

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4882005号
(P4882005)

(45) 発行日 平成24年2月22日(2012.2.22)

(24) 登録日 平成23年12月9日(2011.12.9)

(51) Int.Cl.

F I

G 0 6 F 9/46 (2006.01)

G 0 6 F 9/46 3 5 0

G 0 6 F 9/48 (2006.01)

G 0 6 F 9/46 3 1 1 Z

請求項の数 20 (全 21 頁)

(21) 出願番号 特願2009-540365 (P2009-540365)
 (86) (22) 出願日 平成19年11月13日(2007.11.13)
 (65) 公表番号 特表2010-512577 (P2010-512577A)
 (43) 公表日 平成22年4月22日(2010.4.22)
 (86) 国際出願番号 PCT/US2007/084522
 (87) 国際公開番号 W02008/070410
 (87) 国際公開日 平成20年6月12日(2008.6.12)
 審査請求日 平成22年10月15日(2010.10.15)
 (31) 優先権主張番号 11/635,455
 (32) 優先日 平成18年12月6日(2006.12.6)
 (33) 優先権主張国 米国 (US)

早期審査対象出願

(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
 (74) 復代理人 100115624
 弁理士 濱中 淳宏
 (74) 復代理人 100162950
 弁理士 久下 範子
 (74) 代理人 100077481
 弁理士 谷 義一
 (74) 代理人 100088915
 弁理士 阿部 和夫

最終頁に続く

(54) 【発明の名称】 仮想環境における最適化した割り込み送信

(57) 【特許請求の範囲】

【請求項 1】

仮想コンピューティングシステムに対する割り込みを処理する方法であって、前記仮想コンピューティングシステムは少なくとも1つのプロセッサ、仮想マシンモニター、およびゲストオペレーティングシステムと仮想割り込みコントローラーとを含む少なくとも1つの仮想マシンを備え、前記方法は

仮想割り込みコントローラーが、前記仮想マシンモニターから第1の割り込み要求を受信するステップと、

前記仮想割り込みコントローラーが、前記第1の割り込み要求に対応する割り込みサービスフラグを立て、前記第1の割り込み要求が、ゲストオペレーティングシステムにより自動EOIとしてプログラムされている場合、前記第1の割り込み要求をプロセッサに送信すると、前記割り込みサービスフラグをクリアするステップと

を含むことを特徴とする方法。

【請求項 2】

前記第1の割り込み要求を受信したプロセッサが、前記第1の割り込み要求により規定される動作を実行するステップをさらに含み、前記割り込みサービスフラグをクリアするステップは、前記第1の割り込み要求が規定した動作が完全に実行される前に、前記第1の割り込み要求に対応する前記割り込みサービスフラグをクリアするステップを含むことを特徴とする請求項1に記載の方法。

【請求項 3】

10

20

前記第 1 の割り込み要求を受信したプロセッサが、前記第 1 の割り込み要求が規定した動作を実行するステップと、

前記第 1 の割り込み要求が規定した前記動作を実行する間、前記仮想割り込みコントローラーが、他の割り込み要求を受信するステップと

をさらに含むことを特徴とする請求項 1 に記載の方法。

【請求項 4】

前記第 1 の割り込み要求が規定した動作を実行するステップは、前記仮想マシンモニターにより調停されることを特徴とする請求項 3 に記載の方法。

【請求項 5】

前記仮想コンピューティングシステムは、第 1 のゲストオペレーティングシステムと、第 2 のゲストオペレーティングシステムとを備え、前記仮想マシンモニターは、前記第 1 のゲストオペレーティングシステムを含む第 1 のパーティションと、前記第 2 のゲストオペレーティングシステムを含む第 2 のパーティションを定義し、

前記第 1 の割り込み要求は、前記第 1 のパーティションから前記第 2 のパーティションへの第 1 のパーティション間メッセージに対する要求に対応することを特徴とする請求項 1 に記載の方法。

【請求項 6】

各ゲストオペレーティングシステムが、仮想プロセッサと関連する少なくとも 1 つのメッセージスロットを備え、

前記第 1 のパーティション間メッセージの処理が完了する前に第 2 のパーティション間メッセージが前記第 1 の割り込み要求の仮想プロセッサに関連する前記メッセージスロットに対する待ち行列に入れられる場合、前記仮想マシンモニターが、前記第 1 のパーティション間メッセージに関連するメッセージ保留フラグをセットするステップ

をさらに含むことを特徴とする請求項 5 に記載の方法。

【請求項 7】

前記パーティション間メッセージに関連するメッセージ保留フラグがセットされる場合、前記パーティション間メッセージを受信する仮想プロセッサが、パーティション間メッセージの処理が完了した後にメッセージ終了命令を前記仮想割り込みコントローラーに送信するステップをさらに含むことを特徴とする請求項 6 に記載の方法。

【請求項 8】

仮想コンピューティングシステムに対する割り込みを処理する方法であって、前記仮想コンピューティングシステムは仮想マシンモニター、および仮想割り込みコントローラーとゲストオペレーティングシステムとを含む少なくとも 1 つの仮想マシンを備え、前記方法は

前記仮想マシンモニターが、割り込みを前記ゲストオペレーティングシステムにリダイレクトするステップと、

前記仮想割り込みコントローラーが、前記リダイレクトされたゲストオペレーティングシステムから割り込み終了命令を受信するステップと、

前記仮想割り込みコントローラーが、前記割り込み終了命令は、サービス中としてフラグを立てた最高優先度の割り込みに対応しないことを判定するステップと、

前記仮想割り込みコントローラーが、前記受信した割り込み終了命令に対応する前記割り込みを特定する情報を記録するステップと、

サービス中としてフラグを立てた最高優先度の割り込みに対する割り込み終了命令を受信した後に、前記仮想割り込みコントローラーが、前記割り込み終了命令を処理するステップと

を含むことを特徴とする方法。

【請求項 9】

サービス中の最高優先度の割り込みに対する割り込み終了命令の処理が完了した際に、前記仮想割り込みコントローラーが、前記記録した情報に対応する割り込みに対する割り込み終了命令を処理するステップをさらに含むことを特徴とする請求項 8 に記載の方法。

10

20

30

40

50

【請求項 10】

少なくとも 1 つのプロセッサと、

仮想マシンモニターと、

ゲストオペレーティングシステムおよび仮想割り込みコントローラーを含む仮想マシンと

を備えたシステムであって、

前記仮想割り込みコントローラーは、前記仮想マシンモニターから割り込み要求を受信した後に、前記割り込み要求に対応する割り込みサービスフラグを立て、前記第 1 の割り込み要求が、ゲストオペレーティングシステムにより自動 E O I としてプログラムされている場合、前記第 1 の割り込み要求をプロセッサに送信すると、前記割り込みサービスフラグをクリアすることを特徴とするシステム。

10

【請求項 11】

前記仮想割り込みコントローラーは、第 1 の割り込み要求が規定した動作が完全に実行される前に、前記割り込み要求に対応する割り込みサービスフラグをクリアすることを特徴とする請求項 10 に記載のシステム。

【請求項 12】

前記仮想コンピューティングシステムは、第 1 のゲストオペレーティングシステムと、第 2 のゲストオペレーティングシステムとを備え、前記仮想マシンモニターは、前記第 1 のゲストオペレーティングシステムを含む第 1 のパーティションと、前記第 2 のゲストオペレーティングシステムを含む第 2 のパーティションを定義し、前記第 1 の割り込み要求は、前記第 1 のパーティションから前記第 2 のパーティションへの第 1 のパーティション間メッセージに対する要求に対応し、

20

前記仮想マシンモニターは、記第 1 のパーティションから前記第 2 のパーティションへの第 1 のパーティション間メッセージの処理が完了する前に第 2 のパーティション間メッセージが待ち行列に入れられる場合、前記第 1 のパーティション間メッセージにフラグを立て、前記第 1 のパーティション間メッセージにフラグが立てられている場合に限り、前記パーティション間メッセージを受信するプロセッサが前記第 1 のパーティション間メッセージの処理完了時に前記仮想割り込みコントローラーに信号を送ることを特徴とする請求項 10 に記載のシステム。

