

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5507660号
(P5507660)

(45) 発行日 平成26年5月28日 (2014. 5. 28)

(24) 登録日 平成26年3月28日 (2014. 3. 28)

(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 9/50 (2006. 01)

G 0 6 F 9/46 4 6 2 A

請求項の数 20 (全 32 頁)

(21) 出願番号 特願2012-502133 (P2012-502133)
 (86) (22) 出願日 平成22年3月19日 (2010. 3. 19)
 (65) 公表番号 特表2012-521610 (P2012-521610A)
 (43) 公表日 平成24年9月13日 (2012. 9. 13)
 (86) 国際出願番号 PCT/US2010/028034
 (87) 国際公開番号 W02010/111149
 (87) 国際公開日 平成22年9月30日 (2010. 9. 30)
 審査請求日 平成25年2月5日 (2013. 2. 5)
 (31) 優先権主張番号 12/412, 272
 (32) 優先日 平成21年3月26日 (2009. 3. 26)
 (33) 優先権主張国 米国 (US)

(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) 代理人 100140109
 弁理士 小野 新次郎
 (74) 代理人 100075270
 弁理士 小林 泰
 (74) 代理人 100080137
 弁理士 千葉 昭男
 (74) 代理人 100096013
 弁理士 富田 博行
 (74) 代理人 100153028
 弁理士 上田 忠

最終頁に続く

(54) 【発明の名称】 仮想マシン用非一様仮想メモリーアーキテクチャー

(57) 【特許請求の範囲】

【請求項 1】

仮想マシンをインスタンス化するためのリクエストを受信するステップであって前記リクエストが、前記仮想マシンに関する特性を含んでいるものと、

前記特性に基づいて、前記仮想マシンに関する仮想非一様メモリ・アーキテクチャ (NUMA: Non-Uniform Memory Architecture) ノードトポロジを選択するステップであって前記仮想 NUMA ノードトポロジが、複数の仮想 NUMA ノードを含んでいるものと、

前記仮想マシンを計算機システム上でインスタンス化するステップであって前記仮想マシンが、前記複数の仮想 NUMA ノードを含んでいるものと、

前記複数の仮想 NUMA ノードの内の特定の仮想 NUMA ノードにおけるメモリー圧力に基づいて、前記特定の仮想 NUMA ノードに割り当てられるゲストメモリー容量を調整するステップと、を含む方法。

【請求項 2】

更に、

前記複数の仮想 NUMA ノードの内の第 2 の仮想 NUMA ノードにおけるメモリー圧力が所定の値よりも大きいことを決定するステップと、

前記第 2 の仮想 NUMA ノードを前記計算機システムの第 2 の NUMA ノードへ移動させるステップと、を含む請求項 1 記載の方法。

【請求項 3】

10

20

前記ゲストメモリー容量を調整するステップが更に、
ゲストメモリーの少なくとも1つのメモリーブロックを前記特定の仮想NUMAノードからデコミットするステップと、

前記ゲストメモリーの少なくとも1つのデコミットされたメモリーブロックを第2の仮想NUMAノードへコミットするステップと、を含むことを特徴とする請求項1記載の方法。

【請求項4】

前記ゲストメモリー容量を調整するステップが更に、
前記特定の仮想NUMAノードの少なくとも1つのゲストメモリーブロックがシステムメモリーと分離されていることを決定するステップと、

前記ゲストメモリーの少なくとも1つのメモリーブロックをシステムメモリーの少なくとも1つのメモリーブロック上へマッピングするステップと、を含むことを特徴とする請求項1記載の方法。

【請求項5】

更に、

前記特定の仮想NUMAノードを前記計算機システムの第1のNUMAノード上へマッピングするステップと、

前記特定の仮想NUMAノードを前記計算機システムの第2のNUMAノード上へ移動させるステップと、を含む請求項1記載の方法。

【請求項6】

更に、

仮想プロセッサを特定の前記仮想NUMAノードに追加するステップを含む請求項1記載の方法。

【請求項7】

更に、

前記仮想マシンの仮想プロセッサと、論理プロセッサに割り当てられた前記仮想プロセッサと、NUMAノードに割り当てられ前記論理プロセッサと、及び前記複数の仮想NUMAノードの内の1つに割り当てられた前記仮想プロセッサと、を実行するためのリクエストを受信するステップと、

前記論理プロセッサが前記仮想プロセッサを実行不可能なことを決定するステップと、

前記仮想プロセッサを実行する第2の論理プロセッサを選択するステップであって前記第2の論理プロセッサが、第2のNUMAノードから提供されているもの、を含む請求項1記載の方法。

【請求項8】

計算機システムであって、

仮想マシンを実行するための回路であって前記仮想マシンが、複数の仮想非一様メモリー・アーキテクチャ（NUMA：Non-Uniform Memory Architecture）ノードを含むトポロジーを有しており、前記仮想マシンの前記トポロジーが、前記計算機システムの物理トポロジーから独立して生成されるものと、

前記複数の仮想NUMAノードの内の仮想NUMAノードそれぞれにおけるメモリー圧力を決定するための回路と、

前記複数の仮想NUMAノードの内の仮想NUMAノードそれぞれにおける前記メモリー圧力に基づいて、前記複数の仮想NUMAノードの少なくとも1つに割り当てられるゲストメモリーを調整するための回路と、を含む計算機システム。

【請求項9】

更に、

前記仮想マシンを第2の計算機システムへ送信するための回路を含む請求項8記載の計算機システム。

【請求項10】

更に、

前記複数の仮想NUMAノードの内の第1の仮想NUMAノードを前記計算機システムの第1のNUMAノード上へマッピングするための回路と、

前記複数の仮想NUMAノードの内の第2の仮想NUMAノードを前記計算機システムの前記第1のNUMAノード上へマッピングするための回路と、を含む請求項8記載の計算機システム。

【請求項11】

更に、

前記複数の仮想NUMAノードの内の第1の仮想NUMAノードを前記計算機システムの第1のNUMAノード上へマッピングするための回路と、

10

前記複数の仮想NUMAノードの内の第2の仮想NUMAノードを前記計算機システムの第2のNUMAノード上へマッピングするための回路と、を含む請求項8記載の計算機システム。

【請求項12】

更に、

前記第2の仮想NUMAノードにおけるメモリー圧力が所定の値よりも大きいことを決定するための回路と、

前記第2の仮想NUMAノードを前記計算機システムの第2のNUMAノードへ移動させるための回路と、を含む請求項10記載の計算機システム。

【請求項13】

20

更に、

前記第2の仮想NUMAノードのメモリー圧力が所定の値よりも大きいことを決定するための回路と、

前記仮想マシンの前記第2の仮想NUMAノードを前記計算機システムの前記第1のNUMAノードへ移動させるための回路と、を含む請求項11記載の計算機システム。

【請求項14】

プロセッサ実行可能命令を含む計算機可読記憶媒体であって、

第1の仮想マシンを実行するための命令であって前記仮想マシンが、複数の仮想非一様メモリー・アーキテクチャ(NUMA: Non-Uniform Memory Architecture)ノードを含むトポロジを有しており、前記複数の仮想NUMAノードの内の仮想NUMAノードそれぞれが、仮想プロセッサ及びゲスト物理アドレスを含んでおり、前記仮想マシンの前記トポロジが、計算機システムの物理トポロジから独立して生成されるものと、

30

前記第1の仮想マシンのランタイム実行中に、増設の仮想プロセッサを前記複数の仮想NUMAノードの内の仮想NUMAノードに追加するための命令と、を含む前記計算機可読記憶媒体。

【請求項15】

更に、

前記複数の仮想NUMAノードの内の仮想NUMAノードそれぞれにおけるメモリー圧力を決定するための命令と、

前記複数の仮想NUMAノードの内の仮想NUMAノードそれぞれにおける前記メモリー圧力に基づいて、複数の仮想NUMAノードの内の仮想NUMAノードの少なくとも1つに割り当てられるゲストメモリーを調整するための命令と、を含む請求項14記載の計算機可読記憶媒体。

40

【請求項16】

更に、

前記複数の仮想NUMAノードの内の仮想NUMAノードから仮想プロセッサを移動するための命令を含む請求項14記載の計算機可読記憶媒体。

【請求項17】

更に、

前記複数の仮想NUMAノードに関するNUMA比をゲストオペレーティングシステム

50

へ報告するための命令を含む請求項 1 4 記載の計算機可読記憶媒体。

【請求項 1 8】

更に、

前記仮想マシンを第 1 の NUMA ノードから複数の NUMA ノードへ移動させるための命令を含む請求項 1 4 記載の計算機可読記憶媒体。

【請求項 1 9】

前記ゲストメモリーを調整するための前記命令が更に、

前記複数の仮想 NUMA ノードの内の第 1 の仮想 NUMA ノードの現在の前記メモリー圧力が目標の閾値よりも低いという決定に基づいて、前記第 1 の仮想 NUMA ノードからメモリーをデコミットするための命令を含むことを特徴とする請求項 1 5 記載の計算機可読記憶媒体。

10

【請求項 2 0】

前記ゲストメモリーを調整するための前記命令が更に、

ゲストオペレーティングシステムの現在のメモリー圧力が目標の閾値よりも大きいという決定に基づいて、前記複数の仮想 NUMA ノードの内の第 1 の仮想 NUMA ノードへメモリーをコミットするための命令を含むことを特徴とする請求項 1 5 記載の計算機可読記憶媒体。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

20

本発明は、仮想化技術に関し、具体的には、仮想マシン用のメモリーアーキテクチャーに関する。

【背景技術】

【0 0 0 2】

[0001]仮想化技術は、複数のパーティション間においてハードウェア資源を共有することを可能にし、各パーティションがゲストオペレーティングシステムをホスティングすることを可能にする。通常、仮想マシン技術は、サーバーを統合化し、それらの移植性を増大させるために利用され得る。仮想マシンがより大きくなりそれらの作業負荷が増大するにつれて、それらのある計算機システムから別の計算機システムへ容易に統合及び/又は移動させる機能は、より困難になっている。したがって、より大きな仮想マシンを統合化及び/又は移動させる能力を増大する技法が望まれている。

30

【発明の概要】

【発明が解決しようとする課題】

【0 0 0 3】

本発明の目的は、仮想マシン用の仮想 NUMA アーキテクチャーを達成し、仮想 NUMA ノードにおけるメモリーを調整するための技法を提供することである。

【課題を解決するための手段】

【0 0 0 4】

[0002]本開示の実施形態例は、一方法を記述している。この例において、本方法はこれらに限定しないが、仮想マシンをインスタンス化するためのリクエストを受信するステップであってリクエストが、仮想マシンに関する特性を含んでいるものと、前記特性に基づいて仮想マシンに対する仮想 NUMA ノードトポロジーを選択するステップであって仮想 NUMA ノードトポロジーが、複数の仮想 NUMA ノードを含んでいるものと、計算機システム上の仮想マシンをインスタンス化するステップであって仮想マシンが、複数の仮想 NUMA ノードを含んでいるものと、特定の複数の仮想 NUMA ノードにおけるメモリー圧力に基づいて調整するステップであってゲストメモリー容量が、特定の仮想 NUMA ノードに割り当てられているものと、を含む。前述のものに加えて別の態様が、本開示の一部を形成する請求項、図面、及びテキストに記述される。

40

【0 0 0 5】

[0003]本開示の実施形態例は、一方法を記述している。この例において、本方法はこれ

50

らに限定しないが、仮想マシンを実行するステップであって仮想マシンが、複数の仮想NUMAノードを含むトポロジを有しており、前記仮想マシンのトポロジが、計算機システムの物理トポロジから独立して生成されるものと、複数の仮想NUMAノードそれぞれにおけるメモリー圧力を決定するステップと、複数の仮想NUMAノードそれぞれにおけるメモリー圧力に基づいて調整するステップであってゲストメモリーが、複数の仮想NUMAノードの少なくとも1つに割り当てられているものと、を含む。前述のものに加えて別の態様が、本開示の一部を形成する請求項、図面、及びテキストに記述される。

【0006】

[0004]本開示の実施形態例は、一方法を記述している。この例において、本方法はこれらに限定しないが、第1の仮想マシンを実行するステップであって仮想マシンが、複数の仮想NUMAノードを含むトポロジを有しており、複数の仮想NUMAノードそれぞれが、仮想プロセッサ及びゲスト物理アドレスを含んでおり、仮想マシントポロジが、計算機システムの物理トポロジから独立して生成されるものと、増設の仮想プロセッサを複数の仮想NUMAノードに追加するステップと、を含む。前述のものに加えて別の態様が、本開示の一部を形成する請求項、図面、及びテキストに記述される。

【0007】

[0005]本開示の1つ以上の様々な態様は、本明細書に参照されている本開示の態様に作用するための回路及び/又はプログラムであって回路及び/又はプログラムが、システム設計者の設計選択に従って本明細書に参照されている態様に効果的になるように構成されたハードウェア、ソフトウェア、及び/又はファームウェアの実際の組み合わせのいずれかであり得るものを含むが、これらに制限しないことが当業者によって十分に理解されよう。

【0008】

[0006]前述のものは概要であって、かくして必然的に、詳細の簡素化、一般化、及び省略を含む。当業者は、概要が例示に過ぎず、任意の方法に制限することを意図しないことを十分に理解されよう。

【図面の簡単な説明】

【0009】

【図1】[0007] 本開示の態様が実装され得る計算機システムの例を表している。

【図2】[0008] 本開示の態様を実施するための動作環境を表している。

【図3】[0009] 本開示の態様を実施するための動作環境を表している。

【図4】[0010] 本開示の実施形態においてメモリーがどのように配置され得るか表している。

【図5】[0011] 本開示の態様を実施する動作環境の例を表している。

【図6】[0012] 本開示の態様を実施する動作環境の例を表している。

【図7】[0013] 本開示の態様を実施する動作環境の例を表している。

【図8】[0014] 本開示の態様を例示するブロック図の例を表している。

【図9】[0015] 本開示の態様を実施するための動作手順を表している。

【図10】[0016] 図9の動作手順の代替実施形態を表している。

【図11】[0017] 本開示の態様を実施するための動作手順を表している。

【図12】[0018] 図11の動作手順の代替実施形態を表している。

【図13】[0019] 図12の動作手順の代替実施形態を表している。

【図14】[0020] 図12の動作手順の代替実施形態を表している。

【図15】[0021] 本開示の態様を実施するための動作手順を表している。

【図16】[0022] 図15の動作手順の代替実施形態を表している。

【図17】[0023] 図16の動作手順の代替実施形態を表している。

【発明を実施するための形態】

【0010】

[0024]実施形態は1つ以上の計算機上で実行し得る。図1及び以下の論述は、本開示が実装され得る適切な計算環境の簡潔な概説を提供することを意図している。当業者は、図

