

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5079246号
(P5079246)

(45) 発行日 平成24年11月21日(2012.11.21)

(24) 登録日 平成24年9月7日(2012.9.7)

(51) Int. Cl. F I
G06F 9/46 (2006.01) G O 6 F 9/46 3 5 0
G06F 13/10 (2006.01) G O 6 F 13/10 3 3 0 C

請求項の数 12 (全 20 頁)

(21) 出願番号	特願2006-67701 (P2006-67701)	(73) 特許権者	500046438 マイクロソフト コーポレーション アメリカ合衆国 ワシントン州 9805 2-6399 レッドモンド ワン マイ クロソフト ウェイ
(22) 出願日	平成18年3月13日(2006.3.13)	(74) 代理人	100077481 弁理士 谷 義一
(65) 公開番号	特開2006-252565 (P2006-252565A)	(74) 代理人	100088915 弁理士 阿部 和夫
(43) 公開日	平成18年9月21日(2006.9.21)	(72) 発明者	エリック ピー. トラウト アメリカ合衆国 98052 ワシントン 州 レッドモンド ワン マイクロソフト ウェイ マイクロソフト コーポレーシ ョン内
審査請求日	平成21年2月16日(2009.2.16)		
(31) 優先権主張番号	11/078, 141		
(32) 優先日	平成17年3月11日(2005.3.11)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 仮想マシン環境におけるマルチレベルインターセプト処理のためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

仮想マシン環境内の区分に対するインターセプトを処理するためのシステムであって、前記仮想マシン環境はバーチャライザおよび少なくとも1つの区分を含み、前記システムは、

区分において発生した特定のイベントをインターセプトして、該インターセプトされたイベントを外部モニタに転送するバーチャライザと、

前記バーチャライザから前記インターセプトされたイベントを受け取ると、前記区分のオペレーティングシステムによる仮想プロセッサの制御を一時停止させて、前記インターセプトされたイベントを実行し、該実行中のインターセプトされたイベントに対応するよ
うに前記区分の仮想プロセッサの状態を修正し、前記インターセプトされたイベントの実行が終了すると前記区分の前記オペレーティングシステムに制御を返す、外部モニタと

を備えたことを特徴とするシステム。

【請求項 2】

前記バーチャライザは、前記インターセプトされたイベントを発生させた前記区分の前記オペレーティングシステムの要求に応じて、前記インターセプトされたイベントを前記外部モニタに転送せずに前記バーチャライザ自身で処理することを特徴とする請求項 1 に記載のシステム。

【請求項 3】

前記インターセプトされたイベントは、読取りと、書込みと、ロック付き読取り/修正

／書込み論理演算AND、ORまたはXORと、ロック付き読取り／修正／書込み算術演算ADDまたはSUBと、ロック付き比較演算XCHGおよびCMPXCHGと、HLT、PAUSEまたはCPUID命令と、制御レジスタ(CR)、デバッグレジスタ(DR)およびモデル固有レジスタ(MSR)へのアクセスと、特定の例外と、外部割込みと、I/Oポートへのアクセスと、致命的例外と、のうちのいずれかのイベントであることを特徴とする請求項1に記載のシステム。

【請求項4】

前記インターセプトされたイベントは、プロセッサインターセプトタイプとメモリアンターセプトタイプのいずれかであり、前記プロセッサインターセプトタイプは、特定のプロセッサ命令、レジスタ、例外、または割込みに関連付けられるイベントであり、前記メモリアンターセプトタイプは、特定のメモリ位置にアクセスするイベントであり、前記インターセプトされたイベントのタイプは、前記パーチャライザから独立に前記区分によって指定されることを特徴とする請求項1に記載のシステム。

10

【請求項5】

仮想マシン環境内の区分に対するインターセプトを処理する方法であって、前記仮想マシン環境は、パーチャライザおよび少なくとも1つの区分を備え、前記方法は、

前記パーチャライザが、区分において発生した特定のイベントをインターセプトすることと、

前記パーチャライザが、前記インターセプトされたイベントを特定の区分内で稼働している外部モニタに転送することと、

20

前記外部モニタが、前記区分のオペレーティングシステムによる仮想プロセッサの制御を一時停止させることと、

前記外部モニタが、前記インターセプトされたイベントを実行することと、

前記外部モニタが、実行中の前記インターセプトされたイベントに対応するように前記区分の前記仮想プロセッサの状態を修正することと、

前記外部モニタが、前記インターセプトされたイベントの実行が終了すると前記区分の前記オペレーティングシステムに制御を返すことと

を備えることを特徴とする方法。

【請求項6】

前記パーチャライザは、前記インターセプトされたイベントを発生させた前記区分の前記オペレーティングシステムの要求に応じて、前記インターセプトされたイベントを前記外部モニタに転送せずに前記パーチャライザ自身で処理することを特徴とする請求項5に記載の方法。

30

【請求項7】

前記インターセプトされたイベントは、読取りと、書込みと、ロック付き読取り／修正／書込み論理演算AND、ORまたはXORと、ロック付き読取り／修正／書込み算術演算ADDまたはSUBと、ロック付き比較演算XCHGおよびCMPXCHGと、HLT、PAUSEまたはCPUID命令と、制御レジスタ(CR)、デバッグレジスタ(DR)およびモデル固有レジスタ(MSR)へのアクセスと、特定の例外と、外部割込みと、I/Oポートへのアクセスと、致命的例外と、のうちのいずれかのイベントであることを特徴とする請求項5に記載の方法。

40

【請求項8】

前記インターセプトされたイベントは、プロセッサインターセプトタイプとメモリアンターセプトタイプのいずれかであり、前記プロセッサインターセプトタイプは、特定のプロセッサ命令、レジスタ、例外、または割込みに関連付けられるイベントであり、前記メモリアンターセプトタイプは、特定のメモリ位置にアクセスするイベントであり、前記インターセプトされたイベントのタイプは、前記パーチャライザから独立に前記区分によって指定されることを特徴とする請求項5に記載の方法。

【請求項9】

仮想マシン環境内の区分に対するインターセプトを処理するためのコンピュータ実行可

50

能な命令を記録したコンピュータ読み取り可能な記録媒体であって、前記仮想マシン環境はバーチャライザおよび少なくとも1つの区分を備え、前記コンピュータ実行可能な命令は、

前記バーチャライザが、区分において発生した特定のイベントをインターセプトするための命令と、

前記バーチャライザが、前記インターセプトされたイベントを特定の区分内で稼働している外部モニタに転送するための命令と、

前記外部モニタが、前記区分のオペレーティングシステムによる仮想プロセッサの制御を一時停止するための命令と、

前記外部モニタが、前記インターセプトされたイベントを実行するための命令と、

前記外部モニタが、実行中の前記インターセプトされたイベントに対応するように前記区分の前記仮想プロセッサの状態を修正するための命令と、

前記外部モニタが、前記インターセプトされたイベントの実行が終了すると前記区分の前記オペレーティングシステムに制御を返すための命令と

を備えることを特徴とするコンピュータ読み取り可能な記録媒体。

【請求項10】

前記コンピュータ実行可能な命令は、前記バーチャライザが、前記インターセプトされたイベントを発生させた前記区分の前記オペレーティングシステムの要求に応じて、前記インターセプトされたイベントを前記外部モニタに転送せずに、前記インターセプトされたイベントを前記バーチャライザ自身で処理するための命令をさらに備えることを特徴とする請求項9に記載のコンピュータ読み取り可能な記録媒体。

【請求項11】

前記インターセプトされたイベントは、読取りと、書込みと、ロック付き読取り/修正/書込み論理演算AND、ORまたはXORと、ロック付き読取り/修正/書込み算術演算ADDまたはSUBと、ロック付き比較演算XCHGおよびCMPXCHGと、HLT、PAUSEまたはCPUID命令と、制御レジスタ(CR)、デバッグレジスタ(DR)およびモデル固有レジスタ(MSR)へのアクセスと、特定の例外と、外部割込みと、I/Oポートへのアクセスと、致命的例外と、のうちのいずれかのイベントであることを特徴とする請求項9に記載のコンピュータ読み取り可能な記録媒体。