【請求項 13】

割り込みをゲストオペレーティングシステムにリダイレクトする仮想マシンモニターと

30

ゲストオペレーティングシステムおよび仮想割り込みコントローラーを含む少なくとも 1 つの仮想マシンと

を備えたシステムであって、

前記仮想割り込みコントローラーは、前記リダイレクトされたゲストオペレーティングシステムから割り込み終了命令を受信し、前記割り込み終了命令は、サービス中としてフラグを立てた最高優先度の割り込みに対応しないことを判定し、前記受信した割り込み終了命令に対応する割り込みを特定する情報を記録して、サービス中のより高い優先度の割り込み処理完了後に前記受信した割り込み終了命令を処理可能とすることを特徴とするシステム。

40

【請求項 14】

コンピューターに、仮想コンピューティングシステムとして機能させるためのプログラムを備えたコンピューター読み取り可能な記録媒体であって、前記仮想コンピューティングシステムは、少なくとも 1 つのプロセッサ、仮想マシンモニター、およびゲストオペレーティングシステムと仮想割り込みコントローラーとを含む少なくとも 1 つの仮想マシンを備え、

前記仮想割り込みコントローラーは、前記仮想マシンモニターから割り込みサービス要求を受信した後に、前記割り込みサービス要求に対応する割り込みサービスフラグを立て、前記第 1 の割り込み要求が、ゲストオペレーティングシステムにより自動 E O I として

50

プログラムされている場合、前記割り込みサービス要求をプロセッサに送信すると、前記割り込みサービスフラグをクリアするよう動作す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 15】

前記コンピュータに、割り込みサービス要求が規定した動作を実行させるためのプログラムをさらに備え、

前記仮想割り込みコントローラーは、前記コンピュータが前記割り込みサービス要求が規定した前記動作を実行する間、他の割り込みサービス要求を受信するようさらに動作することを特徴とする請求項 14 に記載のコンピュータ読み取り可能な記録媒体。

【請求項 16】

前記仮想マシンモニターは、第 1 のプロセッサ間メッセージの処理が完了する前に、第 2 のプロセッサ間メッセージが前記第 1 のプロセッサ間メッセージを受信するプロセッサに対する待ち行列に入れられる場合、前記第 1 のプロセッサ間メッセージに関連するメッセージ保留フラグをセットするよう動作することを特徴とする請求項 14 に記載のコンピュータ読み取り可能な記録媒体。

【請求項 17】

第 1 のプロセッサ間メッセージに関連する前記メッセージ保留フラグを、前記第 1 のプロセッサ間メッセージのヘッダ内のビットとして具現化したことを特徴とする請求項 16 に記載のコンピュータ読み取り可能な記録媒体。

【請求項 18】

前記プロセッサ間メッセージを受信するプロセッサは、プロセッサ間メッセージに関連するメッセージ保留フラグがセットされている場合、前記プロセッサ間メッセージの処理が完了した後に仮想割り込みコントローラーに信号を送るよう動作することを特徴とする請求項 16 に記載のコンピュータ読み取り可能な記録媒体。

【請求項 19】

コンピュータに、仮想マシンモニター、およびゲストオペレーティングシステムと仮想割り込みコントローラーとを含む少なくとも 1 つの仮想マシンを備えた仮想コンピューティングシステムとして機能させるためのプログラムを備えたコンピュータ読み取り可能な記録媒体であって、

前記仮想マシンモニターは、割り込みを前記ゲストオペレーティングシステムにリダイレクトするよう動作し、

前記仮想割り込みコントローラーは、前記リダイレクトされたゲストオペレーティングシステムから割り込み終了命令を受信し、前記割り込み終了命令を受信したときに処理中である最高優先度の割り込みに前記割り込み終了命令が対応しない場合、後に処理する前記割り込み終了命令を待ち行列に入れるよう動作す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 20】

前記仮想割り込みコントローラーは、待ち行列に入れた割り込み終了命令に対応する割り込みより優先度が高い割り込みに対応する割り込み終了命令の処理後、前記待ち行列に入れた割り込み終了命令を処理するようさらに動作することを特徴とする請求項 19 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、仮想環境における最適化した割り込み送信に関する。

【背景技術】

【0002】

バーチャルマシン (VM) は、エミュレートされたマシンまたはシステムを提供するためにコンピュータデバイスなど (例えば、ホスト) で動作するソフトウェア構成またはそれと同種のものである。一般に、必ずではないが、VM はアプリケーションなどであり

10

20

30

40

50

、ホスト上で使用して利用アプリケーションなどをインスタンス化し、同時に上記利用アプリケーションを上記ホスト装置から、または上記ホスト上の他のアプリケーションから隔離することができる。1つの一般的な状況では、ホストは配備されている複数のVMを収容することができ、それぞれのVMはホストから利用可能なリソースを用いて何らかの所定の機能を実行する。

【0003】

特に、コンピューターデバイス上にホストされるそれぞれのVMは、仮想的な形ではあるが実質的にコンピューティングマシンであり、従ってその利用アプリケーションおよび外部の両方に対して自身をコンピューティングマシンの様に見せる。例えば、VMおよび/またはその利用アプリケーションは、VMにハードウェアリソースが実際にはなくても、VMのハードウェアリソースにハードウェア要求を発行することができ、実際にそうしている。理解できるように、上記ハードウェア要求は遮断されるか、またはホストにリダイレクトされ、上記ホストがそのハードウェアリソースに基づいて上記ハードウェア要求を処理する。このことを、要求したVMおよび/またはその利用アプリケーションは、一般に全く認識しない。

10

【0004】

一般に、必ずではないが、ホストはそのそれぞれのVMを別々のパーティション、アドレス空間、処理領域、などに展開する。上記ホストは、仮想マシンモニター(「VMM」)などを有する仮想化層を含むことができる。VMMは監視アプリケーションまたは「ハイパーバイザ」として動作し、仮想化層はホストのそれぞれのVMの監視態様を監督および/または管理する。さらにVMMは、それぞれのVMと外部との間を接続しうるものとして動作する。VMMは、自分のアドレス空間で実行される独立したアプリケーションとすることができ、または、直接的に、もしくはデバイスドライバなどのある種のオペレーティングシステム拡張として、VMをホストオペレーティングシステムと、より密接に統合することができる。特にホストのVMMは、ホストのそれぞれのVMおよび/またはその利用アプリケーションからのハードウェア要求を遮断またはリダイレクトすることができ、少なくともその要求処理を支援することができるが、先と同様にこのことを、要求したVMおよび/またはその利用アプリケーションは全く認識しない。

20

【0005】

多くのコンピューティングシステムは複数のプロセッサを含む。マルチプロセッサの仮想マシン環境内のプロセッサは、ゲストモードまたはVMMモードで動作することができる。ゲストモードで実行されるとき、プロセッサは仮想マシンの定義を使用して、仮想マシンのゲストオペレーティングシステムとアプリケーションとを管理し、VMMからの干渉を受けずに引数を変換し、システムリソースを管理する。時々、ゲストオペレーティングシステムまたはアプリケーションが、VMMによる管理が必須のシステムリソースを必要とする場合がある。例えば、VMMがエラー処理、システム障害、または割り込み処理を要求される可能性がある。これらの状況では、プロセッサはVMMモードで動作する。

30

【0006】

近代の処理システムでは割り込みがサポートされ、それによりプロセッサに外部のイベントを通知することができる。例えば、ユーザーがキーボードのキーを押下し、またはネットワークパケットが電線から到着すると、対応する割り込みが生成されてプロセッサに送信される。一般に、割り込みによりプロセッサは現在実行中のものを停止し、現在の実行位置を記録する。それによりプロセッサは割り込み処理後に実行を再開し、特定の割り込みサービスルーチンを実行することができる。

40

【0007】

コンピューティングシステムは、システム内での割り込みの流れを指示および調停する1つまたは複数の割り込みコントローラーを含むことができる。割り込みコントローラーのロジックを別々のハードウェアコンポーネントに具現化すること、プロセッサに統合すること、または仮想化することができる。割り込みコントローラーは、特に、マルチプ

50

ロセッサ環境において割り込み優先度を決定し、割り込みを適切なプロセッサへ送信する責任を有する。仮想環境では、プロセッサおよび割り込みコントローラーを仮想化することができる。これは一般に、仮想マシンモニターのようなソフトウェアと、ハードウェアが提供する仮想化支援機構(virtualization assist)との組合せにより実現される。

【0008】

一般に、割り込みを処理した後、割り込み終了(EOI)命令を介して割り込みコントローラーに通知する。EOI命令は、前の割り込みの処理中に送信が保留された可能性がある他の割り込みを、この時点で送信できることを割り込みコントローラーに伝える。EOI命令は、一般にレジスタからの読み取りまたはレジスタへの書き込みなど、I/OポートまたはメモリーマップドI/Oアクセスを通して割り込みコントローラーに送信される。物理割り込みコントローラーに関しては、EOI命令の処理により何十何百ものサイクルが消費される可能性がある。仮想割り込みコントローラーに関しては、EOI命令の処理により何千ものサイクルが消費される可能性がある。

