

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号
特許第5680554号
(P5680554)

(45) 発行日 平成27年3月4日(2015.3.4)

(24) 登録日 平成27年1月16日(2015.1.16)

(51) Int.Cl.

F I

GO6F 9/48 (2006.01)

GO6F 9/46 (2006.01)

GO6F 9/46 3 1 1 B

GO6F 9/46 3 5 0

GO6F 9/46 3 1 1 F

請求項の数 14 (全 35 頁)

(21) 出願番号	特願2011-548227 (P2011-548227)	(73) 特許権者	591016172
(86) (22) 出願日	平成22年1月26日 (2010.1.26)		アドバンスト・マイクロ・デバイス・
(65) 公表番号	特表2012-515995 (P2012-515995A)		インコーポレイテッド
(43) 公表日	平成24年7月12日 (2012.7.12)		ADVANCED MICRO DEVI
(86) 国際出願番号	PCT/US2010/022111		CES INCORPORATED
(87) 国際公開番号	W02010/085804		アメリカ合衆国、94088-3453
(87) 国際公開日	平成22年7月29日 (2010.7.29)		カリフォルニア州、サニペール、ピー・
審査請求日	平成25年1月28日 (2013.1.28)		オウ・ボックス・3453、ワン・エイ・
(31) 優先権主張番号	61/147,269		エム・ディ・プレイス、メイル・ストップ
(32) 優先日	平成21年1月26日 (2009.1.26)		・68 (番地なし)
(33) 優先権主張国	米国 (US)	(74) 代理人	100108833
(31) 優先権主張番号	12/611,595		弁理士 早川 裕司
(32) 優先日	平成21年11月3日 (2009.11.3)	(74) 代理人	100111615
(33) 優先権主張国	米国 (US)		弁理士 佐野 良太

最終頁に続く

(54) 【発明の名称】 各プロセッサに対して割り込み仮想化を支援するためのゲスト割り込み制御器

(57) 【特許請求の範囲】

【請求項 1】

システム内の複数のゲストの第1のゲストに割り当てられる周辺デバイスをソースとする割り込みをやりとりする割り込みメッセージを、入力/出力メモリ管理ユニット(IOMMU)にて受信することと、

前記IOMMUが、メモリ内の1つ以上の第1のデータ構造にアクセスして前記割り込みを前記第1のゲストにマッピングすることであって、前記1つ以上の第1のデータ構造は他のメモリ位置へのポインタを含むことと、

ゲスト割り込みマネージャが、前記ポインタに応じて、1つ以上の第2のデータ構造をメモリに配置することであって、前記1つ以上の第2のデータ構造は前記第1のゲストに対する割り込み制御器状態を含むことと、

メモリ内の前記1つ以上の第2のデータ構造内に前記割り込みを記録して、前記第1のゲストが実行中であるときに前記割り込みが前記第1のゲストへ受け渡されることを可能にすることとを備えた方法。

【請求項 2】

前記第1のゲストに割り当てられ且つ前記第1のゲスト内での前記割り込みに対する宛先である第1の割り込み制御器へ前記割り込みを伝えることを更に備えた請求項1に記載の方法。

【請求項 3】

割り込み制御器状態は割り込み制御器内の割り込み要求レジスタの状態を含み、前記割

り込みは割り込みベクトルを含み、前記割り込み要求レジスタは前記ベクトルに関連するビット位置を含み、前記割り込みを記録することは前記データ構造内の前記ビット位置におけるビットをセットすることを備えている請求項 1 に記載の方法。

【請求項 4】

前記ビットをセットすることはアトミックである請求項 3 に記載の方法。

【請求項 5】

前記ビットをセットすることはセットされたビットについて前記ビット位置内へアトミックに論理和を取ることを備えている請求項 4 に記載の方法。

【請求項 6】

システム上で実行可能なゲストを対象とする割り込みに対応する割り込みメッセージを受信するように構成されるゲスト割り込みマネージャであって、

前記ゲスト割り込みマネージャは、前記割り込みメッセージが受信されているときに前記ゲストが前記システム内でアクティブでないとしても前記ゲストへの前記割り込みの受け渡しを確実にするために前記割り込みをメモリシステム内のデータ構造内に記録するように構成されており、

前記ゲスト割り込みマネージャは、入力/出力メモリ管理ユニット（IOMMU）からポインタを受信して、前記データ構造を前記メモリシステムに配置するように構成されており、前記ポインタは第 2 のデータ構造から前記 IOMMU によって提供され、前記 IOMMU は、前記第 2 のデータ構造にアクセスして前記割り込みメッセージを前記ゲストにマッピングするように構成されている、ゲスト割り込みマネージャ。

【請求項 7】

前記データ構造は所与のゲストに対して割り込み制御器の状態の少なくとも一部分を記憶するように構成され、前記状態は割り込み要求レジスタの状態を含み、前記ゲスト割り込みマネージャは前記割り込み要求レジスタの前記状態を更新して前記割り込みを記録するように構成される請求項 6 に記載のゲスト割り込みマネージャ。

【請求項 8】

前記割り込み要求レジスタは前記割り込み制御器によってサポートされる各割り込みベクトルに対応するビットを含み、前記割り込みメッセージは前記割り込みの前記割り込みベクトルを含み、前記ゲスト割り込みマネージャは前記割り込みメッセージからの前記割り込みベクトルに対応する前記ビットを更新するように構成される請求項 7 に記載のゲスト割り込みマネージャ。

【請求項 9】

メモリシステムと、

請求項 6 ～ 8 のいずれかに記載のゲスト割り込みマネージャとを備えたシステム。

【請求項 10】

前記割り込みを開始するように構成される周辺デバイスと、

前記 IOMMU とを更に備えた請求項 9 に記載のシステムであって、

前記 IOMMU は前記割り込みを前記メモリ内の 1 つ以上の前記第 2 のデータ構造内のデータに応答する前記ゲストと関連付けるように構成され、前記 IOMMU は前記ゲストと関連している前記割り込みに応答するゲストを識別するゲスト識別子を含む前記割り込みメッセージを送信するように構成されるシステム。

【請求項 11】

前記ゲストに割り当て可能な割り込み制御器を更に備えた請求項 10 に記載のシステムであって、前記割り込み制御器は前記ゲストの前記ゲスト識別子を記憶するように構成され、前記割り込み制御器は受信した割り込みメッセージからの前記ゲスト識別子を前記割り込み制御器内のゲスト識別子と比較するように構成され、前記割り込み制御器は前記ゲスト識別子の比較における整合に応答する前記割り込みを受け入れるように構成されるシステム。

【請求項 12】

前記ゲスト割り込みマネージャは、前記 IOMMU からの前記割り込みメッセージに応

10

20

30

40

50

答し且つ前記データ構造内に記録されている前記割り込みに応答して前記割り込みメッセージを前記割り込み制御器へ送信するように構成される請求項 1 1 に記載のシステム。

【請求項 1 3】

ホストによって制御され且つプロセッサ上で実行可能なゲストを対象とするゲスト割り込みとやりとりするゲスト割り込みメッセージを前記プロセッサに結合されるゲスト割り込み制御器において受信することと、

前記ゲスト割り込み制御器内に記憶されるゲスト識別子を前記ゲスト割り込みメッセージにおける受信したゲスト識別子と整合させることと、

前記ゲスト割り込み制御器内の少なくとも 1 つの宛先識別子をゲストメッセージにおける受信した宛先識別子と整合させることと、

前記ゲスト識別子を整合させること及び前記少なくとも 1 つの宛先識別子を整合させることに応答する前記ゲスト割り込み制御器において前記割り込みを受け入れることとを備えた方法。

【請求項 1 4】

プロセッサ上で実行可能なホストを対象とする第 1 の割り込みとやりとりする第 1 の割り込みメッセージを前記プロセッサに結合される第 1 の割り込み制御器において受信することと、

前記ホストによって制御され且つ前記プロセッサ上で実行可能なゲストを対象とする第 2 の割り込みとやりとりする第 2 の割り込みメッセージを前記プロセッサに結合される請求項 1 3 に記載のゲスト割り込み制御器である第 2 の割り込み制御器において受信することと、

前記第 1 の割り込みメッセージを受信することに応答して前記第 1 の割り込み制御器が前記第 1 の割り込みを前記プロセッサへ受け渡すことと、

前記第 2 の割り込みメッセージを受信することに応答して前記第 2 の割り込み制御器が前記第 2 の割り込みを前記プロセッサへ受け渡すこととを更に備えた請求項 1 3 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

この発明はプロセッサ及び仮想化に関し、更に特定のには仮想マシンゲストに割り込みを受け渡すこと(delivering)に関する。

【背景技術】

【0002】

仮想化は、種々の異なる目的のためにコンピュータシステムにおいて用いられてきた。例えば、仮想化は、「コンテナ(container)」内で特権的ソフトウェア(privileged software)を実行するために用いることができ、コンテナは、仮想マシンを制御する仮想マシンマネージャ(VMM)によって最初に許可されることなしに特権的ソフトウェアが少なくとも何らかの物理マシン状態に直接的にアクセスし且つ/又は少なくとも何らかの物理マシン状態を変化させることを防止するためのものである。そのようなコンテナは、「バグのある(buggy)」又は悪意のあるソフトウェアが物理マシン上で問題を生じさせることを防止し得る。また、仮想化は、2 つ以上の特権的プログラムを同一の物理マシン上で同時に実行することを可能にするために用いられ得る。物理マシンへのアクセスは制御されるので、これらの特権的プログラムが互いに干渉しないようにすることができる。特権的プログラムは、オペレーティングシステムを含むであろうし、そして更に当該ソフトウェアを実行中のハードウェアの完全な制御を有することを期待する他のソフトウェアを含み得る。別の例においては、仮想化は、特権的プログラムによって予定されているハードウェアとは異なるハードウェア上でその特権的プログラムを実行するために用いられ得る。

【0003】

一般に、プロセッサ又はコンピュータシステムの仮想化は、1 つ以上の特権的プログラムに仮想マシン(前述したコンテナ)へのアクセスをもたらすことを含んでいるであろう

10

20

30

40

50

し、特権的プログラムは、仮想マシンを介して完全な制御を有するが、物理マシンの制御はVMMによって保たれる。仮想マシンは、単一のプロセッサ（又は複数のプロセッサ）、メモリ、及び特権的プログラムがそれが実行中のマシン内で見出すことを予定している種々の周辺デバイスを含み得る。仮想マシン要素は、VMMが少なくとも一時的に仮想マシンに割り当てているハードウェアによって実装され得るし、且つ/又はソフトウェア内でエミュレートされ得る。各特権的プログラム（及び幾つかの場合には関連するソフトウェア、例えばオペレーティングシステム上で実行するアプリケーション）は、ここではゲスト(guest)と称される。仮想化は、VMM及びその仮想マシンが実行される物理マシン内で何らかのハードウェア仮想化サポートなしにソフトウェア（例えば上述のVMM）において実装され得る。しかし、何らかのハードウェアサポートが提供されている場合には、仮想化は簡素化されるであろうし且つ/又はより高い性能を達成するであろう。

10

【発明の概要】

【発明が解決しようとする課題】

【0004】

仮想化に伴い生じ得る問題は、割り込み受け渡しの待ち時間である。上述したように、仮想マシンによる使用のために周辺デバイスが割り当てられ得る（仮想マシン内で仮想周辺デバイスとして作用する）。そのような周辺デバイスは、仮想マシン内でソフトウェアによって処理されることになる割り込みを生成することができる。非仮想化環境においては、割り込み処理待ち時間は比較的短い可能性がある。仮想化環境においては、割り込みは、一般的には、VMMによってインターセプトされ(intercepted)、VMMによって処理され、そして何らかのソート(sort)のソフトウェアメカニズムを用いてVMMによって目標とされる仮想マシンへ受け渡され得る。その一方で、割り込み処理待ち時間は、著しく長い可能性がある（例えば100倍長い）。

20

【課題を解決するための手段】

【0005】

1つの実施形態においては、システムは、プロセッサと、プロセッサに結合される第1の割り込み制御器と、プロセッサに結合される第2の割り込み制御器とを備えている。第1の割り込み制御器は、システム内のホストを対象とする第1の割り込みとやりとりする第1の割り込みメッセージを受信することに応答する割り込みに対する信号をプロセッサへ送るように構成される。第2の割り込み制御器は、ホストによって制御され且つプロセッサ上で実行可能なゲストを対象とする第2の割り込みとやりとりする第2の割り込みメッセージを受信することに応答する割り込みに対する信号をプロセッサへ送るように構成される。

30

【図面の簡単な説明】

【0006】

以下の詳細な説明は添付の図面を参照にしこれらをここで簡単に説明する。

【0007】

【図1】図1は仮想化を実装しているコンピュータシステムの1つの実施形態のブロック図である。

【0008】

【図2】図2は図1に示されるホストハードウェアの1つの実施形態のブロック図である。

40

【0009】

【図3】図3はゲストに受け渡されつつある割り込みの1つの実施形態を示すブロック図である。

【0010】

【図4】図4はゲスト進歩的プログラム可能割り込み制御器(APIC)の1つの実施形態を示すブロック図である。

【0011】

【図5】図5はゲストAPIC状態データ構造におけるゲストAPIC状態エントリの1

50

つの実施形態を示すブロック図である。

【0012】

【図6】図6はゲストを目標としている割り込みを受信することに対応する図2に示されるゲスト割り込みマネージャの1つの実施形態の動作を示すフローチャートである。

【0013】

【図7】図7は割り込みメッセージを受信することに対応するゲストAPICの1つの実施形態の動作を示すフローチャートである。

【0014】

【図8】図8はゲストAPIC状態を1つのゲストから他へ変更する仮想マシンモニタの1つの実施形態の動作を示すフローチャートである。

10

【0015】

【図9】図9はゲストAPIC状態エントリ内で割り込み状態を配列する1つの実施形態を示すブロック図である。

【0016】

【図10】図10は割り込みに対してゲストAPIC状態エントリを配置する1つの実施形態を示すブロック図である。

【0017】

【図11】図11は割り込みに対してゲストAPIC状態エントリを配置する別の実施形態を示すブロック図である。

【0018】

20

【図12】図12は割り込みに対してゲストAPIC状態エントリを配置する更に別の実施形態を示すブロック図である。

【0019】

【図13】図13は割り込みに対してゲストAPIC状態エントリを配置するまた更に別の実施形態を示すブロック図である。

【0020】

【図14】図14は図1に示されるホストハードウェアの別の実施形態のブロック図である。

【0021】

【図15】図15はVMMの1つの実施形態を記憶しているコンピュータアクセス可能記憶媒体の1つの実施形態のブロック図である。

30

【発明を実施するための形態】

【0022】

本発明は種々の修正及び代替的な形態を許容する一方で、その特定の実施形態が図面における例によって示されており、そして詳細にここに説明されることになる。しかし、図面及びそれに対する詳細な説明は、開示された特定の形態に本発明を限定することを意図するものではなく、それとは対照的に、添付の特許請求の範囲によって画定されるような本発明の精神及び範囲内にある全ての修正、均等なもの、及び代替案に及ぶことが意図されていることが理解されるべきである。ここで用いられる表題は組織化の目的だけのためのものであり、明細書の範囲を限定することを意味しない。この出願を通して用いられる「かもしれない、あり得る、し得る、であろう、ことがある、可能性がある、してよい、あってよい、することができる(may)」の語は、必須の意味(即ち「なければならない(must)」を意味する)ではなく、許容の意味(即ち「可能性がある」を意味する)である。同様に、「含む(include, includes)」及び「含んでいる(including)」の語は含んでいることを意味し、限定することを意味しない。

40

【0023】

種々のユニット、回路、又は他のコンポーネントが単一又は複数のタスクを実行「するように構成される(configured to)」として説明されることがある。そのような文脈において、「するように構成される」は、動作の間に単一又は複数の当該タスクを実行する「回路を有する(having circuitry that)」ことを一般的に意味する構造の広い記載である

50

。従って、ユニット／回路／コンポーネントは、ユニット／回路／コンポーネントがその時点でオンでないとしても、当該タスクを実行するように構成されていることができる。一般に、「するように構成される」に対応する構造を形成する回路は、当該動作を実装するハードウェア回路を含み得る。同様に、種々のユニット／回路／コンポーネントは、説明の便宜上、単一又は複数のタスクを実行するものとして説明されることがある。そのような説明は、「するように構成される」の句を包含するものとして解釈されるべきである。１つ以上のタスクを実行するように構成されるユニット／回路／コンポーネントを記載することで、当該ユニット／回路／コンポーネントに対して 35 U . S . C . § 112、第 6 段落解釈を行使しないよう明示的に意図するものである。

【 0 0 2 4 】

実施形態においては、コンピュータシステムは、少なくとも 1 つのホスト割り込み制御器と少なくとも 1 つのゲスト割り込み制御器を含む。ホスト割り込み制御器は、ホスト（例えば仮想化された環境における仮想マシンマネージャ又は VMM）によってサービスされることになる割り込みを管理してよい。そのような割り込みの例としては、システム上で実行中のゲストに割り当てられていないコンピュータシステム内のデバイスからの割り込み、VMM がゲストにさすことを望まないシステムレベル割り込み、等を挙げることができる。ゲスト割り込み制御器は、ゲストによってサービスされることになる割り込みを管理してよい。そのような割り込みの例としては、ゲストの仮想マシンに対してデバイスの機能性を提供するようにゲストに割り当てられたそのデバイスによって発行される割り込みを挙げることができる。