【請求項12】

前記インターセプトされたイベントは、プロセッサインターセプトタイプとメモリアンターセプトタイプのいずれかであり、前記プロセッサインターセプトタイプは、特定のプロセッサ命令、レジスタ、例外、または割込みに関連付けられるイベントであり、前記メモリアンターセプトタイプは、特定のメモリ位置にアクセスするイベントであり、前記インターセプトされたイベントのタイプは、前記バーチャライザから独立に前記区分によって指定されることを特徴とする請求項9に記載のコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は一般に、仮想マシン(「プロセッサ仮想化」とも呼ばれる)の分野に関し、より詳細には、本発明は、マルチレベル外部仮想マシンモニタ(VMM: virtual machine monitor)を対象とし、いくつかのインターセプト(intercept)処理機能は、特定の区分(それぞれは仮想マシンすなわち「VM」のインスタンス)で稼働する外部モニタによって、これらの区分化された外部モニタを管理するベースレベルVMMと共に実行される。

【背景技術】

【0002】

コンピュータは、特定のシステム命令セットを実行するように設計された汎用中央処理装置(CPU)を備える。類似するアーキテクチャまたは設計仕様を有する一群のプロセッサは、同じプロセッサファミリのメンバと考えることができる。現在のプロセッサファ

10

20

30

40

50

ミリの例には、アリゾナ州 Phoenix にある Motorola, Inc. 製の Motorola 680X0 プロセッサファミリ、カリフォルニア州 Sunnyvale にある Intel Corporation 製の Intel 80X86 プロセッサファミリ、および、Motorola 製であり、カリフォルニア州 Cupertino にある Apple Computer, Inc. 製コンピュータで使用される PowerPC プロセッサファミリが含まれる。一群のプロセッサのアーキテクチャおよび設計に関する考慮事項が類似するためにこれらのプロセッサが同じファミリ中にある場合があっても、プロセッサは、それらのクロック速度およびその他の性能パラメータにより、ファミリ内で大きく異なることがある。

【0003】

マイクロプロセッサの各ファミリは、そのプロセッサファミリに固有の命令を実行する。プロセッサまたはプロセッサファミリが実行することの可能な命令の総体的なセットは、そのプロセッサの命令セットと呼ばれる。例として、Intel 80X86 プロセッサファミリによって使用される命令セットは、PowerPC プロセッサファミリによって使用される命令セットと互換性がない。Intel 80X86 命令セットは、CISC (Complex Instruction Set Computer) フォーマットに基づく。Motorola PowerPC 命令セットは、RISC (Reduced Instruction Set Computer) フォーマットに基づく。CISC プロセッサは多数の命令を使用し、これらの命令のいくつかは、いくぶん複雑な機能を実行することが可能であるが、実行には、概して多くのクロックサイクルを必要とする。RISC プロセッサは、より少数の利用可能な命令を使用してより単純な機能セットを実行し、これらの機能は、ずっと高いレートで実行される。

【0004】

コンピュータシステム間におけるプロセッサファミリの独自性は通常、コンピュータシステムのハードウェアアーキテクチャの、他の要素間でも非互換性をもたらす。Intel 80X86 プロセッサファミリからのプロセッサを備えて製造されたコンピュータシステムは、PowerPC プロセッサファミリからのプロセッサを備えて製造されたコンピュータシステムのハードウェアアーキテクチャとは異なるハードウェアアーキテクチャを有することになる。プロセッサ命令セットと、コンピュータシステムのハードウェアアーキテクチャとの独自性のため、アプリケーションソフトウェアプログラムは通常、特定のオペレーティングシステムを稼働させる特定のコンピュータシステム上で実行されるように書かれる。

【0005】

(プロセッサ仮想化)

コンピュータ製造業者は、自社の製品ラインに関連するマイクロプロセッサファミリ上で、少なめよりはむしろ多めのアプリケーションが稼働するようにすることによって、自社の市場シェアを最大限にしたいと望む。コンピュータシステム上で稼働可能なオペレーティングシステムおよびアプリケーションプログラムの数を拡大するために、ある技術分野が発展してきたが、この技術分野では、ある種類のCPU (ホストと呼ばれる) を有する所与のコンピュータは、関係ない種類のCPU (ゲストと呼ばれる) の命令をこのホストコンピュータがエミュレートできるようにするエミュレータプログラムを備えることになる。そのため、ホストコンピュータは、所与のゲスト命令に回答して1つまたは複数のホスト命令が呼び出されるようにするアプリケーションを実行することになる。したがって、ホストコンピュータは、それ自体のハードウェアアーキテクチャ用に設計されたソフトウェアと、関係ないハードウェアアーキテクチャを有するコンピュータ用に書かれたソフトウェアとの両方を稼働させることが可能である。より具体的な例として、例えば Apple Computer 製のコンピュータシステムが、PC ベースのコンピュータシステム用に書かれたオペレーティングシステムおよびプログラムを稼働させることができる。また、エミュレータプログラムを使用して、単一のCPU上で複数の非互換オペレーティングシステムを同時に動作させることが可能な場合もある。この構成では、各オペレー

10

20

30

40

50

ティングシステムが相互に互換しないにもかかわらず、エミュレータプログラムは、2つのオペレーティングシステム的一方をホストすることが可能であり、普通なら互換性がないはずのこれらのオペレーティングシステムが同じコンピュータシステム上で同時に稼動できるようにする。

【0006】

ホストコンピュータシステム上でゲストコンピュータシステムがエミュレートされる時、ゲストコンピュータシステムは「仮想マシン」と言われる。というのは、ゲストコンピュータシステムは、ホストコンピュータシステム中で、ある特定ハードウェアアーキテクチャの動作の純粋なソフトウェア表現として存在するだけだからである。エミュレータ、仮想マシン、およびプロセッサエミュレーションという用語は、コンピュータシステム全体のハードウェアアーキテクチャを模倣またはエミュレートする能力を指すために交換可能に使用されることがある。例として、カリフォルニア州 San Mateo にある Connectix Corporation によって生み出された Virtual PC ソフトウェアは、Intel 80X86 Pentium (登録商標) プロセッサと様々なマザーボードコンポーネントおよびカードとを備えたコンピュータ全体をエミュレートする。これらのコンポーネントの動作は、ホストマシン上で稼動している仮想マシン中でエミュレートされる。Power PC プロセッサを有するコンピュータシステムなどのホストコンピュータの、オペレーティングシステムソフトウェアおよびハードウェアアーキテクチャ上で実行されているエミュレータプログラムが、ゲストコンピュータシステム全体の動作を模倣する (mimic)。

【0007】

エミュレータプログラムは、ホストマシンのハードウェアアーキテクチャと、エミュレートされた環境内で稼動するソフトウェアから送られる命令との間のインターチェンジとして働く。このエミュレートされた環境は、仮想マシンモニタ (VMM) によって生み出すことができ、VMM は、ハードウェアのすぐ上で稼動するソフトウェア層であって、VMM が仮想化しているハードウェアと同じインタフェースを公開することによってマシンのリソースすべてを仮想化するソフトウェア層である (これにより VMM は、VMM よりも上で稼動するオペレーティングシステム層に認識されずに (unnotice) いることができる)。この構成では、ホストオペレーティングシステム (HOS) と VMM とが、同じ物理ハードウェア上で並行して稼動することができる。あるいは、エミュレータプログラムは、物理コンピュータハードウェアのすぐ上で稼動して別のハードウェア構成をエミュレートする HOS 自体であってもよい。この実施形態の特定の実装形態では、HOS ソフトウェアは特に、「ハイパーバイザ (hypervisor)」の一実施形態を備えることができる。

【0008】

ハイパーバイザは、HOS のカーネルレベルの近くに存在する制御プログラムであり、コンピュータシステムの1つまたは複数の物理プロセッサを含めたコンピュータシステムハードウェアを、HOS 以外の1つまたは複数の2次オペレーティングシステムが使用できるようにするべく動作する。ハイパーバイザは、2次オペレーティングシステムのための動作環境をエミュレートし、それにより2次オペレーティングシステムは、実際には別のハードウェアおよび/またはオペレーティングシステム環境で動作しており HOS がコンピュータシステムを論理制御しているであろうにもかかわらず、その通例のハードウェアおよび/またはオペレーティングシステム環境で動作しておりそれ自体がコンピュータシステムを論理制御していると認識する (believe)。これは重要なことである。というのは、多くのオペレーティングシステムは、コンピュータシステムのハードウェアを排他的に論理制御しているかのように動作しなければならないものとして機能するからである。したがって、複数のオペレーティングシステムが単一のコンピュータシステム上で同時に機能するためには、各オペレーティングシステムのハイパーバイザは、各オペレーティングシステムがコンピュータシステム全体の排他的制御を有するかのごとく機能するように、他のオペレーティングシステムの存在を覆い隠すよう機能しなければならない