10

20

30

40

50

1の計算機システムがいくつかの実施形態において、計算機システム(200)、(300)、(600)、及び(700)を達成可能なことを十分に理解できよう。これらの実施形態例において、計算機システムは、図1に記述されたコンポーネント及び本開示の態様を例示化するように構成される回路のいくつか又はすべてを含み得る。

【0011】

[0025]本開示を介し使用されている用語「回路」は、ハードウェア割り込みコントローラ、ハードドライブ、ネットワークアダプター、グラフィックスプロセッサ、ハードウェアベースのビデオ/オーディオコーデックのようなハードウェアコンポーネント、及びそのようなハードウェアを操作するために使用されるファームウェア/ソフトウェアを含み得る。同一又は別の実施形態において、用語「回路」は、ファームウェア又は一定の方法で設定される切り換えによって機能(単数又は複数)を実行するように構成されるマイクロプロセッサを含み得る。同一又は別の実施形態例において、用語「回路」は、1つ以上の論理プロセッサ、例えば1つ以上のコアであるマルチコア汎用演算処理装置を含み得る。この例において論理プロセッサ(単数又は複数)は、メモリー、例えば、RAM、ROM、ファームウェア、及び/又は仮想メモリーからロードされる機能(単数又は複数)を実行するように作動可能なロジックを具体化しているソフトウェア命令によって、構成され得る。回路がハードウェア及びソフトウェアの組み合わせを含む実施形態例において、実装者は、ロジックを具体化しているソースコードを書き出し得、その後、論理プロセッサによって処理され得る計算機読み出し可能コードにコンパイルされる。最先端技術がハードウェア、ソフトウェア、又はハードウェア/ソフトウェアの組み合わせの間でほとんど差異がない点へ発展していることを当業者は十分に理解し得るので、機能を達成するハードウェア対ソフトウェアの選択は、単なる設計選択に過ぎない。かくして、当業者は、ソフトウェアプロセスが同等のハードウェア構造に変換され得ることと、ハードウェア構造がそれ自体、同等のソフトウェア処理に変換され得ることと、を十分に理解し得るので、ハードウェア実装対ソフトウェア実装の選択は、設計選択の1つであって実装者に任せられている。

【0012】

[0026]ここで図1を参照すると、例示的な汎用計算システムが表されている。汎用計算システムは、従来の計算機(20)などを含み得、論理プロセッサ(21)、システムメモリー(22)、システムメモリーを含む様々なシステムコンポーネントを論理プロセッサ(21)と接続する及びシステムバス(23)を含んでいる。システムバス(23)は、いくつかのタイプのバス構造のいずれかであり得、メモリーバス、又はメモリーコントローラ、周辺機器用バス、様々なバスアーキテクチャのいずれかを使用するローカルバスを含んでいる。システムメモリーは、読み出し専用メモリー(ROM)(24)及びランダムアクセスメモリー(RAM)(25)を含み得る。始動中など、計算機(20)内部のエLEMENT間に情報を送信する支援をする基本的ルーチンを含んでいる基本入力/出力システム(BIOS)(26)は、ROM(24)にストアされている。計算機(20)は更に、(図示されていない)ハードディスクから読み込むか又はそれに書き出すためのハードディスクドライブ(27)、取り外し可能磁気ディスク(29)から読み出すか又は書き出すための磁気ディスクドライブ(28)、及びCD-ROM又はその他の光学式媒体のような取り外し可能光学式ディスク(31)から読み出すか又はそれに書き込むための光学式ディスクドライブ(30)を含み得る。ハードディスクドライブ(27)、磁気ディスクドライブ(28)、及び光ディスクドライブ(30)がそれぞれ、ハードディスクドライブインターフェース(32)、磁気ディスクドライブインターフェース(33)、及び光学式ドライブインターフェース(34)によってシステムバス(23)に関連付けられるように示されている。ドライブ及びそれらの関連する計算機可読記憶媒体は、計算機(20)に対する命令、データ構造、プログラムモジュール、及びその他のデータの不揮発性計算機可読記憶装置を提供する。本明細書に説明される例示的環境は、ハードディスク、取り外し可能磁気ディスク(29)、及び取り外し可能光学式ディスク(31)を使用しているが、磁気カセット、フラッシュメモリーカード、デジタルビデオ

ディスク、ベルヌーイカートリッジ、ランダムアクセスメモリー（ＲＡＭ）、読み出し専用メモリー（ＲＯＭ）などのような計算機によってアクセス可能なデータをストアし得る別のタイプの計算機可読記憶媒体もまた例示的動作環境において使用され得ることが当業者によって十分に理解される必要がある。通常、実施形態の中にはそのような計算機可読記憶媒体が、本開示の態様を具体化するプロセッサ実行可能命令をストアするために使用され得るものもいくつかある。

【 0 0 1 3 】

[0027]多くのプログラムモジュールは、ハードディスク、磁気ディスク（２９）、光ディスク（３１）、ＲＯＭ（２４）、又はＲＡＭ（２５）上にストアされ得、オペレーティングシステム（３５）、１つ以上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）、及びプログラムデータ（３８）を含んでいる。ユーザーは、キーボード（４０）及びポインティングデバイス（４２）のような入力装置を介し、コマンド及び情報を計算機（２０）へ入力し得る。その他の（図示されていない）入力装置は、マイクロフォン、ジョイスティック、ゲームパッド、衛星放送受信機、スキャナーなどを含み得る。これら及びその他の入力装置は多くの場合、システムバスに接続されるシリアルポートインターフェース（４６）を介し、論理プロセッサ（２１）に接続されるが、しかしながらパラレルポート、ゲームポート又は普遍的なシリアルバス（ＵＳＢ）のような別のインターフェースによって接続される。ディスプレイ（４７）又はその他のタイプの表示装置もビデオアダプター（４８）のようなインターフェースを介しシステムバス（２３）に接続され得る。ディスプレイ（４７）に加えて計算機は典型的に、スピーカー及びプリンターのような別の（図示されていない）周辺出力装置を含む。図１の例示的システムは、ホストアダプター（５５）、小型コンピューター用周辺機器インターフェース（ＳＣＳＩ）バス（５６）、及びＳＣＳＩバス（５６）に接続される外部記憶装置（６２）も含む。

【 0 0 1 4 】

[0028]計算機（２０）は、リモートコンピューター（４９）のような１つ以上のリモートコンピューターとの論理接続を利用するネットワーク環境において作動し得る。リモートコンピューター（４９）は、別の計算機、サーバー、ルーター、ネットワークＰＣ、ピア装置又はその他の一般的ネットワークノードであり得、典型的に、前述した計算機（２０）に関連するエレメントの多く又はすべてを含み得るが、メモリー記憶装置（５０）だけが図１に例示されている。図１に表された論理接続は、ローカルエリアネットワーク（ＬＡＮ）（５１）及び広域ネットワーク（ＷＡＮ）（５２）を含み得る。そのようなネットワーク環境は、オフィス、企業規模コンピューターネットワーク、イントラネット、及びインターネットにおいて一般的である。

【 0 0 1 5 】

[0029]ＬＡＮネットワーク環境において利用されるとき、計算機（２０）は、ネットワークインターフェース又はアダプター（５３）を介しＬＡＮ（５１）へ接続され得る。ＷＡＮネットワーク環境において利用されるとき、計算機（２０）は、典型的に、インターネットのような広域ネットワーク（５２）を介した通信を確立するためのモデム（５４）又はその他の手段を含み得る。内蔵又は外付けがあり得るモデム（５４）が、シリアルポートインターフェース（４６）を介しシステムバス（２３）へ接続され得る。ネットワーク環境において、計算機（２０）又はその一部に関連し表されたプログラムモジュールが、リモートメモリー記憶装置にストアされ得る。示されたネットワーク接続が例示的であって、計算機間において通信リンクを確立する別の手段が使用され得ることを十分に理解されよう。更に、本開示の多くの実施形態がコンピューター化されたシステムに対し特に適切であるように描かれているが、本書における開示をそのような実施形態に限定することは意図されていない。

【 0 0 1 6 】

[0030]ここで図２及び図３を参照すると、それらは計算機システムの高水準のブロック図を表している。図面によって示されているように、計算機システム（２００）は、例え

10

20

30

40

50

ば、記憶装置(208)、ハードドライブ、ネットワークインターフェースコントローラ(NIC)(210)、グラフィックカード(234)、少なくとも1つの論理プロセッサ(212)、ランダムアクセスメモリ(RAM)(214)のような物理ハードウェアデバイスを含み得る。計算機システム(200)は、図1の計算機(20)と同様のコンポーネントも含み得る。1つの論理プロセッサが例示されているが、別の実施形態において、計算機システム(200)は、複数の論理プロセッサ、例えば、プロセッサあたり複数の実行コア及び/又は複数の実行コアをそれぞれ有する複数のプロセッサ、を有し得る。図2の説明を続けると、表されているものは、当技術分野において、仮想マシンモニターとしても参照され得るハイパーバイザー(202)である。表された実施形態において、ハイパーバイザー(202)は、計算機システム(200)のハードウェアへのアクセスを制御し、仲裁するための実行可能命令を含む。概してハイパーバイザー(202)は、子パーティション1から子パーティションN(Nは1より大きい整数)のようなパーティションと呼ばれる実行環境を生成し得る。実施形態において、子パーティションは、ハイパーバイザー(202)によって支援される分離単位と考えられ得、すなわち、子パーティションそれぞれが、ハイパーバイザー(202)及び/又は親パーティションの制御下にある一連のハードウェア資源、例えば、メモリー、デバイス、論理プロセッササイクルなどへマッピングされ得る。実施形態において、ハイパーバイザー(202)は、スタンドアロンソフトウェア製品、オペレーティングシステムの一部、マザーボードのファームウェア内埋め込み、専用集積回路、又はその組み合わせであり得る。

【0017】

[0031]表された例において、計算機システム(200)は、オープンソースコミュニティにおいて、ドメイン0としても考えられ得る親パーティション(204)を含む。親パーティション(204)は、オープンソースコミュニティにおいて、バックエンドドライバとしても知られている仮想化サービスプロバイダー(228)(VSP)(複数)を使用することによって、子パーティション(1~N)を実行しているゲストオペレーティングシステムへリソースを提供するように構成され得る。このアーキテクチャー例において、親パーティション(204)は基本ハードウェアへのアクセスを開閉し得る。概してVSP(228)は、オープンソースコミュニティにおいて、フロントエンドドライバとしても知られている仮想化サービスクライアント(VSC)を介し、ハードウェア資源に対するインターフェースを多重化するために使用され得る。子パーティションそれぞれは、ゲストオペレーティングシステム(220~222)が管理し得、その上で実行するスレッドをスケジューリングし得る仮想プロセッサ(230~232)のような1つ以上の仮想プロセッサを含み得る。通常、仮想プロセッサ(230~232)は、実行可能命令であって、特定のアーキテクチャーを有する物理プロセッサ表現を提供する関連状態情報である。例えば、1つの仮想マシンは、インテルx86プロセッサの特性を有する仮想プロセッサを有し得るが、一方で別の仮想プロセッサは、パワーPCプロセッサの特性を有し得る。この例において、仮想プロセッサは、仮想プロセッサを達成する命令が論理プロセッサによって支援されるように、計算機システムの論理プロセッサへマッピングされ得る。かくして、これらの実施形態例において、複数の仮想プロセッサが、例えば、別の論理プロセッサがハイパーバイザー命令を実行している間、同時に実行し得る。一般的に言えば、図面によって例示されるように、パーティションにおける仮想プロセッサ、様々なVSC、及びメモリーの組み合わせが、仮想マシン(240)又は(242)のような仮想マシンであると考えられ得る。

【0018】

[0032]通常、ゲストオペレーティングシステム(220~222)は、例えば、マイクロソフト(登録商標)、アップル(登録商標)、オープンソースコミュニティから提供されるオペレーティングシステムのような任意のオペレーティングシステムを含み得る。ゲストオペレーティングシステムは、ユーザー/カーネル動作モードを含み得、スケジューラ、メモリーマネージャーなどを含み得るカーネルを有し得る。ゲストオペレーティングシステム(220~222)それぞれが、その上にストアされる電子商取引サーバー、

電子メールサーバーのようなアプリケーションを有し得るファイルシステムと、ゲストオペレーティングシステム自体と、を関連付けている。ゲストオペレーティングシステム(220~222)は、仮想プロセッサ(230~232)上で実行するスレッドをスケジューリングし得、そのようなアプリケーションのインスタンスが達成され得る。

【0019】

[0033]ここで図3を参照すると、それは使用され得る代替アーキテクチャーを例示している。図3は、図2のそれと同様のコンポーネントを表しているが、しかし、この実施形態例において、ハイパーバイザー(202)は仮想化サービスプロバイダー(228)及びデバイスドライバー(224)を含み得、親パーティション(204)は設定ユーティリティ(236)を含み得る。このアーキテクチャーにおいて、ハイパーバイザー(202)は、図2のハイパーバイザー(202)と同一か又は同様の機能を実行し得る。図3のハイパーバイザー(202)は、スタンドアロンのソフトウェア製品、オペレーティングシステムの一部、マザーボードのファームウェア内部に埋め込まれているか、又はハイパーバイザー(202)の一部が専用集積回路によって達成され得る。この例において、親パーティション(204)は、ハイパーバイザー(202)を構成するために使用され得る命令を有し得るが、しかし、ハードウェアアクセスリクエストは、親パーティション(204)へ渡される代わりにハイパーバイザー(202)によって処理され得る。