10

【0009】

仮想マシンモニターの中には、割り込みをパーティション間メッセージングの基礎として使用するものがある。1つのパーティション内で実行されるソフトウェアが同一物理マシン上の第2のパーティション内で実行されるソフトウェアと通信する必要がある場合、パーティション間メッセージングを用いて行うことができる。或るプロセッサがメッセージを送信すると、仮想マシンモニターはメッセージの受信者であるプロセッサに割り込みを送信し、受信プロセッサの割り込みサービスルーチンにメッセージを処理させ、その内容に応答させることができる。

20

【発明の概要】

【0010】

本概要は、以下の発明の詳細な説明でさらに説明されている選択した概念を簡略化した形式で紹介するために提供するものである。本概要は、請求項に記載されている主題事項の重要な特徴または本質的な特徴を特定することを目的とするものでも、請求項に記載されている主題事項の範囲を決定に際して、補助として使用することを目的とするものでもない。

【0011】

本明細書では、マルチプロセッサの仮想環境において割り込みを効率的に処理する機構を説明する。いくつかの実施形態では、ゲストオペレーティングシステムは或る割り込み源を「自動割り込み終了」(自動EOI)としてプログラムすることができる。自動EOIを処理する際、仮想割り込みコントローラーは明示的な割り込み終了(EOI)命令を待たずに、送信された割り込みに対応する割り込みサービスレジスタ内のビットをクリアする。自動EOI割り込みにより、他の割り込みの送信をブロックすることはできない。

30

【0012】

割り込みを使用してパーティション間通信を実施することができる。ゲストオペレーティングシステムがパーティション間メッセージに関連する割り込みを受信すると、ゲストオペレーティングシステムの割り込みサービスルーチンが指定のメッセージスロットからメッセージを読み取り、そのメッセージタイプとペイロードとに基づいて動作を実行する。以下で詳述するように、ゲストオペレーティングシステムは、別のメッセージが同一メッセージスロットに対する待ち行列内にある場合に限り、明示的なメッセージ終了(EOM)命令をメッセージ処理完了時に送信することで、パーティション間メッセージの処理により被るオーバーヘッドを幾らか排除することができる。EOM命令の計算コストは、大まかにEOIの計算コストと等しいが、EOMが送信されるのは、追加のメッセージが待ち行列内にあるという稀なケースに限られる。これにより、パーティション間通信に対する割り込み処理の平均的なコストを大幅に削減することができる。

40

【0013】

割り込みは様々な優先度を有することができる。一般に、高優先度の割り込みは低優先

50

度の割り込み処理に割り込むことができるが、逆は不可能である。仮想環境では、サービス中の最高優先度の割り込みではない割り込みに対して、ゲストオペレーティングシステムがE O I命令を発行することが可能である。本明細書では、仮想割り込みコントローラーが、着信中のE O I命令がサービス中の最高優先度の割り込みに対応するかどうかをチェックする機構を説明する。着信中のE O I命令がサービス中の最高優先度の割り込みに対応しない場合、仮想割り込みコントローラーは、着信中のE O I命令に対する割り込みベクトルを、後にE O Iの発行が必要な割り込みの集合に加える。着信中のE O I命令がサービス中の最高優先度の割り込みに対応する場合、仮想割り込みコントローラーは対応する割り込みに対するE O I命令を処理するだけでなく、後にE O Iを発行するために前もってマークした全ての他の割り込みに対するE O I命令も処理する。

10

【図面の簡単な説明】

【0014】

【図1】コンピューターシステムにおける仮想動作環境に対する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を表すブロック図である。

【図2】仮想化がホストオペレーティングシステムにより（直接またはハイパーバイザを介して）行われる仮想コンピューティングシステムを表すブロック図である。

【図3】ホストオペレーティングシステムと並行して実行される仮想マシンモニターにより仮想化が行われる、代替的な仮想コンピューティングシステムを表すブロック図である。

【図4】ホストオペレーティングシステムとは独立に動作するバーチャライザーにより仮想化が行われる、別の代替的な仮想コンピューティングシステムを表すブロック図である。

20

【図5】割り込みコントローラーを有するコンピューティングシステムの一部を示すブロック図である。

【図6】割り込み要求の処理方法を示すフローチャートである。

【図7】割り込み優先度の例に対するタイムラインを示す図である。

【図8】本明細書の教示に従って自動E O Iを用いた割り込み要求の処理方法を示すフローチャートである。

【図9】本明細書の教示に従ってプロセッサ間メッセージの処理方法を示すフローチャートである。

30

【発明を実施するための形態】

【0015】

或る特定の詳細を、以下の記述および図面で説明して、本発明の様々な実施形態の徹底的な理解を与える。コンピューティング技術およびソフトウェア技術に関連することの多い或る公知な詳細は、様々な実施形態を不必要に不明瞭とすることを回避するため、以下の開示では説明していない。さらに、当業者は、後述する1つまたは複数の詳細がなくとも他の実施形態を実践できることを理解するであろう。最後に、以下の開示におけるステップおよび順序を参照して様々な方法を説明するが、そのような説明は本発明の実施形態を明確に実施するためであり、そのステップおよびステップの順序が本発明を実践するための要件であると捉えるべきではない。

40

【0016】

本明細書で説明されている様々な技術を、ハードウェアもしくはソフトウェア、または必要に応じてその両方の組み合わせと関連させて実施できることは理解されるべきである。従って、方法および装置、またはその或る態様もしくは一部は、有形媒体で具現化したプログラムコード（例えば、命令）の形態を取ることができ、当該有形媒体には、フロッピー（登録商標）ディスク、CD-ROM、ハードドライブ、または任意の他の機械読取可能記憶媒体などがある。プログラムコードをコンピューターなどのマシンにロードして実行すると、そのマシンは本発明を実践する装置となる。プログラム可能なコンピューター上でプログラムコードを実行する場合、コンピューターデバイスは、一般にプロセッサ、プロセッサにより読み取り可能な記憶媒体（揮発性および不揮発性のメモリーおよ

50

び／または記憶要素を含む)を含み、少なくとも1つの入力装置、および少なくとも1つの出力装置を含む。1つまたは複数のプログラムは、例えばAPI、再利用可能コントロールなどを使用して、本発明に関して説明されているプロセスを実施または利用することができる。コンピュータシステムと通信するには、上記プログラムを、高レベルの手続き型言語またはオブジェクト指向プログラミング言語で実施することが望ましい。ただし、必要に応じてプログラムをアセンブリ言語または機械語で実施することができる。いずれにせよ、言語はコンパイル言語またはインタープリタ型言語に変換される言語であってもよく、ハードウェア装置と組み合わせられてもよい。

【0017】

例示的な実施形態は、1つまたは複数のスタンドアロンコンピュータシステムが存在する状況での本発明の態様の利用にすることができるが、本発明はそのようには限定されず、ネットワークまたは分散コンピューティング環境などの任意のコンピューティング環境と関連して実施することができる。さらに、本発明の態様を複数の処理チップまたは装置内で、またはそれらに渡って実施することができ、同様に複数の装置に渡って記憶装置に影響を及ぼすことができる。上記装置には、パーソナルコンピュータ、ネットワークサーバー、ハンドヘルド装置、スーパーコンピュータ、または自動車および航空機などの他のシステムに統合されるコンピュータを含めることができる。

【0018】

概要

仮想マシン環境内で割り込みを効率的に処理する様々な方法およびシステムを説明する。割り込みは、現代のコンピューティングシステムにおいて、一例として、プロセッサに外部イベントを通知すること、マルチプロセッサシステムのプロセッサ間通信を促進することを含む、様々な目的で使用されている。一般に、割り込みは、通常の処理に割り込んで、一時的に制御フローを割り込みサービスルーチン(「ISR」)に回す。コンピューティングシステムの様々な活動が割り込みを引き起こす可能性がある。いくつかの例として、キーボードのキーの押下、ネットワークパケットの受信、およびディスクの読み書きがある。プロセッサ間割り込み(「IPI」)は一種の割り込みであり、それによりマルチプロセッサ環境において或るプロセッサは別のプロセッサに割り込むことができる。IPIをプロセッサ間メッセージングの基礎として使用することができる。