10

20

30

40

50

【0009】

簡単にするために、プロセッサ仮想化プログラム（限定されないがVMMおよびハイパーバイザを含む）を、本明細書では「バーチャライザ（virtualizer）」と総称する。さらに、ハイパーバイザのコンテキスト（状況、態様）で本明細書に開示される本発明の態様はどれも、VMMおよびその他のバーチャライザに対しても等しく有効であり開示されるものとし、またその逆でもある。

【0010】

（インターセプトおよびモニタ）

当業者には知られており理解されるように、バーチャライザ（例えばハイパーバイザやVMM）の主要な機能の1つは、ソフトウェアが区分（VMの個々のインスタンス）中のゲストオペレーティングシステム上で実行されている間に発生したイベントをインターセプトすることである。このコンテキストで、イベントとは、区分の何らかのコンポーネントと、実際には区分の一部ではない物理的または仮想的な何らかのリソースとの間で発生する対話である。例えば、第1の区分で実行されているプログラムが、ある周辺デバイスにデータを送りたい場合があり、この区分のオペレーティングシステムは、このデバイスに対して排他的制御を有すると認識している。しかしこの場合、このプログラムが、その区分のゲストオペレーティングシステムを介してそのデータを送り、ゲストオペレーティングシステムが周辺デバイスと通信しようとする、と、バーチャライザは、同じことをしようとする他の区分による他の試みと共に（すなわちデバイスは実際にはどの特定区分にも専用になっていない）この区分によるこのデバイスへのアクセスを処理するために、この通信をインターセプトする。これらの種類のイベントをインターセプトすることにより、バーチャライザは本質的に、現実には基礎的な物理リソースは実際にいくつかの仮想マシン区分間で共有または区分化されているときに、ゲストOSを騙して、（物理ハードウェア上で実行されているOSが通常そうするように）マシンのリソースすべてを所有していると思わせる。これに関して、バーチャライザは、このようなイベントをインターセプトすることと、これらが発生したときにインターセプトにตอบสนองすることとの両方を担う。残念ながら、当業者にとっては容易に明らかな理由でより単純なバーチャライザが望まれるときに、各区分の外で動作する単一のバーチャライザにこの種の機能を集中化させる場合、バーチャライザを非常に複雑にする必要がある。

【発明の開示】

【発明が解決しようとする課題】

【0011】

したがって、より単純でありながらインターセプト機能をなお提供するバーチャライザ-インターセプトモデルが、当技術分野で必要とされている。本発明は、かかる解決策のひとつを提供する。

【課題を解決するための手段】

【0012】

本発明の様々な実施形態は、インターセプト関係の機能の大部分を（各区分の外に存在する）ベースレベルのバーチャライザから除去して、その代わりにこの機能の大部分を各区分に直接に組み込むように設計された、マルチレベルバーチャライザを対象とする。いくつかの実施形態では、区分内で稼動し特定のインターセプトイベントにตอบสนองする「外部モニタ」によって、いくつかのインターセプト処理機能が実行され、ベースレベルバーチャライザは、これらの外部モニタを各区分内にインストールし、その後、単一区分と区分間との両方のインターセプトイベントについて、外部モニタを管理する。この分散型のインターセプト処理手法により、複雑さのずっと低いバーチャライザが可能になり、インターセプト機能が各区分中にまで引き上げられ、各区分中では、各外部モニタが、その区分中の対応するゲストオペレーティングシステムのリソースを使用してインターセプトイベントを解決する。

【0013】

前述の概要、ならびに好ましい実施形態に関する後続の詳細な記述は、添付の図面と共に読めばよりよく理解される。本発明を例示するために、図面には本発明の例示的な構造が示してある。しかし、本発明は、開示される特定の方法及び手段に限定されない。

【発明を実施するための最良の形態】

【0014】

本発明の主題について、法定の要件を満たすために具体的に述べる。しかし、この記述自体は、本特許の範囲を限定するものではない。そうではなく、発明者は、特許請求される本主題が他の方法でも実施されて、他の現在または将来の技術に関連する、本明細書に述べるステップに類似する種々のステップまたはステップの組合せを含めることができることを企図している。さらに、用語「ステップ」は、本明細書では、利用される方法の種々の要素を意味するのに使用される場合があるが、この用語は、個々のステップの順序が明示的に述べられていない限り、本明細書に開示される様々な複数のステップ内またはステップ間の任意の特定の順序を暗示するものと解釈されるべきではない。

【0015】

(コンピュータ環境)

本発明の多くの実施形態は、コンピュータ上で実行することができる。図1および後続の考察は、本発明を実施することのできる適切なコンピューティング環境に関する簡単で一般的な記述を提供するものである。必須ではないが、本発明を、クライアントワークステーションやサーバなどのコンピュータによって実行される、プログラムモジュールなどのコンピュータ実行可能命令の一般的なコンテキストで述べる。一般に、プログラムモジュールは、特定のタスクを実行するか特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、コンポーネント、データ構造などを含む。さらに、当業者なら理解するであろうが、本発明は、ハンドヘルドデバイス、マルチプロセッサシステム、マイクロプロセッサベースのまたはプログラム可能な民生用電子機器、ネットワークPC、ミニコンピュータ、メインフレームコンピュータなどを含めた他の種類のコンピュータシステム構成でも実施することができる。本発明は分散コンピューティング環境で実施することもでき、その場合、通信ネットワークを介してリンクされたりリモート処理デバイスによってタスクが実行される。分散コンピューティング環境では、プログラムモジュールは、ローカルとリモートとの両方のメモリ記憶デバイスに位置することができる。

【0016】

図1に示すように、例示的な汎用コンピューティングシステムは、従来型のパーソナルコンピュータ20などを備え、コンピュータ20は、処理ユニット21と、システムメモリ22と、システムメモリを含む様々なシステムコンポーネントを処理ユニット21に結合するシステムバス23とを備える。システムバス23は、様々なバスアーキテクチャのいずれかを使用した、メモリバスまたはメモリコントローラ、周辺バスおよびローカルバスを含む、いくつかの種類のバス構造のうちのいずれかとすることができる。システムメモリは、ROM(Read Only Memory)24およびRAM(Random Access Memory)25を含む。ROM24には、起動中などにパーソナルコンピュータ20内の要素間で情報を転送するのを支援する基本ルーチンを含むBIOS(basic input/output system)26を記憶することができる。パーソナルコンピュータ20はまた、ハードディスク(図示せず)に対して読み書きするためのハードディスクドライブ27と、取外し可能な磁気ディスク29に対して読み書きするための磁気ディスクドライブ28と、CD-ROMやその他の光媒体など取外し可能な光ディスク31に対して読み書きするための光ディスクドライブ30を備えることもできる。ハードディスクドライブ27、磁気ディスクドライブ28、および光ディスクドライブ30は、それぞれハードディスクドライブインタフェース32、磁気ディスクドライブインタフェース33、および光ドライブインタフェース34によってシステムバス23に接続される。ドライブおよびそれらに関連するコンピュータ読み取り可能な媒体は、コンピュータ読み取り可能な命令、データ構造、プログラムモジュール、およびその他のデータの揮発性記憶域をパーソナルコンピュータ20に提供する。本明細書に述べる例

10

20

30

40

50

示的な環境は、ハードディスク、取外し可能な磁気ディスク29、および取外し可能な光ディスク31を利用しているが、磁気カセット、フラッシュメモ리카ード、デジタルビデオディスク、ベルヌーイカートリッジ、RAM、ROMなど、コンピュータによってアクセス可能なデータを記憶可能な他の種類のコンピュータ読み取り可能な媒体をこの例示的な動作環境で使用してもよいことを、当業者は理解されたい。