【 0 0 2 5 】

システム内のハードウェアは、ホスト割り込み制御器又はゲスト割り込み制御器のいずれかへ割り込みを送信してよい。代替的には、割り込みは、両方の割り込み制御器へ送信されてよく、そしてこれらの割り込み制御器の一方は、それがホスト割り込みであるかゲスト割り込みであるかに基いて割り込みを受け入れてよい。各割り込み制御器は、プロセッサに結合されていてよく、また割り込みを受け渡すためにそのプロセッサと通信してよい。幾つかの実施形態においては、割り込みをゲストへ受け渡すためにゲスト割り込み制御器を設けることは、ホスト割り込み待ち時間と概ね等しい待ち時間を伴うゲスト割り込み受け渡しをもたらし得る。

【 0 0 2 6 】

実施形態においては、ゲスト割り込み制御器は、ホスト割り込み制御器の複製であってよい。即ち、ゲスト割り込み制御器は、ホスト割り込み制御器と同一のハードウェアを含んでいてよい。別の実施形態においては、ゲスト割り込み制御器は、ホスト割り込み制御器の一部のみを複製していてよい。ゲスト割り込み制御器によって複製されていない部分は、VMM インターセプト及び VMM によるエミュレーション(emulation)を介して仮想化されてよい。

[仮想化概説]

【 0 0 2 7 】

図 1 は仮想化を実装しているコンピュータシステム 5 の 1 つの実施形態のブロック図を示している。図 1 の実施形態においては、多重ゲスト 10 A ~ 10 N が示されている。ゲスト 10 A は、ゲストオペレーティングシステム (OS) 12 と、ゲスト OS 12 上で実行される 1 つ以上のアプリケーション 14 A ~ 14 N とを含む。ゲスト 10 N は特権的コード 16 を含む。ゲスト 10 A ~ 10 N は、仮想マシンマネージャ (VMM) 18 によって管理される。VMM 18 及びゲスト 10 A ~ 10 N はホストハードウェア 20 上で実行され、ホストハードウェア 20 は、コンピュータシステム 5 内に含まれる物理的ハードウェアを備えていてよい。1 つの実施形態においては、VMM 18 は一連の仮想マシン制御ブロック (VMCB) 22 を維持していてよい。各ゲスト 10 A ~ 10 N に対して 1 つの VMCB 22 があってよい。図 1 においては例示のために VMCB 22 が VMM 18 の一部として図示されているが、VMCB 22 は、メモリ内に記憶されていてよく且つ / 又はホストハードウェア 20 内のディスクドライブのような不揮発性媒体上に記憶されていて

よい。

【 0 0 2 8 】

ホストハードウェア 20 は、一般的にはコンピュータシステム 5 内に含まれるハードウェアの全てを含む。種々の実施形態において、ホストハードウェア 20 は、1 つ以上のプロセッサ、周辺デバイス、及びこれらのコンポーネントを結合するために用いられる他の回路を含んでいてよい。具体的には、ホストハードウェア 20 は、1 つ以上のホスト割り込み制御器、1 つ以上のゲスト割り込み制御器、及び / 又は 1 つ以上のゲスト割り込みマネージャを含んでいてよい。例えば、パーソナルコンピュータ (P C) 型システムは、複数のプロセッサを結合するノースブリッジ (Northbridge) と、メモリと、進歩的なグラフィックポート (advanced graphic port) (A G P) インタフェースを用いるグラフィクスデバイスとを含んでいてよい。加えて、ノースブリッジは、周辺コンポーネントインタフェース (P C I) バスのような周辺バスに結合してよく、周辺バスには種々の周辺コンポーネントが直接的に又は間接的に結合されていてよい。レガシー (legacy) 機能性を提供し且つ / 又はレガシーハードウェアに結合するために、P C I バスに結合されるサウスブリッジ (Southbridge) もまた含まれていてよい。種々の実装において、ゲスト割り込みマネージャは、ノースブリッジ内、サウスブリッジ内、又はインタフェースの 1 つ上のデバイス内に実装されていてよい。ホスト割り込み制御器及びゲスト割り込み制御器は、各プロセッサに対して実装されてよく、又はプロセッサの群によって共有されてよい。他の実施形態においては、種々のハードウェアコンポーネントをリンクするために、他の回路が用いられてよい。例えば、ハイパートランスポート (HyperTransport) (商標) (H T) リンクが複数のノードをリンクするために用いられてよく、ノードの各々は、1 つ以上のプロセッサ、ホストブリッジ、及びメモリ制御器を含んでいてよい。各ノードもノースブリッジを含んでいてよく、そのノースブリッジは、ゲスト割り込みマネージャ並びに / 又はホスト割り込み制御器及びゲスト割り込み制御器を含んでいてよい。代替的には、ホストブリッジがゲスト割り込みマネージャ並びに / 又はホスト割り込み制御器及びゲスト割り込み制御器を含んでいてもよい。ホストブリッジは、H T リンクを介してデイジーチェーン形式 (daisy chain fashion) で周辺デバイスと結合するために用いられてよい。任意の所望の回路 / ホストハードウェア構造が用いられてよい。

【 0 0 2 9 】

V M M 1 8 は、ゲスト 1 0 A ~ 1 0 N の各々に対する仮想化を提供するように構成されてよく、そしてホストハードウェア 20 へのゲスト 1 0 A ~ 1 0 N のアクセスを制御してよい。V M M 1 8 はまた、ホストハードウェア 20 上での実行のためにゲスト 1 0 A ~ 1 0 N をスケジューリングすることに関与してよい。V M M 1 8 は、ホストハードウェア 20 内で仮想化のために提供されるハードウェアサポートを用いるように構成されてよい。例えば、イベントをインターセプトしそして処理のためにゲストから V M M 1 8 へ退出するためのハードウェアを含めて、プロセッサは仮想化のためのハードウェアサポートを提供してよい。ゲスト割り込みマネージャ及び / 又はゲスト割り込み制御器は、仮想化をサポートするために設けられるハードウェアであってもよい。

【 0 0 3 0 】

幾つかの実施形態においては、V M M 1 8 は、ホストハードウェア 20 上で実行される「薄い (thin) 」スタンドアローンソフトウェアプログラムとして実装されてよく、そしてゲスト 1 0 A ~ 1 0 N に対して仮想化を提供する。そのような V M M 実装は、しばしば「ハイパーバイザ (hypervisor) 」と称されることがある。他の実施形態においては、V M M 1 8 はホスト O S 内へ統合されていてよく又はホスト O S 上で実行されてよい。そのような実施形態においては、V M M 1 8 は、ホスト O S 内の任意のドライバ、システム B I O S によって提供されるプラットフォームシステム管理モード (S M M) コード、等を含めてホスト O S に依存するであろう。従って、ホスト O S コンポーネント (及びプラットフォーム S M M コードのような種々の低レベルコンポーネント) は、ホストハードウェア 20 上で直接的に実行され、また V M M 1 8 によっては仮想化されない。V M M 1 8 及びホスト O S (含まれる場合には) は、1 つの実施形態においては、一緒にホストと称される

ことがある。一般に、ホストは、使用に際してホストハードウェア 20 の直接制御内にある任意のコードを含んでいてよい。例えば、ホストは、VMM 18 であってよく、ホスト OS を伴う VMM 18 であってよく、あるいはホスト OS 単独（例えば非仮想化環境において）であってよい。

【0031】

種々の実施形態において、VMM 18 は、完全仮想化(full virtualization)、準仮想化(paravirtualization)、又は両方をサポートしてよい。更に、幾つかの実施形態においては、VMM 18 は、準仮想化されている複数のゲスト及び完全仮想化されている複数のゲストを同時に実行してよい。

【0032】

完全仮想化では、ゲスト 10A ~ 10N は、仮想化が発生中であることに気付かない。各ゲスト 10A ~ 10N は、その仮想マシン内に連続的なゼロベースのメモリを有してよく、そして VMM 18 は、シャドウ(shadow)ページテーブル又は入れ子にされた(nested)ページテーブルを用いてホスト物理アドレス空間へのアクセスを制御してよい。シャドウページテーブルは、ゲスト仮想アドレスからホスト物理アドレスへ再マッピングしてよく（事実上は、ゲスト 10A ~ 10N 内のメモリ管理ソフトウェアによって割り当てられたゲスト「物理アドレス」のホスト物理アドレスへの再マッピング）、一方、入れ子にされたページテーブルは、入力としてのゲスト物理アドレスを受信し、そしてホスト物理アドレスへマッピングしてよい。シャドウページテーブル又は入れ子にされたページテーブルを各ゲスト 10A ~ 10N に対して用いて、VMM 18 は、ゲストがホストハードウェア 20 内の他のゲストの物理メモリへアクセスしないことを確実にし得る。

【0033】

準仮想化では、ゲスト 10A ~ 10N は、少なくとも部分的には VM アウェア(VM-aware)であってよい。そのようなゲスト 10A ~ 10N は、VMM 18 とメモリページについて交渉してよく、従って、ゲスト物理アドレスをホスト物理アドレスへ再マッピングすることは必要ないかもしれない。1つの実施形態では、準仮想化において、ゲスト 10A ~ 10N は、ホストハードウェア 20 内の周辺デバイスと直接的に相互作用することが許可されてよい。任意の所与の時間において、周辺デバイスは、単一又は複数のゲスト 10A ~ 10N によって「所有されて(owned)」よい。1つの実装においては、例えば、周辺デバイスは、その周辺デバイスをその時点で所有している 1つ以上のゲスト 10A ~ 10N と共に保護ドメイン内へマッピングされてよい。周辺デバイスを所有しているゲストのみが、その周辺デバイスと直接的に相互作用することができる。保護ドメイン内のデバイスが別の保護ドメイン内のゲストに割り当てられているページに対して読み出し/書き込みすることを防止するための保護メカニズムがあってもよい。

【0034】

前述したように、VMM 18 は各ゲスト 10A ~ 10N に対して VMCB 22 を維持してよい。VMCB 22 は、一般的には、対応するゲスト 10A ~ 10N に対して VMM 18 によって割り当てられる記憶区域内に記憶されるデータ構造を備えていてよい。1つの実施形態においては、VMCB 22 はメモリのページを備えていてよいが、他の実施形態は、より大きな若しくはより小さなメモリ区域を用いてよく且つ/又は不揮発性ストレージのような他の媒体上での記憶を用いてよい。1つの実施形態においては、VMCB 22 はゲストのプロセッサ状態を含んでいてよく、ゲストのプロセッサ状態は、ゲストの実行をスケジューリングされるときにホストハードウェア 20 におけるプロセッサ内へロードされてよく、またゲストが退出するときに（当該スケジューリングされた時間を完了したこと、又はゲストを退出するためにプロセッサが検出している 1つ以上のインターセプトのいずれかに起因して）、VMCB 22 へ記憶し戻されてよい。幾つかの実施形態においては、VMCB 22 に対応しているゲストへ制御を移す命令（「仮想マシン実行(Virtual Machine Run) (VMRUN)」命令）を介してプロセッサ状態の一部分のみがロードされ、そして他の望まれる状態は、VMRUN 命令を実行するのに先立ち VMM 18 によってロードされてよい。同様に、そのような実施形態においては、プロセッサ状態の一部分

10

20

30

40

50

のみが、ゲスト退出に際してプロセッサによってV M C B 2 2へ記憶されてよく、またV M M 1 8は、必要に応じて任意の追加的な状態を記憶することに関与してよい。他の実施形態においては、V M C B 2 2は、プロセッサ状態が記憶されている別のメモリ区域へのポインタを含んでいてよい。また、1つの実施形態においては、2つ以上の退出メカニズムが規定されてよい。1つの実施形態においては、記憶されている状態の量及びロードされる状態のロケーションは、どの退出メカニズムが選択されているかに依存するであろう。

【 0 0 3 5 】

1つの実施形態においては、V M M 1 8はまた、V M M 1 8に対応しているプロセッサ状態を記憶するように割り当てられるメモリの区域を有していてよい。V M R U Nが実行される場合、V M M 1 8に対応しているプロセッサ状態は、その区域内に保存されてよい。ゲストがV M M 1 8へ退出すると、その区域からのプロセッサ状態がその区域からリロードされてV M M 1 8が実行を継続することを可能にし得る。1つの実装においては、例えばプロセッサは、V M M 1 8保存区域のアドレスを記憶するためにレジスタ（例えば特定モデル向けレジスタ又はM S R）を実装していてよい。

【 0 0 3 6 】

また、V M C B 2 2は、ゲストに対して有効にされているインターセプトイベントを識別するインターセプト構成と、有効にされたインターセプトイベントが検出される場合にゲストを退出するためのメカニズムとを含んでいてよい。1つの実施形態においては、インターセプト構成は、一連のインターセプト表示と、プロセッサがサポートしている各インターセプトイベントに対する1つの表示とを含んでいてよい。インターセプト表示は、対応するイベントをプロセッサがインターセプトすべきか否か（又は別の見方をすればインターセプトが有効にされているか否か）を表示してよい。ここで用いられているように、イベントがゲスト内で「インターセプトされる(intercepted)」のは、そのイベントがそのゲスト内で発生すべきでプロセッサがそのイベントの処理のためにそのゲストを退出する場合である。1つの実施形態においては、インターセプト構成は、2つの退出メカニズムのどちらが用いられるかを表示する第2の一連の表示を含んでいてよい。他の実施形態は2つを超える退出メカニズムを規定してよい。別の実施形態においては、インターセプト構成は、第1の退出メカニズムがイベントに対して用いられるべきか否かを表示する、インターセプトイベント毎の一連のインターセプト表示と、第2の退出メカニズムがイベントに対して用いられるべきか否かを表示する、インターセプトイベント毎の第2の一連のインターセプト表示とを備えていてよい。

【 0 0 3 7 】

一般に、退出メカニズムは、ゲスト実行から退出する（通常は再起動可能なやり方で）と共に他のコードの実行を開始するためにプロセッサによって実行される動作を規定してよい。1つの実施形態においては、1つの退出メカニズムは、少量のプロセッサ状態を保存することと、ミニバイザ(Minivisor)に対する状態をロードすることとを含んでいてよい。ミニバイザは、ゲスト物理アドレス空間内で実行されてよく、そして比較的単純なインターセプト処理を実行してよい。大量のプロセッサ状態を保存すると共にV M Mのプロセッサ状態をロードする別の退出メカニズムがV M Mへ退出してよい。従って、インターセプトイベントは、イベントに応じて異なる命令コードによって処理されてよい。また、比較的単純なインターセプト処理は、実行するのにより少ない時間を費やすであろう「より軽量な(lighter weight)」退出メカニズムを通して処理されてよく、これにより幾つかの実施形態においては性能が改善され得る。より複雑な処理は、「より重たい(heavier weight)」メカニズムが退出のために用いられた後で、V M M内で実行されてよい。従って、この実施形態においては、V M M 1 8は、ゲスト1 0 A ~ 1 0 Nが内部的に処理することをV M M 1 8が望まないイベントをインターセプトするようにプロセッサを構成してよく、またV M M 1 8は、退出メカニズムを使用すべきプロセッサを構成してもよい。イベントとしては、命令（即ち、命令を実行する代わりに命令をインターセプトする）、割り込み、例外、及び/又はゲスト実行の間に発生し得る他の所望のイベントを挙げることが

10

20

30

40

50

できる。

【 0 0 3 8 】

1つの実施形態においては、V M C B 2 2は、V M C B 2 2をロードするに際して特定の動作をプロセッサに実行させることができる他の制御ビットを更に含んでいてよい。例えば、制御ビットは、プロセッサ内にT L Bをフラッシュする(flush)するための命令を含んでいてよい。他の制御ビットは、ゲストに対する実行環境(例えば割り込み処理モード、ゲストに対するアドレス空間識別子、等)を指定してよい。更に他の制御ビットは、何故ゲストが退出したか等を説明する退出コードとやりとりするために用いられてよい。

【 0 0 3 9 】

一般的に、「ゲスト」は、コンピュータシステム5内での実行のために仮想化されることになる任意の1つ以上のソフトウェアプログラムを備えていてよい。ゲストは、特権的なモードにおいて実行される少なくとも何らかのコードを含んでいてよく、従ってそれが実行中のコンピュータシステムにわたって完全な制御を有することを期待する。前述したように、ゲスト1 0 Aは、ゲストがその内部にゲストO S 1 2を含む例である。ゲストO S 1 2は任意のO Sであってよく、例えば、マイクロソフトコープ(レッドモンド、ワシントン)から入手可能なウィンドウズO Sの任意のもの、I B Mコーポレーション(アーモック、ニューヨーク)から入手可能なリナックス、A I Xのような任意のユニックス型オペレーティングシステム、サンマイクロシステムズインク(サンタクララ、カリフォルニア)から入手可能なソラリス、ヒューレットパッカードカンパニー(ボロアルト、カリフォルニア)から入手可能なH P - U X、等であってよい。ゲスト1 0 Nは、非O S特権的

【 0 0 4 0 】

尚、ここで用いられる場合における1 0 Nのような参照番号内の「N」の文字は、当該参照番号を有する任意の数の要素を総称的に表示することを意味している(例えば任意の数のゲスト1 0 A ~ 1 0 Nは1つのゲストを含む)。また、文字「N」を用いる異なる参照番号(例えば1 0 Nと1 4 N)は、特に断りのない限り、同様の数の異なる要素が設けられていることを表示するようには意図されていない(例えば、ゲスト1 0 A ~ 1 0 Nの数は、アプリケーション1 4 A ~ 1 4 Nの数とは異なっていてよい)。