【0020】

[0034]ここで図4を参照すると、それは仮想マシンを含む実施形態において、メモリーがどのように配置され得るか例示している。例えば、計算機システム(200)のような計算機システムは、メモリーアドレスを有しているRAM(214)を有し得る。システム物理メモリーアドレスを仮想マシンへ報告する代わりに、ハイパーバイザー(202)は、システム物理アドレス、例えば、ゲスト物理アドレス(GPA)に関する別のアドレスをゲストオペレーティングシステムのメモリーマネージャーへ提示し得る。ゲストオペレーティングシステムがその後、ゲスト物理アドレスを操作し得、ハイパーバイザー(202)がGPA及びSPAによる関係を維持する。図面によって示されるように、実施形態において、GPA及びSPAは、メモリーブロックの中へ配置され得る。概してメモリーブロックは、1つ以上のメモリーページを含み得る。GPAとSPAとの間の関係は、「Enhanced Shadow Page Table Algorithms」と題する米国特許出願No. 11/128,665に記載されているようなシャドウページテーブルによって維持され得、その内容をすべて参照として本明細書に組み込む。作動中、ゲストオペレーティングシステムがブロック1のGPAにデータをストアしたとき、データは、実際にはシステム上のブロック6のような異なるSPAにストアされ得る。

【0021】

[0035]手短に述べると図5は、本開示の態様を実施するための動作環境を表している。例えば、多くの計算機システム(504~510)がデータセンター(500)に共に接続され得る(4つの計算機システムが表されているが、当業者はデータセンター(500)が、より多いか又はより少ない計算機システムを含み得ることを十分に理解できよう)。表されている計算機システムは、異なるトポロジを有し得、その上、それらは、異なる特性、例えば、異なるRAM数、異なるRAM速度、異なる論理プロセッサ数、及び/又は異なる速度を持つ論理プロセッサを有し得る。

【0022】

[0036]管理システム(502)は、図1の計算機システム(20)、及び/又は計算機システム(200)、(300)、(600)、又は(700)と同様のコンポーネントを有し得る。すなわち、実施形態において、管理システム(502)は、図6又は図7に関し後述される対象項目を含む計算機システムであり得る。

【0023】

[0037]図面の概要を続けると、図6は、完全対称型マルチプロセッシングトポロジ(SMP)又は「フラット」トポロジを有する計算機システム(600)を表している。通常、SMPは、単一の共有メモリーに接続される複数のプロセッサを含むコンピュー

ターアーキテクチャーである。この手続きにおいては、メモリーコントローラー（６０２）が、メモリーへのデータフロー及びメモリーからのデータフローを管理し得る。メモリーアクセスは、論理プロセッサ（２１２Ａ～Ｆ）それぞれに対し一様であり得、論理プロセッサそれぞれは、メモリー範囲全体、すなわち、システム物理アドレス（６２２～６３２）をアクセスし得る。このトポロジは、比較的少ない数のプロセッサを用いた計算機システムに対し十分に動作するが、しかしながら計算機システムは多くのプロセッサを含んでいて、すべてが共有メモリーバスへのアクセスを求めて競合し、システム性能が低下し得る。その上、計算機システムの複雑さが、著しく増加し、次々に１プロセッサあたりの価格をつり上げる。

【００２４】

10

[0038]手短に述べると計算機システム（６００）は、計算機（２００）又は（３００）と同一か又は同様のコンポーネントを含み得る。図面によって示されるように、計算機システム（６００）は、ＲＡＭ（２１４）へのアクセスを開閉するメモリーコントローラー（６０２）を介し連結された複数の論理プロセッサ（２１２Ａ～２１２Ｆ）を有し得る（６つの論理プロセッサが表されているが、計算機システムはより多いか又は少ないものを有し得る）。前述のものと同様に、論理プロセッサ（２１２Ａ～２１２Ｆ）それぞれは、異なる特性、例えば、クロック速度、キャッシュサイズなどを有し得る。この手続きにおいては、メモリーコントローラー（６０２）が、ＲＡＭ（２１４）へのデータフロー及びＲＡＭ（２１４）からのデータフローを管理し得る。

【００２５】

20

[0039]ハイパーバイザー（２０２）がインスタンス化され得、それが計算機システム（６００）のハードウェアを制御し得る。ハイパーバイザー（２０２）が１つ以上の仮想マシン（２４０～２４２）を管理し得、それぞれは、仮想ＮＵＭＡノード（６０６～６１２）のような仮想ＮＵＭＡノードを有し得る。仮想ＮＵＭＡノード（６０６～６１２）が使用され得、ゲストアプリケーション又はゲストオペレーティングシステム（２２０）及び（２２２）のようなゲストオペレーティングシステムへ仮想トポロジを報告することによって、仮想マシンのリソースを統合化し得る。図面によって示されるように、仮想ＮＵＭＡノード（６０６～６１２）それぞれは、１つ以上の仮想プロセッサ（２３０Ａ～Ｄ、２３２Ａ～Ｄ）、及びゲスト物理アドレス（６１４～６１６）及び（６１８～６２０）を有し得る。通常、ハイパーバイザー（２０２）は、１つ以上の論理プロセッサを有する仮想ＮＵＭＡノード（６０６～６１２）それぞれと、ＲＡＭ（２１４）からのシステム物理アドレスと、を支援し得る。すなわち、ハイパーバイザー（２０２）は、仮想プロセッサスレッドを実行するために使用され得る理想的なプロセッサとして１つ以上の論理プロセッサを設定し得る。

30

【００２６】

[0040]手短に述べると図７は、ＮＵＭＡノード（７０２～７０６）を含むトポロジを有する計算機システム（７００）を表している。ＮＵＭＡノードを有する計算機システムは、通常、より小さな計算機システムから作り上げられた計算機として考えられ得る。この例において、ＮＵＭＡノードそれぞれ（６０６～６１２）は、１つ以上の論理プロセッサ及びローカルメモリーを含み得る。ＮＵＭＡノードの内部メモリーは、ローカルメモリーであると考えられ、別のＮＵＭＡノードの内部メモリーは、ノード内部のプロセッサだけが同一のメモリーバスと接続され得るので、リモートメモリーであると考えられる。ＮＵＭＡノードは、キャッシュ貫性ドメイン相互接続によって相互に接続され、１つのＮＵＭＡノード内のプロセッサは、首尾一貫した方法で別のＮＵＭＡノード内のメモリーをアクセス可能にする。かくして、システム物理アドレス（６２２～６３２）は、プロセッサそれぞれに関し一様である。すなわち言い換えると、システム物理アドレス２０，０００は計算機システムにおいて、すべてのプロセッサに対し同一である。相違は、いくつかのプロセッサに関し、メモリーアドレス２０，０００が、例えば、それらのＮＵＭＡノード内部のローカルメモリーアドレスであって、別のプロセッサに対するメモリーアドレス２０，０００が、例えば、それらのＮＵＭＡノードの外部、リモートであ

40

50

ることである。通常、ローカルメモリーは、リモートメモリーよりも高速にアクセスされ得、ローカルアクセス時間対リモートアクセス時間の間の関係はNUMA比として参照される。NUMA比1対2は、特定のリモートシステム物理アドレスをアクセスするプロセッササイクル数が、ローカルシステム物理アドレスよりも2倍かかることを意味する。NUMAは、任意の1つのメモリーバス上のプロセッサ数を制限することによって、SMPシステムによってもたせられるボトルネックを軽減し、通常、同一の論理プロセッサ数を有するSMP計算機システムよりも高価にならない。

【0027】

[0041]計算機システム(700)は、計算機(200)又は(300)と同一か又は同様のコンポーネントを含み得る。図面によって示されるように、この動作環境において、計算機システム(700)は、相互接続(708)によって接続された3つのNUMAノード(702~706)を含む(けれども計算機は多いか又は少ないものを有し得る)。図面によって例示されるように、NUMAノードそれぞれの内部のプロセッサ数は可変であり得、ノードそれぞれはそれ自身のRAMを有し得る。

【0028】

[0042]図7と同様にハイパーバイザー(202)は、計算機システム(700)のハードウェアを制御し得る。ゲストオペレーティングシステム又はモノリスアプリケーションがブートしたとき、それが前述したものと同様の仮想マシン(240)及び(242)のトポロジを検出し得る。仮想NUMAノード(606~612)それぞれが、仮想プロセッサスレッドを実行するために使用され得る同一のNUMAノードから1つ以上の理想的なプロセッサ及びシステム物理アドレスを割り当てられ得る。

【0029】

[0043]計算機システム(600)及び(700)が2つの仮想マシン(240)及び(242)を含んでいるように表されているが、別の実施形態において、それらは、より多いか又はより少ない仮想マシンを実行し得る。その上、仮想マシンそれぞれが2つの仮想NUMAノードを有しているように表されているが、別の実施形態において、仮想マシンは、より多いか又はより少ない仮想NUMAノードを有し得る。仮想NUMAノードも2つの仮想プロセッサを有するように表されているが、別の実施形態において、仮想NUMAノードは、より多いか又はより少ない仮想プロセッサを有し得る。更に、仮想NUMAノードそれぞれは、別の仮想NUMAノードと異なるトポロジを有し得、例えば、ある仮想NUMAノードは4つの仮想プロセッサ及び8ギガバイトRAMを有し得るが、別の仮想NUMAノードは2つの仮想プロセッサ及び4ギガバイトRAMとを有し得る。

【0030】

[0044]図8は、本開示の態様において使用され得る環境のブロック図を表している。図面によって示されるように、ダイナミックメモリー仮想化サービスプロバイダー(DMVSP)(802)として知られ得る、仮想マシンに割り当てられるメモリーを管理するコンポーネントが例示されていて、仮想NUMAノードにアクセス可能なメモリー容量を調整するために使用され得る。図面によって示されるように、DMVSP(802)は、仮想化サービスクライアント、すなわち、ダイナミックメモリー仮想化サービスクライアント(DMVSC)(804)及び/又は(806)として知られ得る1つ以上のブルーニングドライバーに関連付けられ得る(仮想NUMAノードあたり1つのDMVSCが表されているが、別の実施形態において、1パーティションあたり1つのDMVSCが使用され得る)。概してDMVSC(804)及び/又は(806)は、DMVSP(802)によって使用され得る仮想NUMAノードのメモリーを調整するための情報を提供し得、DMVSCそれぞれも、それが関連付けられる仮想NUMAノードからメモリーをコミット及びデコミットする支援をし得る。DMVSC(804)、(806)、及びDMVSP(802)は、「Partition Bus」と題する米国特許出願No. 11/128,647に記載されている仮想化バスを介し通信し得、その内容全体をすべて参照として組み込む。その上、更にDMVSC及びDMVSPの態様は「Dynamic Vi

「Virtual Machine Memory Management」と題する米国特許出願No. 12/345,469に記載されていて、その内容全体をすべて参照として組み込む。

【0031】

[0045]図8の説明を続けると、本システムは、ワーカプロセス(812)を含み得、子パーティション(単数又は複数)を管理可能に表されている。ワーカプロセス(812)は、メモリーを子パーティションに割り当て得る仮想基盤ドライバー(VID)(810)と連動して作動し得る。例えば、VID(810)は、ゲスト物理アドレスとシステム物理アドレスとの間の関係確立及び解除し得る。図8は、メモリーマネージャー(808)を含み得るゲストオペレーティングシステム(220)のような、ゲストオペレーティングシステムを含み得るパーティションも表している。通常、メモリーマネージャー(808)は、それらのリクエスト時、メモリーをアプリケーションに割り当て、それがもはやアプリケーションによって必要とされないとき、メモリーを解放し得る。

10

【0032】

[0046]以下は、プロセスの実装を表している一連の流れ図である。理解を容易にするため、流れ図は、最初の流れ図が「大画像」の観点を介する実装を表し、その後の流れ図が更なる追加及び/又は詳細を提供するように統合化されている。更に当業者は、点線によって表された動作手順が、任意に考えられることを十分に理解されよう。

【0033】

[0047]ここで図9に移ると、それは動作(900~910)を含む本開示の態様を実施するための動作手順を表している。動作手順は、動作(900)で開始し、動作(902)は、仮想マシンをインスタンス化するためのリクエストを受信するステップであってリクエストが、仮想マシンに関する特性を含んでいるもの、を例示している。例えば、図6又は図7を参照すると、ハイパーバイザー(202)は、仮想マシン(240)などの仮想マシンを作成するためのリクエストを受信し得る例えば、リクエストは、管理システム(502)の図2又は図3の親パーティション(204)などから受信され得る、リクエストは、新しい仮想マシンに対するものであり得るか又はそれは、前にセーブされた仮想マシンをインスタンス化するためのリクエストであり得る。仮想マシン(240)が新しい仮想マシンであるとき、仮想マシンの特性、例えば、仮想マシンに割り当てられるRAM数、仮想プロセッサ数、又は仮想マシンが有する必要があるI/O装置のタイプ、が例えば管理者によって設定され得る。

20

30

【0034】

[0048]図9の説明を続けると、動作(904)は、特性に基づいて仮想マシン用の仮想NUMAノードトポロジ選択をするステップであって仮想NUMAノードトポロジが、複数の仮想NUMAノードを含んでいるもの、を示している。例えば、親パーティション(204)におけるプロセス(及び/又はハイパーバイザー(202))は、受信された特性に基づいて仮想マシン(240)に関するトポロジを決定し得る。例えば、親パーティション(204)は、仮想NUMAノード(606)などの仮想NUMAノードに関するデフォルトサイズを識別する情報を含み得る。親パーティション(204)におけるプロセスは、仮想マシン(240)に対する仮想NUMA数を決定するためのデフォルトサイズ及び所望の特性を記述している情報を使用し得る。特定の例において、所望される特性は、10ギガバイトのRAMを有する6つのプロセッサ仮想マシンであり得る。仮想NUMAノードのデフォルトサイズが2つの仮想プロセッサ及び4ギガバイトのRAMを含む場合、管理システム(502)は、仮想マシン(240)が3つの仮想NUMAノードを含むことを示す構成ファイルを生成し得る。

40

【0035】

[0049]実施形態において、デフォルトの仮想NUMAノードサイズは、管理者によるか又は管理システム(502)によって設定され得る。図5に移ると、管理システム(502)は、データセンター(500)における計算機システム(504~510)の物理トポロジを識別する情報、例えば、計算機システム(504~510)それぞれが(もし

50

あれば)いくつかのNUMAノードを有しているか、RAM速度、計算機システム(504~510)それぞれがどのくらいの容量RAMを有しているか、RAMがどのように配列されているか、プロセッサ速度、それぞれのプロセッサがいくつかのコアを有しているかなどを識別する情報、を取得し得る1つ以上のプログラムを実行し得る。