【0017】

ハードディスク、磁気ディスク29、光ディスク31、ROM24、またはRAM25には、いくつかのプログラムモジュールを記憶することができ、これらのプログラムモジュールにはオペレーティングシステム35、1つまたは複数のアプリケーションプログラム36、その他のプログラムモジュール37、およびプログラムデータ38が含まれる。ユーザは、キーボード40やポインティングデバイス42などの入力デバイスを介して、パーソナルコンピュータ20にコマンドおよび情報を入力することができる。その他の入力デバイス(図示せず)には、マイクロホン、ジョイスティック、ゲームパッド、衛星受信アンテナ、スキャナなどを含めることができる。これらおよびその他の入力デバイスは、システムバスに結合されたシリアルポートインタフェース46を介して処理ユニット21に接続されることが多いが、パラレルポート、ゲームポート、またはUSB(Universal Serial Bus)など、他のインタフェースによって接続されてもよい。モニタ47または他の種類の表示デバイスも、ビデオアダプタ48などのインタフェースを介してシステムバス23に接続される。モニタ47に加えて、パーソナルコンピュータは通常、スピーカやプリンタなど、その他の周辺出力デバイス(図示せず)も備える。図1の例示的なシステムはまた、ホストアダプタ55、SCSI(Small Computer System Interface)バス56、および、SCSIバス56に接続された外部記憶デバイス62も備える。

【0018】

パーソナルコンピュータ20は、リモートコンピュータ49など、1つまたは複数のリモートコンピュータへの論理接続を使用して、ネットワーク化された環境で動作することができる。リモートコンピュータ49は、別のパーソナルコンピュータ、サーバ、ルータ、ネットワークPC、ピアデバイス、またはその他の一般的なネットワークノードであってよく、通常はパーソナルコンピュータ20に関して上述した要素の多くまたはすべてを備えるが、図1にはメモリ記憶デバイス50だけが示してある。図1に示す論理接続は、LAN(Local Area Network)51およびWAN(Wide Area Network)52を含む。このようなネットワーキング環境は、オフィス、企業規模のコンピュータネットワーク、イントラネット、およびインターネットでよく見られる。

【0019】

LANネットワーキング環境で使用されるときは、パーソナルコンピュータ20は、ネットワークインタフェースまたはアダプタ53を介してLAN51に接続される。WANネットワーキング環境で使用されるときは、パーソナルコンピュータ20は通常、インターネットなどのワイドエリアネットワーク52を介した通信を確立するためのモデム54または他の手段を備える。モデム54は内蔵でも外付けでもよく、シリアルポートインタフェース46を介してシステムバス23に接続される。ネットワーク化された環境では、パーソナルコンピュータ20に関して示したプログラムモジュールまたはその一部を、リモートのメモリ記憶デバイスに記憶することができる。図示のネットワーク接続は例であり、コンピュータ間で通信リンクを確立するための他の手段を使用してもよいことは理解されるであろう。さらに、本発明の多くの実施形態は、コンピュータ化されたシステムに特によく適合すると考えられるが、本明細書中のいずれのものも、本発明をかかると限定することを意図していない。

【0020】

(仮想マシン)

概念的に見ると、コンピュータシステムは一般に、基礎的なハードウェア層上で稼動す

10

20

30

40

50

る1つまたは複数のソフトウェア層を備えている。この階層化は、抽象化の理由で行われる。所与のソフトウェア層のためのインタフェースを定義することにより、この層を、それよりも上にある他の層によって様々に実装することが可能である。適切に設計されたコンピュータシステムでは、各層は、その直下の層だけについて識別している（また、直下の層だけに依拠する）。これにより、ある層または「スタック」（隣接する複数の層）を、この層またはスタックよりも上の層に悪影響を与えずに置き換えることができる。例えば、ソフトウェアアプリケーション（上方の層）は通常、より低いレベルのオペレーティングシステム（下方の層）に依拠して、ファイルを何らかの形の永久記憶装置に書き込むが、これらのアプリケーションは、データをフロッピー（登録商標）ディスクに書き込むか、ハードドライブに書き込むか、ネットワークフォルダに書き込むかの違いを理解する必要はない。この下方の層がファイル書込み用の新しいオペレーティングシステムコンポーネントで置き換えられても、上方の層のソフトウェアアプリケーションの動作は影響を受けない。

【0021】

階層化されたソフトウェアのフレキシビリティは、仮想マシン（VM）が、実際には別のソフトウェア層である仮想ハードウェア層を提示することを可能にする。このようにして、VMは、VMよりも上のソフトウェア層に対して、これらのソフトウェア層がそれら自体の個別コンピュータシステム上で稼動しているという錯覚を生み出すことが可能であり、したがってVMは、複数の「ゲストシステム」が単一の「ホストシステム」上で同時に稼動できるようにすることが可能である。

【0022】

図2は、コンピュータシステム中のエミュレートされた動作環境での、ハードウェアおよびソフトウェアアーキテクチャの論理階層化を表す図である。図では、エミュレーションプログラム94が、物理ハードウェアアーキテクチャ92の上で直接的または間接的に稼動する。エミュレーションプログラム94は、(a)ホストオペレーティングシステムと並行して稼動する仮想マシンモニタ、(b)ネイティブエミュレーション能力を有する特殊化されたホストオペレーティングシステム、または(c)ハイパーバイザコンポーネントを備えるホストオペレーティングシステムとすることができ、(c)の場合はハイパーバイザコンポーネントがエミュレーションを実行する。エミュレーションプログラム94は、ゲストハードウェアアーキテクチャ96をエミュレートする（ゲストハードウェアアーキテクチャ96は破線で示してあり、それにより、このコンポーネントが「仮想マシン」であること、すなわち実際に存在するのではなくエミュレーションプログラム94によってエミュレートされたハードウェアであることを示す）。ゲストハードウェアアーキテクチャ96上では、ゲストオペレーティングシステム98が実行され、ゲストオペレーティングシステム98上では、ソフトウェアアプリケーション100が稼動する。図2のエミュレートされた動作環境では、またエミュレーションプログラム94の作用により、ソフトウェアアプリケーション100は、ホストオペレーティングシステムおよびハードウェアアーキテクチャ92との互換性を概して持たないオペレーティングシステム上で稼動するように設計されているにもかかわらず、コンピュータシステム90中で稼動することが可能である。

【0023】

図3Aに、物理コンピュータハードウェア102のすぐ上で稼動するホストオペレーティングシステムソフトウェア層104を含む、仮想化されたコンピューティングシステムを示す。ホストオペレーティングシステム（ホストOS）104は、ホストOSがエミュレートしている（すなわち「仮想化」している）ハードウェアと同じインタフェースを公開することによって、物理コンピュータハードウェア102のリソースへのアクセスを提供し、これによりホストOSは、ホストOSよりも上で稼動するオペレーティングシステム層に認識されずにいることができる。エミュレーションを実行するためには、ホストオペレーティングシステム104は、ネイティブエミュレーション能力を備えた特別設計のオペレーティングシステムとすることができ、あるいは、エミュレーションを実行するた

10

20

30

40

50

めの組込みハイパーバイザコンポーネントを備えた標準的なオペレーティングシステムとすることもできる。

【0024】

再び図3Aを参照すると、ホストOS104の上には、2つの仮想マシン（VM）実装形態VMA108およびVMB110があり、VMA108は例えば、仮想化されたIntel386プロセッサとすることができ、VMB110は例えば、Motorola680X0ファミリのプロセッサのうちの1つの仮想化バージョンとすることができる。各VM108および110の上には、それぞれゲストオペレーティングシステム（ゲストOS）A112およびB114がある。ゲストOSA112の上では、2つのアプリケーションすなわちアプリケーションA1116およびアプリケーションA2118が稼動しており、ゲストOSB114の上では、アプリケーションB1120が稼動している。

10

【0025】

図3Aに関して、以下のことに留意するのが重要である。すなわち、（破線で示す）VMA108およびVMB110は、仮想化されたコンピュータハードウェアの表現であり、これらはソフトウェア構造としてのみ存在し、これらは専用ソフトウェアコードが存在することによって実行可能になり（possible）、この専用ソフトウェアコードは、VMA108およびVMB110をそれぞれゲストOSA112およびゲストOSB114に提示するだけでなく、ゲストOSA112およびゲストOSB114が実際の物理コンピュータハードウェア102と間接的に対話するのに必要なすべてのソフトウェアの