[ホストハードウェア]

【 0 0 4 1 】

次に図2を参照すると、ホストハードウェア2 0の1つの実施形態を表すブロック図が示されている。図示される実施形態においては、ホストハードウェア2 0は、複数のプロセッサ3 0 A ~ 3 0 Bと、それぞれのホスト進歩的プログラム可能割り込み制御器(host Advanced Programmable Interrupt Controllers)(h A P I C)3 2 A ~ 3 2 Bと、それぞれのゲストA P I C(g A P I C)3 4 A ~ 3 4 Bと、随意的な付加的g A P I C3 4 C ~ 3 4 Dと、ブリッジ3 6(ゲスト割り込みマネージャ3 8、入力/出力(I/O)メモリ管理ユニット(I O M M U)4 0、及びメモリ制御器4 2を含む)と、複数のインタフェース回路(I F)4 4 A ~ 4 4 Cと、メモリインタフェース回路(M I F)4 6と、I O A P I C5 0を含んでいてよい随意的なブリッジ4 8と、周辺機器5 2 A ~ 5 2 B(そのうちの幾つかはI O A P I C5 4のようなI O A P I Cを含んでいてよい)と、メモリ5 6とを含む。プロセッサ3 0 A ~ 3 0 Bは、図2に示されるように、ブリッジ3 6と結合され、またそれぞれのh A P I C3 2 A ~ 3 2 B及びg A P I C3 4 A ~ 3 4 Dと結合されている。h A P I C3 2 A ~ 3 2 B及びg A P I C3 4 A ~ 3 4 Dはブリッジ3 6と結合され、ブリッジ3 6は、インタフェース回路4 4 A ~ 4 4 C及びメモリインタフェース回路4 6と結合されている。メモリインタフェース回路4 6はメモリ5 6と結合され、インタフェース回路4 4 Aはブリッジ4 8と結合され、ブリッジ4 8は周辺機器5 2 A ~ 5 2 Bと結合されている。

【 0 0 4 2 】

図示される実施形態においては、各プロセッサ3 0 A ~ 3 0 Bは、関連するh A P I C3 2 A ~ 3 2 B及び少なくとも1つの関連するg A P I C3 4 A ~ 3 4 Dを有している。

この実施形態においては、割り込みは、インテルコーポレーション（サンタクララ、カリフォルニア）（Intel Corporation (Santa Clara, CA)）によって記述されている A P I C 規格に従ってホストハードウェア 20 内で通信されてよい。その規格においては、各プロセッサは、プロセッサそれ自身、他のプロセッサ、内部 A P I C 割り込みソース、及び周辺機器に関連する I O A P I C からの割り込みを受信する関連するローカル A P I C を有している。ローカル A P I C は、係属中の複数の割り込みを優先付けし、割り込みがプロセッサ上で進行中である別の割り込みよりも高い優先度である場合且つ / 又は割り込みがプロセッサのその時点のタスクよりも優先度が高い場合には、その割り込みをプロセッサへ送信する。

【 0 0 4 3 】

10

図 2 の実施形態においては、h A P I C 3 2 A ~ 3 2 B は、プロセッサのホスト割り込み（即ち、ホストによって処理されるべき割り込み）に対するローカル A P I C であってよく、また g A P I C 3 4 A ~ 3 4 D は、プロセッサのゲスト割り込み（即ち、それぞれのプロセッサ 3 0 A ~ 3 0 B 上でアクティブなゲストによって処理されるべき割り込み）に対するローカル A P I C であってよい。ゲストは、そのゲストがその時点で当該プロセッサ上で実行中である場合（例えば、V M R U N 命令がそのゲストに対するそのプロセッサ上で実行完了しており、且つゲスト退出が発生しなかった場合）、又はそのゲストは既に退出しており且つ V M M 1 8 が実行中であるが、そのゲストはそのプロセッサ上で再度実行されることが期待されている場合に、そのプロセッサ上でアクティブであり得る。

【 0 0 4 4 】

20

V M M 1 8 がプロセッサ 3 0 A ~ 3 0 B 上のゲストをスケジューリングする場合、V M M 1 8 は、当該プロセッサ 3 0 A ~ 3 0 B の g A P I C 3 4 A ~ 3 4 D をそのゲストに対応する g A P I C 状態と共にロードしてよい。具体的には、所与のゲストは多重仮想 C P U (v C P U) を有していてよい。V M M 1 8 は、プロセッサ 3 0 A ~ 3 0 B 上で実行するゲストの v C P U をスケジューリングしてよく、そしてそのゲストの仮想マシン内の当該 v C P U に対する割り込み状態と共に g A P I C 3 4 A ~ 3 4 D をロードしてよい。また、ゲストがアクティブである間に信号を送られるゲスト（及び v C P U ）を対象としている任意の割り込みが、g A P I C 3 4 A ~ 3 4 D によって捕獲されてよい。g A P I C 3 4 A ~ 3 4 D は、前述した A P I C 規格に従ってゲストを割り込ませてよい。

【 0 0 4 5 】

30

所与のプロセッサ 3 0 A ~ 3 0 B に対する h A P I C 3 2 A ~ 3 2 B 及び g A P I C 3 4 A ~ 3 4 D は、そのプロセッサに対する任意のインタフェースを有していてよい。例えば、複数のローカル A P I C とそれらのそれぞれのプロセッサとの間で任意のインタフェースが用いられてよい。各 A P I C は、割り込みがサービスのために受け渡されている最中のプロセッサへ独立に信号を送るように構成されてよい。プロセッサがゲストを実行中であり且つゲスト割り込みが信号を送られている場合、プロセッサは、ゲストコードに割り込み且つゲストの仮想マシン内で正しい割り込みハンドラの実行を開始するように構成されてよい。従って、実施形態においては、ゲスト割り込みは、ホスト内の割り込みの受け渡しと同様の待ち時間を伴って受け渡され得る。プロセッサがゲストを実行中であり且つ h A P I C が割り込みに信号を送る場合には、プロセッサは、ゲストから V M M 1 8 へ退出してホスト割り込みを処理するように構成されてよい。プロセッサがゲストを実行中でない場合には、g A P I C によって信号を送られている割り込みは、ゲストが再度実行されるまでプロセッサによってマスキングされてよい。プロセッサがゲストを実行中でなく且つ h A P I C が割り込み内へ信号を送る場合には、プロセッサは、ホスト実行に割り込み且つホスト割り込みハンドラへ分岐するように構成されてよい。

40

【 0 0 4 6 】

1 つの実施形態においては、2 つ以上の g A P I C 3 4 A ~ 3 4 D がプロセッサ 3 0 A ~ 3 0 B 毎に含まれていてよい。各 g A P I C 3 4 A ~ 3 4 D は、異なるゲスト / v C P U に対応する A P I C 状態を記憶してよい。そのような実施形態においては、各 g A P I C 3 4 A ~ 3 4 D は、ゲスト割り込みの信号をプロセッサへ送るときに各 g A P I C 3 4

50

A ~ 3 4 D がどのゲストに対応しているかを識別するように構成されてよい（又はどのゲストがその時点で各 g A P I C 3 4 A ~ 3 4 D に割り当てられているかを識別する内部レジスタをプロセッサ 3 0 A ~ 3 0 B が有してよい）。V M M 1 8 が実行中でない場合にゲスト割り込みをマスキングすると同様に、異なるゲストがその時点で実行中である場合には、プロセッサはゲスト割り込みをマスキングしてよい。代替的には、各 g A P I C 3 4 A ~ 3 4 D は、対応するゲストがスケジューリングされる場合に V M M 1 8 によってアクティブにセットされ得るアクティブ表示を含んでいてよく、また g A P I C 3 4 A ~ 3 4 D は、そのゲスト割り込みへ、対応するゲストがアクティブの場合にのみ信号を送るように構成されてよい。プロセッサ 3 0 A ~ 3 0 B 毎に 2 つ以上の g A P I C 3 4 A ~ 3 4 D を含むことにより、多重ゲストがプロセッサ上で時間内に実行されるようにスケジューリングされている場合に、g A P I C 状態移動 (state movement) の量を低減することができる。例えば、プロセッサ 3 0 A ~ 3 0 B 毎に N 個の g A P I C 3 4 A ~ 3 4 D がある場合（N は 0 より大きい整数）、N 個までの異なるゲストは、任意のゲストに対して g A P I C 状態が保存されることが必要になるであろう前に、実行のためにスケジューリングされてよい。プロセッサ 3 0 A ~ 3 0 B 毎に 2 つ以上の g A P I C 3 4 A ~ 3 4 D を実装する幾つかの実施形態においては、g A P I C 3 4 A ~ 3 4 D は、割り込みメッセージが適切に受け入れられ且つ / 又は記録されることを確実にする追加的な状態を含んでいてよい。例えば、g A P I C 3 4 A ~ 3 4 D は、対応する仮想マシンが対応するプロセッサ 3 0 A ~ 3 0 B 上でその時点で実行中であるか否かを識別する「現在実行中 (currently running)」表示を含んでいてよい（V M M 実行に対するサスペンション (suspension) にあるのとは対照的に、又は別の仮想マシンが実行中である間とは対照的に）。仮想マシンが実行中であることを現在実行中表示が示している場合には、g A P I C は割り込みメッセージを受け入れてよい。仮想マシンが実行中でないことを現在実行中表示が示している場合には、g A P I C は割り込みが受け入れられない旨の信号を送ってよい。代替的には、g A P I C は、割り込みが受け入れられない旨の信号を g A P I C が送ることになるか否かを示す追加的な表示を含んでいてよい。そのような実施形態においては、現在実行中表示が現在実行中でないと示し且つ割り込みが受け入れられない旨の信号を g A P I C が送ることになることを非容認表示が示している場合には、割り込みが受け入れられなかった旨の信号を g A P I C が送ってよい。このような機能性は、実行中でないゲストに対して割り込みが受信されていることを検出するために用いられてよく、割り込みの対象にされているゲストをスケジューリングするために用いられてよい。

【 0 0 4 7 】

g A P I C 3 4 A ~ 3 4 D は、h A P I C 3 2 A ~ 3 2 B 内に含まれるハードウェアの少なくとも一部分を含んでいてよく、また、そのハードウェアの全てを含んでいてよい（例えば h A P I C 3 2 A ~ 3 2 B の複製であってよい）。g A P I C 3 4 A ~ 3 4 D は、g A P I C 3 4 A ~ 3 4 D がどのゲストに割り当てられることになるかを識別するために、A P I C 状態に加えてゲスト識別子 (I D) と共にプログラム可能であってよい。ゲストが多重 v C P U を含む場合には、物理 A P I C _ I D 及び論理 A P I C _ I D がゲスト内の v C P U を識別し得る。1 つの実施形態においては、ゲスト I D は、周辺デバイスに対して I O M M U 4 0 によってサポートされるドメイン I D と同じであってよい。代替的には、ゲスト I D は別個に管理される資源であってよい。いずれの場合においても、V M M 1 8 は、ゲスト I D をゲストに割り当ててよく、そして各ゲストに対して g A P I C 3 4 A ~ 3 4 D が適切にプログラムされることを確実にし得る。v C P U 及び / 又は g A P I C 及び / 又は当該対は、ここではより簡潔にゲスト内の割り込みの宛先と称されることがある。宛先は最終的には、割り込みをサービスすることになる v C P U であってよいが、対応する g A P I C もまた、対応するプロセッサに関連しており且つ割り込みを記録することから、宛先と考えられてよい。

【 0 0 4 8 】

g A P I C 3 4 A ~ 3 4 D 及び h A P I C 3 2 A ~ 3 2 B は、割り込みを受信するためにブリッジ 3 6 と結合されている。g A P I C 3 4 A ~ 3 4 D 及び h A P I C 3 2 A ~ 3

10

20

30

40

50

2 Bへ割り込みを運ぶために、任意のインタフェースが用いられ得る。例えば、A P I C 割り込み運搬のために実装される任意のインタフェースが用いられ得る。1つの実施形態においては、割り込みメッセージを運ぶために、プロセッサ30 A ~ 30 Bへの/プロセッサ30 A ~ 30 Bからの他の動作とやりとりするために用いられるのと同じ通信メカニズム（例えば、プロセッサ30 A ~ 30 Bによって開始されるメモリ読み出し/書き込み動作、キャッシュコヒーレンシメンテナンステナンスに対するプローブ、等）が用いられ得る。別の方法としては、g A P I C 34 A ~ 34 Dとh A P I C 32 A ~ 32 Bの結合が、プロセッサ30 A ~ 30 Bのブリッジ36への結合と共有されてよい。代替的には、例えばg A P I C 34 A ~ 34 D及びh A P I C 32 A ~ 32 BがA P I C「3線インタフェース(3 wire interface)」を用いている場合には、プロセッサ30 A ~ 30 Bは、ブリッジ36への別個のパスを有してよい。割り込みメッセージは、送信されている割り込み及び割り込みの宛先を識別する任意のインタフェース上での任意の通信であってよい。例えば、割り込みは関連する割り込みベクトルを有してよく、そして割り込みベクトルは割り込みメッセージの一部であってよい。割り込みメッセージはまた、ゲストID及び宛先ID（例えば論理A P I C __ ID又は物理A P I C __ ID）を含んでよい。

【0049】

h A P I C 32 A ~ 32 BはローカルA P I Cと同様であってよい。例えば、h A P I Cはホスト割り込みのためには用いられないから、h A P I C 32 A ~ 32 Bは、ゲスト識別のための追加的なハードウェアを含んでいなくてよい。代替的には、h A P I C 32 A ~ 32 Bは追加的なハードウェアを含んでよいが、追加的なハードウェアは、h A P I C 32 A ~ 32 Bがホスト割り込みのためのものであることを表示するようにプログラムされてよい。ブリッジ36によってh A P I C 32 A ~ 32 B及びg A P I C 34 A ~ 34 Dへ送信される割り込みメッセージは、ゲスト割り込みをホスト割り込みに対抗するものとして識別してよく、またゲスト割り込みに対するゲストIDを含んでよい（又はゼロ若しくは全バイナリ1のような予約されたゲストIDを用いてホスト割り込みを表示してよい）。h A P I C 32 A ~ 32 Bは、ホスト割り込みとして識別される割り込みを受け入れるように構成されてよく（ホスト割り込みの物理A P I C __ ID又は論理A P I C __ IDが対応するh A P I C __ IDに整合する場合）、またg A P I C 34 A ~ 34 Dは、それらのそれぞれのゲストに対するゲスト割り込みを受け入れるように構成されてよい（ゲストIDが整合する場合、及びゲスト割り込みの物理A P I C __ ID又は論理A P I C __ IDが対応するg A P I C __ IDに整合する場合）。

【0050】

g A P I C 34 A ~ 34 Dはアクティブゲストに対する割り込みを管理してよいが、幾つかのゲストは非アクティブかもしれない（且つ又はゲスト割り込みによって対象とされ得る非アクティブv C P Uを有しているかもしれない）。1つの実施形態においては、ゲスト割り込みマネージャ38は、非アクティブゲストに対するゲスト割り込み状態を維持し且つアクティブゲストに対するg A P I Cが確実にそれらの割り込みを受信するように構成されてよい。

【0051】

特に、1つの実施形態においては、ゲスト割り込みマネージャ38は分散型割り込み受け渡しスキームを採用してよく、そのスキームにおいては、ゲスト割り込みマネージャ38は、ブリッジ36内で受信される各ゲスト割り込みを記録するように構成されてよく、またゲスト割り込みを各g A P I C 34 A ~ 34 Dへ送信するように構成されてよい。g A P I C 34 A ~ 34 Dが割り込みを受け入れる場合には、ゲスト割り込みによって対象とされるゲストはアクティブである。どのg A P I C 34 A ~ 34 Dも割り込みを受け入れない場合には、ゲスト割り込みによって対象とされるゲストは非アクティブである。

【0052】

図示される実施形態においては、ゲスト割り込みマネージャ38は、メモリ56におけるg A P I C状態データ構造58内のシステム5において規定されるゲストに対するg A P I C状態を維持するように構成されてよい。g A P I C状態データ構造58は、システ

10

20

30

40

50

ムにおいて規定される各 g A P I C に対する g A P I C 状態エントリ（例えばシステム内の各ゲスト 1 0 A ~ 1 0 N における各 v C P U に対する 1 つのエントリ）を含んでいてよい。g A P I C は、それがシステム内のアクティブゲスト又は非アクティブゲストのいずれかに関連する場合に、システムにおいて規定されてよい。従って、ゲスト割り込みを受信することに応答して、ゲスト割り込みマネージャ 3 8 は、当該割り込みによって対象とされるゲスト / v C P U に対して g A P I C 状態データ構造 5 8 内の g A P I C 状態を更新するように構成されてよい。ゲスト割り込みマネージャ 3 8 は、1 つの実施形態においては、ゲストがアクティブか否かに関係なく g A P I C 状態を更新するように構成されてよい。1 つの対象より多くを有するマルチキャスト割り込み及びブロードキャスト割り込みに対しては、ゲスト割り込みマネージャ 3 8 は、各割り込み宛先に対して g A P I C 状態データ構造 5 8 内の g A P I C 状態を更新するように構成されてよい。代替的には、ゲスト割り込みマネージャ 3 8 は、これらの多重宛先割り込みに対して V M M 1 8 を頼るように構成されてよい。ゲスト割り込みマネージャ 3 8 は、そのような実施形態においては、V M M 1 8 にアクセス可能なメモリロケーション内に割り込みを記録するように構成されてよく、また V M M 1 8 へ信号を送ってメッセージを処理するように構成されてよい。