【0036】

[0050]通常、仮想NUMAノードサイズが、データセンター(500)の仮想マシンの動作に影響する。例えば、仮想NUMAノードサイズは、例えば、メモリー及び/又はプロセッサにおいて増加するにつれて、仮想NUMAノードの移植性が減少する。すなわち言い換えると、大きな仮想NUMAノードが、仮想マシンを移動させることがより困難とし得る。このことは、仮想NUMAノードが、仮想NUMAノードを達成するために「フラット」リソースを有するNUMAノード又は計算機システムのどちらか一方に割り当てられる必要があるため生じる。例えば、仮想NUMAノードあまりに大きな、例えば、それが非常に大きなRAMか又はあまりに多くの仮想プロセッサを有する場合、データセンター(500)において、より小さいNUMAノードに適合させることが不可能であって、かくして、仮想マシンを移動させる能力を制限する。その上、より大きな仮想NUMAノードが単に、より小さな複数のNUMAノードに割り当てられた場合、仮想マシンノードの性能は、ローカルメモリアクセス時間とリモートメモリアクセス時間との間に存在する相違による理由によって減少する。

【0037】

[0051]他方では、仮想NUMAノードのサイズが減少するにつれて、ゲストオペレーティングシステムの性能は悪影響を及ぼされ得る。この非効率性は、ゲストオペレーティングシステムがアプリケーションを分離しようとし、それが単一の仮想NUMAノードに対する自身の実行であるために生じ得る。ゲストオペレーティングシステムは、この場合、抑制され、性能が低下する。

【0038】

[0052]したがって、実施形態において、管理システム(502)は、データセンター(500)に対する最適な仮想NUMAノードサイズを決定することによって、移植性と効率性との間のバランスを決めることができる。例えば、実施形態において、管理システム(502)の論理プロセッサがプログラムを実行し得、データセンターにおける平均NUMAノードサイズ、例えば平均論理プロセッサ数、平均RAM数などを決定し得、システムにおいて仮想NUMAノードのサイズを平均NUMAノードと同一か又はそれよりも小さくするように設定し得る。別の実施形態において、プログラムは、データセンター(500)における仮想NUMAノードサイズを最小NUMAノードよりもわずかに小さく設定するように構成され得る。実施形態において、仮想NUMAノードサイズが平均サイズ又は最小サイズよりもわずかに小さいように設定され得、計算機システムが大量にコミットされた場合、2つ以上の仮想NUMAノードが単一のNUMAノードへ割り当てられ得る。特定の例において、最小のNUMAノードが4つの論理プロセッサ及び8ギガバイトRAMを有する場合、仮想NUMAノードサイズは、例えば、2つの仮想プロセッサ及び4ギガバイトRAMに設定され得る。

【0039】

[0053]動作(906)は、計算機システム上に仮想マシンをインスタンス化するステップであって仮想マシンが、複数の仮想NUMAノードを含んでいるもの、を示している。実施形態において、ハイパーバイザー(202)が論理プロセッサによって実行され得、複数の仮想NUMAノードを有する仮想マシンがインスタンス化され得る。

そして図6及び/又は図7を参照すると、例えば、仮想NUMAノード(606~608)を有する仮想マシン(240)は、計算機システム(600)又は(700)によって達成され得る。すなわち、VID(810)は、RAMからのシステム物理アドレスを有する仮想マシン(240)のゲスト物理アドレスと、1つ以上の論理プロセッサを有する仮想プロセッサと、を支援し得る。例えば、ゲスト物理アドレスブロック(614)がシステム物理アドレスブロック(622)を用いて支援され得、ゲスト物理アドレスブ

10

20

30

40

50

ロック(616)がシステム物理アドレスブロック(624)によって支援され得る。ハイパーバイザースレッドがその後、仮想プロセッサを支援している論理プロセッサ上にスケジューリングされ得、仮想プロセッサの命令指標が実行され得る。図6及び図7によって示されるように、仮想マシンそれぞれのトポロジが、基本ハードウェアトポロジから独立して生成され得る。すなわち、仮想マシンのトポロジそれぞれは、それを達成する計算機システムの基本物理トポロジから分離されている。

【0040】

[0054]実施形態において、仮想マシンBIOS又はブートファームウェアは、仮想マシンのトポロジを記述し得、例えば、それは仮想NUMAノード、任意の仮想NUMAノードサイズ、及び仮想NUMAノードに関するNUMA比をモノリスアプリケーションのゲストオペレーティングシステムに対し有しているか否か説明し得る。データ構造は処理され得、ゲストOS(220)又はアプリケーション及びそれが、仮想NUMAノードの存在利点を取り入れるために、OS又はアプリケーションによって使用され得る。例えば、ゲストオペレーティングシステム(220)は、アプリケーションの実行がローカルのみであるように、NUMA非認識アプリケーションスレッドを仮想NUMAノードと一体化しようとし得る。別の例において、SQLサーバーのようなデータベース管理プログラムは、ロックを仮想NUMAノードへローカルに割り当て得、データベースが、仮想NUMAノード全域の読み出し/書き出しリクエストを分割し得る。更に別の例において、ゲストオペレーティングシステム(220)は、仮想マシンに仮想NUMAノードそれぞれに関するページプールを生成し得る。

【0041】

[0055]図9の説明を続けると、動作(908)は、特定の仮想NUMAノードに割り当てられるゲストメモリー容量を、特定の複数の仮想NUMAノードにおけるメモリー圧力に基づいて調整するステップを示している。例えば、論理プロセッサ、例えば、図6又は図7の論理プロセス(212A)は、DMVSP(802)の命令指標を実行し得、仮想NUMAノード(606)のような仮想NUMAノードが、利用可能なゲスト物理アドレス量を調整し得る。すなわち、DMVSP(802)が実行され得、メモリーが仮想NUMAノードが経験している圧力に基づいてコミット又はデコミットされ得る。

【0042】

[0056]実施形態において、仮想NUMAノード(606~608)それぞれが、利用可能なメモリー容量によって、ゲストオペレーティングシステム(220)の性能がどのくらい影響されるかをメモリー圧力が識別し得る。この情報は、DMVSC(804)及び/又は(806)のような例えば、1つ以上のDMVSCによってゲストオペレーティングシステム(220)のランタイム中に算出され得、DMVSP(802)へ送信され得る。例えば、メモリー圧力は、仮想NUMAノードにおける異なるメモリー圧力レベルを識別し得る一連の値によって提示され得る。仮想NUMAノードにおけるリソースが、より圧迫されるにつれて、すなわち、仮想NUMAノード上の現在の作業負荷を効率的に実行するために要求されるメモリー容量が増加するにつれて、DMVSC(804)が値を改定し得、この情報をDMVSP(802)へ伝達し得る。

【0043】

[0057]実施形態において、メモリー圧力情報が、ゲストオペレーティングシステム(220)から受信された情報からDMVSC(804)によって算出され得る。例えば、DMVSC(804)が、仮想NUMAノード(606)に関するオペレーティングシステムのページング情報をメモリーマネージャー(808)から受信するように構成され得る。ゲストオペレーティングシステムのページングレートが、メモリーマネージャー(808)及びキャッシュマネージャによって公開される2つのカウンター、すなわち、ページングレート及びキャッシュ回転速度を介しモニターされ得る。

【0044】

[0058]同一又は別の実施形態において、DMVSC(804)が、メモリーマネージャー(808)からノード(606)を仮想NUMAに関連付けられた物理メモリー通知を

受信し、この情報を使用し得、仮想NUMAノード(606)のメモリー圧力を計算し得る。例えば、メモリーマネージャー(808)は、仮想NUMAノード(606)に関連付けられたゲストオペレーティングシステム(220)における活動に基づいて、高いメモリーの通知及び低いメモリーの通知を出力し得る。メモリーマネージャー(808)は、低メモリー閾値(LMT)及び高メモリー閾値(HMT)に基づいてこれらの通知を起動し得る。特定の実施形態例において、低メモリーリソース通知イベントを信号送出する有効メモリーのデフォルトレベルは、およそ4GBあたり約32MB、最大64MBであり得る。例えば、高メモリーリソースの通知イベントを信号送出するデフォルトレベルは、デフォルトの低メモリーの値の3倍であり得る。2つの間の中間メモリーの有効レベルは、高メモリーの閾値レベルと低メモリーの閾値レベルとの間隔を分割することによって決定され得る。当業者は、これらの値が例示的であり得、変更が本開示の趣旨から逸脱せずに実施され得ることを十分に理解できよう。

10

【0045】

[0059]これらの通知が、別のものと一緒にDMVSC(804)によって使用され得、仮想NUMAノード(606)のメモリー圧力を算出し得る。例えば、レベルそれぞれは、例えば0~4の値に関連付けられ得、別の任意のパフォーマンスカウンターを考慮に入れた場合、それらの値も関連付けられ得る。パフォーマンスカウンターそれぞれに関する値はその後、仮想NUMAノード(606)の現在のメモリー圧力を算出し得るために使用され得る。特定の例において、メモリー圧力が、より大きい又は小さいパフォーマンスカウンター値を取ることによって算出され得る。別の例において、パフォーマンスカウンターの平均値が、メモリー圧力として使用され得る。更に別の実施形態において、より洗練されたアルゴリズムが使用され得、計算において、以前のパフォーマンスカウンターを考慮に入れ、相対重量に作用するスカラーをパフォーマンスカウンターそれぞれに割り当て、メモリー圧力を算出し得る。

20

【0046】

[0060]メモリーをコミットするための決定が実行されたとき、DMVSP(802)は、様々な技法を使用し得、その1つがホット(hot)追加動作である。例えば、オペレーティングシステムの中には、物理メモリー範囲が、システムのリブートを必要とせずに実行中のオペレーティングシステムに追加されることを可能にするホット追加を支援するものもいくつかある。すなわち、メモリーマネージャー(808)は、実行中のシステムへメモリーの動的追加を支援するように構成され得る。ホット追加の実施形態において、DMVSC(804)が、メモリーマネージャー(808)のホット追加インターフェースをアクセスするように構成され得、DMVSC(804)が、ホット追加されたGPA及びそれらがその仮想NUMAノードに関連付けられていることを記述しているメッセージをゲストオペレーティングシステム(220)へ送信し得る。メモリーマネージャー(808)がその後、ゲストオペレーティングシステム(220)が利用可能な新しいメモリー、仮想NUMAノード(606)上で実行しているドライバー、アプリケーション、又は別の任意のプロセスを生成し得る。例えば、DMVSC(804)は、V1D(810)がGPAとSPAとの間に関係を生成した後、DMVSP(802)からホット追加メモリーアドレスを受信し得る。

30

40

【0047】

[0061]同様に、ホット移動動作が使用され得、仮想NUMAノード(606)のような仮想NUMAノードからメモリーアドレスをデコミットし得る。例えば、DMVSC(804)は、メモリーがホット移動されたことを示すメッセージをゲストオペレーティングシステム(220)へ送信し得る。DMVSC(804)は、メモリーマネージャー(808)が、仮想NUMAノード(606)からGPAブロックを移動用に提供するようにリクエストし得る。この例において、DMVSC(804)はその後、メモリーマネージャー(808)の移動APIを呼び出し得、ゲストオペレーティングシステム(220)からGPAを移動し得る。ホット移動が使用される実施形態において、移動されるメモリーは、ゲストの現在のコミットに対し数えられず、メモリーマネージャー(808)は、

50

オペレーティングシステムによってメモリーを移動するために使用される同様の技法を使用し、内部カウンターを調整し得、マザーボードから物理的に移動される。

【 0 0 4 8 】

[0062]別の実施形態において、メモリーは、バルーニング技法を使用することによって、仮想NUMAノードヘデコミットされ得る。すなわち、メモリーは、仮想NUMA(606)ノードにおけるゲスト物理アドレスを、それらを支援している物理アドレスから分離することによってデコミットされ得る。例えば、論理プロセッサ(212B)は、DMVSC(804)の命令指標を実行し得、メモリーマネージャー(808)が、DMVSC(804)によって、例えば1つ以上のメモリーブロックを使用するための一定のメモリー容量を予約するリクエストメッセージをメモリーマネージャー(808)へ送信し得る。メモリーマネージャー(808)は、DMVSC(804)及びDMVSC(804)の内部で排他的に使用するためにメモリーをロックし得、メモリーのGPAをDMVSP(802)へ送信し得る。この例において、DMVSP(802)は、GPAをVID(810)へ送信し得、VID(810)は、これらのGPAに関するエントリーをシャドウページテーブル内のSPAへ移動し得る。この例において、メモリーマネージャー(808)は、GPAがまだ利用可能であるが、しかしながら実際には、GPAがもはやシステム物理アドレスによって支援されないことを識別する情報を含み得る。この例において、メモリーマネージャー(808)は、ロックされたGPAを使用しないでそれらを支援するSPAが再び割り当てられ得る。

10

【 0 0 4 9 】

20

[0063]分離されたゲスト物理アドレスは、物理アドレスに再度関連付けられ得る。この例において、メモリーページをコミットするためのリクエストがVID(810)によって受信され得、VID(810)は、リクエストを満たしSPAを取得し得、アドレスの範囲をDMVSP(802)へ送信し得る。実施形態において、VID(810)は、システム効率性を増大するために隣接するSPAの範囲を取得するように構成され得る。この例において、VID(810)は、仮想NUMAノード(606)に関連付けられるDMVSC(804)によって、排他的に使用するためにロックされるGPAを有するゲストオペレーティングシステム(220)を決定し得る。VID(810)は、ロックされたGPAとSPAとの間に関係を生成し得、メッセージをDMVSP(802)へ送信し得る。DMVSP(802)はその後、メッセージをDMVSC(804)へ送信し得、DMVSC(804)は、GPAがアンロックされ、仮想NUMAノード(606)に関連付けられたメモリーマネージャー(808)のメモリーブールへ返却され得ることを示すメモリーマネージャー(808)へメッセージを送信し得る。

30

【 0 0 5 0 】

[0064]実施形態において、VID(810)は、ホット追加技法、又はGPAが膨張するか否かによってバルーニング技法を使用するかどちらか一方を決定し得る。例えば、VID(810)がNUMAノード(606)ヘコミットするSPAを受信したとき、それは、任意のGPAがDMVSC(804)によってロックされるか否か決定し得る。ロックされたGPAが存在する例において、VID(810)は、それがメモリーをホット追加する前にSPAを用いてそれらを支援し得る。メモリーが仮想NUMAノード(606)ヘコミットされる前に、それは0になり得、それに関連するキャッシュラインが、セキュリティのためにフラッシュされ得る。メモリーを0にすることによって、1つのパーティションに以前関付けられたメモリーのコンテンツが別のパーティションヘリークされない。