【0019】

コンピューティングシステムは、一般にシステム内の割り込みの流れを指示および調停する1つまたは複数の割り込みコントローラーを含む。割り込みコントローラーは着信した割り込みの優先付け、それらの割り込みをマルチプロセッサ内の適切なプロセッサに割り当てる責任を有する。割り込みコントローラーをハードウェアで実現することができ、従って割り込みコントローラーは別々のコンポーネントとすることができ、またはプロセッサに統合することができる。割り込みコントローラーを仮想化してもよい。割り込みコントローラーの仮想化は、一般にソフトウェアと、ハードウェアが与える仮想化支援機構とを組み合わせることで実現される。ソフトウェアは、物理割り込みコントローラーと同一の基本機能を行う仮想マシンモニターの一部であってもよい。

【0020】

一般に、それぞれの割り込み源は、指定の割り込み優先度を有する。一実施例として、これらの優先度に0から255の番号を付与することができる。255は最高優先度であり、0は最低優先度である。高優先度の割り込みは、低優先度の割り込みを処理する割り込みサービスルーチンに割り込むことができるが、低優先度の割り込みは高優先度の割り込みに割り込むことはできない。割り込みサービスルーチンが実行を完了すると、それを実行していたプロセッサは一般にEOI命令を発行して、割り込みコントローラーに割り込みの処理が完了したこと、および保留された低優先度の割り込みをこの時点で送信できることを知らせる。

【0021】

10

20

30

40

50

仮想割り込みコントローラーでは、E O I 命令は、物理割り込みコントローラーが E O I 命令に応答して行うものと同じの動作を行うソフトウェアで実施される。これには一般に、E O I ポートまたはレジスタへのアクセスを遮断してソフトウェアハンドラを呼び出すことが含まれる。遮断およびソフトウェアハンドラを組み合わせると、一般に E O I 命令の処理に何千何万ものサイクルが必要となり、仮想環境における割り込みサービスルーチンの動作に大量のオーバーヘッドが加わる。

【 0 0 2 2 】

本明細書で説明する方法およびシステムにより、割り込みを効率的に処理する機構が実現される。多くの場合において、E O I 命令を省略し、割り込み送信に関する仮想化オーバーヘッドを大幅に削減することができる。プロセッサ間メッセージングに用いる I P I の場合、丁度処理したメッセージを含むスロットに対して第 2 のメッセージが既に待ち行列内にあるときに限り、メッセージ E O I の終了を送信する必要がある。いくつかの場合、物理割り込みに対して選択的に E O I を送信することができるが、これはその物理割り込みがサービス中の最高優先度の割り込みであるかどうかに関係である。

【 0 0 2 3 】

一般的な仮想化

オペレーティングシステムおよびプロセッサの命令セットにおける多様性のため、ソフトウェアの相互運用性が低下する可能性がある。高レベル言語およびオペレーティングシステム双方におけるメモリーおよび I / O 抽象化により、ハードウェアリソースへの依存性を幾分削減することができるが、幾分かの依存性は残ることになる。多くのオペレーティングシステムが特定のシステムアーキテクチャに対して開発され、ハードウェアリソースを直接管理するように設計されている。これにより、利用可能なソフトウェアおよびオペレーティングシステムに関してコンピューターシステムの柔軟性が限定される可能性があり、特にシステムが複数ユーザーにより共有されるとき、セキュリティおよび障害隔離に悪影響を及ぼす可能性がある。

【 0 0 2 4 】

仮想化により、セキュリティおよび信頼性を高めつつ柔軟性を向上させる機構が実現される。プロセッサ、メモリー、および I / O 装置は、仮想化できるサブシステムの例である。サブシステムを仮想化すると、仮想インターフェースと、その仮想インターフェースを通して利用可能な仮想リソースとが、仮想化を実施した実システムのインターフェースとリソースとにマッピングされる。仮想化を、サブシステムだけでなくマシン全体に適用することができる。仮想マシンのアーキテクチャは、実マシン上のソフトウェア層で実施される。

【 0 0 2 5 】

概念的観点からは、コンピューターシステムは、一般に基本的なハードウェア層で実行される 1 つまたは複数のソフトウェア層を含む。この階層化は抽象化のために行われている。所与のソフトウェア層に対してインターフェースを定義することで、その層をそれより上の他の層と別々に実施することができる。良く設計されたコンピューターシステムでは、それぞれの層は直下の層に関してのみ知っている（および、直下の層にのみ依存する）。これにより層または「スタック」（複数の隣接層）を、上記層またはスタックより上の層に悪影響を及ぼすことなく置換することができる。例えば、ソフトウェアアプリケーション（上層）は、一般に低レベルのオペレーティングシステム（低層）に依存して何らかの形の永久記憶装置にファイルを書き込むが、これらのアプリケーションはフロッピー（登録商標）ディスク、ハードドライブ、またはネットワークフォルダへデータを書き込むことの違いを理解する必要はない。この低層を、ファイル書き込み用の新しいオペレーティングシステムのコンポーネントで置換しても、上層のソフトウェアアプリケーションの動作は、影響を受けないままである。

【 0 0 2 6 】

階層化したソフトウェアの柔軟性により、V M (virtual machine) は、実際は別のソフトウェア層にある仮想ハードウェア層を表現することができる。このように、V M はそ

10

20

30

40

50

れより上のソフトウェア層に対して、自分がプライベートなコンピューターシステム上で実行されているかのように錯覚させることができ、従って、VMにより複数の「ゲストシステム」を単一の「ホストシステム」上で並列に実行することができる。

【0027】

図1は、コンピューターシステム内の仮想環境に対する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を表す図である。図1で、仮想化プログラム110は、物理ハードウェアアーキテクチャ112上で直接的または間接的に実行される。仮想化プログラム110は、(a)ホストオペレーティングシステムと並行に実行される仮想マシンモニターとすることもでき、または(b)ハイパーバイザコンポーネントを有するホストオペレーティングシステムとすることができる。この場合、ハイパーバイザコンポーネントが仮想化を行う。仮想マシンモニターという用語は、様々な種類の仮想化プログラムのうち何れかに対する一般的な用語として用いる。仮想化プログラム110は、(このコンポーネントがパーティションまたは「仮想マシン」である事実を表すため点線で示してある)ゲストハードウェアアーキテクチャ108を仮想化する。即ち、ハードウェアは実際には存在しないが、その代わりに仮想化プログラム110により仮想化される。ゲストオペレーティングシステム106を、ゲストハードウェアアーキテクチャ108上で実行し、ソフトウェアアプリケーション104を、ゲストオペレーティングシステム106上で実行することができる。図1の仮想化した動作環境では、ソフトウェアアプリケーション104が、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112と一般に非互換であるオペレーティングシステム上で実行するように設計されていても、ソフトウェアアプリケーション104をコンピューターシステム102内で実行することができる。

【0028】

次に、図2は、物理コンピューターハードウェア202の直上で実行されるホストオペレーティングシステム(ホストOS)のソフトウェア層204を含む仮想コンピューティングシステムを示す。ここで、ホストOS 204は、オペレーティングシステムA 212およびオペレーティングシステムB 214によりそれぞれ使用されるパーティションA 208およびパーティションB 210にインターフェースを公開することで、物理コンピューターハードウェア202のリソースへのアクセスを提供する。これによりホストOS 204は、その上で実行されるオペレーティングシステム層212および214に注目されずに済む。繰り返しになるが、仮想化を行うため、ホストOS 204はネイティブな仮想化機能を持つ特殊設計のオペレーティングシステムとすることができ、あるいは仮想化を行うための組み込みのハイパーバイザコンポーネントを持つ標準的なオペレーティングシステムとすることができる(図示せず)。

【0029】

再度図2を参照する。ホストOS 204の上には、2つのパーティション、パーティションA 208およびパーティションB 210がある。パーティションA 208は、例えば仮想のインテル386プロセッサとすることができ、パーティションB 210は、例えばモトローラ680X0系プロセッサの1つを仮想化したものとすることができる。それぞれのパーティション208およびパーティション210内には、それぞれゲストオペレーティングシステム(ゲストOS)A 212およびゲストOS B 214がある。ゲストOS A 212上では2つのアプリケーション、アプリケーションA 1216およびアプリケーションA 2218が実行され、ゲストOS B 214上ではアプリケーションB 1220が実行されている。