20

【0026】

図3Bに、代替の仮想化されたコンピューティングシステムを示す。このコンピューティングシステムでは、ホストオペレーティングシステム104と並行して稼動する仮想マシンモニタ（VMM）104'によってエミュレーションが実行される。いくつかの実施形態では、VMMは、ホストオペレーティングシステム104の上で稼動し、およびホストオペレーティングシステム104を介してのみコンピュータハードウェアと対話する、アプリケーションとすることができる。他の実施形態では、図3Bに示すように、VMMは代わりに、部分的に独立したソフトウェアシステムを構成することができ、そのソフトウェアシステムはあるレベルではホストオペレーティングシステム104を介して間接的にコンピュータハードウェア102と対話するが、他のレベルでは（ホストオペレーティングシステムがコンピュータハードウェアと直接に対話するのと同様の方式で）コンピュータハードウェア102と直接に対話する。他の実施形態では、VMMは、完全に独立したソフトウェアシステムを構成することができ、このソフトウェアシステムは、（コンピュータハードウェア102の使用を調整して競合を回避することなどに関する限りにおいては、やはりホストオペレーティングシステム104と対話するが）ホストオペレーティングシステム104を利用せずに、（ホストオペレーティングシステムがコンピュータハードウェアと直接に対話するのと同様の方式で）すべてのレベルでコンピュータハードウェア102と直接に対話する。

30

【0027】

VMMを実装するためのこれらの変形はすべて、本明細書に述べる本発明の代替実施形態を形成することが予期され、本明細書中のどんなものも、本発明をいずれか特定のVMM構成に限定するものと解釈されるべきではない。加えて、（おそらくハードウェアエミュレーションのシナリオで、）VMA108および/またはVMB110をそれぞれ介したアプリケーション116、118、120の間の対話に関するどんな言及も、実際にはアプリケーション116、118、120とVMMとの間の対話であると解釈されるべきである。同様に、（おそらくコンピュータ命令をコンピュータハードウェア102に対して直接的または間接的に実行するための）アプリケーションVMA108および/またはVMB110と、ホストオペレーティングシステム104および/またはコンピュータハードウェア102との間の対話に関するどんな言及も、実際には、VMMと、適宜ホスト

40

50

オペレーティングシステム 104 またはコンピュータハードウェア 102 との間の対話であると解釈されるべきである。

【0028】

(イベントインターセプトの概観)

本明細書で先に論じたように、バーチャライザは、コンピュータシステム上で稼動する薄いソフトウェア層であり、1つまたは複数の抽象仮想マシン (VM) インスタンス (それぞれは区分とも呼ばれる) を生み出すことを可能にする。これらの各インスタンスは通常、それ自体のメモリと、1つまたは複数のプロセッサと、I/O デバイスと (これらは、実際の物理リソースに対応するか、または実際にはどんな物理的意味においても存在しない仮想化されたリソースに対応することができる) を備えた、実際のコンピュータのよう

10

に振舞う。仮想化を実施するために、ほとんどのバーチャライザは、ゲストソフトウェアが区分中で実行される間に発生する特定のイベントをインターセプトする必要がある。 (ゲストソフトウェアは、特定の仮想マシン内で稼動するソフトウェアであり、通常はオペレーティングシステムおよび1つまたは複数のアプリケーションからなる。) 特定のイベントをインターセプトすることは、仮想化において重要である。というのは、そうすることによってバーチャライザは、実際には物理マシン上のすべてのリソースは共有または区分化されている (すなわち細分されて特定の仮想マシンに割り当てられている) ときに、

20

ゲストを騙して、これらのリソースすべてを所有していると思わせることができるからである。例えば、4つのプロセッサと、プロセッサの総数を報告するレジスタとを有し、さらに、システム上で稼動する4つの仮想マシン (区分) を有するコンピュータシステム

30

であって、各仮想マシンが単一のプロセッサに割り当てられているコンピュータシステムを考へてみる。この例で、これらの各仮想マシンのゲストソフトウェアが「プロセッサ数」レジスタを読み取るとき、通常は値「4」を読み取ることになる (これは実際には物理プロセッサの総数であり、したがって物理レジスタ中の値である) 。しかしこの場合、バーチャライザは、各ゲストOSを騙して、1つのプロセッサシステム上で稼動していると認識させる必要がある。これを達成するために、バーチャライザは、ゲストが「プロセッサ数」レジスタから読み取るときに「インターセプト」イベントを生成するように、各仮想マシン (区分) の仮想プロセッサを構成する。このインターセプトは、VMがレジスタを直接に読み取るのを防止し、制御をゲストからバーチャライザに移行させる。次いでバーチャライザは、レジスタ中に実際に存在する正常値「4」を、バーチャライザが報告

40

したい値、この場合なら「1」で、オーバーライド (上書き : `override`) することができる。

【0029】

当業者には周知であり容易に理解されるように、一般にインターセプトは、プロセッサインターセプトとメモリインターセプトという2つの範疇に分けることが可能である。プロセッサインターセプトは、特定のプロセッサ命令、レジスタ、例外、または割込みに結び付けられるイベントである。メモリインターセプトは、指定のメモリ位置に対する読み取りまたは書き込みをプロセッサに行わせるアクションに特有であり、明示的なアクセス (例えば、メモリから読み取る「MOVE」命令) 、または暗黙的なアクセス (例えば、プロセッサがTLB (`translation look-aside buffer`) ミス

40

に回答してページテーブルを検索 (`walk`) するときにプロセッサによって生成される参照) が含まれる。

【0030】

従来、(各区分の外に存在する) バーチャライザは、インターセプトを「インストールする」ことと、インターセプトが発生したときにそれに「応答する」こととの両方を担う。インターセプトをインストールする技法は、当該のプロセッサアーキテクチャに大きく依存する。ほとんどのプロセッサは、いくつかのイベント時に「トラップ」する方法を提供する。いくつかの場合では、インターセプトに対する制御は、細粒または粗粒とすることが可能である。後者の場合、バーチャライザは、単一の粗粒インターセプトをインストールし、次いで、インターセプトしたい特定のイベントを「フィルタにかけて選

50

び出し」

、他のインターセプトイベントは無視する必要がある場合がある。プロセッサが特定のインターセプトタイプに対する細粒制御を提供する場合は、バーチャライザは、欲しいインターセプトだけを要求することが可能である。インターセプト報告は一般にかなりのオーバーヘッド処理コストを追加し、これはシステム性能を低下させるので、理想的には、インターセプトは、バーチャライザが標準的なプロセッサ挙動をオーバーライドする必要があるときにだけ行われるべきである。先の4プロセッサ4VMコンピュータシステムの例に再び言及するが、ユーザが4プロセッサの物理システム上で単一の4プロセッサ仮想マシンを構成した場合は、「プロセッサ数」レジスタへのアクセスをインターセプトする必要はなく、したがって、適切に書かれたバーチャライザなら、この場合にインターセプトを回避することになり、それにより、不要なオーバーヘッドコストおよび性能打撃(hit)を回避することになる。

10

【0031】

インターセプトの必要を低減するための別の手法は、「シャドウレジスタ(Shadow Register)」の一般概念である。シャドウレジスタは、仮想マシン内で動作しているときにアクセスされる、実際のレジスタのコピーである。再び我々の先の例に言及するが、システムがシャドウ「プロセッサ数」レジスタを実装した場合、仮想化された環境の外で動作するプロセス(バーチャライザ自体を含む)は、このレジスタにアクセスして実際の値を読み取ることになるが、仮想化された環境の中で稼動するゲストプロセスの場合は、このレジスタにアクセスするとシャドウ値を読み取ることになり、バーチャライザは、シャドウレジスタが適切なプロセッサ数を報告するようにプログラム済みである。この情報をあらかじめプログラムしておくことにより、バーチャライザは、このレジスタアクセスに対するインターセプトを要求する必要はなく、その代わりに、全般的に、各区分に実際の値ではなくシャドウ値を読み取らせる。しかし、この手法は、コストのかかるインターセプトイベントを回避するが、プロセッサがさらに複雑になること(通常は追加のレジスタ)を必要とし、その結果、(任意のアクションを実行可能な汎用インターセプトハンドラと比較して)フレキシビリティがいくらか減少する。