【 0 0 5 3 】

幾つかの実施形態においては、ゲスト割り込みマネージャ 3 8 は、ゲスト I D 及び / 又はゲスト割り込みメッセージ内の他の情報に直接的に応答して、g A P I C 状態エントリを g A P I C 状態データ構造 5 8 内に配置するように構成されてよい。他の実施形態においては、g A P I C 状態データ構造 5 8 における柔軟性を提供するために且つ / 又はメモリ空間を節約するために、ゲスト割り込みマネージャ 3 8 は、g A P I C 状態マッピングテーブル 6 0 を用いて g A P I C 状態エントリを g A P I C 状態データ構造 5 8 内に配置するように構成されてよい。g A P I C 状態データ構造 5 8 の種々の実施形態及びマッピングテーブル 6 0（幾つかの実施形態に対して）が図 1 0 ~ 1 3 に示されており、また更に詳細に後で論じられる。従って、ゲスト割り込みに応答して、ゲスト割り込みマネージャ 3 8 は、g A P I C 状態マッピングテーブル 6 0 を用いて g A P I C 状態エントリを配置するように、そして g A P I C 状態エントリを更新して割り込みを記録するように構成されてよい。

【 0 0 5 4 】

1 つの実施形態では、g A P I C 状態データ構造 5 8 は、g A P I C 状態のサブセット (subset) を記憶してよい。サブセットは、ハードウェア 2 0（例えば I O M M U 4 0 と併用されるゲスト割り込みマネージャ 3 8）によって追跡される g A P I C 状態であってよい。より特定的には、サブセットは、対応するゲストが非アクティブである間に変化し得る g A P I C 状態の一部分であってよい。例えば、1 つの実施形態では、周辺機器 5 2 A ~ 5 2 B は、対応するゲストが非アクティブである間に割り込みを信号で伝えてよく、これにより対応する割り込み要求は g A P I C によって捕獲され得る。割り込み要求は g A P I C 状態データ構造 5 8 内で追跡されてよい。他の g A P I C 状態は、どの割り込みがプロセッサによって稼動中 (in-service) であるか、プロセッサのタスク優先度、等を追跡してよい。これらの値は、ゲストがアクティブである場合にのみ変化し得る。実施形態においては、ゲストが非アクティブである場合に変化しないであろう g A P I C 状態は、図 2 に V M M 管理 g A P I C 状態データ構造 6 8 として示される 1 つ以上の他のデータ構造を用いて V M M 1 8 によって追跡されてよい。V M M 1 8 は、V M M 管理 g A P I C 状態データ構造 6 8 とシステム内でアクティブ化しているゲスト及び非アクティブ化しているゲストの一部としての g A P I C 3 4 A ~ 3 4 D との間で、状態を転送してよい。

【 0 0 5 5 】

図示される実施形態においては、g A P I C 状態マッピングテーブル 6 0 及び g A P I C 状態データ構造 5 8 はメモリ 5 6 内に記憶されているとして示されているが、一方又は両方の部分部分は、ゲスト割り込みマネージャ 3 8 及び / 又はブリッジ 3 6 へアクセス可能なキャッシュによってキャッシュされてよい。加えて又は代替的には、1 つ以上の g A P I C 状態エントリのための専用のメモリがブリッジ 3 6 内に実装されてよい。専用のメ

モリは、g A P I C 3 4 A ~ 3 4 D 内へ及び g A P I C 3 4 A ~ 3 4 D から外へ速やかに切り換えられ得る一連の「高速の(fast)」g A P I C 状態を記憶してよい。他の g A P I C 状態は、メモリ 5 6 内でより低速にアクセス可能であってよい。幾つかの実施形態においては、高速の g A P I C 状態切り換えはゲスト割り込みマネージャ 3 8 によって処理されてよい一方で、より低速の g A P I C 状態切り換えは V M M 1 8 によって処理されてよい。

【 0 0 5 6 】

A P I C 割り込みメカニズムにおいては、各プロセッサ(そのローカル A P I C を介して)は、物理 A P I C __ I D 及び論理 A P I C __ I D を有してよい。物理 A P I C __ I D は A P I C __ I D レジスタ内に記憶される。物理 A P I C __ I D は、物理受け渡しモード割り込み(physical delivery mode interrupt)によって表示される物理 A P I C __ I D と 1 対 1 を原則として整合する。論理 A P I C __ I D は、論理宛先レジスタとしてローカル A P I C 内に記憶される。論理 A P I C __ I D はクラスタ I D 及びローカル A P I C __ I D を有しており、ここでローカル A P I C __ I D はワンホットベクトル(one-hot vector)である。論理受け渡しモード割り込みは、割り込みをクラスタ内の 1 つ以上のローカル A P I C へ受け渡すために、ワンホットベクトル内に任意の一連のビットを含んでいてよい。従って、論理 A P I C __ I D を整合させることは、クラスタ I D を比較することと、ローカル A P I C 内のワンホットベクトルの一連のビットと同じ位置でのローカル A P I C __ I D ベクトル内の一連のビットを検出することとを含んでいてよい。別の方法では、論理受け渡しモード割り込みにおけるローカル A P I C __ I D ベクトルが、ローカル A P I C のローカル A P I C __ I D ベクトルと論理積を取られて(logically ANDed)よく、結果が非ゼロであり且つクラスタ I D が整合する場合には、ローカル A P I C は論理割り込みの対象である。論理 A P I C __ I D は、より簡潔にここでは論理 I D と称されることがあり、また同様に物理 A P I C __ I D は、より簡潔にここでは物理 I D と称されることがある。割り込みに関連して与えられる I D (論理又は物理)は、割り込みの宛先 I D と称されることがある。割り込みに対する対応する受け渡しモードは、割り込みの宛先 I D を識別することができる。

【 0 0 5 7 】

g A P I C 3 4 A ~ 3 4 D は、物理受け渡しモード及び論理受け渡しモードの両方を同様にサポートしてよい。上で強調したようなモードに従う割り込みメッセージ内で A P I C __ I D を整合させることに加えて、g A P I C 3 4 A ~ 3 4 D は、割り込みメッセージ内のゲスト I D を g A P I C 内のゲスト I D と整合させてよい。

【 0 0 5 8 】

I O M M U 4 0 は、I / O が開始したメモリ動作(例えば、周辺機器 5 2 A ~ 5 2 B をソースとする又は周辺機器 5 2 A ~ 5 2 B に代わる D M A 制御器によるメモリ読み出し/書き込み動作)のために仮想アドレスから物理アドレスへのマッピングを実行するように構成されてよい。翻訳動作の一部として、I O M M U 4 0 は、デバイステーブル 6 2 及び随意的に割り込みリダイレクトテーブル 6 4 へアクセスするように構成されてよい。デバイステーブル 6 2 は、各周辺機器 5 2 A ~ 5 2 B に対するエントリを含んでいてよい(また、複数の周辺機器が結合される周辺機器インタフェース上の 2 つ以上の識別子を含む周辺機器に対する多重エントリを含んでいてよい)。デバイステーブル 6 2 は、メモリ読み出し/書き込み動作のメモリアドレスを翻訳するための I / O ページテーブル(図示せず)へのページテーブルポインタを含んでいてよく、また割り込みリダイレクトテーブル 6 4 へのポインタを含んでいてよい。幾つかの実施形態においては、デバイステーブル 6 2 は、ゲストに割り当てられる周辺機器に対するゲスト I D を記憶してよい。1 つの実施形態においては、ゲスト I D は、I O M M U 4 0 においてデバイスアクセス保護のために用いられるドメイン I D と同じであってよい。代替的には、ゲスト I D は別々に割り当てられてよい。実施形態においては、デバイステーブル 6 2 はまた、g A P I C 状態マッピングテーブル 6 0 (用いられている場合)へのポインタ、又は g A P I C 状態データ構造 5 8 へのポインタを記憶してよい。別の実施形態においては、ゲスト I D 及び/又はテーブ

ル 6 0 / データ構造 5 8 へのポインタは、割り込みリダイレクトテーブル 6 4 内に記憶されてよい。割り込みリダイレクトテーブル 6 4 は、割り込みをその元の宛先及び / 又は割り込みベクトルから新たな宛先及び / 又は割り込みベクトルへリダイレクトするために用いられてよい。この開示の残りにおける簡潔さのために、ゲスト ID がデバイステーブル 6 2 からのドメイン ID であり且つマッピングテーブル 6 0 及び / 又は g A P I C 状態データ構造 5 8 へのポインタがデバイステーブル 6 2 内に記憶されている実施形態が用いられることになる。しかし、この開示の残りにおける実施形態は、上で議論したように多くの場合に修正されてよい。

【 0 0 5 9 】

他の実施形態においては、ゲスト割り込みマネージャ 3 8 は設けられていなくてよい。そのような構成は、例えば、ゲストが 1 つのプロセッサ 3 0 A ~ 3 0 B から他へ移動させられる (are migrated) 場合において V M M 1 8 がデバイステーブル 6 2 及び / 又は割り込みリダイレクトテーブル 6 4 を更新するとき、及び非アクティブゲストに代わってプロセッサ 3 0 A ~ 3 0 B が割り込みを受信するように設けられている場合 (所望に応じて、メモリ 5 6 内の g A P I C 状態データ構造 5 8 を更新し且つ / 又は割り込みをサービスするために) に可能である。

【 0 0 6 0 】

メモリ制御器 4 2 は、プロセッサ 3 0 A ~ 3 0 B によって発行されるメモリ動作 (例えば命令フェッチ、ロード / ストアデータアクセス、翻訳のためのプロセッサページテーブルアクセス、等)、ゲスト割り込みマネージャ 3 8 からのメモリ動作 (例えば g A P I C 状態データ構造 5 8 及び / 又は g A P I C 状態マッピングテーブル 6 0 を読み出し / 更新するための)、I O M M U 4 0 (例えば I / O ページテーブル、デバイステーブル 6 2、及び割り込みリダイレクトテーブル 6 4 にアクセスするための)、及びインタフェース回路 4 4 A ~ 4 4 C から受信するメモリ動作を受信するように結合されていてよい (幾つかの実施形態においては)。メモリ制御器 4 2 は、メモリ動作をオーダし、そしてメモリ 5 6 と通信してそれらメモリ動作を実行するように構成されてよい。メモリインタフェース回路 4 6 は、メモリ 5 6 への物理レベルアクセスを実行してよい。

【 0 0 6 1 】

メモリ 5 6 は任意の種類メモリを備えていてよい。例えば、メモリ 5 6 は、ダイナミックランダムアクセスメモリ (D R A M)、例えば同期 D R A M (S D R A M)、ダブルデータレート (D D R、D D R 2、D D R 3、等) S D R A M、R A M B U S _ D R A M、スタティック R A M、等を備えていてよい。メモリ 5 6 は、多重メモリチップを備えた 1 つ以上のメモリモジュール、例えばシングルインラインメモリモジュール (S I M M)、デュアルインラインメモリモジュール (D I M M) 等を含んでいてよい。

【 0 0 6 2 】

この実施形態では、ゲスト割り込みマネージャ 3 8、I O M M U 4 0、及びメモリ制御器 4 2 を含むことに加えて、ブリッジ 3 6 はまた、プロセッサ 3 0 A ~ 3 0 B、h A P I C 3 2 A ~ 3 2 B、g A P I C 3 4 A ~ 3 4 D、及びインタフェース回路 4 4 A ~ 4 4 C に結合される回路の間で通信する他の通信機能を含んでいてよい。例えば、図示される実施形態においては、別のブリッジ 4 8 がインタフェース回路 4 4 A に結合されていてよく、そしてブリッジ 4 8 は、インタフェース回路 4 4 A によって用いられるプロトコルと周辺機器 5 2 A ~ 5 2 B によって用いられるプロトコルとの間の通信をブリッジするように構成されていてよい。1 つの実施形態においては、インタフェース回路 4 4 A ~ 4 4 C は、例えば上述した H T インタフェースを実装していてよく、またブリッジ 4 8 は、H T から別のインタフェース、例えば P C I エクスプレス (P C I Express) (P C I e) インタフェースへブリッジしてよい。そのような実施形態においては、周辺機器 5 2 A ~ 5 2 B は P C I e デバイスであってよい。ブリッジ 4 8 はまた、他のインタフェースへブリッジするように構成されていてよく、あるいは別のブリッジがブリッジ 4 8 に結合されて他のインタフェースへブリッジしてもよい。任意の単一又は複数の周辺機器インタフェースが用いられてよい。加えて、周辺機器 5 2 A ~ 5 2 B は、H T インタフェースと直接的に結合

10

20

30

40

50

するように構成されるHT周辺機器を備えていてよい。そのような周辺機器は、ブリッジ48を必要としないであろう。

【0063】

1つの実施形態においては、ブリッジ48及び/又は1つ以上の周辺機器52A~52Bは、IOAPIC(図2においては50及び54)を含んでいてよい。IOAPICは、周辺機器からの割り込み要求を受信することと、割り込み要求をhAPIC32A~32B及びゲスト割り込みマネージャ38へ送信するために割り込みメッセージを形成することとに参与していてよい(gAPIC34A~34Dへの送信及び/又はメモリ内のgAPIC状態データ構造58内への記録のために)。

【0064】

上述したように、1つの実施形態においては、インタフェース回路44A~44CはHTインタフェース上で通信するように構成されてよい。インタフェース回路44A~44Cは、HTを用いて周辺デバイス/ブリッジと通信するように構成されてよい。また、幾つかの実施形態においては、インタフェース回路44A~44Cは、プロセッサを伴う他のノード、hAPIC、gAPIC、等と結合されるように構成されてよい。そのような実施形態においては、ブリッジ36は、前述した回路に加えてコヒーレンス管理回路を含んでいてよい。

【0065】

プロセッサ30A~30Bは、任意の命令セットアーキテクチャを実装していてよく、そして命令セットアーキテクチャにおいて定義される命令を実行するように構成されてよい。プロセッサ30A~30Bは、任意のマイクロアーキテクチャ、例えば、スーパーパイプラインにされたもの(superpipelined)、スーパースケラ(superscalar)、及び/又はこれらの組み合わせ、順序内(in-order)又は順序外(out-of-order)の実行、投機的実行(speculative execution)、等を含んでいてよい。プロセッサ30A~30Bは、所望に応じてマイクロコーディング技術を実装していてもよいし実装していなくてもよい。

【0066】

周辺機器52A~52Bは任意の種類の周辺デバイスを備えていてよい。周辺機器52A~52Bは、記憶デバイス、例えば、磁気的な、ソリッドステートの、又は光学的なディスクドライブ、フラッシュメモリのような不揮発性メモリデバイス等を含んでいてよい。周辺機器52A~52Bは、I/Oデバイス、例えば、ユーザI/Oデバイス(キーボード、マウス、ディスプレイ、音声入力、等)、ネットワークデバイス、ユニバーサルシリアルバス(USB)又はファイヤワイヤ(Firewire)のような外部インタフェースデバイス等を含んでいてよい。

【0067】

図示される実施形態においては、プロセッサ30A~30B、ブリッジ36、hAPIC32A~32B、gAPIC34A~34D、インタフェース回路44A~44C、及びメモリインタフェース回路46は、集積回路66として単一の半導体基板上に集積化されてよい。他の実施形態は、所望に応じて異なる量の集積化及び個別回路を実装してよい。尚、図2においては種々の数のコンポーネント、例えばプロセッサ、hAPIC、gAPIC、インタフェース回路、周辺機器、ブリッジ等が図示されているが、他の実施形態は、所望に応じて1つ以上の任意の数の各コンポーネントを実装してよい。

【0068】

他の実施形態においては、IOMMU40及びゲスト割り込みマネージャ38の位置は変わってよい。例えば、一方又は両方は、ブリッジ48内、周辺機器52A~52B内、ブリッジに結合される別のブリッジ内、等にあってよい。

【0069】

図示される実施形態においては、図2に示されるように、各gAPIC34A~34D及びhAPIC32A~32Bは、特定のプロセッサ30A~30Bに付随している。従って、この実施形態においては、所与の割り込み制御器は対応するプロセッサ30A~30Bに専用である。より具体的には、図2においては、hAPIC32A並びにgAPIC

10

20

30

40

50

C 3 4 A 及び 3 4 C はプロセッサ 3 0 A の専用であり、そして h A P I C 3 2 B 並びに g A P I C 3 4 B 及び 3 4 D はプロセッサ 3 0 B の専用である。割り込み制御器は、対応するプロセッサへ任意の方法で割り込みの信号を送ってよい。一般的に、信号を送ること(signaling)は、割り込みが必要とされていることを表示してよい。信号を送ることは割り込みベクトルを含んでいてよく、即ち割り込みベクトルは割り込みが受け渡された後に実行されるソフトウェアによって読まれてよい。割り込みを受け渡すことは、実施形態においては、プロセッサに信号を送ること及びプロセッサが割り込みを受け入れることを参照してよい。割り込みをサービスすることは、割り込みベクトルに関連する割り込みサービスルーチンを実行して、割り込みしているデバイスによって必要とされる動作を実行することを参照してよい。