【0030】

図2に関して、(点線で示す)パーティションA 208とパーティションB 214とは、ソフトウェア構造としてのみ存在する仮想コンピューターハードウェア表現であることに留意することが重要である。これらは、特殊な仮想化ソフトウェアを実行することで可能となるものである。当該ソフトウェアは、パーティションA 208およびパーティションB 210をそれぞれゲストOS A 212およびゲストOS B 214に

提示するだけでなく、ゲストOS A 212およびゲストOS B 214が実物理コンピュータハードウェア202と間接的に相互作用するために必要な、全てのソフトウェアステップも実行する。物理コンピュータハードウェア202は、単一プロセッサ環境の場合は単一のCPU (central processing unit) 222を、マルチプロセッサ環境の場合は複数のCPU 222、224、226を含むことができる。

【0031】

図3は、ホストオペレーティングシステム306と並行に実行される仮想マシンマネージャ304により仮想化が行われる、代替的な仮想コンピューティングシステムを示す。或る場合、仮想マシンマネージャ304は、ホストオペレーティングシステム306上で実行され、ホストオペレーティングシステム306を通してのみコンピュータハードウェア302と相互作用するアプリケーションとすることができる。図3に示すような他の場合、仮想マシンマネージャ304は、代わりに部分的に独立したソフトウェアシステムを含むことができる。当該ソフトウェアシステムは、いくつかのレベルではホストオペレーティングシステム306を介してコンピュータハードウェア302と間接的に相互作用するが、他のレベルでは、仮想マシンマネージャ304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相互作用するのと同様に)コンピュータハードウェア302と直接相互作用する。さらに他の場合、仮想マシンマネージャ304は、完全に独立したソフトウェアシステムを含むことができる。当該ソフトウェアシステムは、ホストオペレーティングシステム306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相互作用するのと同様に)全てのレベルでコンピュータハードウェア302と直接相互作用する(しかしながら、コンピュータハードウェア302の使用を調整し、衝突を回避するなどのため、ホストオペレーティングシステム306とも相互作用する)。

【0032】

図3に示す実施例では、2つのパーティション、パーティションA 308およびパーティションB 310は、概念的に仮想マシンマネージャ304の上にある。それぞれのパーティション308およびパーティション310内には、それぞれゲストオペレーティングシステム(ゲストOS) A 312およびゲストOS B 314がある。ゲストOS A 312上では2つのアプリケーション、アプリケーションA1 316およびアプリケーションA2 318が実行され、ゲストOS B 314上ではアプリケーションB1 320が実行されている。物理コンピュータハードウェア302は、単一プロセッサ環境の場合は単一のCPU (central processing unit) 322を、マルチプロセッサ環境の場合は複数のCPU 322、324、326を含むことができる。

【0033】

図4は、ハイパーバイザ404により仮想化が行われる別の代替的な仮想コンピューティングシステムを示す。ハイパーバイザ404は、ホストオペレーティングシステムを用いずにコンピュータハードウェア402と直接相互作用できる独立したソフトウェアシステムを含む。物理コンピュータハードウェア402は、単一プロセッサ環境の場合は単一のCPU (central processing unit) 422を、マルチプロセッサ環境の場合は複数のCPU 422、424、426を含むことができる。

【0034】

図4に示す実施例では、2つのパーティション、パーティションA 408およびパーティションB 410は、概念的に仮想マシンマネージャ404の上にある。それぞれのパーティション408およびパーティション410内には、それぞれゲストオペレーティングシステム(ゲストOS) A 412およびゲストOS B 414がある。ゲストOS A 412上では2つのアプリケーション、アプリケーションA1 416およびアプリケーションA2 418が実行され、ゲストOS B 414上ではアプリケーションB1 420が実行されている。ゲストOS A 412はホストOSのサービスを提

供する。物理コンピューターハードウェア402は、単一プロセッサ環境の場合は単一のCPU (central processing unit) 422を、マルチプロセッサ環境の場合は複数のCPU 422、424、426を含むことができる。

【0035】

上述のパーティションを実施するこれらの変形の全ては例示的な実施例に過ぎず、本開示の発明を特定の仮想化の態様に限定するものと解釈すべきではない。

【0036】

一般的な割り込み処理

図5は、割り込みコントローラーを有するマルチプロセッサコンピューティングシステムの一部分の実施例を示すブロック図である。任意の数の装置502、504、506、508が、割り込み要求源の役割を果たすことができる。装置502、504、506、508は、例えばキーボード、ディスクドライブ、ネットワークカードなどの物理装置とすることができ、または仮想装置とすることができ、割り込み要求を、プロセッサ510、512、514の何れかにより生成してもよい。

【0037】

割り込みコントローラー516は、割り込み要求の処理を調停および指示する。割り込みコントローラー516は、プログラム可能な割り込みコントローラー(「PIC」)または高度なプログラム可能な割り込みコントローラー(「APIC」)などの物理装置とすることができ、あるいは、割り込みコントローラー516を仮想化してもよく、この場合、その機能はVMM内のソフトウェアハンドラなどのソフトウェアにより実行される。

【0038】

大部分の割り込みコントローラーは、要求されたサービス中の割り込み要求を追跡する。これは2つのビットベクトルを用いてしばしば行われ、それぞれのビットは個々の割り込み源を表す。第1のビットベクトルは割り込み要求レジスタ518と呼ばれ、第2のビットベクトルは割り込みサービスレジスタ520と呼ばれる。割り込みコントローラー516が割り込みの要求を受信すると、割り込み要求レジスタ518内の対応ビットをセットする。割り込みコントローラー516が割り込みをプロセッサ510、512、または514に送信すると、割り込みコントローラー516は割り込み要求レジスタ518内の対応ビットをクリアして、割り込みサービスレジスタ520内の対応ビットをセットする。割り込みコントローラー516がEOIを受信すると、割り込みコントローラー516は、対応する割り込みがもはや行われていないことを知り、従って割り込みサービスレジスタ520内の対応ビットをクリアする。この時点で、割り込みコントローラー516は割り込み要求レジスタ518をスキャンして、まだサービスされていない、要求された最高優先度の割り込みを決定する。上記割り込みの優先度がサービス中の最高優先度の割り込みより高い場合、割り込みコントローラーはサービス中の低優先度の割り込みの割り込みサービスルーチンに割り込む。

【0039】

従来の割り込み要求のライフサイクルを図6に示す。装置またはソフトウェアは割り込み要求(「IRQ」)をアサートすることで処理を開始する(602)。割り込み要求を多様な情報源により生成することができる。限定するためではなく例として、情報源にはキーボード、マウス、音源カード、モデム、通信ポート、時間測定装置、およびソフトウェア命令を含めることができる。「レベルトリガ」割り込みでは、割り込みを知らせたい装置は、割り込み要求の電圧を「アクティブ」と定義される規定のレベルまで上げ、割り込みが処理されるまでそのレベルに保つ。「エッジトリガ」割り込みでは、割り込み要求は割り込み要求ラインのレベル遷移により送信され、割り込みを送信したい装置は割り込み要求ラインにパルスを駆動して、その後ラインを休止状態に戻す。

【0040】

割り込みコントローラーがIRQを検出すると(604)、割り込みコントローラーは、IRQが現在サービス中のどの割り込みよりも高い優先度を持つかどうかを判定する(

606)。この判定は、割り込みサービスレジスタ520(図5)を検証することで可能である。IRQの検出時により優先度が高い割り込みがサービス中である場合、割り込みコントローラーは割り込み要求レジスタ内の対応ビットのフラグを立て、保留されたIRQを記録する(608)。要求された割り込みがサービス中のどの割り込みよりも高い優先度を有する場合、割り込みコントローラーは割り込みサービスレジスタ内の対応ビットのフラグを立て(610)、対応する割り込みサービスルーチンを実行するよう適切なプロセッサに知らせる(612)。割り込みサービスルーチンの実行完了後(614)、プロセッサは、一般にI/Oポートまたはメモリーマップドレジスタを読み書きすることで、割り込みが処理されたことを割り込み装置に示すことができる(616)。プロセッサは、割り込みコントローラーにEOI命令を送信することで、割り込みが処理されたことを知らせる(618)。当該EOI命令は、一般にレジスタ読み書きなどの、I/OポートまたはメモリーマップドI/Oアクセスを通して送信される。割り込みコントローラーはEOI命令を処理し、割り込みサービスレジスタ内の対応するフラグをクリアする(620)。EOI命令は、割り込みコントローラーに、保留された低優先度の割り込みがこの時点で送信できることを伝える。例えば、高優先度のキーボード割り込みが処理されている間にネットワークパケットが到着して割り込みを引き起こす場合に、割り込みコントローラーにより高優先度のキーボード割り込みが完全に処理されるまで、ネットワークパケットの割り込み要求を保留しておくことができる。EOI命令の処理には、物理割り込みコントローラーに対して何十または何百ものサイクルを要し、仮想割り込みコントローラーに対しては何千ものサイクルを要する可能性がある。