40

【 0 0 5 1 】

[0065]ここで図10に移ると、それは図9の追加的動作(1010~1020)を含んでいる動作手順の代替実施形態を表している。動作(1010)は、第2の複数の仮想NUMAノードにおけるメモリー圧力が、所定の値よりも大きいことを決定するステップと、第2のNUMAノードを計算機システムの第2の仮想NUMAノードヘ移動させるステップと、を示している。そして図7に移ると、例えば、実施形態において、第2の仮想N

50

UMA ノード (6 0 8) におけるメモリー圧力が増大し得る。すなわち、メモリー圧力の指標値が、仮想 NUMA ノード (6 0 8) が圧迫されていることを示す DMVSP (8 0 2) によって、受信され得る。この例において、仮想マシン (2 4 0) 又は個々の仮想 NUMA ノード (6 0 8) は、目標の圧力値を有し得、現在の圧力値が管理者によって設定された目標値よりも大きい場合がある。目標の圧力値は、DMVSP (8 0 2) によってアクセスされ得るデータ構造でストアされ得る。実行中の仮想マシン又は仮想 NUMA ノードの現在の圧力値は、その後受信され得る。DMVSP (8 0 2) は、実行中の仮想マシン又は仮想 NUMA ノードリストを介し連続的にステップし得、メモリー圧力値を目標の値に減少させるためにメモリーをコミットし得、圧力を目標の値に増大するためにメモリーをデコミットし得る。

10

【 0 0 5 2 】

[0066]例において、DMVSP (8 0 2) は、NUMA ノード仮想 NUMA ノード (6 0 6) 及び (6 0 8) を現在ホスティングしていて、例えば、NUMA ノード (7 0 2) がその仮想 NUMA ノード双方に対する目標メモリー圧力値を取得するためのメモリーを十分に割り当てることが不可能なことを決定するように構成され得る。この例において、DMVSP (8 0 2) は、ハイパーバイザー (2 0 2) へ信号を送信するように構成され得、ハイパーバイザー (2 0 2) は、仮想 NUMA ノードの 1 つを NUMA ノード (7 0 2) 上へ移動しようと試みるように構成され得る。ハイパーバイザー (2 0 2) は、NUMA ノード (7 0 2 ~ 7 0 6) の現在の作業負荷をチェックし得、例えば、NUMA ノード (7 0 4) が、仮想 NUMA ノードをホスティングし得、それに十分なリソースを割り

20

【 0 0 5 3 】

[0067]図 1 0 の説明を続けると、動作 (1 0 1 2) は、ゲストメモリーの少なくとも 1 つのメモリーブロックを特定の仮想 NUMA ノードからデコミットするステップと、デコミットされたゲストメモリーの少なくとも 1 つのメモリーブロックを第 2 の仮想 NUMA ノードへコミットするステップと、を例示している。例えば、DMVSP (8 0 2) は、例えば、メモリーを仮想 NUMA ノード (6 0 6) からデコミットし、メモリーを仮想 NUMA ノード (6 0 8) へコミットするように構成され得る。この例において、仮想 NUMA ノード (6 0 6) 及び (6 0 8) は、単一の NUMA ノード又は「フラット」アーキテクチャーによって支援され得る。この実施形態例において、DIMIVSP (8 0 2) は、例えば、仮想 NUMA ノード (6 0 8) へコミットされ得る利用可能な有効メモリーが存在しないとき、仮想 NUMA ノード (6 0 6) からメモリーをフリーにしようと試み得る。別の例において、DMVSP (8 0 2) は、例えば、メモリーを仮想 NUMA ノード (6 1 0) からデコミットし、メモリーを仮想 NUMA ノード (6 0 8) へコミットするように構成され得る。すなわち、メモリーがある仮想マシンから取得され得、別の仮想

30

40

【 0 0 5 4 】

[0068]そして図 6 を参照すると、特定の例において、仮想 NUMA ノード (6 0 6) 及び (6 0 8) が計算機システム (6 0 0) のリソースへマッピングされ得る。この例において、DMVSP (8 0 2) は、別の仮想 NUMA ノードを、例えば仮想マシン (2 4 0) における低優先順位の仮想 NUMA ノード又は例えば、最低優先順位の仮想マシンを開始するメモリー優先順位でチェックし得る。例えば、目標の閾値よりも小さなメモリー圧力値を有する仮想 NUMA ノード (6 0 6) のような仮想 NUMA ノードが検出された場合、DMVSP (8 0 2) は、メモリーをデコミット開始し得、仮想 NUMA ノード (6 0 6) からメモリーを移動し得る。デコミットが完了した後、コミット動作が開始され得

50

、メモリーは、仮想NUMAノード(608)へホット追加され得るか又は膨張したゲスト物理アドレスがシステム物理アドレスに再度関連付けられ得る。

【0055】

[0069]そして図7を参照すると、特定の例において、DMVSP(802)は、例えば、メモリーの優先順位で同一のNUMAノード(702)によって支援される別の仮想NUMAノードをチェックし得る。例えば、仮想NUMAノード(608)と同一のNUMAノード上の仮想NUMAノードが、目標の閾値よりも小さなメモリー圧力値を有していることが検出された場合、DMVSP(802)はメモリーをデコミット開始し得る。デコミットが完了した後にコミット動作が開始され得、メモリーが仮想NUMAノード(608)へホット追加され得るか又は膨張されたゲスト物理アドレスがシステム物理アドレスと再度関連付けられ得る。

10

【0056】

[0070]図10の説明を続けると、動作(1014)は、特定の仮想NUMAノードの少なくとも1つのゲストメモリーブロックが、システムメモリーと分離されることを決定するステップと、ゲストメモリーの少なくとも1つのメモリーブロックをシステムメモリーの少なくとも1つのメモリーブロック上へマッピングするステップと、を表している。例えば、実施形態において、DMVSP(802)が論理プロセッサによって実行され得、SPA(624)を用いて仮想NUMAノード(606)におけるGPAを支援する決定が実行され得る。例えば、GPAがDMVSC(804)によって予約され得、SPAが別の仮想NUMAノード又は親パーティション(204)のどちらか一方へ再度割り当てられ得る。この例において、メモリーページをコミットするためのリクエストがVID(810)によって受信され得、VID(810)がリクエストを満足するSPAを取得し得、アドレス範囲をDMVSP(802)へ送信し得る。実施形態において、VID(810)は、システム効率を増大するために隣接したSPAの範囲を取得するように構成され得る。NUMA実施形態において、VID(810)は、仮想NUMAノード(606)を実行中の同一のNUMAノードから隣接したSPAの範囲を取得するように構成され得る。VID(810)は、ロックされたGPAとSPAとの間に関係を生成し得、メッセージをDMVSP(802)へ送信し得る。DMVSP(802)はその後、メッセージをDMVSC(804)へ送信し得、DMVSC(804)は、GPAがアンロックされ得、仮想NUMAノード(606)と関連付けられたメモリープールへ返却され得ることを示すメッセージをメモリーマネージャー(808)へ送信し得る。

20

30

【0057】

[0071]図10の説明を続けると、動作(1016)は、特定の仮想NUMAノードを計算機システムの第1のNUMAノード上へマッピングするステップと、特定の仮想NUMAノードを計算機システムの第2のNUMAノード上へ移動させるステップと、を例示している。そして図7を参照すると、例えば、ゲストオペレーティングシステム(220)は、NUMAノード(702)及び(704)のような少なくとも2つのNUMAノード全域に拡大され得る。そして図7を参照すると、例えば、ハイパーバイザー(202)は、仮想NUMAノード(606)及び(608)が、NUMAノード(702)上で実行するようにスケジューリングし得る。この例において、ハイパーバイザー(202)は、NUMAノード(702)が圧迫されていることを示す信号を受信し得る。例えば、ゲストオペレーティングシステム(220)は、仮想NUMAノード(606)及び(608)がメモリー上において低いことを示す信号を生成し得る。この例において、ハイパーバイザー(202)は、NUMAノード(702)上に仮想NUMAノード(608)を移動することによって、圧迫されているNUMAノードの作業負荷を減少させるように構成され得る。

40

【0058】

[0072]図10の説明を続けると、動作(1018)は、仮想プロセッサを特定の仮想NUMAノードへ追加するステップを例示している。例えば、実施形態において、仮想プロセッサ(230B)のような仮想プロセッサが、仮想マシン(240)のランタイ

50

ム実行中に例えば、プロセッサのホット追加動作を使用して追加され得る。すなわち、仮想NUMAノード(606)は、ある地点において単一の仮想プロセッサ(230A)だけを有し得、その後、別のプロセッサが追加される。実施形態において、新たに追加されたプロセッサが、仮想プロセッサ(230A)を支援しているプロセッサへ割り当てられ得るか、又は別の論理プロセッサが、仮想プロセッサ(230B)のスレッドへ割り当てられ得、実行され得る。NUMA実施形態において、別の論理プロセッサが使用され、仮想プロセッサ(230B)を支援している場合、それが仮想NUMAノード(606)において別の仮想プロセッサを支援している同一のNUMAノード(702)から割り当てられ得る。

【0059】

[0073]図10の説明を続けると、動作(1020)は、仮想マシンの仮想プロセッサと、論理プロセッサに割り当てられた仮想プロセッサと、NUMAノードに割り当てられた論理プロセッサと、仮想NUMAノードに割り当てられた仮想プロセッサと、を実行するためのリクエストを受信するステップと、論理プロセッサが仮想プロセッサを実行不可能であることを決定するステップと、仮想プロセッサを実行するための第2の論理プロセッサを選択するステップであって第2の論理プロセッサが、第2のNUMAノードから提供されているものと、を例示している。そして図7を参照すると、例えば、実施形態において、ハイパーバイザー(202)は、仮想プロセッサ(230A)から仮想プロセッサスレッドを実行するためのリクエストを受信し得、理想プロセッサ(212A)、例えば仮想プロセッサ(230A)を支援するプロセッサ上にスレッドをスケジューリングしようと試み得る。この例において、ハイパーバイザー(202)は、論理プロセッサ(212A)が過度にコミットされていて、仮想プロセッサスレッドを実行不可能なことを検出し得る。この場合、ハイパーバイザー(202)が実行され得、それが、仮想プロセッサスレッドを実行するための別の論理プロセッサを選択し得る。例えば、ハイパーバイザー(202)は、同一のNUMAノード上において異なる論理プロセッサを選択しようと試み得る。例えば、NUMAノードが過度にコミットされている場合、ハイパーバイザー(202)は、仮想プロセッサ(230A)、例えば論理プロセッサ(212E)を実行するためのリモートプロセッサを選択するように構成され得る。この例において、リモートノード上のスレッドを待つか又はスケジューリングすることに関する決定が、NUMAノード(704)に関連するNUMA比を使用し実行され得る。NUMA比が低い場合、かつ、理想的なプロセッサを待つ予測時間が長い場合、NUMAノード(704)上にスレッドをスケジューリングする決定が実行され得る。他方では、NUMA比が高い場合、かつ、待つ予測時間が短い場合、待機する決定が実行され得る。

【0060】

[0074]ここで図11に移ると、それは動作(1100)、(1102)、(1104)、及び(1106)を含んでいる本開示の態様を実施するための動作手順を表している。動作(1100)が動作手順を開始し、動作(1102)は、仮想マシンを実行するステップであって仮想マシンが、複数の仮想NUMAノードを含むトポロジーを有しているものと、仮想マシンのトポロジーが、計算機システムの物理トポロジーから独立して生成されていることを示している。例えば、ハイパーバイザー(202)は、複数の仮想NUMAノードを有する仮想マシンを実行し得る。図6によって示されるような仮想NUMAノード(606)及び(608)を含む仮想マシン(240)が作成され得る。仮想NUMAノードはそれぞれ、1つ以上の仮想プロセッサ(230A~D)と、ゲスト物理アドレス(614)及び(616)とを有し得る。この実施形態において、仮想NUMAノード(606)及び(608)が基本ハードウェアトポロジーから独立して生成され得る。すなわち、仮想マシンのトポロジーは、図6及び図7によって表されるような基本ハードウェアと関係しない。かくして、この実施形態において、仮想マシンのトポロジーそれぞれは、それを達成する計算機システムの基本物理トポロジーから分離されている。

【0061】

[0075]図11の説明を続けると、動作(1104)は、複数の仮想NUMAノードそれぞれにおけるメモリ圧力を決定するステップを例示している。そして図8を参照すると、例えば、仮想NUMAノード(606)及び(608)それぞれに対するメモリ圧力が取得され得、例えば、生成され及び/又は受信され得る。メモリ圧力情報は、ゲストの実績が仮想NUMAノード(606~608)それぞれが利用可能なメモリ容量によってどのように作用されるか識別し得る。この情報は、例えば、DMVSC(804)及び/又は(806)のような1つ以上のDMVSCによってゲストオペレーティングシステム(220)のランタイム中に算出され得、DMVSP(802)へ送信され得る。すなわち、特定の実施形態において、論理プロセッサは、DMVSC(804)又は(806)の命令指標を実行し得、仮想NUMAノードそれぞれに対するメモリ圧力情報を生成し得る。この情報はその後、例えば、DMVSP(802)へ送信され得る。

10

【0062】

[0076]実施形態例において、メモリ圧力情報は、0~4の範囲の一連の値を含み得、それぞれの値は、ゲストOSが、仮想NUMAノード(606~608)のリソースによる経験をしている異なるレベルのメモリ圧力を識別し得る。ゲストオペレーティングシステムが、より圧迫されるにつれて、すなわち、現在の作業負荷を効率的に実行するために必要なメモリ容量が増加するにつれて、DMVSC(804)及び(806)は、それらの値を改め、この情報をDMVSP(802)へ伝達し得る。

【0063】

[0077]図11の説明を続けると、動作(1106)は、複数の仮想NUMAノードの少なくとも1つに割り当てられるゲストメモリを複数の仮想NUMAノードそれぞれにおけるメモリ圧力に基づいて調整するステップを示している。動作(1206)を含む実施形態において、論理プロセッサ(212A)は、DMVSP(802)の命令指標を実行し得、例えば、仮想NUMAノード(606)におけるゲスト物理アドレス量を調整し得る。すなわち、DMVSP(802)は、仮想NUMAノード(606)におけるメモリ圧力に基づいてメモリをコミット又はデコミットし得、例えばプロセスが、圧迫されている仮想NUMAノード(606)に割り当てられた場合、メモリがコミットされ得る。