10

【 0 0 7 0 】

次に図 3 を参照すると、1 つの実施形態に対する周辺機器から g A P I C への割り込みの進行を表すブロック図が示されてる。他のプロセッサからの割り込み(プロセッサ間割り込み(interprocessor interrupts)、又は I P I)もまた、ゲスト割り込みマネージャ 3 8 へ送信されてよく、そしてその時点以降は図 3 と同様に取り扱われてよい。代替的には、I P I を開始しているプロセッサからの I P I を受信する g A P I C は、更新をゲスト割り込みマネージャ 3 8 へ送信してよく(受信しているゲストが非アクティブである場合にそのゲストに対して g A P I C 状態を更新するために)、また I P I (ゲスト I D を含む)を他の g A P I C へ送信してもよい。

【 0 0 7 1 】

20

図示される実施形態においては、周辺機器 5 2 A は、割り込みが望まれていることを決定する。周辺機器 5 2 A 内の I O A P I C 5 4 (図 2 参照)は、周辺機器 5 2 A の動作に応答して割り込みメッセージを生成してよい。具体的には、I O A P I C 5 4 は、望まれている割り込み(例えば、周辺機器 5 2 A によって必要とされるサービスに基いて、周辺機器 5 2 A が多重機能を実装している場合に特定の機能が割り込みの信号を送ること、等)に対応する割り込みベクトルを生成してよい。割り込みベクトルは、割り込み通信の一部であり、そして割り込みソースを識別し、割り込みを優先付ける、等のためにソフトウェアによって用いられてよい。幾つかの場合には、割り込みベクトルは I O M M U 4 0 によって再マッピングされてよく、そのため図 3 においては、割り込みベクトルは「元のベクトル(original vector)」として示されている。周辺機器 5 2 A は、割り込みメッセージを I O M M U 4 0 へ送信してよい(矢印 A)。この実施形態においては、割り込みは、例えば P C I e 規格において定義されるようなメッセージが信号で送られた割り込み(message-signalled interrupt) (M S I) の形態で送信されてよい。他の実施形態は、割り込みを任意の望ましい方法で送信してよい。一般的に、送信は、割り込み、その受け渡しモード(例えば論理又は物理)、及び割り込みの宛先 I D (D e s t I D) を識別してよい。

30

【 0 0 7 2 】

I O M M U 4 0 は M S I を受信してよい。M S I は周辺機器の識別子を含んでいる。例えば、P C I プログラミングモデルを実装しているインタフェースは、各デバイスをバス番号及び当該バスのデバイス番号で識別することができる(階層的な且つ/又は並列な形態にあるシステム内に多重 P C I インタフェースが存在することを可能にする)。デバイスは多重「機能(functions)」を有していてよく、多重機能は、物理デバイス上の別々の複数の仮想デバイスであってよく、あるいはデバイス上での複数の動作の分割(partitioning)であってよい。識別子はまた、機能番号を含んでいてよい。従って、この実施形態においては、識別子は、バスデバイス機能(Bus-Device-Function)又は B D F と称されることがある。I O M M U 4 0 は、B D F を用いてデバイステーブル 6 2 内へ索引付ける(index)ことができ(矢印 B)、また周辺機器 5 2 A に対応するデバイステーブルエントリを識別することができる。エントリは、幾つかの実施形態においては、ゲスト I D と、g A P I C 状態マッピングテーブル 6 0 又は g A P I C 状態データ構造 5 8 へのポインタとを含んでいてよい(矢印 C)。この実施形態においては、デバイステーブルエントリはまた

40

50

、デバイスに対応する割り込みリダイレクトテーブル 6 4 を識別し得る割り込みリダイレクトテーブルポインタ (I R T P) を含んでいてよい (矢印 C 1) 。割り込みリダイレクトテーブル 6 4 は、元の割り込みベクトルによって索引付けられてよく、そして出力ベクトル及び宛先 I D (D e s t I D 、例えば論理又は物理 A P I C _ I D) を割り込みに対して提供してよい (矢印 C 2) 。

【 0 0 7 3 】

図 3 は M S I がベクトル 4 2 、ゲスト I D 9 9 へ再マッピングされる例を示している。再マッピングはゲスト I D を付加することを含んでいてよく、そしてベクトルはまた、割り込みリダイレクトテーブル 6 4 が用いられている場合に变化させられてよい。そうでない場合には、M S I からの元の割り込みベクトルは、割り込みメッセージ内で提供されてよい。割り込みベクトル 4 2 及びゲスト I D 9 9 の具体例が用いられている図 3 における点は、大括弧、即ち [] で囲まれて図示されている。

【 0 0 7 4 】

I O M M U 4 0 は、ゲスト I D (例えばこの例では 9 9) を含めて割り込みメッセージをゲスト割り込みマネージャ 3 8 へ送信してよい。割り込みメッセージはまた、割り込みベクトル (例えばこの例では 4 2) 及び宛先 I D を含む。割り込みメッセージはまた、g A P I C 状態マッピングテーブル 6 0 又は g A P I C 状態データ構造 5 8 へのポインタを含んでいてよい (矢印 D) 。

【 0 0 7 5 】

g A P I C 状態マッピングテーブル 6 0 を実装している実施形態においては、ゲスト割り込みマネージャ 3 8 は、ポインタ並びに場合によってはゲスト I D 及び / 又は宛先 I D のような他の情報を用いて g A P I C 状態ポインタを g A P I C 状態マッピングテーブル 6 0 内に配置してよい (矢印 E 1 、そしてゲスト割り込みマネージャ 3 8 への戻りのポインタは矢印 E 2 で示されている) 。 g A P I C 状態ポインタは、g A P I C 状態データ構造 5 8 内の g A P I C 状態エントリを識別してよく、またゲスト割り込みマネージャ 3 8 は、g A P I C 状態ポインタを用いて g A P I C 状態データ構造 5 8 内での g A P I C 状態更新を実行してよい (矢印 E) 。この実施形態においては、g A P I C 状態更新は、ベクトル 4 2 に対応する割り込み要求レジスタ内のビットをセットしてよい。割り込み要求レジスタ (I R R) は、図 4 に関して後で更に詳細に説明される。

【 0 0 7 6 】

1 つの実施形態においては、g A P I C 状態 5 8 への更新はアトミック (atomic) であってよい。例えば、ゲスト割り込みマネージャ 3 8 は、g A P I C 状態エントリにおける割り込み要求レジスタのその時点での状態内へセットされつつある割り込み要求についてアトミックに論理和を取る (atomically ORs) アトミック O R トランザクションを生成してよい。アトミックな動作は、動作が多重ステップとして実装されている場合であっても、一単位として効果的に実行される動作であってよい。アトミックに更新されつつあるロケーションへアクセスしようとしているオブザーバは、アトミックな更新に先立つ値又はアトミックな更新の後の値のいずれかを受信するが、中間値は受信しないであろう。アトミックに更新されつつあるロケーションへアクセスしようとしているオブザーバは、その更新をアトミックな動作の前又はアトミックな動作が完了した後のいずれかに実行するが、アトミックな動作の間には実行しない。この実施形態はアトミック O R を実装しているが、他の実施形態は、より一般的なアトミックな更新動作を実装してよい。例えば、アトミックな更新は、修正されるべきでない対象のビットを識別する A N D マスクを含んでいてよく、またどのビットが論理和を取られるべきかを識別する O R マスクを含んでいてよい。他の実装もまた可能である。例えば、比較及び交換実装 (compare and swap implementation) が用いられてよく、この実装においては、メモリロケーションからの元の値が読まれ、そして比較及び交換動作が元の値及び論理和を取られた新たな値に対して実行される。比較が失敗した場合には、プロセスは繰り返されてよい (新たな元の値を読み、そして比較及び交換を実行する) 。必要であれば、ループのフェイルアウト (fail out) に対してバックオフ (backoff) 及び / 又はタイムアウトのメカニズムが用いられてよい。

10

20

30

40

50

【 0 0 7 7 】

ゲスト割り込みマネージャ 38 はまた、割り込みベクトル、ゲスト ID、及び宛先 ID を含む割り込みメッセージを g A P I C 3 4 A ~ 3 4 D へブロードキャストしてよい (矢印 F)。g A P I C の 1 つ (図 3 では g A P I C 3 4 A) は、ゲスト ID 9 9 と宛先 ID に整合する論理又は物理 A P I C _ I D とを有していてよく、従って、g A P I C 3 4 A は、それが割り込みメッセージを受け入れたことを表示する承認 (acknowledgement) (A c k) で割り込みメッセージに応答してよい (矢印 G)。g A P I C 3 4 A はまた、その割り込み要求レジスタを更新して、この実施形態ではベクトル 4 2 に対応するビットをセットしてよい。割り込みが任意の進行中の割り込み (もしあれば) 及び / 又はプロセッサのタスク優先度よりも高い場合には、g A P I C 3 4 A はまた、割り込みの信号をプロセッサ 3 0 A へ送ってよい。他の g A P I C 3 4 B ~ 3 4 D は、ブロードキャスト割り込みメッセージに
10 応答してよいが、それらは割り込みの対象ではないから、受け入れを承認しなくてよい (矢印 H)。論理割り込みに対しては、論理割り込みが多重対象を識別している場合には、2 つ以上の承認があってよい。

【 0 0 7 8 】

上述のメカニズムを用いると、ゲスト割り込みマネージャ 38 は、どの g A P I C 3 4 A ~ 3 4 D がどのゲストに割り当てられているかを「知っている (aware) 」必要がない。どの g A P I C 3 4 A ~ 3 4 D がどのゲストに割り当てられているか、及びどれが対象となる g A P I C に対してのみ割り込みを送信するかをゲスト割り込みマネージャ 38 が追跡する他の実施形態が検討される。ゲスト割り込みマネージャ 38 は、g A P I C を自動的に追跡してよく、あるいは g A P I C が他のゲストに再割り当てされる毎に V M M 1 8 によってプログラムされてよい。そのような実施形態においては、ゲスト割り込みマネージャ 38 は、対象となる g A P I C にのみ割り込みメッセージを送信してよい。
20

【 0 0 7 9 】

割り込みの h A P I C 3 2 A ~ 3 2 B への送信は、通常の A P I C の方法で実行されてよい。具体的には、割り込みは、ゲスト割り込みマネージャ 38 による場合には動作させられなくてよいが、実施形態においては、他の点については図 3 の動作と同様であってよい。

【 0 0 8 0 】

尚、ゲスト割り込みマネージャ 38 はここではブロックとして図示されているが、ゲスト割り込みマネージャ 38 を実装している回路が配布されてもよい。例えば、実施形態においては、ポインタを受信し、g A P I C 状態マッピングテーブル 6 0 を随意的に処理し、そして g A P I C 状態データ構造 5 8 に対する更新を生じさせるゲスト割り込みマネージャ 38 の一部分が I O M M U 4 0 内に含まれていてよく、この場合、I O M M U 4 0 は、g A P I C 状態データ構造 5 8 に対するアトミック O R 及び g A P I C 3 4 A ~ 3 4 D へ送信されるべき割り込みメッセージを送信する。1 つ以上の物理ロケーションにおけるゲスト割り込みマネージャ 38 の任意の実装が用いられてよい。
30

【 0 0 8 1 】

次に図 4 を参照すると、g A P I C 3 4 A の 1 つの実施形態のブロック図が示されている。他の g A P I C 3 4 B ~ 3 4 D は同様であってよい。図 4 の実施形態においては、g A P I C 3 4 A は、割り込み要求レジスタ (I R R) 7 0、割り込みサービスレジスタ (I S R) 7 2、トリガモードレジスタ (T M R) 7 4、タスク優先度レジスタ (T P R) 7 6、制御ユニット 7 8、物理 ID レジスタ 8 0、論理 ID レジスタ 8 2、ゲスト ID レジスタ 8 4、及び随意的な他の A P I C 状態 8 6 を含む。制御ユニット 7 8 は、I R R 7 0、I S R 7 2、T M R 7 4、T P R 7 6、物理 ID レジスタ 8 0、論理 ID レジスタ 8 2、ゲスト ID レジスタ 8 4、及び他の A P I C 状態 8 6 と結合される。加えて、制御ユニット 7 8 は、割り込みを受信するためにゲスト割り込みマネージャ 38 と通信するように結合され、またプロセッサ 3 0 A と通信するためにプロセッサインタフェースに結合される。
40

【 0 0 8 2 】

ゲスト割り込みマネージャ 38 からの割り込みメッセージを受信することに応答して、制御ユニット 78 は、g A P I C 3 4 A に対応するゲストを割り込みが対象としている場合に、割り込みを I R R 7 0 内に書き込むように構成されてよい。I R R 内の割り込み要求の位置は、割り込みベクトルに対応する。I R R は「固定された(fixed)」割り込みを追跡してよい。他の割り込み種類としては、非マスク可能割り込み(non-maskable interrupt) (N M I)、システム管理割り込み(S M I)、レガシ外部割り込み(legacy external interrupt) (e x t I N T)、等を挙げることができる。これらの割り込みは、他の A P I C 状態 86 の一部として取り扱われてよい。

【0083】

1つの実施形態においては、割り込みメッセージはまた、各割り込みに対するトリガモード(レベル又はエッジ)を含んでいてよい。T M R 7 4 は、どちらのトリガモードを割り込みに適用するかの表示を記憶してよい。例えば、エッジトリガされる割り込みは T M R 7 4 内のバイナリ 0 によって表されてよく、またレベルトリガされるものはバイナリ 1 によって表されてよい。他の実施形態においては、エッジトリガされる割り込みのみが g A P I C 3 4 A 内でサポートされてよく、また T M R 7 4 (及び g A P I C 状態データ構造 58 内のそのコピー)は除かれてよい。別の実施形態においては、T M R 7 4 は、仮想レベル敏感な割り込みを V M M 1 8 が記録することを可能にするために再度目的を持たされてよい。V M M 1 8 は、対応する割り込みベクトルに対してプロセッサ 30 A によって割り込みの終点の信号が送られた場合に、プロセッサ 30 A が V M M 1 8 から退出することを表示する種々のビットを T M R 7 4 内にセットしてよい。

【0084】

固定された割り込みに対しては、g A P I C 3 4 A は、割り込み要求がプロセッサへ受け渡されるべきかを決定するために、割り込み要求及びインサース割り込みを優先付けるように構成されてよい。概して、最も高い優先度の割り込み要求が最も高い優先度のインサース割り込み(プロセッサがそのソフトウェア実行に割り込んだ結果、割り込みに対応する割り込みハンドラを実行した場合には、割り込みはインサースである)よりも高い優先度である場合には、制御ユニット 78 は、要求された割り込みをプロセッサ 30 A へ受け渡すように構成されてよい。また、T P R 7 6 は、プロセッサ 30 A によって受け入れられている最も低い優先度レベルの割り込みを確立するために、ソフトウェアによってプログラムされてよい。制御ユニット 78 は、最も高い優先度の割り込み要求を、それが最も高い優先度のインサース割り込みよりも高い優先度である場合、及びそれが T P R 7 6 において表示されている優先度よりも高い場合に、受け渡すように構成されてよい。

【0085】

プロセッサ 30 A が割り込みを取る場合には、プロセッサは、割り込み承認コマンドで g A P I C 3 4 A に対して応答してよい。制御ユニット 78 は、最も高い優先度の割り込み要求を I R R 7 0 から取り除くと共に当該割り込みをインサースとして I S R 7 2 内に記録するように構成されてよい。I S R 内の割り込みに対応するインサース表示の位置は、当該割り込みの割り込みベクトルに対応してよい。プロセッサ 30 A は、単一又は複数の割り込みサービスルーチンを実行して割り込みをサービスしてよい。割り込みサービスルーチンは、割り込みサービスが完了した旨の信号を送るための、g A P I C 3 4 A への割り込み終点(end of interrupt) (E O I) コマンドで終了してよい。制御ユニット 78 は、E O I コマンドに応答して最も高い優先度のインサース割り込みを I S R 7 2 から除くように構成されてよい。

【0086】

I R R 7 0、I S R 7 2、及び T M R 7 4 の各々は、g A P I C 3 4 A によってサポートされる各割り込みベクトルに対応するロケーションを有している。図示される実施形態においては、ベクトル 0 乃至 255 がサポートされている。割り込みベクトル番号はまた、他の割り込みとのその相対的優先順位を示してよい(例えば、より大きなベクトル番号はそれより小さなベクトル番号よりも高い優先度であり、他の実施形態においては反対で

ある)。各割り込みベクトルに対して、IRR70は、当該割り込みベクトルで割り込みが要求されているか否かを表示する割り込み要求ビットを記憶する。例えば、表示は、セットされている場合に要求を表し且つクリアの場合に要求がないことを表すビットであってよい。同様に、各割り込みベクトルに対して、ISR72は、割り込みが当該割り込みベクトルに対してインサースビスであるか否かを示すインサースビスビットを記憶する（例えば、セットの場合にインサースビスを表し且つクリアの場合にインサースビスでないことを表す）。各割り込みベクトルに対してTMR74はトリガモードを記憶する。IRR70、ISR72、及びTMR74の各々に対して、レジスタ内のビットロケーションは、割り込みに対応する割り込みベクトル番号と等しい。