20

【0032】

より具体的には、従来のバーチャライザは、本質的にモノリシックであった。すなわち、インターセプトをインストールすることと、インターセプトされるイベントが発生したときにそれに応答することとの両方を担っていた。したがって、これらのバーチャライザは今まで、必然的に複雑で厄介であった。図4は、プロセッサがサポートする手段が何であれ(例えば例外、デフォルト、トラップなど)それを介してインターセプトをインストールすることと、インターセプトされるイベントすべてに応答することとの両方を行う、従来のバーチャライザのモノリシックな性質を示すブロック図である。図では、物理コンピュータハードウェア402の上で動作するバーチャライザ404が、2つの仮想マシン環境VMA408およびVMB410を仮想化する。ゲストオペレーティングシステム(GOS)A412がVMA408中で実行され、GOSB414がVMB410中で実行される。ソフトウェアアプリケーションA1416およびA2418が、VMA408上のGOSA412中で実行され、ソフトウェアアプリケーションB1420が、VMB410上のGOSB414中で実行される。この実施形態では、バーチャライザはハイパーバイザとして示されており(ただし、その他のバーチャライザも予期され使用できる)、VMA408は「1次区分」である。すなわち、VMAのゲストオペレーティングシステム(GOS)Aはハイパーバイザによって利用され、普通ならホストオペレーティングシステムから提供されるであろう機能を提供する(当業者には理解されるであろうが、ホストオペレーティングシステムは、この特定の仮想化構成では不要であり存在しない)。

30

40

【0033】

インターセプトに関して、バーチャライザ404(やはりこの場合もハイパーバイザ)は、完全なインターセプト機能450を備え、これによりバーチャライザ404の複雑さは非常に高まる。したがって、この構成では、第1の区分(例えばVMA408)中のイ

50

イベントがインターセプトを誘因 (t r i g g e r) し、プロセッサに、第1の区分で稼動するゲストOS (例えばゲストOS A 4 1 2) からパーチャライザ 4 0 4 に制御を移行させる。次いでパーチャライザ 4 0 4 は、インターセプトの原因を決定してそれに応答し、完了時に制御をゲストOS A 4 1 2 に返す。

【 0 0 3 4 】

(マルチレベルインターセプトイベント)

本発明の様々な実施形態は、パーチャライザとインターセプト処理方法とのためのマルチレベル設計を対象とし、このマルチレベル設計では、インターセプトに関係する複雑さの大部分をベースレベルのパーチャライザ (仮想化された環境の外で実行されるパーチャライザコンポーネント) の外に移して、各区分 (仮想マシン) 内にある外部モニタ中に置くことができる。この構成では、(いくつかの実施形態では汎用ハイパーバイザに統合することができる) ベースレベルのパーチャライザは単純なままであり、インターセプト処理のいくらかまたはすべては、各区分内すなわちゲストレベルで稼動する1つまたは複数の外部モニタによって、または、インターセプトイベントを生成した区分と同じ区分または異なる区分で実行される。

10

【 0 0 3 5 】

図5は、本発明のいくつかの実施形態に関するマルチレベルインターセプト手法を示すブロック図である。図では、物理コンピュータハードウェア 5 0 2 の上で動作するパーチャライザ 5 0 4 が、2つの仮想マシン環境 V M A 5 0 8 および V M B 5 1 0 を仮想化する。G O S A 5 1 2 が V M A 5 0 8 中で実行され、G O S B 5 1 4 が V M B 5 1 0 中で実行される。ソフトウェアアプリケーション A 1 5 1 6 および A 2 5 1 8 が、V M A 5 0 8 上の G O S A 5 1 2 中で実行され、ソフトウェアアプリケーション B 1 5 2 0 が、V M B 5 1 0 上の G O S B 5 1 4 中で実行される。この実施形態でも、パーチャライザはやはりハイパーバイザとして示されており (ただし、その他のパーチャライザも予期され使用できる)、V M A 5 0 8 は「1次区分」である。すなわち、V M A のゲストオペレーティングシステム (G O S) A はハイパーバイザによって利用され、普通ならホストオペレーティングシステムから提供されるであろう機能を提供する。

20

【 0 0 3 6 】

しかし、インターセプトに関して、パーチャライザ 5 0 4 (やはりこの場合もハイパーバイザ) は、最小限のインターセプト機能 5 5 0 だけを備え、その代わりに、インターセプト機能の大部分は、各区分 5 0 8 および 5 1 0 中の1つまたは複数の外部モニタ (「E M」) 5 6 2 および 5 6 4 に組み込まれている (ここでは各区分の V M 仮想化に組み込まれているのが示されているが、代替の実施形態では、区分中の V M とゲストOS とアプリケーションとに対する相対的な E M の位置は、異なってもよい)。

30

【 0 0 3 7 】

この構成では、図5と6の両方を参照すると (後者の図は、マルチレベル方法がインターセプトイベントを処理する方法を示すプロセスフロー図である)、ステップ 6 0 2 で、第1の区分 (例えば V M A 5 0 8) 中のイベントがインターセプトを誘因し、ステップ 6 0 4 で、プロセッサに、第1の区分中で稼動するゲストOS (例えばゲストOS A 5 1 2) からパーチャライザ 5 0 4 に制御を移行させる。ステップ 6 0 6 で、次いでパーチャライザ 5 0 4 は、それ自体のデフォルトの (かつ単純な) 処理機構を使用してインターセプトを処理するか、あるいはインターセプトを「外部モニタ」 (例えば E M 5 6 2) に渡すかを決定する。例えばパーチャライザは、インターセプトのソースがデフォルト処理を要求した場合に、(ステップ 6 0 8 でインターセプトをそれ自体で処理する) 前者を選択することができる、あるいは反対に、例えばデフォルト挙動を使用しないようインターセプトのソースが要求した場合に、(ステップ 6 1 0 を介してインターセプトを E M に転送する) 後者を選択することができる。パーチャライザがインターセプトを処理するために外部モニタに転送しようとする場合、パーチャライザはまず、ステップ 6 1 0 で、インターセプトを処理することになる特定の外部モニタに信号を送り、ステップ 6 1 2 で、この E M が2つの方法のうち的一方でインターセプトに応答できるようにする。2つの方法は次の

40

50

とおりである。(a)(EMが特定のインターセプトイベントを処理することができないときなど)EMは、ステップ614で、インターセプトをバーチャライザに返し、インターセプトを通常のデフォルト方式で処理するようバーチャライザに要求する。あるいは(b)EMは、ステップ616でインターセプトを処理する(かつ、いくつかの実施形態では、バーチャライザによるデフォルトアクションが不要であることをバーチャライザに通知する)。インターセプト処理が完了すると、ステップ618で、制御はインターセプトのソース(例えばゲストOSA412)に返される。

【0038】

実施の際、外部モニタを各区分ごと、各仮想プロセッサごと、各インターセプトクラスごとにバーチャライザに登録することが可能であり、あるいは、登録の代わりに、インターセプトを処理する必要があるときにバーチャライザがEMをポーリングすることができる(すなわち、1つの外部モニタがインターセプトを処理すると決定するまで、順番に呼び出すことができる)。本発明のいくつかの実施形態では、外部モニタは区分ごとにのみ登録することができ、それにより、所与の区分に関するすべての要求されたインターセプトは、指定の外部モニタに向けられることになる。しかし、他の実施形態では、複数の区分が、様々な区分で稼動する代替の外部モニタを有することができ、これにより、各外部モニタが各区分に対してわずかに異なる挙動(例えば異なるバージョンの外部モニタ)を有することができる。

【0039】

いくつかの実施形態では、循環依存およびデッドロック状況を防止するために、区分はそれ自体のインターセプトを処理することが許されないものとして行うことができる。したがって、これらの実施形態では、各区分の外部モニタは別の区分内で稼動する必要がある。1次区分を有するハイパーバイザベースのシステムでは、当業者には理解されるように、このことは、少なくとも1つの区分(1次区分)は外部モニタを有することが不可能であることを意味し、したがって、このような実施形態の1次区分は、バーチャライザ(ハイパーバイザ)によって提供されるデフォルトのインターセプト処理だけに依拠することが可能である。他方、(ホストオペレーティングシステムを有する)ホストされるシステムに関するいくつかの代替の実施形態では、循環依存およびデッドロック状況に関する同様の懸念は実際のホストオペレーティングシステムには当てはまらないので、外部モニタはホストOSおよびすべての区分内で稼動することが可能である。