【0041】

図7は、割り込み優先度の一般的な概念を示す例であり、限定を意図するものではない。時刻 t_1 に優先度10の割り込み源が、割り込みを要求すると仮定する(702)。割り込みコントローラーはプロセッサに割り込み、当該プロセッサがその割り込み源に関連する割り込みサービスルーチンを呼び出す(704)。ここで、優先度10の割り込みに対するISRが処理を完了する前の時刻 t_2 に、優先度200の割り込み源が、同一プロセッサ宛ての割り込みを要求すると仮定する(706)。割り込みコントローラーは再度プロセッサに割り込み、優先度200の割り込みに対するISR708の実行を開始し、優先度10の割り込みに対するISRを一時停止する(710)。優先度200のISR708が完了する前の時刻 t_3 に、第3の割り込み源が、優先度50で割り込みを要求すると仮定する(712)。割り込みコントローラーは、プロセッサが優先度200のISRの実行を時刻 t_4 で完了するまで、この割り込みの送信を保留する(714)。優先度50のISRは、優先度200のISR708が完了した後に呼び出される(716)。優先度50のISR716が完了すると、時刻 t_5 で優先度10のISRを再開し(718)、時刻 t_6 で完了する。

【0042】

仮想割り込みコントローラー

仮想環境では、プロセッサおよび割り込みコントローラーを仮想化することができる。これは、ソフトウェア(例えば、VMM)とハードウェアが提供する仮想化支援機構との組合せにより行われる。一般的な構成では、EOI命令はVMMによりエミュレートされる。これは、EOIポートまたはレジスタへのアクセスを遮断することで行われる。この遮断によりVMM内で、EOIに応答して物理割り込みコントローラーと同一の機能を実行するソフトウェアハンドラが呼び出される。遮断およびソフトウェアハンドラの組合せにより、何千何万ものサイクルが必要とされる可能性がある。これにより、仮想環境内での実行時にISRに対して大量のオーバーヘッドが加わる。

【0043】

一般にVMMは割り込みを受け入れ、仮想割り込みとしてゲストオペレーティングシステムにリダイレクトする。割り込みを様々な情報源から生成することができ、当該情報源には限定するためではなく例として、物理ハードウェア装置、ハードウェア装置をエミュレートするパーティション、メッセージを送りたい、またはイベントを別のパーティショ

ンに送信したいパーティション、またはパーティションに信号を送りたいVMMが含まれる。VMMは、一般に割り込みを受け入れた後に、EOI命令を物理割り込みコントローラーに発行する。レベルトリガ割り込みに対して、ゲストオペレーティングシステム内のISRがEOI命令を仮想割り込みコントローラーに対して実行および発行するまで、EOI命令を発行するのは、一般に安全ではない。APICのような或る物理割り込みコントローラーは、サービス中の最高優先度の割り込みに対してのみEOIを発行することを可能とする。しかしながら或る状況では、VMMは、サービス中の最高優先度の割り込みではない割り込みに、選択的にEOIを発行する必要がある場合がある。例として、最初の割り込みの優先度が低い2つのレベルトリガ割り込みが次々と到着する際に生ずることを考える。VMMは第1の割り込みを受入れ、ゲストオペレーティングシステムにリダイレクトする。ゲスト内のISRがEOI命令を発行する前に、第2の割り込みが到着する。その後、ゲストオペレーティングシステムが、第1の割り込みでEOI命令を発行する。この場合、VMMは第1の割り込みでEOIを発行することはできない。なぜならば、最高優先度の割り込みが既にサービス中であり、EOI命令を物理割り込みコントローラーに発行すると、その最高優先度の割り込みでEOIを発行することになるからである。

【0044】

本明細書の開示によると、ゲストオペレーティングシステムはいくつかの割り込み源を「自動EOI」としてプログラムすることができる。割り込み源を自動EOIとマークすると、従来の割り込み優先度の振る舞いが修正される。自動EOI割り込みにより、他の割り込みの送信はブロックされない。従って、他の自動EOI割り込みを含む任意の他の割り込みがその関連するISRの実行に割り込み可能という点で、実質的に自動EOI割り込みは、最低優先度の割り込みと同様に振舞う。

【0045】

自動EOI割り込みが送信されると、割り込みサービスレジスタ内の自動EOI割り込みに関連するビットは、即座にクリアされる。効果的に、自動EOI割り込みが送信された時点で、仮想割り込みコントローラーがEOIを自動生成する。自動EOI割り込みを用いて、割り込み源が自身を管理し、前の割り込みが処理されたことを知るまで後続の割り込みを要求しないことが望ましい。そうでなければ、それぞれの後続の割り込みが前のISRに割り込み、プロセッサのスタックをオーバーフローさせる可能性がある。

【0046】

一実施形態では、自動EOIのプロパティを、SINT (synthetic interrupt source) に関連する仮想レジスタ内で指定する。仮想レジスタのフォーマットは以下の通りである。

【0047】

【表1】

ビット	説明	属性
63:18	予約済み (RsvdP) (値は変更不可)	読み取り／書き込み
17	自動EOI 割り込み送信時に暗黙的なEOI発行を行う場合にセット	読み取り／書き込み
16	SINTをマスクする場合にセット	読み取り／書き込み
15:8	予約済み (RsvdP) (値は変更不可)	読み取り／書き込み
7:0	ベクトル	読み取り／書き込み

【0048】

仮想プロセッサの生成時に、全てのSINTレジスタのデフォルト値は、0x00000000000000010000である。このように、全てのSINTはデフォルトでマスクされる。ゲストは、適切なベクトルでプログラミングし第16ビットをクリアすることでそれらのマスクを取らなければならない。

【 0 0 4 9 】

自動 E O I フラグは、割り込みを仮想プロセッサに送信するとき暗黙的な E O I 発行を VMM が行わなければならないことを示す。さらに、VMM は、仮想割り込みコントローラー内のサービス中のレジスタにある対応フラグを自動的にクリアする。ゲストがこの振る舞いを可能とする場合、ゲストはその割り込みサービスルーチン内で明示的な E O I 発行を行ってはならない。

【 0 0 5 0 】

図 8 は、本明細書の開示に従う自動 E O I 割り込み要求のライフサイクルを示す。装置またはソフトウェアは、I R Q をアサートすることで処理を開始する (8 0 2)。割り込みコントローラーは、I R Q を検出すると (8 0 4)、I R Q が、現在サービス中とマークされているどの割り込みよりも高い優先度を持つかどうかを判定する (8 0 6)。この判定は、割り込みサービスレジスタ 5 2 0 (図 5) を検証することで可能である。I R Q の検出時に高優先度の割り込みがサービス中とマークされている場合、割り込みコントローラーは割り込み要求レジスタ内の対応ビットのフラグを立て、その保留中の要求を記録する (8 0 8)。要求された割り込みがサービス中とマークされたどの割り込みよりも高い優先度を有する場合、割り込みコントローラーは適切なプロセッサに、対応する割り込みサービスルーチンを実行するよう伝え (8 1 0)、割り込みサービスレジスタ内の対応ビットをクリアする (8 1 2)。プロセッサが割り込みサービスルーチンの実行を完了した後 (8 1 4)、一般に I / O ポートまたはメモリーマップドレジスタを読み書きすることで、プロセッサは割り込み装置に割り込みが処理されたことを伝えることができる (8 1 6)。

【 0 0 5 1 】

従って、自動 E O I 割り込み要求により、プロセッサが明示的な E O I 命令を発行する必要はない。割り込みをプロセッサに送信したとき割り込みサービスレジスタ内の対応ビットはクリアされているので (8 1 2)、自動 E O I 割り込みにより他の割り込みの送信がブロックされることはない。割り込みの送信時にその割り込みに対して効率的に E O I が発行されるので、E O I 命令の処理に通常必要な計算サイクルは使用されない。