【0087】

図示される実施形態においては、複数の割り込みは、係属中の割り込み要求がプロセッサへ受け渡されるべきかを決定するための優先度レベルを割り当てられている群内へ配置される。例えば、割り込みベクトル0乃至15は優先度レベル0を割り当てられ、割り込みベクトル16乃至31は優先度レベル1を割り当てられ、優先度レベル15での割り込みベクトル240乃至255まで同様である。この実施形態においては、優先度レベル番号が増大することは、優先度レベルが高くなることを示す。制御ユニット78は要求優先度レベルを計算することができ、要求優先度レベルは、少なくとも1つの割り込み要求がIRR70内に係属中である最も高い優先度レベルである。制御ユニット78はまた、インサースビス優先度レベルを計算することができ、インサースビス優先度レベルは、少なくとも1つの割り込みがISR72内でインサースビスとして表示されている最も高い優先度レベルである。制御ユニット78は、要求優先度レベルがインサースビス優先度レベルを超過し且つTPR76内で表示される優先度レベルを超過する場合に、割り込みを受け渡してよい。尚、図示される実施形態においては、16個の優先度レベル群内で256個の割り込みベクトルがサポートされるが、他の実施形態においては、それより多い若しくは少ない割り込みベクトル及び/又はそれより多い若しくは少ない優先度レベル群がサポートされてよい。

【0088】

物理IDレジスタ80及び論理IDレジスタ82は、gAPIC34Aに割り当てられる物理APIC__ID及び論理APIC__IDをそれぞれ記憶してよい。ゲストIDレジスタ84は、gAPIC34Aに割り当てられるゲストIDを記憶してよい。従って、制御ユニット78は、割り込みのゲストIDがゲストIDレジスタ84内のゲストIDに整合する場合であって、割り込みが物理であり且つ割り込みにおけるAPIC__IDが物理IDレジスタ80内の物理IDに整合するか、あるいは割り込みが論理であり且つ割り込みにおけるAPIC__IDが論理IDレジスタ82内の論理IDに整合するかのいずれかの場合に、ゲスト割り込みマネージャ38からの割り込みを受け入れるように構成されてよい。

【0089】

他のAPIC状態86は、内部的に生成される割り込み、タイマ、ローカルベクトルテーブル等を含んでいてよい。種々の実施形態においては、他のAPIC状態86の幾つか又は全ては、gAPIC34A内に含まれていてよく、又はVMM18へのインターセプト及び状態のVMM18エミュレーションと共に仮想化されてよい。

【0090】

hAPIC32A~32Bは、それらがゲストIDレジスタを含んでいなくてよいことを除いてgAPIC34Aと同様であってよい。代替的には、hAPIC32A~32B及びgAPIC34A~34Dは、同一ハードウェアの例(instances)であってよく(gAPIC34A~34Bが全てのAPIC状態を実装している場合)、またゲストIDレジスタは、ゲストIDが有効であるか否かを表示する有効ビットを含んでいてよく、あるいはゲストIDレジスタはhAPICを表示するためにゼロにプログラムされてよい。

【0091】

次に図5を参照すると、gAPIC状態エントリ90の1つの実施形態及びVMM管理

10

20

30

40

50

の g A P I C 状態エントリ 9 2 の 1 つの実施形態のブロック図が示されている。図 5 の図示は状態の論理ビューであってよい。メモリ内の状態の実際の配列は、幾つかの実施形態に対して図 9、12、及び 13 に示されるように変化し得る。

【0092】

一般的に、g A P I C 状態エントリ 9 0 は、g A P I C 状態に対応するゲストがアクティブでない間に変化し得る g A P I C を少なくとも含んでいてよい。本実施形態においては、周辺デバイスが割り込みの信号をゲストへ送ってよく、ゲストは I R R 状態を変化させてよい。しかし、I S R 状態は、ゲストにおける v C P U が割り込みを受け入れる場合にのみ変化してよく、その割り込みは、ゲストがアクティブでない場合には生じないであろう。同様に、T P R が v C P U によって変化させられ、従って T P R は、ゲストがアクティブでない間には変化しないであろう。V M M 1 8 は、V M M 管理の g A P I C 状態エントリ 9 2 におけるそのような状態の保存及び復元を管理することができる。

10

【0093】

従って、図 4 と同様の g A P I C 3 4 A の実施形態に対しては、g A P I C 状態エントリ 9 0 は I R R 7 0 の状態を含んでいてよい。V M M 管理の g A P I C 状態エントリ 9 2 は、I S R 7 2、T M R 7 4、T P R 7 6 の状態、及び種々の他の A P I C 状態 8 6 を含んでいてよい。V M M 管理の g A P I C 状態エントリ 9 2 はまた、ゲスト I D 並びに論理及び物理 I D を記憶してよく、あるいはこれらはエントリ 9 2 を選択することにおいて固有のものであってよい（即ち、V M M 1 8 は、それらの値を用いてデータ構造 6 8 からエントリ 9 2 を選択してよい）。

20

【0094】

次に図 6 を参照すると、ゲストに対して I O M M U 4 0 から割り込みメッセージを受信することに応答するゲスト割り込みマネージャ 3 8 の 1 つの実施形態の動作を表すフローチャートが示されている。理解の容易化のために複数のブロックが特定の順序で示されているが、他の順序が用いられてもよい。これらのブロックは、ゲスト割り込みマネージャ 3 8 内の組み合わせ論理において並列に実行されてよい。これらのブロック、これらのブロックの組み合わせ、及び / 又はフローチャートは、多重クロックサイクルでパイプライン化されてよい。一般的に、ゲスト割り込みマネージャ 3 8 は、図 6 に示される動作を実装するように構成されてよい。

【0095】

30

幾つかの実施形態においては、割り込みメッセージの処理は、割り込みが論理であるか物理であるかに応じて（即ち、割り込みの受け渡しモードが論理であるか物理であるかに応じて）変化してよい。例えば図 11 の実施形態においては、論理割り込み及び物理割り込みに対して異なるテーブルが読まれる。図 12 及び 13 においては、論理テーブル及び物理テーブルはメモリ内で隣接していてもよいが、論理割り込みに対して論理テーブルを配置するベースアドレスポインタにはオフセットが付加されてよく、また物理割り込みに対してはオフセットは付加される必要はない。従って、ゲスト割り込みマネージャ 3 8 は、割り込みが論理であるかあるいは物理であるかを決定するように構成されてよい（判断ブロック 100）。他の実施形態は、受け渡しモードに基づいて変化しなくてよく、そして判断ブロック 100 は除かれてよい（また、以下に議論されるブロードキャスト又はより多くの宛先に対するチェックは、両方に対する 1 つのチェックにまとめられてよい）。

40

【0096】

割り込みが論理である場合（判断ブロック 100、「イエス」行程）、ゲスト割り込みマネージャ 3 8 は、論理割り込みから g A P I C 状態データ構造 5 8 内の対応する g A P I C 状態エントリ 9 0 へのマッピングを決定するように構成されてよい（ブロック 102）。図 10 ~ 13 に示されるように、種々の実施形態は異なるマッピングを実装してよく、従って決定は変化し得る。ゲスト割り込みマネージャ 3 8 は、g A P I C 状態エントリ 9 0 内で代表される I R R 内の割り込みベクトルに対応するビットをセットするように構成されてよい（ブロック 104）。論理割り込みは多重宛先を有していてもよい（例えば、クラスタ内の宛先は、1 つ以上のセットビットを有していてもよいビットベクトルである）

50

。論理割り込みがより多くの宛先を含む場合（判断ブロック106、「イエス」行程）、ゲスト割り込みマネージャ38は、各追加的な宛先に対してブロック102及び104を繰り返すように構成されてよい。代替的には、図12の実施形態においては、後で更に詳細に説明されるように、論理宛先ビットベクトルは、1つの動作におけるgAPIC状態エントリへ書き込まれてよい。ゲスト割り込みマネージャ38は、割り込みメッセージをgAPIC34A~34Dへ送信するように構成されてよい（ブロック108）。

【0097】

割り込みが物理である場合（判断ブロック100、「ノー」行程）、ゲスト割り込みマネージャ38は、物理割り込みからgAPIC状態データ構造58内の対応するgAPIC状態エントリ90へのマッピングを決定するように構成されてよい（ブロック110）。図10~13に示されるように、種々の実施形態は異なるマッピングを実装してよく、従って決定は変化し得る。ゲスト割り込みマネージャ38は、gAPIC状態エントリ90内で代表されるIRR内の割り込みベクトルに対応するビットをセットするように構成されてよい（ブロック112）。物理割り込みは、ブロードキャスト又は単一の宛先であってよい。物理割り込みがブロードキャストである場合（判断ブロック114、「イエス」行程）、ゲスト割り込みマネージャ38は、ゲストの仮想マシン（例えば各vCPU）内の各宛先に対してブロック110及び112を繰り返すように構成されてよい。代替的には、図12の実施形態においては、後で更に詳細に説明されるように、ブロードキャストは、1つの動作におけるgAPIC状態エントリ内に記録されてよい。ゲスト割り込みマネージャ38は、割り込みメッセージをgAPIC34A~34Dへ送信するように構成されてよい（ブロック108）。

【0098】

gAPIC状態エントリ90内で代表されるIRR内でビットをセットすることは、アトミックOR動作として実行されてよく、アトミックOR動作においては、セットされたビットはメモリロケーション内の他のIRRビット内へ論理和を取られる。アトミックOR動作の実際の実装は、ロックされた読み出し/修正/書き込み動作から、ORを1つの動作として実行するように規定された特殊目的回路まで、変化してよい。

【0099】

上述したように、他の実施形態においては、比較及び交換動作が実行されてよい。

【0100】

別の実施形態においては、2つ以上の宛先を伴う論理割り込み及び、ブロードキャストである物理割り込みが、VMM18にアクセス可能なデータ構造（例えばイベントキュー(event queue)）内に割り込みを記録することによって、ゲスト割り込みマネージャ38により取り扱われてよい。ゲスト割り込みマネージャ38はまた、VMM18にイベントを知らせるために、VMM18へ信号を送るように構成されてよい（例えば、プロセッサ30A~30Bの1つにおける仮想マシンからの退出を生じさせる）。代替的には、ゲスト割り込みマネージャ38は、VMM18へ周期的にのみ信号を送ってよく（例えば、Nミリ秒毎に1回且つ/又はイベントキューにおける高ウォーターマーク(high watermark)で1回）、そしてVMM18は、周期的にイベントキューをチェックしてよい他、信号を送ることがサポートしたであろうよりも迅速に任意のイベントをサービスし得る。1つの実施形態においては、イベントキューは、ゲスト割り込みマネージャ38に代えてIOMMU40によって管理されてよい。

【0101】

次に図7を参照すると、ゲスト割り込みマネージャ38からの割り込みメッセージを受信することに対応するgAPIC34A~34Dの1つの実施形態の動作を表すフローチャートが示されている。理解の容易化のために複数のブロックが特定の順序で示されているが、他の順序が用いられてもよい。これらのブロックは、gAPICの組み合わせ論理において並列に実行されてよい。これらのブロック、これらのブロックの組み合わせ、及び/又はフローチャートは、多重クロックサイクルでパイプライン化されてよい。一般的に、gAPICは、図7に示される動作を実装するように構成されてよい。

【 0 1 0 2 】

1つの実施形態においては、g A P I Cは、そのゲストID（ゲストIDレジスタ84内の、図4参照）をゼロにセットすることによって非アクティブ化される。従って、割り込みメッセージを受信することに応答して、g A P I CゲストIDがゼロである場合（判断ブロック120、「イエス」行程）、g A P I Cは非アクティブであり割り込みを処理しなくてよい。他の実施形態は他の方法（例えばレジスタ内のアクティブビット）でg A P I Cを非アクティブ化してよく、そして判断ブロック120は、g A P I Cアクティブ／非アクティブに対するチェックに従って修正されてよい。

【 0 1 0 3 】

g A P I CゲストIDが非ゼロである場合、g A P I Cは、ゲストIDを、受信した割り込みのゲストIDと比較する他、受信した宛先IDをレジスタ80及び82内のそれぞれ論理ID及び物理ID（図4参照）と比較するように構成されてよい。g A P I CのゲストIDが受信したゲストIDと整合しない場合（判断ブロック122、「ノー」行程）、g A P I Cはその時点で異なるゲストに割り当てられており、従ってg A P I Cは割り込みによって対象とされていない。g A P I Cは、割り込みの非承認で応答するように構成されてよい（ブロック124）。非承認は、g A P I Cが割り込みを受信したが、割り込みは対応するプロセッサで対象とされておらず、従って受け入れられなかったとg A P I Cが判断したことを表示してよい。同様に、g A P I CのゲストIDは受信したゲストIDと整合するが、割り込みは論理であり且つg A P I Cの論理IDと整合せず、あるいは割り込みは物理であり、単一の宛先であり、且つg A P I C物理アドレスと整合しない場合（判断ブロック126及び128、「ノー」行程）、g A P I Cは、割り込みの非承認で応答するように構成されてよい（ブロック124）。

【 0 1 0 4 】

論理割り込みを整合させることは、一般的に、論理IDのクラスタID部分を均等のために比較することと、g A P I Cの論理IDレジスタ内のセットビットもまた割り込みから受信された論理IDの宛先部分においてセットされていることを検出することを含んでいてよい。割り込みの論理IDの宛先部分における他のビットもまた、2つ以上の宛先がある場合にセットされてよい。物理IDは、ゲストIDが整合する限りブロードキャスト物理割り込みが整合として取り扱われ得ることを除いて、均等のために比較されてよい。

【 0 1 0 5 】

割り込みが論理であり且つ論理IDと整合する場合（判断ブロック126、「イエス」行程）、又は割り込みが物理であり且つ物理IDと整合し若しくはブロードキャストである場合（判断ブロック128、「イエス」行程）、g A P I C34は、対応するプロセッサ30A～30Bへの提示のためにg A P I Cが割り込みを受け入れていることを表示する承認でゲスト割り込みマネージャ38に応答するように構成されてよい（ブロック130）。g A P I Cはまた、IRRレジスタ70を更新して割り込みメッセージ内の割り込みベクトルに対応する割り込み要求ビットをセットするように構成されてよい（ブロック132）。g A P I Cは、任意のインサースervice割り込みに関する割り込みの優先度及び／又はタスク優先度レジスタを再評価するように構成されてよく（ブロック134）、また再評価に基づいてプロセッサへ割り込みの信号を送るように構成されてよい（ブロック136）。即ち、g A P I Cは、割り込みの優先度がインサースervice割り込みよりも高い優先度であり且つタスク優先度レジスタよりも高い優先度である場合に、割り込みの信号を送るように構成されてよい。

【 0 1 0 6 】

次に図8を参照すると、g A P I C状態を1つのゲストから他へ変化させるVMM18の1つの実施形態の動作を表すフローチャートが示されている。即ち、図8のフローチャートは、g A P I C34A～34Dを1つのゲスト／v C P Uから他のゲスト又は同一ゲスト内の他のv C P Uへ再割り当てすることを提示してよい。理解の容易化のために複数のブロックが特定の順序で示されているが、他の順序が用いられてもよい。一般的に

10

20

30

40

50

、VMM 18は、システム5上で実行されるときに図8に示される動作を実装する命令を含んでよい。

【0107】

VMM 18は、gAPIC状態データ構造58内の「古いゲスト(old guest)」(gAPICから非アクティブ化されているゲスト)に対応するgAPIC状態エントリ90のロケーションを決定してよい(ブロック140)。gAPIC状態エントリ90内のデータは、それがgAPICによって修正されてしまっていることがあるという理由で、「新鮮でない(stale)」と考えられる。例えば、IRRビットは、対応するプロセッサへ割り込みを受け渡すことに応答してリセットされてしまっているかもしれない。従って、VMM 18はgAPIC状態エントリ90内のIRRをゼロにしてよい(ブロック142)。VMM 18は、ゲストIDレジスタ84、論理IDレジスタ82、及び物理IDレジスタ80をクリアしてよい(ブロック144)。この動作により、ゲストID、論理ID、及び物理IDは割り込みメッセージに整合しなくなることになるので、gAPICが任意の追加的な割り込みを受け入れることを妨げることができる。レジスタ80~84がクリアされ(ブロック144)且つIRR状態がgAPIC状態エントリ90へ書き込まれる前に、割り込みが送信され得ていることが可能である。従って、割り込み状態の損失を回避するために、IRR70からgAPIC状態エントリ90内へIRR状態についてアトミックに論理和を取ってよい(ブロック146)。VMM 18はまた、古いゲストに関連するVMM管理のgAPIC状態エントリ92へ他のgAPIC状態を書き込んでよい。

【0108】

VMM 18は、gAPICのIRR、ISR、及びTMRのレジスタ70、72、及び74をクリアして古いゲストの割り込み状態を解除してよい(ブロック150)。VMM 18は、gAPICに割り当てられつつあるゲストに対する新たなゲストID、論理ID、及び物理IDをそれぞれゲストIDレジスタ84、論理IDレジスタ82、及び物理IDレジスタ80に書き込んでよい(ブロック152)。ブロック152が一旦実行されると、gAPICはそのゲストのための割り込みを受け入れることを開始してよい。VMM 18は、「新たなゲスト(new guest)」(gAPICにおいてアクティブ化されつつあるゲスト)に対するgAPIC状態エントリ90を決定してよく(ブロック154)、そしてgAPIC状態エントリ90からIRR状態を読み出してよい(ブロック156)。レジスタ80~84のプログラミングは、gAPICに割り込みを受け入れることを開始させるであろうから、VMM 18がエントリを読み出した後にgAPIC状態エントリ90内に記録されていたIRR内の割り込みをgAPICが受け入れてしまっていることが可能である。従って、VMM 18は、IRRレジスタ70内へIRR状態についてアトミックに論理和を取ってよい。即ち、gAPICは、IRRレジスタ70に対するアトミックなOR動作をサポートしてよい(ブロック158)。VMM 18は、新たなゲストに対してVMM管理のgAPIC状態エントリ92から他の状態を読み出してよく(ブロック160)、そしてその状態をgAPICへ書き込んでよい(ブロック162)。尚、ブロック160及び162は、ブロック150の後の任意の他の点で実行されてもよい。