【0040】

(代替の実施形態)

本明細書で先に述べたように、インターセプトは、メモリインターセプトとプロセッサインターセプトという2つのクラスに分けることが可能である。とりわけ、メモリインターセプトは、メモリマッピングレジスタ、ROM、プレーナVRAMなど、通常のRAMのように働かないメモリの領域を仮想化するのに使用され、限定はされないが次のことのために使用することができる。すなわち、読取り(1から16バイト)と、書込み(1から16バイト)と、ロック付き読取り/修正/書込み論理演算AND、OR、XOR(1、2、4、または8バイト)と、ロック付き読取り/修正/書込み算術演算ADD、SUB(1、2、4、または8バイト)と、ロック付き比較演算XCHG、CMPXCHG(1、2、4、または8バイト)である。対照的に、とりわけ、プロセッサインターセプトは、次のことのために使用することができる。すなわち、HLT、PAUSE、およびCPUID命令と、制御レジスタ(CR)、デバッグレジスタ(DR)、モデル固有レジスタ(MSR)へのアクセスと、特定の例外と、外部割込みと、I/Oポートへのアクセス(IN/OUT命令)と、致命的例外(例えばトリプルフォールト)である。

【0041】

加えて、本発明のいくつかの実施形態では、ゲスト命令の途中でインターセプトを外部モニタに報告することができる。例えば、単一の命令が、関連するインターセプトを有する2つのメモリアドレスにアクセスする場合がある。この場合、ゲスト命令を完了することができる前に、両方のインターセプトを処理する必要がある。これらのインターセプト

10

20

30

40

50

は、連続的に検出および処理され、いくつかの実施形態では、バーチャライザは命令完了を実施することを担う。例えば、CPUIDは、レジスタEAX、EDX、EBX、ECX中の、プロセッサに関する情報を返し、バーチャライザがCPUIDインターセプトを外部モニタに報告するために選択できる方法は2つある。すなわち、(a)外部モニタは、CPUID命令を完了し、逆に適切なレジスタに値を書き込み、CPUID命令を通過してEIP(命令ポインタ)をインクリメントするか、あるいは(b)外部モニタは、CPUID命令から提供して欲しい値を返す。前者の場合では、バーチャライザは、命令間のインターセプトを配信している。外部モニタから見ると、仮想化されたプロセッサはまだCPUID命令を実行しておらず、外部モニタがインターセプトの処理を終えたとき、CPUID命令の実行は完了することになる。他方、後者の場合では、バーチャライザは、逆にこれらの値を適切なレジスタに書き込み、CPUID命令の最後を通過してEIPをインクリメントすることを担う。この方法は、バーチャライザが命令中インターセプトを配信しているとき、より単純な外部モニタおよびよりよい性能をもたらし、外部モニタから見ると、仮想化されたプロセッサはCPUID命令を実行し始めている。外部モニタがインターセプトの処理を終えたとき、バーチャライザはCPUID命令を完了することになる。加えて、いくつかの実施形態は、可能な場合は常に命令中インターセプトをサポートする。というのは、より単純な外部モニタインタフェースがもたらされるとともに、例外に対するインターセプトや、命令境界で必然的に発生する割込みなど、命令間インターセプトも可能だからである。

【0042】

さらに、前に論じたように、ベースレベルモニタは、各インターセプトタイプにつき、「デフォルトの」(通常は単純な)処理を定義する。例えば、CPUID命令はデフォルトで、物理プロセッサが通常返すことになると同じCPU情報を返し、致命的例外(例えばx86アーキテクチャに対するトリプルフォールト)はデフォルトで、区分を終了させる。しかし、いくつかの実施形態では、外部モニタは、基礎をなすバーチャライザを呼び出して、どのインターセプトをオーバーライドしたいかを指定することにより、特定のインターセプトクラスのデフォルトの挙動をオーバーライドすることができ、オーバーライドされないインターセプトは、引き続きデフォルト方式で処理される。

【0043】

ほとんどのバーチャライザはまた、明確に定義された機構を介してゲストソフトウェアがバーチャライザと通信できるようにもする。本発明のいくつかの実施形態では、これは、総合MSR(モデル固有レジスタ)を使用することによって行うことができる。外部モニタがこれらのMSRアクセスの挙動をオーバーライドできるようにすることによって、外部モニタは、バーチャライザが存在しないように見せることが可能である。バーチャライザはまた、異なるバージョンのハイパーバイザまたはVMMの存在をシミュレートすることも可能である。これは、再帰的な仮想化を可能にする。すなわち、ハイパーバイザまたはVMMは、別のハイパーバイザまたはVMMの最上部で稼動する区分内で稼動することができる。これは、ハイパーバイザ/VMMの新バージョンのプロトタイプを作成すること、古い方のバージョンとの後方互換性を提供すること、または第三者の実装との互換性を提供することのために有用となることが可能である。

【0044】

ベースレベルのバーチャライザと外部モニタとの両方を含む設計では、インターセプトが発生したときに外部モニタに信号を送り、任意選択でインターセプト関係のパラメータを渡す方法が必要である。外部モニタは、今は他のコードの実行に使用されているかもしれない別の区分で稼動しているので、いくつかの実施形態では、この信号を配信するための非同期機構が必要である。一実装形態では、この信号は、総合割込みとして配信することができる(本明細書で先に挙げた相互参照の特許出願を参照されたい)。従来のプロセッサ割込み機構は、特定の割込みソースに回答して割込みハンドラを実行することを必要とするが、総合割込みコントローラは、この機構を拡張して、割込みハンドラにパラメータを渡せるようにする。インターセプトの場合、これらのパラメータは、保留中のインタ

ーセプトに関する情報を示し、外部モニタは、この情報を使用して効率的にインターセプトを処理することが可能である。

【0045】

(結び)

本明細書に述べた様々なシステム、方法、および技法は、ハードウェアまたはソフトウェア、あるいは適切ならこの両方の組合せで実施することができる。したがって、本発明の方法および装置、またはそのいくつかの態様または部分は、フロッピー（登録商標）ディスク、CD-ROM、ハードドライブ、またはその他の任意のマシン読み取り可能な記憶媒体など、有形媒体に組み入れられたプログラムコード（すなわち命令）の形をとることができ、プログラムコードがコンピュータなどのマシンにロードされマシンによって実行されると、このマシンは本発明を実施するための装置になる。プログラムコードがプログラム可能なコンピュータ上で実行される場合、コンピュータは一般に、プロセッサと、プロセッサによって読み取り可能な記憶媒体（揮発性と不揮発性のメモリおよび/または記憶素子を含む）と、少なくとも1つの入力デバイスと、少なくとも1つの出力デバイスを備えることになる。1つまたは複数のプログラムは、コンピュータシステムと通信するために高レベルのプロシージャ指向またはオブジェクト指向プログラミング言語で実現されることが好ましい。しかし、1つまたは複数のプログラムは、望むならアセンブリ言語またはマシン言語で実現することが可能である。いずれの場合でも、言語はコンパイルまたは解釈される言語とすることができ、ハードウェア実装形態と組み合わせることができる。

10

20

【0046】

本発明の方法および装置はまた、電気ワイヤリングまたはケーブリング、光ファイバ、あるいはその他の任意の形の伝送など、何らかの伝送媒体を介して伝送されるプログラムコードの形で具体化することもでき、プログラムコードが、EPROM、ゲートアレイ、プログラマブルロジックデバイス（PLD）、クライアントコンピュータ、ビデオレコーダなどのマシンによって受け取られてロードおよび実行されると、このマシンは本発明を実施するための装置になる。汎用プロセッサ上で実施されるとき、プログラムコードは、プロセッサと組み合わさって、本発明の索引付け機能を実施するように動作する独特の装置を提供する。