20

【0064】

[0078]実施形態において、DMVSP(802)によって構成される論理プロセッサ(212)がメモリをコミットするか又はデコミットする決定をしたとき、それがメモリーブロックベース単位でそのように実行し得る。例えば、DMVSP(802)は、メモリーブロックをコミット又はデコミットし得、メモリーステータスがどのように変化するかチェックし得る。メモリーステータスが変化していない場合、DMVSP(802)は、別のメモリーブロックをコミット又はデコミットし得る。

30

【0065】

[0079]ここで図12に移ると、それは動作(1208~1216)を含んでいる図11の動作手順1(100)の代替実施形態を表している。図面によって例示されるように動作(1208)は、仮想マシンを第2の計算機システムへ送信するステップを示している。そして図5を参照すると、例えば、実施形態において、仮想マシンの状態が1つ以上の構成ファイルにセーブされ得、別の計算機システムへ、例えば、計算機(504)から(506)へ送信され得る。計算機システム(506)のハイパーバイザーは、ファイルを読み出し得るか又はファイル仮想マシンをインスタンス化し得る。

40

【0066】

[0080]仮想マシンのトポロジーは、仮想マシンへ移動させてリストアする能力に作用する。具体的には、基本ハードウェアトポロジーを検出可能にする決定及び仮想NUMAノードサイズが、仮想マシンがどのくらい十分に実行するか及びそれが容易に移動され得るか否かに影響する。例えば、仮想NUMAノードサイズは、移動させる機能に作用し、例えば、仮想NUMAノードサイズが増加するにつれて仮想NUMAノードの移植性が減少し、仮想NUMAノードサイズが減少するにつれて仮想マシンの性能も低下する。加える

50

と、基本計算機のトポロジーを検出し得る仮想マシンは、NUMA認識オペレーティングシステム及びアプリケーションが、ブート時間に仮想マシンが検出する第1のトポロジーに基づいて自らを最適化するという事実のために容易に移動され得ず、これらの最適化は、将来、仮想マシンが移動され得る計算機上において十分に実行し得ない。かくして、仮想NUMAノードをゲストオペレーティングシステムに公開することによって、それがブートしたとき、オペレーティングシステムがNUMAノードを使用するように最適化され得る。仮想NUMAノードを正確にサイジングすることによって、仮想マシンが、データセンター(500)における多くの様々な計算機システムに対し最適化され得る。

【0067】

[0081]そして図6を参照すると、例えば、仮想マシン(240)は2つ以上の仮想NUMAノード(606)及び(608)を含み得る。ハイパーバイザー(202)は、論理プロセッサ(212A~D)を用いて仮想プロセッサ(230A~D)を支援し得る。ゲストオペレーティングシステム(220)がブートしたとき、それが仮想NUMAノード(606)及び(608)を検出し得、仮想NUMAノード(606)及び(608)を使用するプロセスをスケジューリング及び実行を最適化するように構成され得る。しばらくして仮想マシン(240)は、図7によって表されたものと同様の物理トポロジーを有する計算機システムへ移動され得る。図7のハイパーバイザー(202)は、論理プロセッサ(212A及びB)を用いて仮想プロセッサ(230A及びB)を支援し得、論理プロセッサ(212E及びF)を用いて仮想プロセッサ(230C及びD)を支援し得る。ゲストオペレーティングシステム(220)は、基本計算機トポロジーがSMPからNUMAへ変更した時でも図6の計算機システム上で実行したときと同一の方法で作動し続け得る。

【0068】

[0082]図12の説明を続けると、動作(1210)は、第1の複数の仮想NUMAノードを計算機システムの第1のNUMAノードへマッピングするステップと、第2の複数の仮想NUMAノードを計算機システムの第1のNUMAノードへマッピングするステップと、を示している。そして図7を参照すると、例えば、論理プロセッサ(212A)のような論理プロセッサは、ハイパーバイザー命令を実行し得、仮想NUMAノード(606)及び(608)のような仮想NUMAノードをNUMAノード(702)のようなNUMAノードへ一体化し得る。より具体的には、論理プロセッサは、システム物理アドレスを用いてNUMAノード(702)のRAM(214)からゲスト物理アドレス(614)及び(616)を支援し得、論理プロセッサ(212A~D)を用いて仮想プロセッサ(230A~D)を支援し得る。

【0069】

[0083]図12の説明を続けると、動作(1212)は、第1の複数の仮想NUMAノードを計算機システムの第1のNUMAノード上へマッピングするステップと、第2の複数の仮想NUMAノードを計算機システムの第2のNUMAノード上へマッピングするステップと、を示している。そして図7を参照すると、例えば、論理プロセッサ(212A)のような論理プロセッサは、ハイパーバイザー命令を実行し得、仮想NUMAノード(606)をNUMAノード(702)へ割り当て得、仮想NUMAノード(608)をNUMAノード(704)へ割り当て得る。この例において、ハイパーバイザー(202)が実行されるとき、ハイパーバイザースケジューラーは、仮想プロセッサ(230A~B)からのスレッドを論理プロセッサ(212A~D)上にスケジューリングし得、仮想プロセッサ(230C又はD)からのスレッドを論理プロセッサ(212E~G)上にスケジューリングし得る。

【0070】

[0084]ここで図13に移ると、それは図12の動作手順の代替実施形態を例示していて、第2の仮想NUMAノードにおけるメモリー圧力が、所定の値よりも大きいことを決定するステップと、第2の仮想NUMAノードを計算機システムの第2のNUMAノードに移動させるステップと、を示している動作(1314)を含んでいる。そして図7に移る

と、例えば、実施形態において、第2の仮想NUMAノード(608)におけるメモリー圧力が増大し得る。すなわち、メモリー圧力の値の指標は、仮想NUMAノード(608)が圧迫されていることを示すDMVSP(802)を受信し得る。この例において、仮想NUMAノード(606~612)及び/又は仮想マシン(240~242)それぞれは、目標の圧力値を有し得、仮想NUMAノード(608)に対する現在の圧力値が管理者によって設定された目標値よりも大きい場合がある。目標圧力値は、DMVSP(802)によってアクセスされ得るデータ構造でストアされ得る。現在実行中の仮想マシン又は仮想NUMAノードの圧力値がその後、受信され得る。DMVSP(802)は、実行中の仮想マシン又は仮想NUMAノードの一覧を介し順次にステップし、メモリー圧力値を目標値に減少させるようにメモリーをコミットし、目標値に対する圧力を増やすようにメモリーをデコミットする。

10

【0071】

[0085]ここで図14に移ると、それは図12の動作手順の代替実施形態を例示していて、第2の仮想マシンメモリー圧力が、所定の値よりも大きいことを決定するステップと、仮想マシンの第2の仮想NUMAノードを計算機システムの第1のNUMAノードへ移動させるステップと、を示している動作(1416)を含んでいる。実施形態において、少なくとも2つの仮想マシン、例えば、仮想マシン(240)及び(242)が実行され得る。

この例において、仮想マシン(240)の仮想NUMAノードが、例えば、図7のNUMAノード(702)及び(704)へマッピングされ得、仮想マシン(242)の仮想NUMAノードが、例えば、NUMAノード(706)へマッピングされ得る。この例において、仮想マシン(240)及び(242)それぞれ及び/又は仮想NUMAノード(606~612)それぞれは目標の圧力値を有し得、DMVSP(802)によってアクセスされ得るデータ構造でストアされ得る。この例において、第2の仮想マシンのメモリー圧力(242)は、仮想マシン、例えば、多くの読み出し又は書き出しリクエストを受信した仮想マシン(242)における活動のために増大し得、その値がDMVSP(802)によって受信され得る。DMVSP(802)が、実行中の仮想マシン又は仮想NUMAノードの現在の圧力値を受信し得、実行中の仮想マシン又は仮想NUMAノードのリストを介し連続的にステップし得、メモリー圧力を軽減するために、メモリーが仮想マシン(242)にコミットされ得るか否か決定し得る。

20

30

【0072】

[0086]圧力が、メモリーをコミット又はデコミットすることによって軽減できない状況において、DMVSP(802)は、ハイパーバイザー(202)へ信号を送信するように構成され得、ハイパーバイザー(202)は、メモリー圧力を軽減するために、計算機システムのリソースを再度割り当てしようと試み得る。例えば、ハイパーバイザー(202)は、NUMAノード(702~706)の現在の作業負荷をチェックし得、例えば、NUMAノード(702)が、仮想マシン(240)から仮想NUMAノードをホスティングし得、仮想NUMAノード(608)をNUMAノード(702)へ再度割り当てし得ることを決定できる。すなわち、ハイパーバイザー(202)は、VID(810)と連動してゲスト物理アドレス(616)をシステム物理アドレス(712)へ再度マッピングし得、論理プロセッサ(212A及びD)を仮想プロセッサ(230C及びD)に対する理想的なプロセッサとして設定し得る。その後、ハイパーバイザーは、仮想NUMAノード(610)をNUMAノード(704)に再度マッピングし得、仮想マシン(242)の仮想NUMAノード(610~612)それぞれの内部のメモリーを、そのメモリー圧力を減少させるために調整し得る。

40

【0073】

[0087]ここで図15に移ると、それは動作(1500)、(1502)、及び(1504)を含んでいる本開示の態様を実施するための動作手順を例示している。動作(1500)が動作手順を開始し、動作(1502)は、第1の仮想マシンを実行するステップであって仮想マシンが、複数の仮想NUMAノードを含むトポロジーを有しているものと、

50

複数の仮想NUMAノードそれぞれが、仮想プロセッサ及びゲスト物理アドレスを含んでいるものと、仮想マシンのトポロジーが、計算機システムの物理トポロジーから独立して生成されていることを示している。例えば、図7のハイパーバイザー(202)は、複数の仮想NUMAノード(608~610)を有している仮想マシン(240)を実行し得る。仮想NUMAノード(606)及び(608)はそれぞれ、1つ以上の仮想プロセッサ(230A~D)と、ゲスト物理アドレス(614)及び(616)と、を有し得る。この実施形態において、仮想NUMAノード(606)及び(608)が計算機システム(700)のリソースへマッピングされ得る。例えば、論理プロセッサ(212A及び212B)は仮想プロセッサ(230A及び230B)に対する理想的なプロセッサとして設定され得、ゲスト物理アドレス(614)がシステム物理アドレス(710)によって支援され得る。同様に論理プロセッサ(212E)及び(212F)が仮想プロセッサ(230C)及び(230D)に対する理想的なプロセッサとして設定され得、ゲスト物理アドレス(616)がシステム物理アドレス(714)へマッピングされ得る。この実施形態において、仮想NUMAノード(606)及び(608)は、基本ハードウェアトポロジーから独立している。すなわち、仮想マシンのトポロジーは、図6及び図7によって表されるような基本ハードウェアと関係しない。かくして、この実施形態において、仮想マシンのトポロジーは、それを達成する計算機システムの基本物理トポロジーから分離されている。

10

【0074】

[0088]図15の説明を続けると、動作(1504)は、増設の仮想プロセッサを複数の仮想NUMAノードに追加するステップを示している。例えば、実施形態において、増設の仮想プロセッサが、例えば、仮想NUMAノード(606)のような仮想NUMAノードに追加され得る。

20

この例において、仮想プロセッサのような仮想プロセッサが、例えば、プロセッサのホット追加動作を使用し仮想マシン(240)のランタイム実行中に追加され得る。実施形態において、新しく追加された仮想プロセッサが仮想プロセッサ(230A)を支援しているプロセッサへ割り当てられ得るか、又は別の論理プロセッサが新しく追加された仮想プロセッサのスレッドを実行する理想的なプロセッサとして設定され得る。NUMA実施形態において、論理プロセッサは、仮想NUMAノード(606)を支援している同一のNUMAノード(702)から割り当てられ得る。

30

【0075】

[0089]ここで図16に移ると、それは付加的な動作(1606~1612)を含んでいる図15によって表されている動作手順の代替実施形態を示している。動作(1606)は、複数の仮想NUMAノードそれぞれにおけるメモリー圧力を決定するステップと、複数の仮想NUMAノードそれぞれにおけるメモリー圧力に基づいて複数の仮想NUMAノードの少なくとも1つに割り当てられるゲストメモリーを調整するステップと、を示している。図8を参照すると、仮想NUMAノード(606)及び(608)のそれぞれに対するメモリー圧力が取得され得、例えば生成され及び/又は受信され得る。メモリー圧力情報は、仮想NUMAノードそれぞれが利用可能なメモリー容量によってゲスト性能がどのように作用されるか識別し得る。DMVSC(804)及び(806)は、例えば、メモリーマネージャー(808)から物理メモリー通知及び/又はゲストオペレーティングシステムのページング情報を受信し、それを利用し、仮想NUMAノード(606)及び(608)それぞれのメモリー圧力を計算するように構成され得る。

40

【0076】

[0090]本例を続けると、論理プロセッサ(212A)は例えば、DMVSP(802)の命令指標を実行し得、例えば、仮想NUMAノード(606)におけるゲスト物理アドレス量を調整し得る。すなわち、DMVSP(802)は、ゲストOSがリソースのために仮想NUMAノード(606)において経験しているメモリー圧力に基づいて、メモリーをコミット又はデコミットする。

【0077】

50

[0091]図16の説明を続けると、動作(1608)は、複数の仮想NUMAノードから仮想プロセッサを移動するステップを例示している。(1706)を含む実施形態において、ハイパーバイザー(202)が論理プロセッサによって実行され得、仮想プロセッサが仮想NUMAノード(606)から移動され得る。例えば、ハイパーバイザー(202)は、ゲストオペレーティングシステム(220)のホット移動APIをアクセスし得、例えば、仮想プロセッサ(230B)を仮想NUMAノード(608)から移動し得る。