【 0 0 5 2 】

仮想割り込みコントローラーを使用することで、物理割り込みコントローラーがそのような機能をサポートしないときに、VMM は物理割り込みに対して選択的に E O I を発行することができる。これは、保留 E O I のリスト、即ち、後に E O I を発行する必要がある割り込みを保持することで実現することができる。例えば、ゲストオペレーティングシステムが E O I 命令を発行すると、VMM は、E O I が発行されている割り込みが本当に物理割り込みコントローラー内でサービス中の最高優先度の割り込みであるかどうかをチェックすることができる。最高優先度の割り込みでない場合、VMM は、単に割り込みを保留 E O I のリストに追加する。他方、E O I を発行されている割り込みがサービス中の最高優先度の割り込みである場合、VMM は現在の割り込みに E O I を発行するだけでなく、保留 E O I のリストにある他の割り込みにも E O I を発行する。

【 0 0 5 3 】

パーティション間メッセージング

VMM の中には、割り込みをパーティション間メッセージングの基礎として用いるものもある。パーティションは VMM が置いた分離境界であり、仮想マシンに対する「コンテナ」である。或るパーティション内で実行されるソフトウェアが同一マシン上の第 2 のパーティション内で実行されているソフトウェアと通信する必要がある場合、パーティション間メッセージを用いて通信を行うことができる。これらのメッセージは、一般に少量のペイロードを含む。例えば、或る公知のハイパーバイザの場合、メッセージは、最大 2 4 0 バイトのメッセージペイロード、および 1 6 バイトのヘッダを含むことができる。メッセージが送信されると、宛先パーティションに関連する仮想プロセッサが実行可能となるまで、そのメッセージはハイパーバイザにより待ち行列に入れられる。その時点で、ハイパーバイザは割り込みをその仮想プロセッサに送信することができる。これにより、

対応するISRが呼び出される。ISRはメッセージの読み取り、およびその内容に対処する責任を有する。上述のように、ISRは、一般に割り込みの処理後にその割り込みに対して「EOI発行」しなければならない。この場合、EOIはメッセージの読み取り後に送信される。パーティション間メッセージングは、可能な限り高速でなければならない。従来のEOI機構では、好ましくないオーバヘッドを加え、恐らく何万ものサイクルを使用しないと、メッセージが読み出されたこと、および後続のメッセージならびに低優先度の保留された割り込みが送信可能であることを仮想割り込みコントローラーに通知することができない。

【0054】

本明細書の開示によると、パーティション間メッセージが処理されたことを伝えるための、従来のEOI機構のオーバヘッドを回避することができる。一実施形態では、メッセージスロットをそれぞれのSINTに与え、メッセージのレイアウトを以下の表1で記述するデータ構造により定義する。

表 1

```
typedef struct
```

```
{
    UINT8 MessagePending:1;
    UINT8 Reserved:7;
} HV_MESSAGE_FLAGS;
```

```
typedef struct
```

```
{
    HV_MESSAGE_TYPE MessageType;
    UINT8PayloadSize;
    HV_MESSAGE_FLAGS MessageFlags
    UINT8Reserved[2];
    union
    {
        HV_PARTITION_ID Sender;
        HV_PORT_ID Port;
    };
} HV_MESSAGE_HEADER;
```

```
#define HV_MESSAGE_MAX_PAYLOAD_BYTE_COUNT 240
```

```
#define HV_MESSAGE_MAX_PAYLOAD_QWORD_COUNT 30
```

```
typedef struct
```

```
{
    HV_MESSAGE_HEADER Header;
    UINT64
    Payload[HV_MESSAGE_MAX_PAYLOAD_QWORD_COUNT];
} HV_MESSAGE;
```

【0055】

図9は、本明細書の開示に従ってプロセッサ間メッセージを処理する実施形態を示すフローチャートである。送信プロセッサは、自動EOIと指定した特定のSINTに対応するプロセッサ間メッセージを送る(902)。VMMはメッセージをメッセージ待ち行列に追加し(904)、特定のSINTに対応するメッセージスロットが空かどうかを判定する(906)。前のメッセージがメッセージスロット内にまだ存在する場合、VMMはスロットにあるメッセージのヘッダ内のメッセージ保留ビットをセットする(908)。メッセージスロットが空である場合、VMMはメッセージをメッセージスロットに

10

20

30

40

50

コピーし(910)、特定のSINTに関連する割り込みを受信プロセッサに送信する(912)。受信プロセッサ上で実行されているゲストOSがSINTに関連する割り込みを受信すると、そのISRはメッセージを対応するメッセージスロットから読み取り、メッセージタイプおよびペイロードに基づいて動作を行う(914)。メッセージの処理が完了すると、ISRはメッセージタイプをクリアする(916)。例えば、表1に定義したデータ構造によると、ISRは特定の値をHV_MESSAGE_TYPEに書き込むことで、メッセージタイプをクリアすることができる。ISRは次いで、丁度処理したメッセージのメッセージ保留ビットを調べる(918)。メッセージ保留ビットがセットされていない場合、これはメッセージスロットに対する待ち行列にメッセージがないことを意味し、ISRが必要とする動作はない(920)。これが最もよくある場合である。特に、ISRがEOI命令を送信する必要はなく、従って大量の計算オーバーヘッドが回避される。

10

【0056】

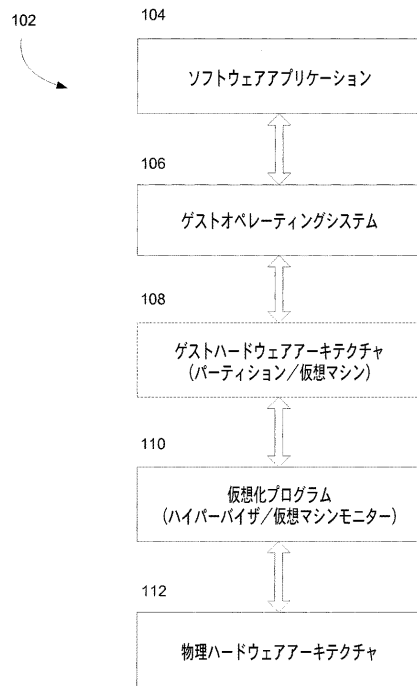
メッセージ保留ビットがセットされている場合、ISRは仮想割り込みコントローラにメッセージ終了(「EOM」)命令を送信し(922)、待ち行列に入れたメッセージの送信を再試行するようVMMに伝える。EOM命令の計算コストは大まかにEOI命令と同じであるが、EOMは、追加のメッセージがメッセージスロットに対する待ち行列に入れられる稀な場合にのみ送信される。従って、プロセッサ間メッセージを処理する平均的なコストは大幅に削減される。

【0057】

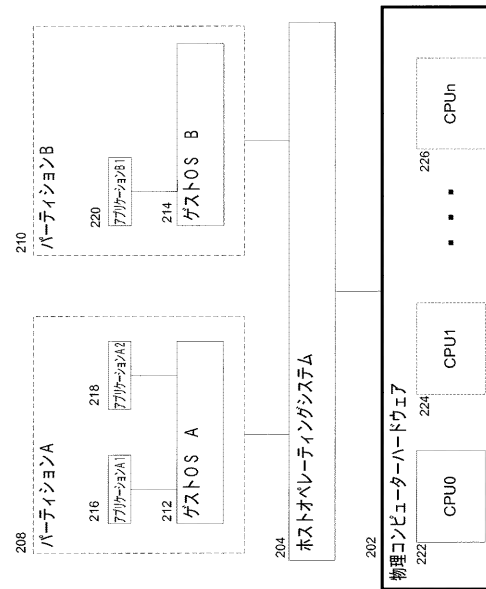
20

本開示を様々な実施形態に関連して説明し、様々な図面で示したが、本開示から逸脱することなく、同様な態様を使用し、または開示した実施形態の記載された態様に対して修正および追加を行って本開示の同一機能を実行することができることは理解される。例えば、本開示の様々な態様において、仮想環境における割り込み処理の動作効率を向上させる機構を開示した。しかしながら、これら説明した態様と同等な他の機構も、本明細書の教示により考慮されている。従って、本開示はどの態様にも限定されず、添付の特許請求の範囲に従って広範囲に解釈されるべきである。