【0047】

様々な図の好ましい実施形態に関して本発明を述べたが、本発明を逸脱することなく、その他の類似の実施形態を使用して、または述べた実施形態に修正および追加を施して、本発明と同じ機能を実施することもできることを理解されたい。例えば、本発明の例示的な実施形態は、パーソナルコンピュータの機能をエミュレートするデジタルデバイスのコンテキストで述べたが、本発明は、本明細書に述べたようなデジタルデバイスに限定されず、ゲーム用コンソール、ハンドヘルドコンピュータ、ポータブルコンピュータなど、有線であれ無線であれ既存のまたは新たに出現する任意の数のコンピューティングデバイスまたは環境にも適用することができ、通信ネットワークを介して接続されネットワークを介して対話する任意の数のこのようなコンピューティングデバイスに適用することができることは、当業者なら理解するであろう。さらに、特に無線ネットワーク化デバイスの数が急増し続けているので、本明細書では、ハンドヘルドデバイスオペレーティングシステムやその他の特定用途向けハードウェア/ソフトウェアインタフェースシステムを含めて、様々なコンピュータプラットフォームが企図されることを強調しておくべきである。したがって、本発明は、いずれか単一の実施形態に限定されるべきではなく、添付の特許請求の範囲による幅および範囲で解釈されるべきである。

30

40

【0048】

最後に、本明細書に記載の開示した実施形態は、他のプロセッサアーキテクチャ、コンピュータベースのシステム、またはシステム仮想化で使用されるように適合させることもでき、そのような実施形態も、本明細書で行った開示によって明白に予期され、したがって本発明は、本明細書に述べた特定の実施形態に限定されるべきではなく、最も広範に解

50

積されるべきである。同様に、プロセッサ仮想化以外の目的で総合命令を使用することも、本明細書で行った開示によって明白に予期され、プロセッサ仮想化以外のコンテキストにおけるこのようなどんな総合命令利用も、本明細書で行った開示の中に最も広範に読み込まれるべきである。

【図面の簡単な説明】

【0049】

【図1】本発明の態様を組み込むことのできるコンピュータシステムを表すブロック図である。

【図2】コンピュータシステム中のエミュレートされた動作環境での、ハードウェアおよびソフトウェアアーキテクチャの論理階層化を示す図である。

【図3A】エミュレーションがホストオペレーティングシステムによって（直接にまたはハイパーバイザを介して）実行される、仮想化されたコンピューティングシステムを示す図である。

【図3B】ホストオペレーティングシステムと並行して稼動する仮想マシンモニタによってエミュレーションが実行される、代替の仮想化されたコンピューティングシステムを示す図である。

【図4】プロセッサがサポートする手段が何であれ（例えば例外、デフォルト、トラップなど）それを介してインターセプトをインストールすることと、インターセプトされたイベントすべてに応答することとの両方を行う、従来のパーチャライザのモノリシックな性質を示すブロック図である。

【図5】本発明のいくつかの実施形態に関するマルチレベルインターセプト手法を示すブロック図である。

【図6】マルチレベル方法がインターセプトイベントを処理する方法を示すプロセスフロー図である。

【符号の説明】

【0050】

92 物理ハードウェアアーキテクチャ

94 エミュレーションプログラム

96 ゲストハードウェアアーキテクチャ（仮想マシン）

98 ゲストオペレーティングシステム

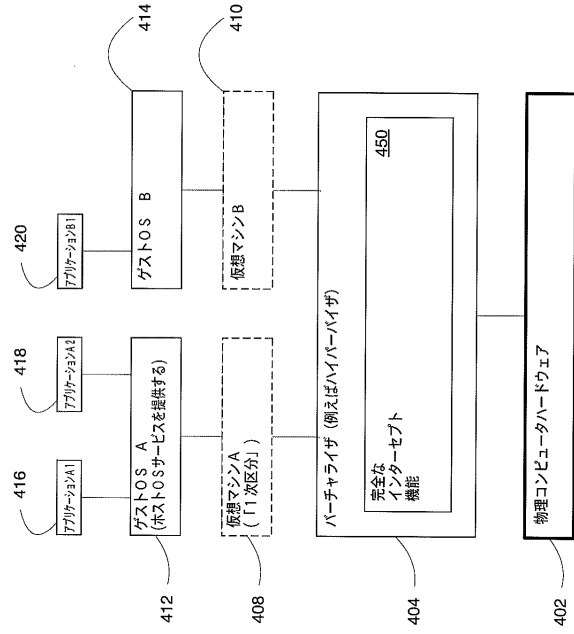
100 ソフトウェアアプリケーション

10

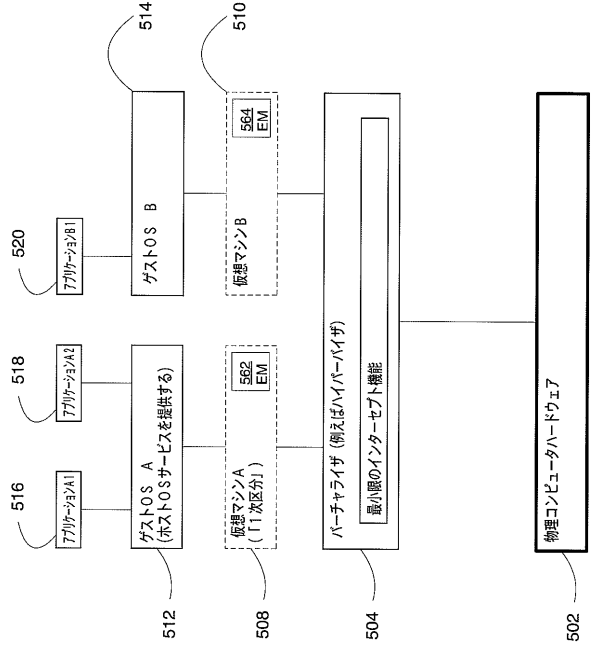
20

30

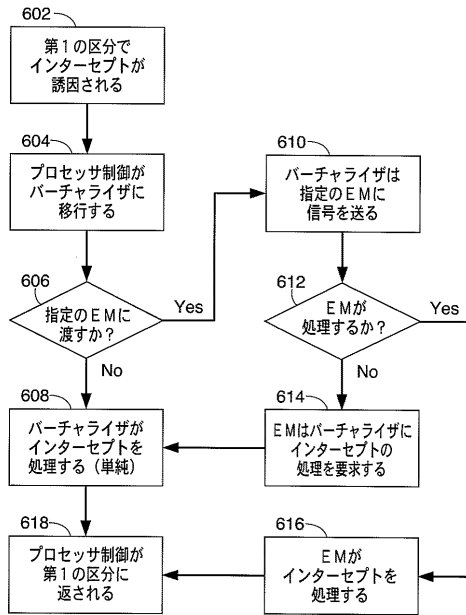
【図4】



【図5】



【図6】



フロントページの続き

- (72)発明者 レネ アントニオ ベガ
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 ジョイ ガングリー
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内

審査官 久保 光宏

- (56)参考文献 米国特許出願公開第2004/0117532 (US, A1)
特開平2 - 19937 (JP, A)
国際公開第2004/095272 (WO, A1)
特表2008 - 508598 (JP, A)
特開平2 - 239334 (JP, A)
特開平3 - 100833 (JP, A)
特開2001 - 318797 (JP, A)
特表2007 - 505402 (JP, A)
特開昭62 - 221041 (JP, A)
特開平9 - 282196 (JP, A)
特開平9 - 81401 (JP, A)
Barham, P., et.al., "Xen and the art of Virtualization", Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP '03), 2003年12月, p.164-177, ISBN :1-58113-757-5
梅野英典 (外4名), 「仮想計算機システムの高性能化方式」, 情報処理, 日本, 社団法人情報処理学会, 1990年12月15日, Vol.31, No.12, 第1665 ~ 1680頁, ISSN:0447-8053
梅野英典 (外2名), 「仮想計算機システムにおける論理プロセッサをスケジュールする新方式の開発と評価」, 情報処理学会論文誌, 日本, 社団法人情報処理学会, 2003年 3月15日, Vol.44, No.3, 第868 ~ 882頁, ISSN:0387-5806
「性能強化が進み本格的実用期に入った仮想計算機システム」, 日経コンピュータ, 日本, 日経マグローヒル社, 1984年 7月23日, 1984年7月23日号 (No.74), 第63 ~ 82頁, ISSN :0285-4619

- (58)調査した分野(Int.Cl., DB名)
G06F9/46 - 9/54,
G06F13/10 - 13/14,
CSDB (日本国特許庁)