【0078】

[0092]ここで動作(1610)に移ると、それは複数の仮想NUMAノードに関するNUMA比をゲストオペレーティングシステムへ報告するステップを示している。例えば、実施形態において、ハイパーバイザー(202)は、仮想NUMAノード(606~608)に関するNUMA比を生成し得、この情報が図6又は図7のどちらか一方のゲストオペレーティングシステム(220)へ報告され得る。実施形態において、ハイパーバイザー(202)は、仮想NUMAノードに関するNUMA比を示しているデータ構造を仮想マシンのファームウェアテーブルに生成し得、ゲストオペレーティングシステム(220)がブートしたとき、ゲストが、テーブルを読み出し得、情報を利用し得、スケジューリング決定を実行し得る。例えば、ゲストオペレーティングシステム又はNUMA認識アプリケーションが、NUMA比を利用し得、リモートNUMAノードからのリソースを使用するか否か決定し得る。例えば、オペレーティングシステムは、実行される準備をしている保留中のスレッドを有し得る。この例において、オペレーティングシステムも、フリーになる理想的なプロセッサを一定の時間、待つように構成され得るか、さもなければ、それは、所定よりも小さなNUMA比を有するスレッドをリモートプロセッサ上にスケジューリングする。この場合、スケジューラーが受け入れる時間はNUMA比に依存している。

【0079】

[0093]ここで動作(1612)に移ると、それは仮想マシンを第1のNUMAノードから複数のNUMAノードへ移動させるステップを示している。例えば、実施形態において、ハイパーバイザー命令は論理プロセッサによって実行され得、仮想マシン(240)は複数のNUMAノード(704)及び(706)へマッピングされ得る。この例において、計算機システム(700)は乱用され得、例えば仮想マシン(242)は、リソースの大部分を使用している場合があって、そのため、仮想マシン(242)が計算機(700)から移動される。この状況において、ハイパーバイザー(202)は、計算機システム(700)のリソースを再割り当てし得、仮想マシン(240)をNUMAノード(704)及び(706)へ再マッピングし得る。

【0080】

[0094]ここで図17に移ると、それは図16の動作手順の動作(1714)及び(1716)を含む代替実施形態を表している。実施形態において、動作(1606)は、仮想NUMAノードの現在のメモリー圧力が、目標の閾値よりも低いという決定に基づいて第1の複数の仮想NUMAノードからメモリーをデコミットするステップを示している動作(1714)を含み得る。例えば、実施形態において、DMVSP(802)は、仮想NUMAノード(608)のメモリー圧力が、目標範囲よりも低いとき、仮想NUMAノード(606)からメモリーをデコミットするように構成され得る。例えば、実施形態において、DMVSP(802)は、例えば、コミットされ得る親パーティションにおいてメモリーも利用可能でないとき、及び仮想NUMAノード(608)が受け入れられない圧力を経験しているとき、仮想NUMAノード(606)からメモリーをフリーにしようと試み得る。メモリーが仮想NUMAノード(606)からデコミットされた場合、非同期メッセージが、メモリーをデコミットする指示するためのDMVSC(804)が送信され得る。ゲストOS(220)に関連するDMVSC(804)が応答したとき、それは仮想NUMAノード(606)内部の新しいメモリー圧力を示し得る。場合によっては、メモリー圧力が、メモリーの移動動作に応答し増大され得るものもある。

【 0 0 8 1 】

[0095]図 1 7 の説明を続けると、動作 (1 7 1 6) は、ゲストオペレーティングシステムの現在のメモリー圧力が目標の閾値よりも大きいという決定に基づいて、第 1 の複数の仮想 NUMA ノードへメモリーをコミットするステップを示している。例えば、実施形態において、DMVSP (8 0 2) は、仮想 NUMA ノード (6 0 6) のメモリー圧力が目標範囲よりも大きいとき、メモリーを仮想 NUMA ノード (6 0 6) へコミットするように構成され得る。この場合、メモリーが利用可能な場合、それが仮想 NUMA ノード (6 0 6) に割り当てられ得る。すなわち、DMVSP (8 0 2) は、利用可能なメモリー容量によって仮想 NUMA ノード (6 0 6) の性能がどのように作用されるか識別するメモリー圧力情報を取得し得、メモリーを仮想 NUMA ノード (6 0 6) へ追加し得る。特定の例において、メモリー圧力情報は数値であり得る。この例において、DMVSP (8 0 2) は、現在のメモリー圧力値を仮想 NUMA ノード (6 0 6) に対する最小値を示す情報テーブルと比較し得、仮想 NUMA ノード (6 0 6) のメモリー圧力が最小と等しくなるまでメモリーを調整し得る。例えば、管理者は、仮想 NUMA ノード (6 0 6) 及び (6 0 8) に対する最小値を有する不可欠なアプリケーションを実行するゲストオペレーティングシステムを構成し得る。

10

【 0 0 8 2 】

[0096]前述の詳細な説明は、例及び / 又は動作図面を介し、システム及び / 又はプロセスの様々な実施形態を詳細に説明している。そのようなブロック図及び / 又は例が 1 つ以上の機能及び / 又は動作を含む場合、そのようなブロック図又は例の中の機能及び / 又は動作それぞれは、様々な範囲のハードウェア、ソフトウェア、ファームウェア、又は実際にはその任意の組み合わせによって個別に及び / 又はまとめて実装され得ることが当業者によって理解されよう。

20

【 0 0 8 3 】

[0097]本明細書に記述された本対象事項の具体的な態様が示され、説明されているが、当業者にとって、本明細書の教示に基づく変更及び修正が本明細書に記述された対象項目及びそのより広い態様から離れずに実行され得ることは明らかであって、したがって、添付の請求項は、それらの変更及び修正すべてが本明細書に記述された対象項目の本当の趣旨及び範囲内にあるようにその範囲内に包含される。

30

【 符号の説明 】

【 0 0 8 4 】

- 2 0 計算機システム
- 2 1 論理プロセッサ
- 2 2 システムメモリー
- 2 3 システムバス
- 2 4 読み出し専用メモリー (R O M)
- 2 5 ランダムアクセスメモリー (R A M)
- 2 6 基本入力 / 出力システム (B I O S)
- 2 7 ハードディスクドライブ
- 2 8 磁気ディスクドライブ
- 2 9 取り外し可能磁気ディスク
- 3 0 光学式ディスクドライブ
- 3 1 取り外し可能光学式ディスク
- 3 2 ハードディスクドライブインターフェース
- 3 3 磁気ディスクドライブインターフェース
- 3 4 光学式ドライブインターフェース
- 3 5 オペレーティングシステム
- 3 6 アプリケーションプログラム
- 3 7 その他のプログラムモジュール
- 3 8 プログラムデータ

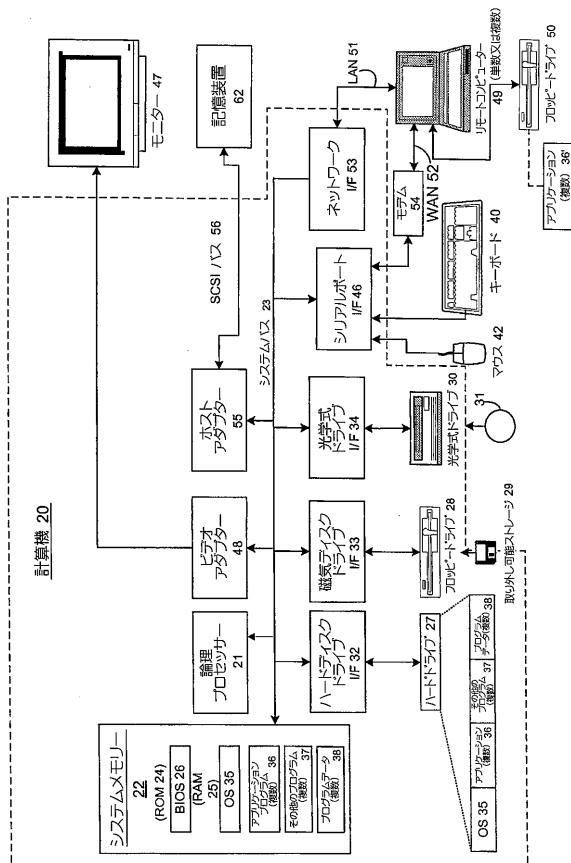
40

50

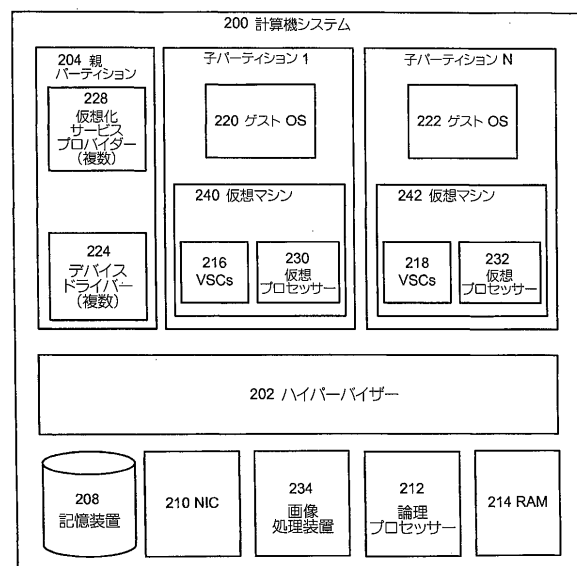
4 0	キーボード	
4 2	ポインティングデバイス	
4 6	シリアルポートインターフェース	
4 7	ディスプレイ	
4 8	ビデオアダプター	
4 9	リモートコンピューター	
5 0	メモリー記憶装置	
5 1	ローカルエリアネットワーク (L A N)	
5 2	広域ネットワーク (W A N)	
5 3	アダプター	10
5 4	モデム	
5 5	ホストアダプター	
5 6	小型コンピューター用周辺機器インターフェース (S C S I) バス	
6 2	外部記憶装置	
2 0 0	計算機システム	
2 0 2	ハイパーバイザー	
2 0 4	親パーティション	
2 0 8	記憶装置	
2 1 0	ネットワークインターフェースコントローラー (N I C)	
2 1 2	論理プロセッサ	20
2 1 4	ランダムアクセスメモリー (R A M)	
2 1 6	仮想化サービスクライアント (V S C)	
2 1 8	仮想化サービスクライアント (V S C)	
2 2 0	ゲストオペレーティングシステム	
2 2 2	ゲストオペレーティングシステム	
2 2 4	デバイスドライバ	
2 2 8	仮想化サービスプロバイダー	
2 3 0	仮想プロセッサ	
2 3 2	仮想プロセッサ	
2 3 4	画像処理装置	30
2 3 6	設定ユーティリティ	
2 4 0	仮想マシン	
2 4 2	仮想マシン	
3 0 0	計算機システム	
5 0 0	データセンター	
5 0 2	管理システム	
5 0 4	計算機システム	
5 0 6	計算機システム	
5 0 8	計算機システム	
5 1 0	計算機システム	40
6 0 0	計算機システム	
6 0 2	メモリーコントローラー	
6 0 6	仮想 N U M A ノード	
6 0 8	仮想 N U M A ノード	
6 1 0	仮想 N U M A ノード	
6 1 2	仮想 N U M A ノード	
6 1 4	ゲスト物理アドレス	
6 1 6	ゲスト物理アドレス	
6 1 8	ゲスト物理アドレス	
6 2 0	ゲスト物理アドレス	50

- | | |
|-------|---------------------------------|
| 6 2 2 | システム物理アドレスブロック |
| 6 2 4 | システム物理アドレスブロック |
| 6 2 6 | システム物理アドレスブロック |
| 6 2 8 | システム物理アドレスブロック |
| 6 3 0 | システム物理アドレスブロック |
| 6 3 2 | システム物理アドレスブロック |
| 7 0 0 | 計算機システム |
| 7 0 2 | NUMA ノード |
| 7 0 4 | NUMA ノード |
| 7 0 6 | NUMA ノード |
| 7 0 8 | 相互接続 |
| 8 0 2 | ダイナミックメモリー仮想化サービスプロバイダー (DMVSP) |
| 8 0 4 | ダイナミックメモリー仮想化サービスクライアント (DMVSC) |
| 8 0 6 | ダイナミックメモリー仮想化サービスクライアント DMVSC |
| 8 0 8 | ゲストオペレーティングシステムメモリーマネージャー |
| 8 1 0 | 仮想基盤ドライバ (VID) |
| 8 1 2 | ワーカプロセス |