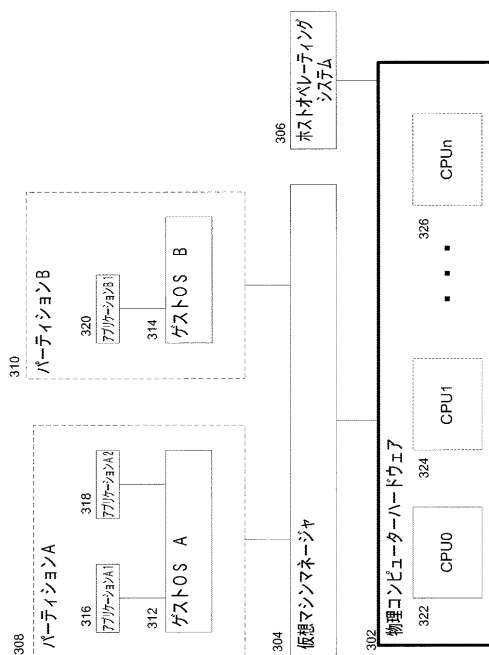
【図 1】



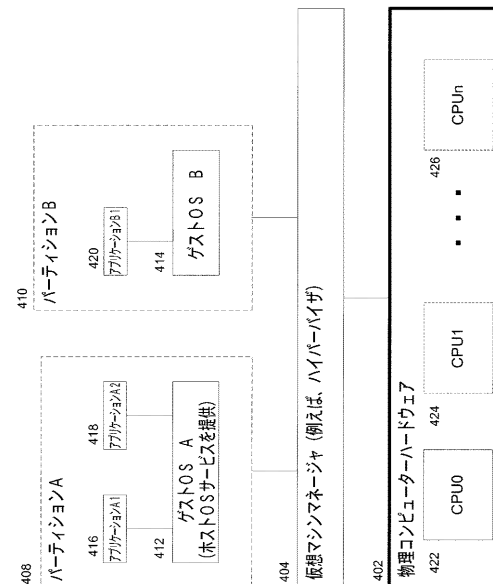
【図 2】



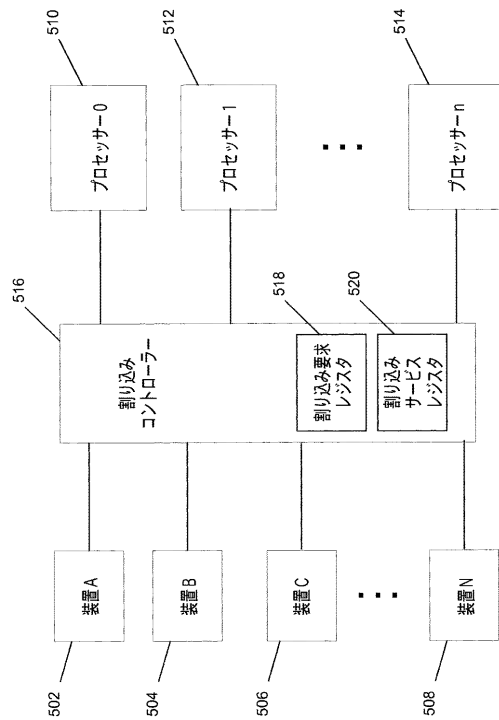
【図 3】



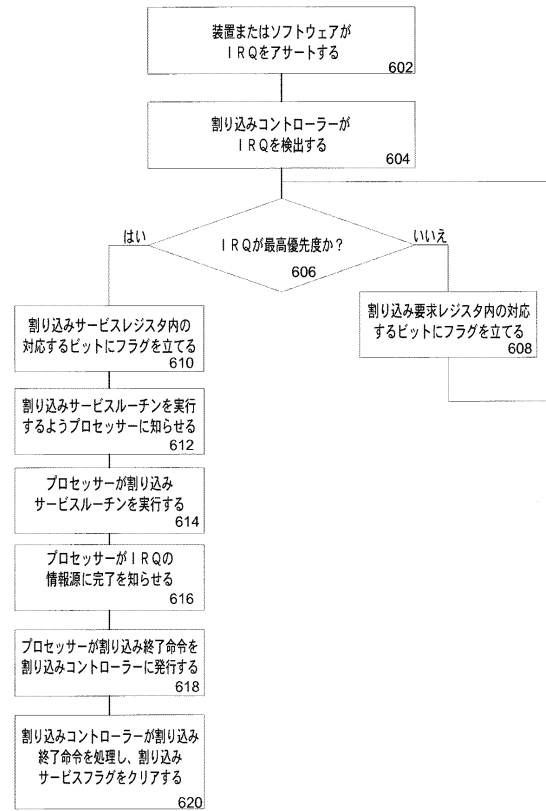
【図 4】



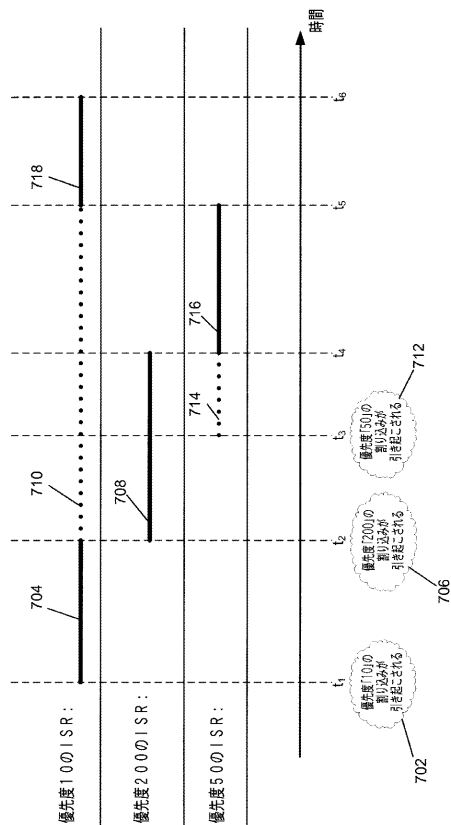
【図 5】



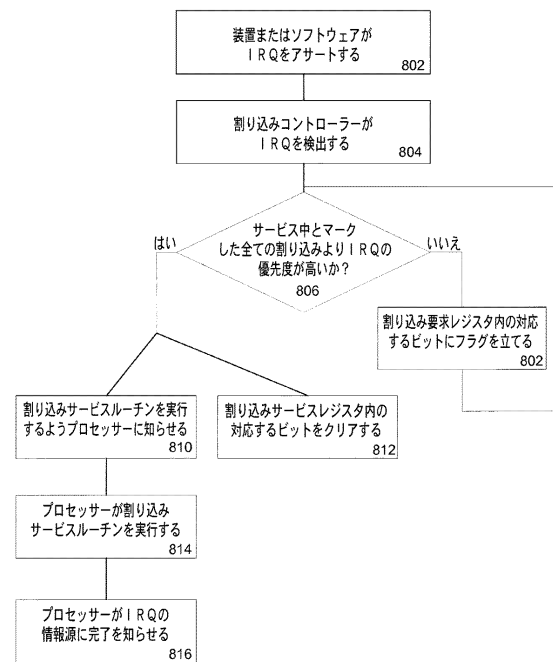
【図 6】



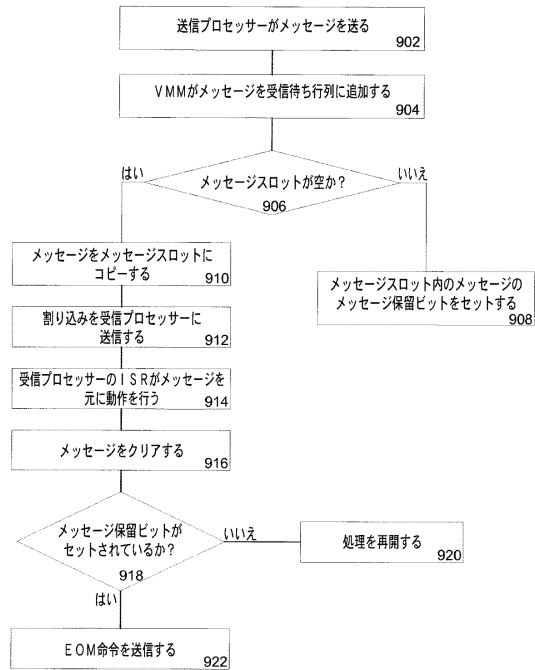
【図 7】



【図 8】



【図 9】



フロントページの続き

- (72)発明者 エリック ビー．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 インターナショナル パテント内
- (72)発明者 シューヴァブラタ 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 インターナショナル パテント内
- (72)発明者 ルネ 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 インターナショナル パテント内

審査官 鈴木 修治

- (56)参考文献 特開平05-040643(JP,A)
特開平01-093830(JP,A)
特開平06-083642(JP,A)
特開平09-097184(JP,A)
特開平01-191234(JP,A)
特開平03-065734(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 9/46
G06F 9/48