【0109】

ブロック140~148は、概してgAPICからのゲストを非アクティブ化するための動作を代表してよく、一方ブロック150~162は、概してgAPIC内のゲストをアクティブ化するための動作を代表してよい。従って、図8に水平な破線によって示されるように、gAPIC内のゲストを非アクティブ化することのみをVMM 18が望む場合には、水平な破線よりも上方のブロックが実行されてよい。gAPIC内のゲストをアクティブ化することのみをVMM 18が望む場合には、水平な破線よりも下方のブロックが実行されてよい。

【0110】

次に図9を参照すると、実施形態に対するgAPIC状態エントリ90内のgAPIC状態の1つの例示的な配列170を表すブロック図が示されている。図9の実施形態においては、IRRの各ビットは異なるバイト内に記憶される。例えば、図9におけるIRR

10

20

30

40

50

ビット0、即ちIRR0は、メモリ内の一連の連続するバイトのバイト0内に記憶され、IRR1はバイト1内に記憶され、バイト255内に記憶されているIRR255まで同様に続く。図示される実施形態においては、IRRビットはバイトのビット0内に記憶されているが、任意のビット位置が用いられてよい。バイト内の他のビットは、図示される実施形態では無関係(don't cares)(DC)である。各ビットを別々のバイト(メモリアクセスの最小単位)内に記憶することによって、各ビットは、他のビットに影響を与えることなしに個別に書き込まれ得る。従って、ビットは書き込みを介してバイトへセットされてよく、これはアトミックな動作である。セットビットをバイト内のIRRビット位置へ書き込むこと及び他のバイトを更新しないことによって、IRRビットのアトミックORが結果であり得る。他の実施形態においては、アトミックORは他の方法において達成されてよく、またIRR状態のビットは他の方法において記憶されてよい。

10

【0111】

次に図10を参照すると、gAPIC状態エントリ90を配置する1つの実施形態のブロック図が示されている。図示される実施形態においては、デバイステーブル62及び割り込みリダイレクトテーブル64が示される他、gAPIC状態マッピングテーブル60の実施形態が示されている。実施形態においては、割り込みを送信した周辺機器のBDFがデバイステーブル62内への索引(index)として用いられ、そしてエントリは、BDFが割り当てられているゲストに対するゲストIDを含んでいてよい。またこの実施形態においては、エントリは割り込みリダイレクトテーブルポインタ(IRTP)を含んでおり、IRTPは割り込みリダイレクトテーブル64のベースを指す。割り込みリダイレクトテーブル64内への索引は、割り込みのための割り込み識別子である。割り込み識別子は、割り込みベクトルを含んでいてよく、また物理又は論理のいずれかである割り込みの受け渡しモード(Delmode)を含んでいてもよい。選択されたエントリは、新たなベクトル及び宛先ID(DestID)を含んでいてよい。割り込みリダイレクトテーブル64を用いない実施形態においては、周辺機器によって供給される割り込みベクトル及び宛先IDは、gAPIC状態マッピングテーブル60を直接的に索引付けるために用いられる。

20

【0112】

gAPIC状態マッピングテーブル60は、gAPIC状態マッピングテーブルベースアドレスを介してメモリ内に配置されてよい。ベースアドレスは、種々の実施形態において、全てのゲストに対して同じであってよく、ゲスト固有であってよく、あるいはデバイステーブル62内に記憶されていてよい。図10においては、ベースアドレスは、一連の階層テーブルの最も高いレベル(L3)を識別し、階層テーブルは、より低いレベルのテーブル(例えばL2及び、L2を指していないL3からのポインタによって表示される同様のテーブル)へのポインタを記憶していてよい。L2テーブルは更に低いレベルのテーブル(L1)へのポインタを記憶していてよく、そのテーブルは、gAPIC状態データ構造58内のgAPIC状態エントリ90へのポインタを記憶していてよい。他の実施形態は、図10に示される3レベルよりも多い又は少ないレベルを含め、階層内の任意の数のレベルを用いてよい。

30

【0113】

gAPIC状態マッピングテーブル60内の各レベルL3~L1内への索引は、デバイステーブル62からのゲストIDを連結させることから形成される値の一部分、周辺機器からの又は割り込みリダイレクトテーブル64からの割り込みベクトル、及び周辺機器からの又は割り込みリダイレクトテーブル64からの宛先IDであってよい。レベルL3~L1への索引は、連結された値の全てのビットを消費してよく、従って、ゲストID、ベクトル、及び宛先IDの各組み合わせは、それ自身の固有のポインタをgAPIC状態マッピングテーブル60内に有していてよい。しかし、幾つかのポインタは同一のgAPIC状態エントリ90を指し示すことがある(例えば、1つの実施形態においては、同一のgAPICの論理ID及び物理IDが同一のgAPIC状態90へのポインタを有していることがある)。

40

50

【 0 1 1 4 】

次に図 1 1 を参照すると、g A P I C 状態エントリ 9 0 を配置する他の実施形態のブロック図が示されている。図示される実施形態においては、デバイステーブル 6 2 及び割り込みリダイレクトテーブル 6 4 が示されている他、g A P I C 状態マッピングテーブル 6 0 の実施形態が示されている。実施形態においては、割り込みを送信した周辺機器の B D F がデバイステーブル 6 2 内への索引として用いられ、そしてエントリは、B D F が割り当てられているゲストに対するゲスト I D を含んでいてよい。またこの実施形態においては、エントリは割り込みリダイレクトテーブルポインタ (I R T P) を含んでおり、I R T P は割り込みリダイレクトテーブル 6 4 のベースを指す。デバイステーブル 6 2 は更に、g A P I C 状態マッピングテーブル 6 0 内のテーブルへの 1 つ以上のポインタを含んでいてよい。具体的には、ゲスト物理テーブルへのポインタ及びゲスト論理テーブルへの別のポインタが記憶されていてよい。ゲスト物理テーブルは、物理宛先 I D を g A P I C 状態エントリ 9 0 へマッピングしてよい。即ち、ゲスト物理テーブルは、宛先 I D によって索引付けられてよく、そしてポインタを g A P I C 状態エントリ 9 0 へ記憶してよい。同様に、ゲスト論理テーブルは、論理宛先 I D を g A P I C 状態エントリ 9 0 へマッピングしてよい。

10

【 0 1 1 5 】

割り込みリダイレクトテーブル 6 4 内への索引は、割り込みのための割り込み識別子である。割り込み識別子は、割り込みベクトルを含んでいてよく、また物理又は論理のいずれかである受け渡しモード (D e l m o d e) を含んでいてもよい。選択されたエントリは、新たなベクトル及び宛先 I D (D e s t I D) を含んでいてよい。割り込みリダイレクトテーブル 6 4 を用いない実施形態においては、周辺機器によって供給される割り込みベクトル及び宛先 I D は、g A P I C 状態マッピングテーブル 6 0 を直接的に索引付けるために用いられてよい。

20

【 0 1 1 6 】

次に図 1 2 を参照すると、g A P I C 状態エントリ 9 0 を配置する別の実施形態のブロック図が示されている。この実施形態では、g A P I C 状態マッピングテーブル 6 0 はない。図 1 0 ~ 1 1 の実施形態と同様に、割り込みを送信した周辺機器の B D F がデバイステーブル 6 2 内への索引として用いられ、そしてエントリは、B D F が割り当てられているゲストに対するゲスト I D を含んでいてよく、また随意的に割り込みリダイレクトテーブルポインタ (I R T P) を含んでおり、I R T P は割り込みリダイレクトテーブル 6 4 のベースを指す。デバイステーブル 6 2 は更に、g A P I C 状態データ構造 5 8 内のテーブルのベースへの少なくとも 1 つのポインタ (P t r) を含んでいてよい。図示される実施形態においては、テーブルは、ゲスト物理セクション 1 8 0 及びゲスト論理セクション 1 8 2 を含んでいる。セクション 1 8 0 及び 1 8 2 は、図面における明瞭さのために図 1 2 においてこれらの間の空白と共に図示されているが、セクション 1 8 0 及び 1 8 2 はメモリ内で隣接していてよい。即ち、ゲスト物理セクション 1 8 0 の先頭は、ゲスト論理セクション 1 8 2 の最後尾と隣接していてよい。デバイステーブル 6 2 は更に、ゲスト論理部分 1 8 2 の先頭を表示する論理限界 (logical limit) (L L i m) フィールドを含んでいてよい。他の実施形態においては、ゲスト物理部分 1 8 0 及びゲスト論理部分 1 8 2 は隣接していなくてよく、また個々のポインタがデバイステーブル 6 2 エントリ内に記憶されてそれぞれゲスト物理部分 1 8 0 及びゲスト論理部分 1 8 2 を表示してよい。

30

40

【 0 1 1 7 】

図 1 2 の実施形態においては、ゲスト物理部分 1 8 0 は割り込みベクトルによって索引付けられてよい (周辺機器からの、又は割り込みリダイレクトテーブル 6 4 からのいずれか) 。ゲスト物理部分 1 8 0 内の各エントリは、ゲスト物理マシン内でサポートされている宛先 I D (例えば、図 1 2 において 0 乃至 6 3 の数字が付されている 6 4 個までの宛先) に対応するビットベクトルを備えていてよい。物理割り込みに応答して、ゲスト割り込みマネージャ 3 8 は、宛先 I D に対応する割り込みベクトルに対するエントリ内のビットをセットするように構成されてよい。ブロードキャスト割り込みに対しては、ゲスト割り

50

込みマネージャ 38 は、仮想マシン内の vCPU の数までの割り込みベクトルに対応するエントリ内の各ビットをセットするように構成されてよい。

【0118】

ゲスト論理部分 182 は、論理 ID のクラスタ部分及びベクトルによって索引付けられてよい。クラスタ部分は索引の最も有意なビットであってよく、従って、ゲスト論理部分 182 は、各論理クラスタに対応する複数のクラスタ部分（図 12 におけるクラスタ 0 乃至クラスタ N）へと分割される。各クラスタの範囲内で複数のエントリが割り込みベクトルによって配列され、各エントリは論理 ID のベクトル部分に対応するビットベクトルを記憶している。図示される実施形態においては、16 までの宛先が 1 つのクラスタ内に含まれてよい（例えば論理 ID のビットベクトル部分は 16 ビットであってよい）。論理割り込みに応答して、ゲスト割り込みマネージャ 38 は、論理 ID のビットベクトル部分の、割り込みベクトルに対応するエントリのコンテンツとの論理和を取るように構成されてよい。

10

【0119】

従って、図 12 の実施形態は、ブロードキャスト物理割り込みの記録及び gAPIC 状態データ構造 58 に対する単一の更新内に多重宛先を有する論理割り込みの記録をサポートすることができる。gAPIC に対する gAPIC 状態エントリは、gAPIC の物理 ID に対応するゲスト物理部分 180 の列を備えていてよく、その列は gAPIC の論理 ID によって表示されるクラスタからの列との論理和を取られ、その列は gAPIC の論理 ID のビットベクトル部分内のセットビットによって識別される。gAPIC 内のゲストを非アクティブ化することに対応して gAPIC 状態データ構造 58 を更新することは、ゲストに対応する列の 1 つをゼロにすることと、他の列へ IRR を書き込むこととを含んでいてよい。

20

【0120】

図 13 は gAPIC 状態エントリ 90 を配置する別の実施形態である。図 13 の実施形態は、ゲスト物理部分 180 及びゲスト論理部分 182 の配列が異なること以外は図 12 の実施形態と同様である。各エントリは IRR に対応し、従って各エントリは各割り込みベクトルに対するビットを含む。ゲスト物理部分 180 は、割り込みの物理 ID によって索引付けられ、そしてゲスト論理部分 182 は、割り込みの論理 ID によって索引付けられる。割り込みベクトルに対応する IRR ビットは、割り込みの受け渡しモードに応じて論理部分 182 又は物理部分 180 のいずれかにおいてセットされる。ゲスト/vCPU に対する gAPIC 状態は、そのゲスト/vCPU に割り当てられた物理 ID に対応するゲスト物理部分 180 からの行と、そのゲスト/vCPU に割り当てられた論理 ID に対応するゲスト論理部分 182 からの行との論理和である。

30

【0121】

次に図 14 を参照すると、ホストハードウェア 20 の別の実施形態のブロック図が示されている。図示される実施形態には、図 2 における集積回路 66 と同様の 2 つの集積回路 66A ~ 66B が含まれる。従って、図示されるように、各集積回路は、集積回路 66A 内の 34A ~ 34D 及び集積回路 66B 内の 34E ~ 34G のような gAPIC を含んでいてよい。各集積回路 66A ~ 66B は、それぞれのゲスト割り込みマネージャ 28A ~ 28B 及び IOMMU（図 14 には図示せず）を含んでいてよい。集積回路 66A ~ 66B の少なくとも一方はメモリ 56A ~ 56B に結合され、また両集積回路 66A ~ 66B は随機的にメモリを含んでいてよい。集積回路 66A ~ 66B は、図示される実施形態においては、インタフェース回路 44C 及び 44D を介して結合される。他の実施形態においては、3 つ以上の集積回路 66A ~ 66B が設けられてよく、また種々の集積回路が任意の望ましい方法で相互接続されてよい。

40

【0122】

1 つの実施形態においては、各ゲスト割り込みマネージャ 28A ~ 28B が有効にされて、そして同一の集積回路内の gAPIC 34A ~ 34G で対象にされる割り込みメッセージを管理してよい。従って、ゲスト割り込みマネージャ 28A ~ 28B は、ゲスト割り

50

込み受け渡しに対して拡張性のある解決法を提供することができる。ゲスト割り込みマネージャ 28A ~ 28B によって用いられるデータ構造は、1つのメモリ(例えばメモリ 56A)内に記憶されていてよく、あるいは各ゲスト割り込みマネージャ 28A ~ 28B は、それ自身のデータ構造をそれ自身のメモリ 56A ~ 56B 内に有してよい。データ構造へのアクセスに対しては何らかの競合(contention)があり得るが、多くの場合には、周辺機器が特定のゲスト(集積回路 66A ~ 66B の一方におけるプロセッサ上で実行中の)に割り当てられ、従って実際の競合の量は比較的小さいであろう。

【0123】

別の実施形態においては、ゲスト割り込みマネージャ 28A ~ 28B の一方が有効にされて、そしてシステム内の g A P I C 34A ~ 34G に対するゲスト割り込み受け渡しを実行してよい。そのような実施形態は、集積回路 66A ~ 66B の間での相互接続を介して、より多くの割り込み関連のトラフィックを経験するであろうが、ゲスト割り込み管理に対する主要目的の概念的な簡潔さを提供することもできる。

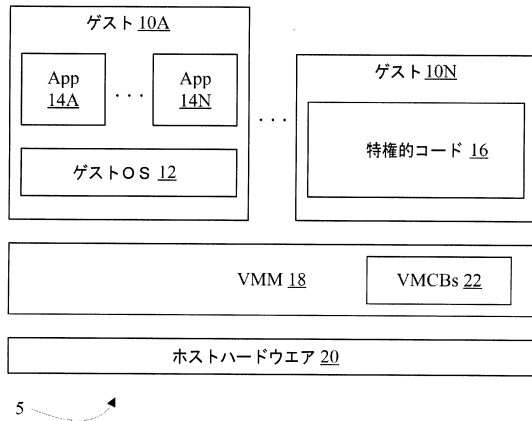
【0124】

次に図 15 を参照すると、コンピュータアクセス可能記憶媒体 200 のブロック図が示されている。一般的に言って、コンピュータアクセス可能記憶媒体は、命令及び/又はデータをコンピュータへ供給するための使用の間にコンピュータによってアクセス可能な任意の記憶媒体を含んでいてよい。コンピュータアクセス可能記憶媒体の例としては、磁気媒体又は光学媒体のような記憶媒体、例えばディスク(固定の又は取り外し可能な)、テープ、CD-ROM、又はDVD-ROM、CD-R、CD-RW、DVD-R、DVD-RW、及び/又はブルーレイ(Blu-Ray)ディスクを挙げることができる。記憶媒体の更なる例としては、揮発性の又は不揮発性のメモリ媒体、例えば、RAM(例えば同期ダイナミックRAM(SDRAM)、ラムバスDRAM(RDRAM)、スタティックRAM(SRAM)、等)、ROM、フラッシュメモリ、ユニバーサルシリアルバス(USB)インタフェースのような周辺機器インタフェース又は任意の他のインタフェースを介してアクセス可能な不揮発性メモリ(例えばフラッシュメモリ)、等を挙げることができる。記憶媒体は、微小電気機械的システム(MEMS)の他、ネットワーク及び/又はワイヤレスリンクのような通信媒体を介してアクセス可能な記憶媒体を含む。図 15 におけるコンピュータアクセス可能記憶媒体 200 はVMM 18 を記憶していてよく、VMM 18 は、図 8 のフローチャート及び/又はVMM 18 に割り当てられる任意の機能性を実装して

