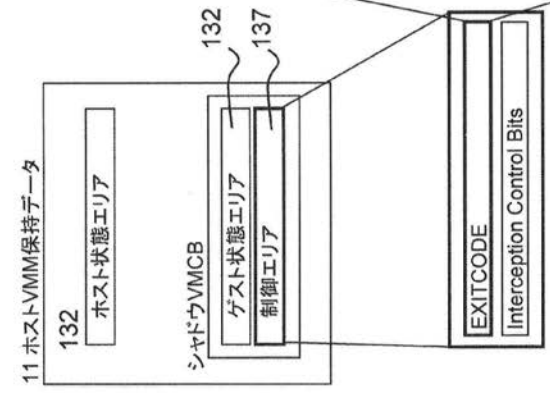
(*) ゲストVMCB中のVM実行制御領域(VM-EXECUTION CONTROL FIELD)等の設定値に依存して、通知条件が可変。

【 図 1 1 】

AMD-64(AMD-V) CPUにおける#VMEXIT要因一覧

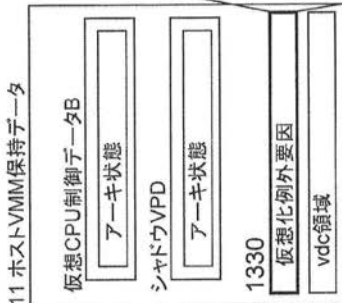
EXITCODE	Name	判定	通知条件	
0-15	VMEXIT_CR[0-15] READ	VMRUN命令発行起因の#VMEXIT	ゲストVMIOB設定による(*2)	
16-31	VMEXIT_CR[0-15] WRITE		ゲストVMIOB設定による(*2)	
32-47	VMEXIT_DR[0-15] READ		ゲストVMIOB設定による(*2)	
48-63	VMEXIT_DR[0-15] WRITE		ゲストVMIOB設定による(*2)	
64-95	VMEXIT_EXCP[0-31]		ゲストVMIOB設定による(*2)	
96	VMEXIT_INTR		ゲストVMIOB設定による(*2)	
97	VMEXIT_NMI		ゲストVMIOB設定による(*2)	
98	VMEXIT_SMI		ゲストVMIOB設定による(*2)	
99	VMEXIT_INIT		ゲストVMIOB設定による(*2)	
100	VMEXIT_VINTR		ゲストVMIOB設定による(*2)	
::	::		::	
128	VMEXIT_VMRUN	VMEXIT発生時の#VMEXIT	ゲストVMIOB設定による(*2)	
129	VMEXIT_VMMCALL		ゲストVMIOB設定による(*2)	
130	VMEXIT_VMLoad		ゲストVMIOB設定による(*2)	
131	VMEXIT_VMSAVE		ゲストVMIOB設定による(*2)	
::	::			::
136	VMEXIT_ICEBP		VMEXIT発生時の#VMEXIT	ゲストVMIOB設定による(*2)
1024	VMEXIT_NPF			○
-1	VMEXIT_INVALID			○

(*2) ゲストVMIOB中のInterception Control Bits等の設定値に依存して、通知条件が可変。



【 図 1 5 】

1331	Handoff Cause	Description	判定 PAL_VPS_RESUMEプログラージャ呼び出し契機 の仮想化例外	1334 ゲストVMMIに対する仮想化例外通知 条件
1		Due to MOV-to-AR instruction.		○
2		Due to MOV-to-AR-imm instruction.		○
3		Due to MOV-from-AR instruction.		○
4		Due to MOV-to-CR instruction.		ゲストvdc領域設定による(*4) ゲストvdc領域設定による(*4)
5		Due to MOV-from-CR instruction.		○
6		Due to MOV-to-PSR instruction.		○
7		Due to MOV-from-PSR instruction.		○
8		Due to itc.d instruction.		○
9		Due to itc.i instruction.		○
10		Due to MOV-to-RR instruction.		○
11		Due to MOV-to-DBR instruction.		ゲストvdc領域設定による(*4)
12		Due to MOV-to-IBR instruction.		ゲストvdc領域設定による(*4)
13		Due to MOV-to-PKR instruction.		○
14		Due to MOV-to-PMC instruction.		ゲストvdc領域設定による(*4)
15		Due to MOV-to-PMD instruction.		ゲストvdc領域設定による(*4)
16		Due to itr.d instruction.		○
17		Due to itr.i instruction.		○
18		Due to MOV-from-RR instruction.		○
19		Due to MOV-from-DBR instruction.		ゲストvdc領域設定による(*4)
20		Due to MOV-from-IBR instruction.		ゲストvdc領域設定による(*4)
21		Due to MOV-from-PKR instruction.		○
22		Due to MOV-from-PMC instruction.		○
23		Due to MOV-from-CPUID instruction.		○
24		Due to ssm instruction.		○
25		Due to rsm instruction.		○
26		Due to ptc.l instruction.		○
27		Due to ptc.g instruction.		○
28		Due to ptc.ga instruction.		○
29		Due to ptr.d instruction.		○
30		Due to ptr.i instruction.		○
31		Due to thash instruction.		○
32		Due to ttag instruction.		○
33		Due to tpa instruction.		○
34		Due to tak instruction.		○
35		Due to ptc.e instruction.		○
36		Due to cover instruction.		○
37		Due to rfi instruction.		○
38		Due to bsw.0 instruction.		○
39		Due to bsw.1 instruction.		○ (*3)
40		Due to vmsw instruction.		○ ゲストvdc領域設定による(*4)



(*3) 仮想CPU上で稼動するPAL_VPS_RESUMEプログラージャコード内でvmsw命令を発行してホストVMMIへ通知する場合に有効。
 (*4) ゲストvdc領域で指定される仮想化例外発生抑制条件によって、仮想化例外の通知条件が可変。
 ゲストvdc領域は、VT-i spec定義の Virtualization Acceleration Control (vac) Fieldの情報を保持するデータ構造。詳細は VT-i spec参照。

IPF(VT-i) CPUにおける仮想化例外要因一覧

フロントページの続き

(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 吉田 美彦

(56)参考文献 特開2007-035045(JP,A)

特表2007-513405(JP,A)

特開2007-183951(JP,A)

特開2006-099332(JP,A)

米国特許出願公開第2007/0006228(US,A1)

伊藤 宏通, CPUの仮想化技術をXenで活用する VMXドメイン徹底研究, Linux
WORLD 第5巻 第4号, 日本, (株)IDGジャパン, 2006年 4月 1日, 第5巻
, P.136~141

(58)調査した分野(Int.Cl., DB名)

G06F 9/46