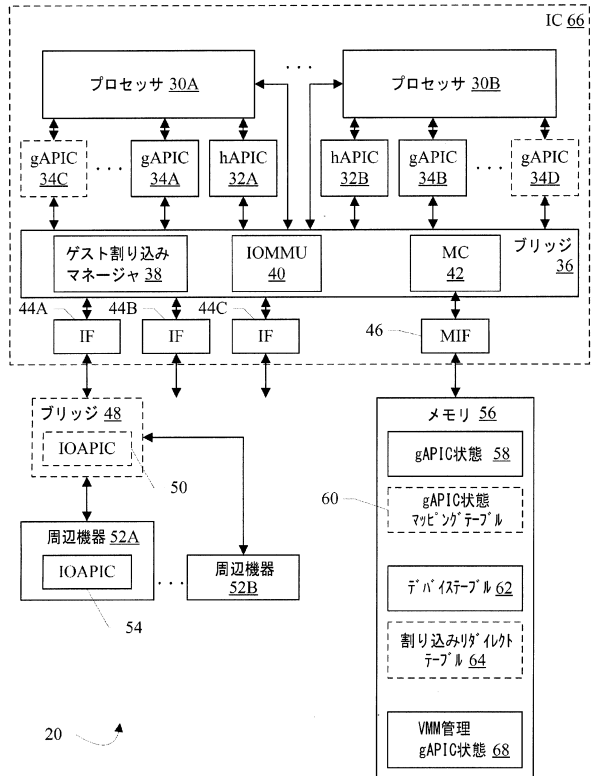
【0125】

上述の開示が完全に理解されるならば、多くの変形及び修正が当業者に明らかになるであろう。以下の特許請求の範囲は、全てのそのような変形及び修正を包含するものと解釈されることが意図されている。

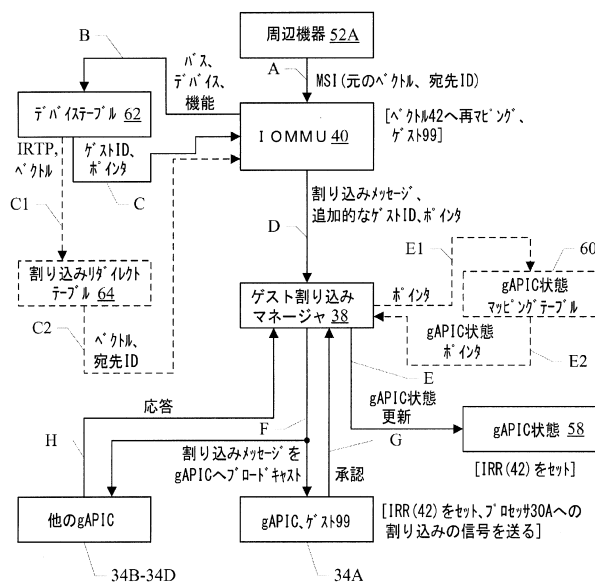
【図 1】



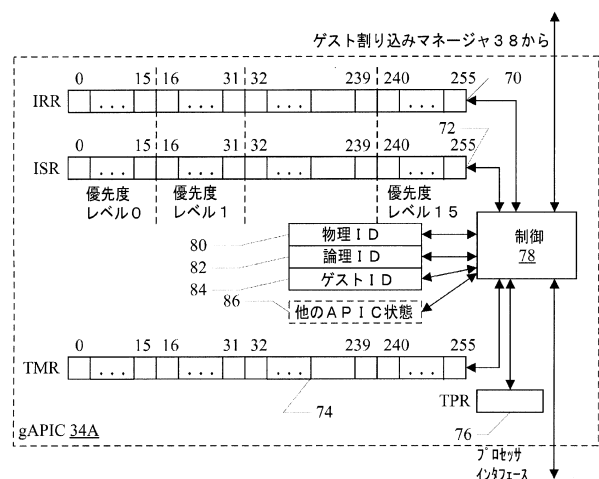
【図 2】



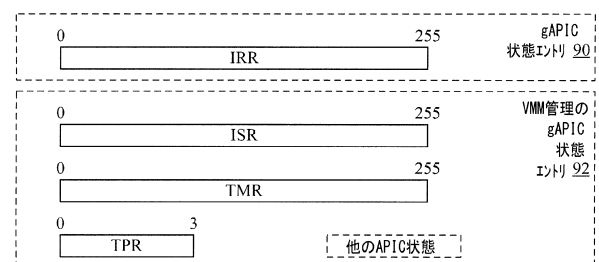
【図 3】



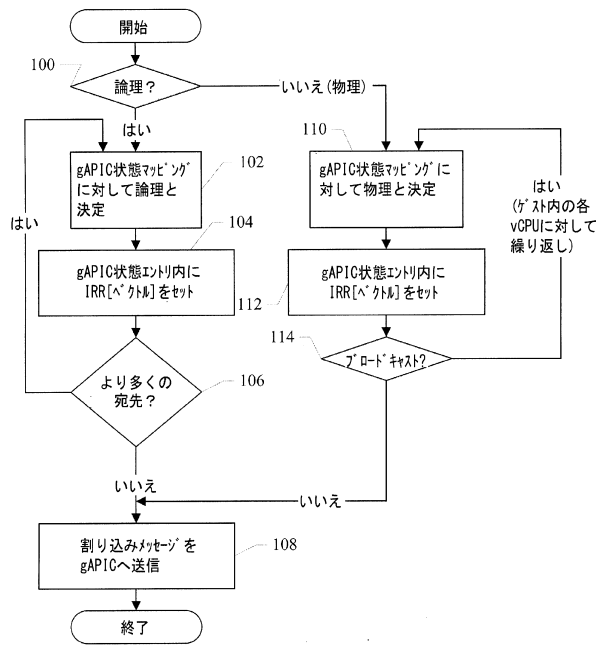
【図 4】



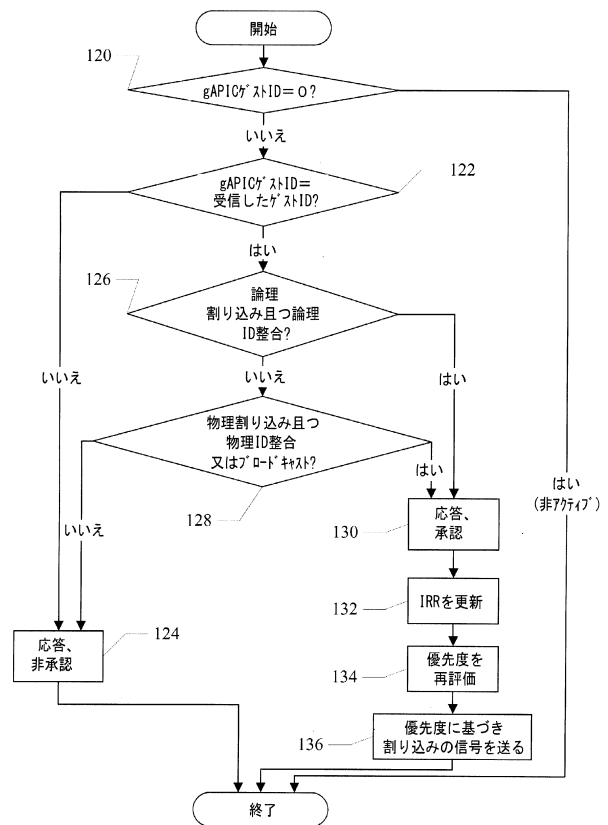
【図 5】



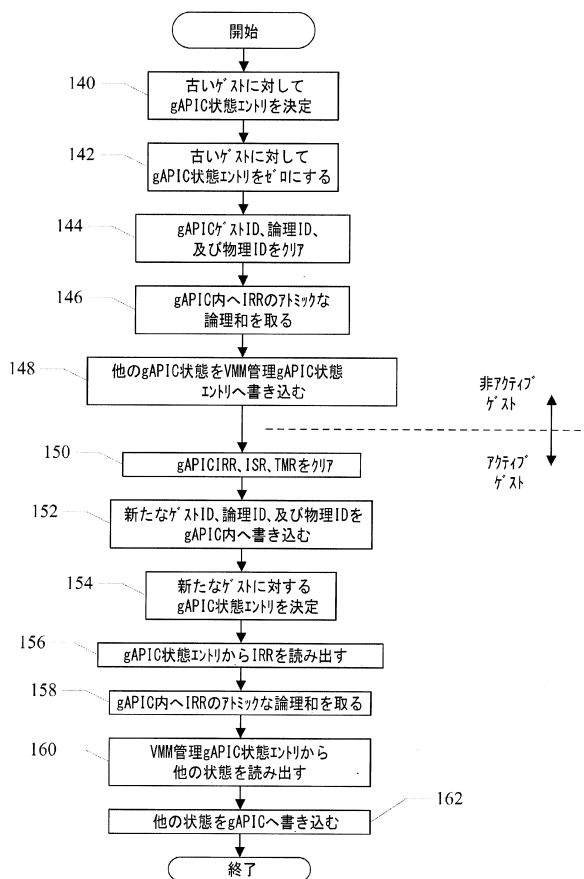
【図 6】



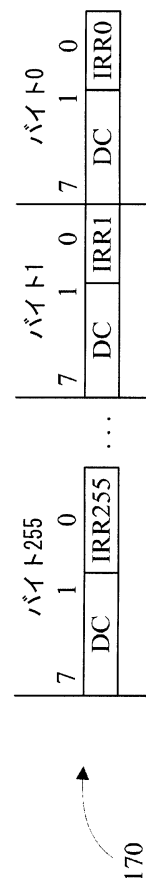
【図 7】



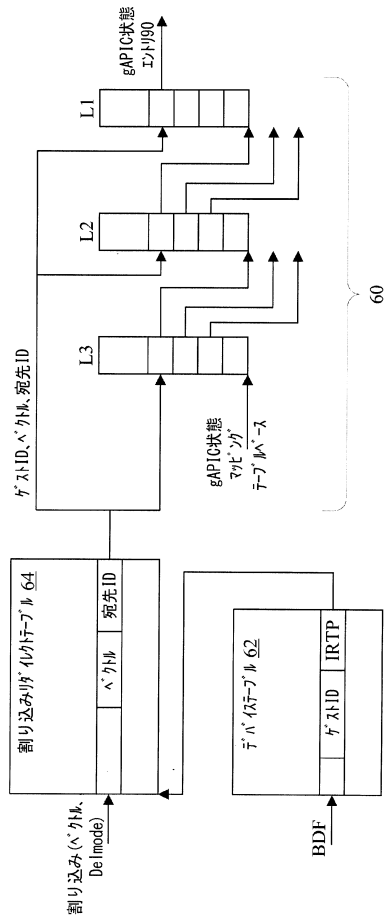
【図 8】



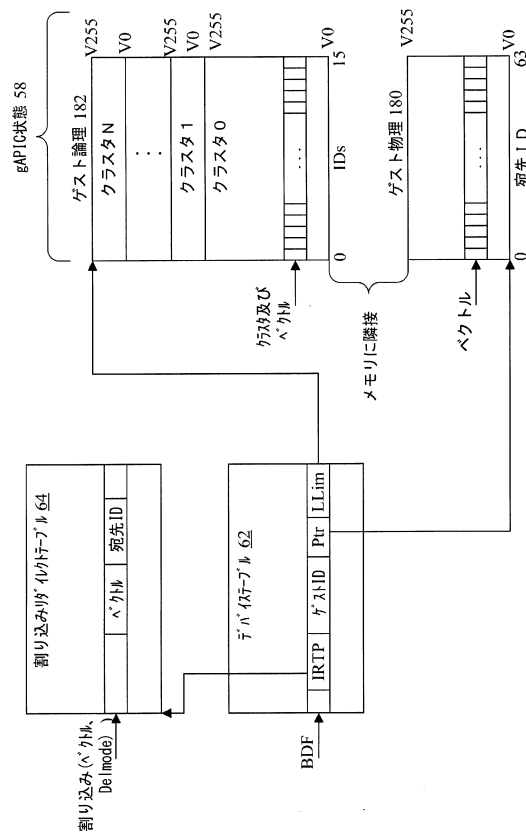
【図 9】



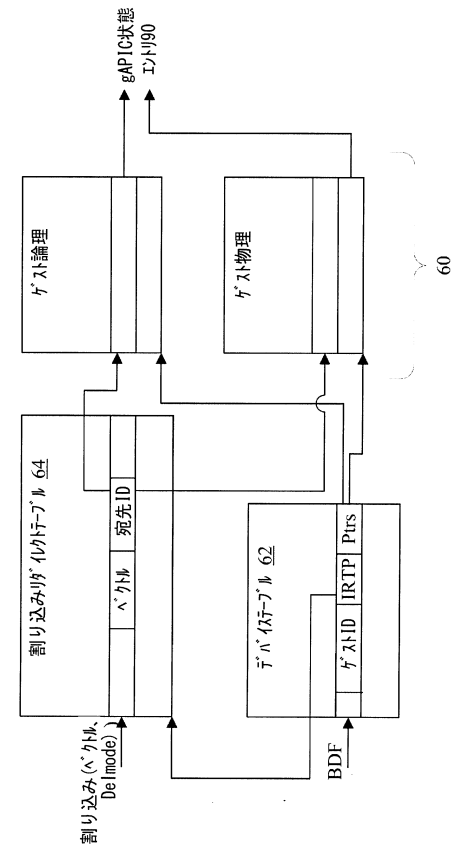
【 図 1 0 】



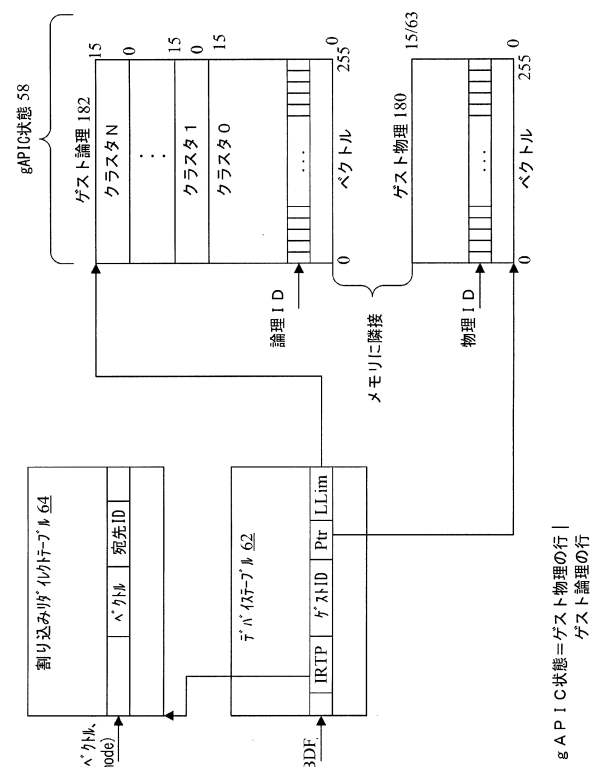
【 図 1 2 】



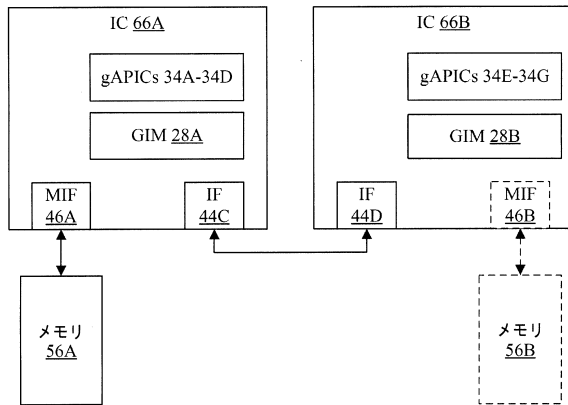
【 図 1 1 】



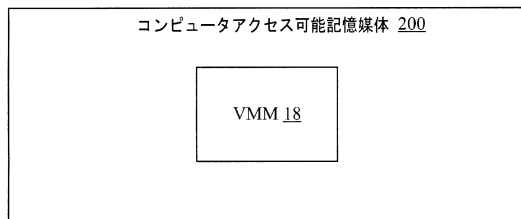
【 図 1 3 】



【図 14】



【図 15】



フロントページの続き

(31)優先権主張番号 12/611,607

(32)優先日 平成21年11月3日(2009.11.3)

(33)優先権主張国 米国(US)

(31)優先権主張番号 12/611,622

(32)優先日 平成21年11月3日(2009.11.3)

(33)優先権主張国 米国(US)

(74)代理人 100162156

弁理士 村雨 圭介

(72)発明者 ベンジャミン シー・セレブリン

アメリカ合衆国、カリフォルニア州 94086、サニーベイル、マホガニー レイン 771

(72)発明者 ドナルド ダブリュー・マコーレー

アメリカ合衆国、テキサス州 78734、レイクウェイ、エッジウォーター コーヴ 102

(72)発明者 ジョン エフ・ウィーデルヒルン

アメリカ合衆国、カリフォルニア州 95125、サン ホセ、ルビノ ドライブ #115 3110

(72)発明者 エリザベス エム・クーパー

アメリカ合衆国、カリフォルニア州 95033、ロス ガトス、ファーヴ リッジ ロード 18764

(72)発明者 マーク ディー・ヒュンメル

アメリカ合衆国 02038 マサチューセッツ州、フランクリン、スチュワート ストリート 68

審査官 大塚 俊範

(56)参考文献 米国特許第07209994(US, B1)

特開昭61-036850(JP, A)

特表2008-535099(JP, A)

米国特許出願公開第2004/0117532(US, A1)

(58)調査した分野(Int.Cl., DB名)

G06F 9/46 - 9/54