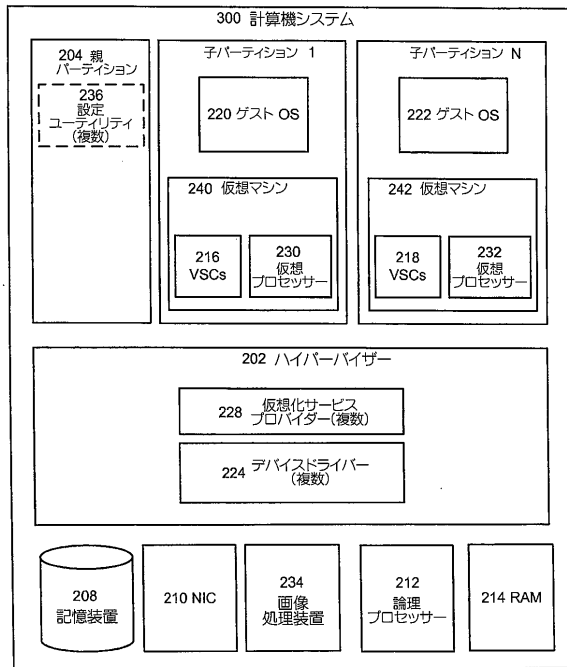
【圖 1】



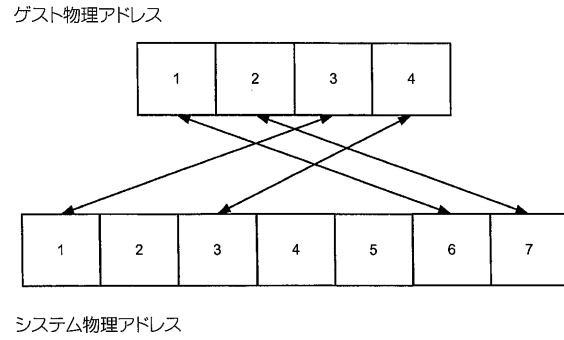
【圖 2】



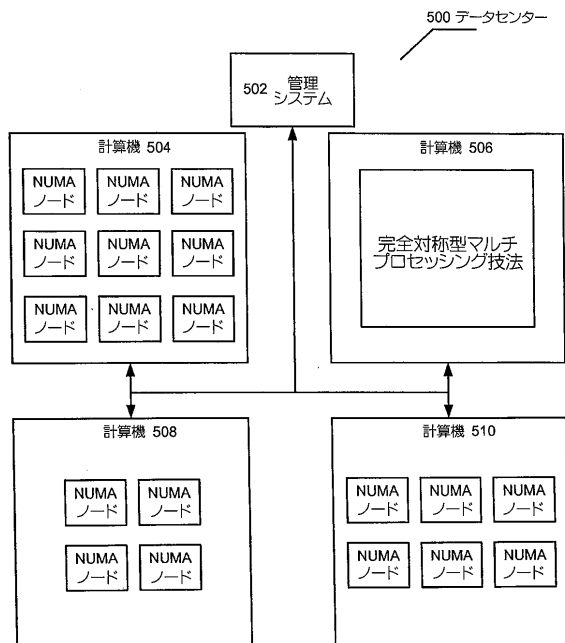
【図 3】



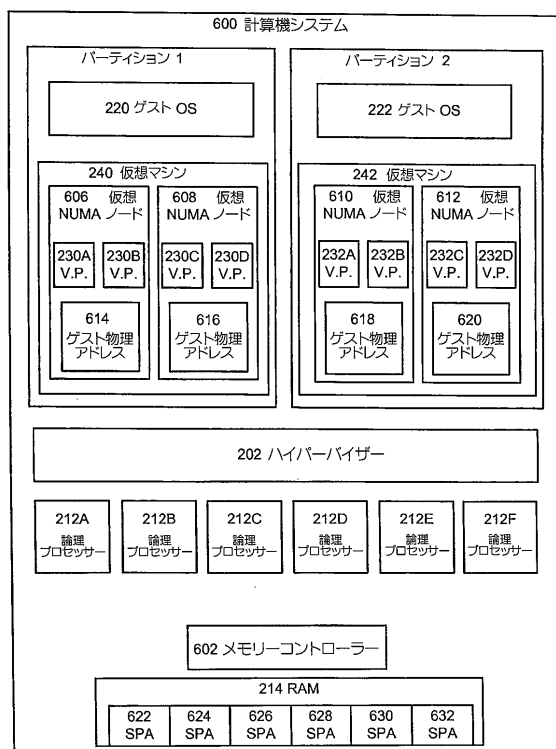
【図 4】



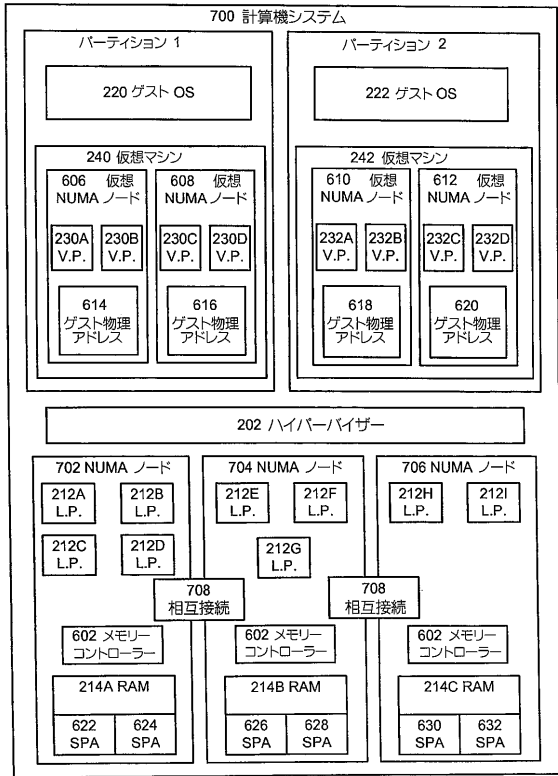
【図 5】



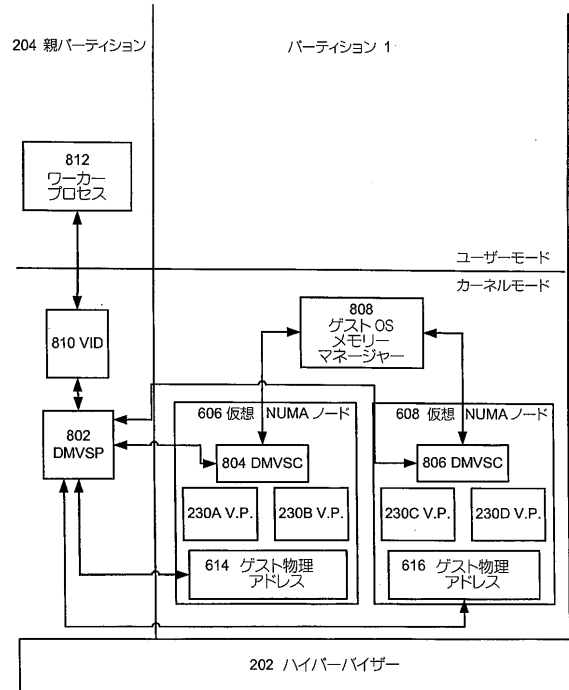
【図 6】



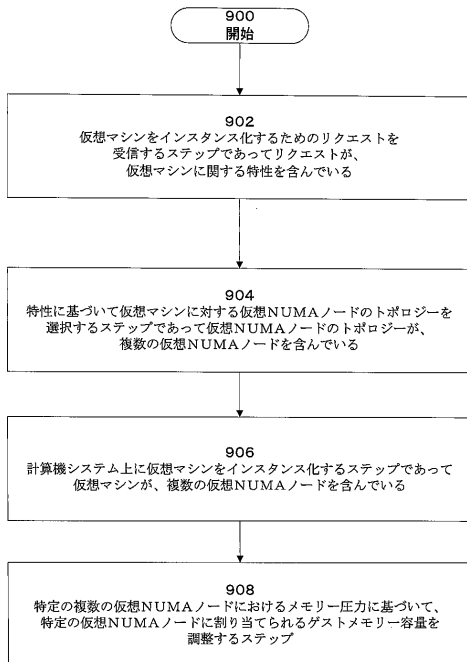
【図 7】



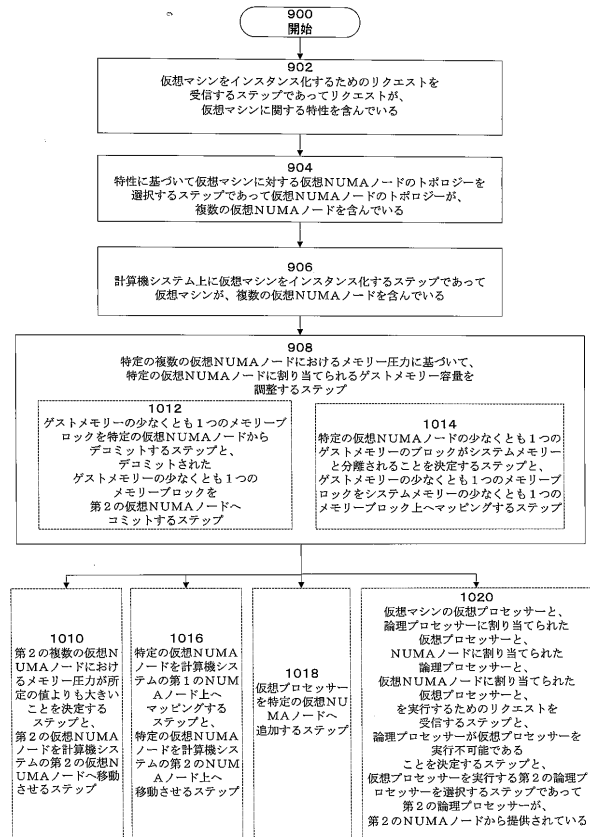
【図 8】



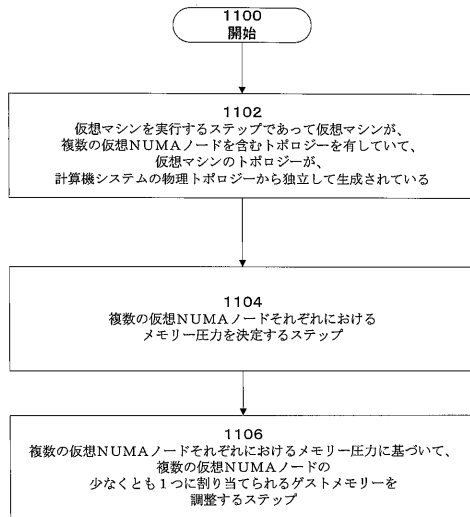
【図 9】



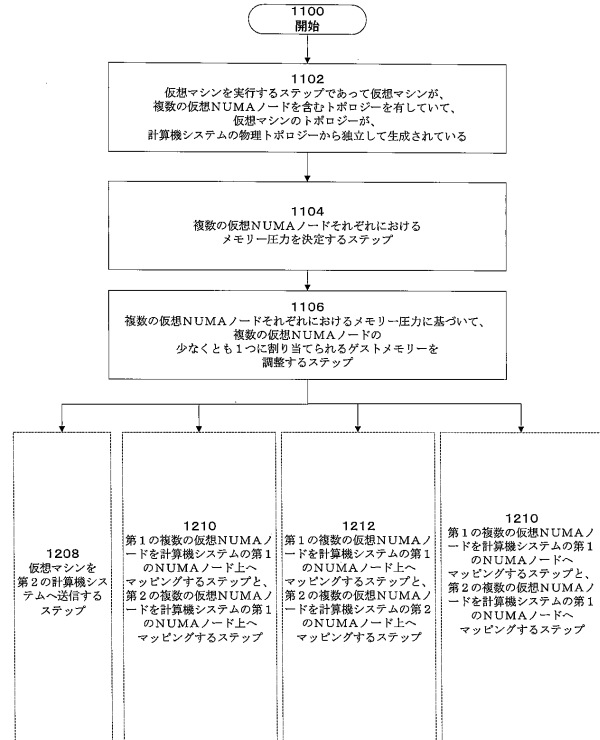
【図 10】



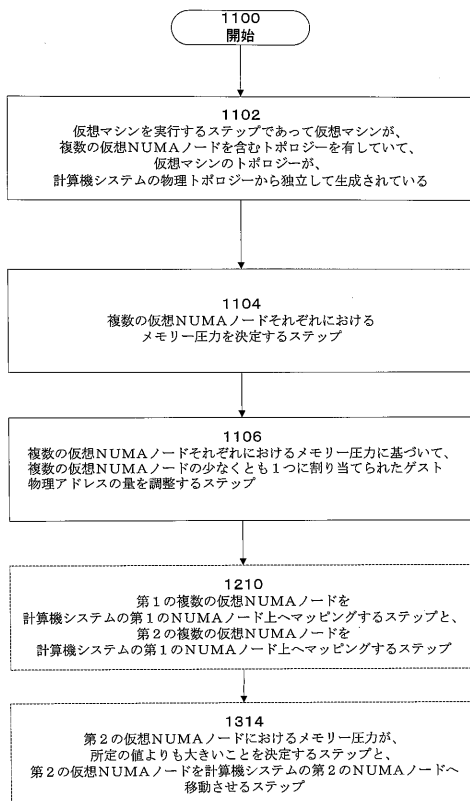
【 図 1 1 】



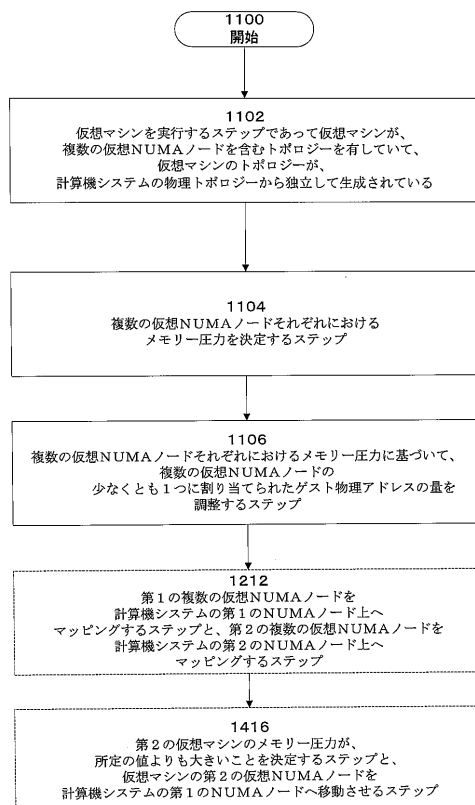
【 図 1 2 】



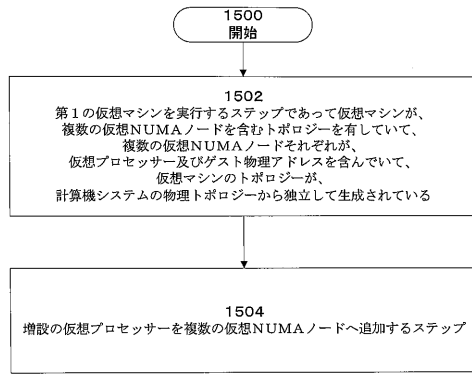
【 図 1 3 】



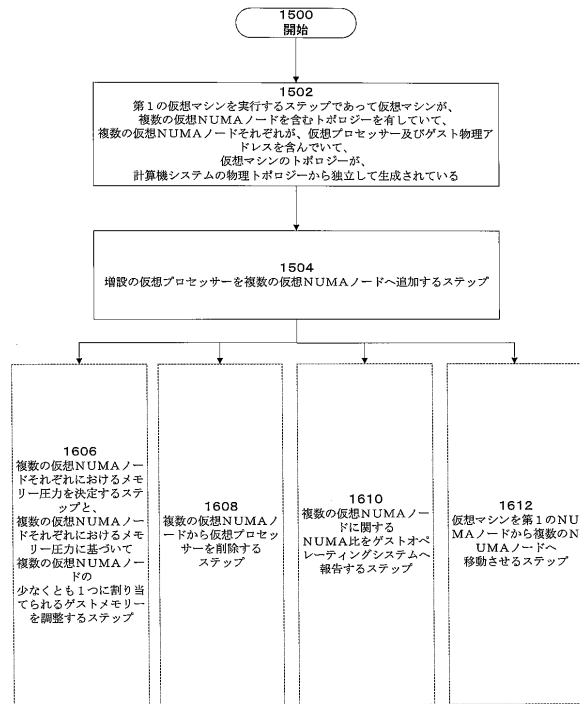
【 図 1 4 】



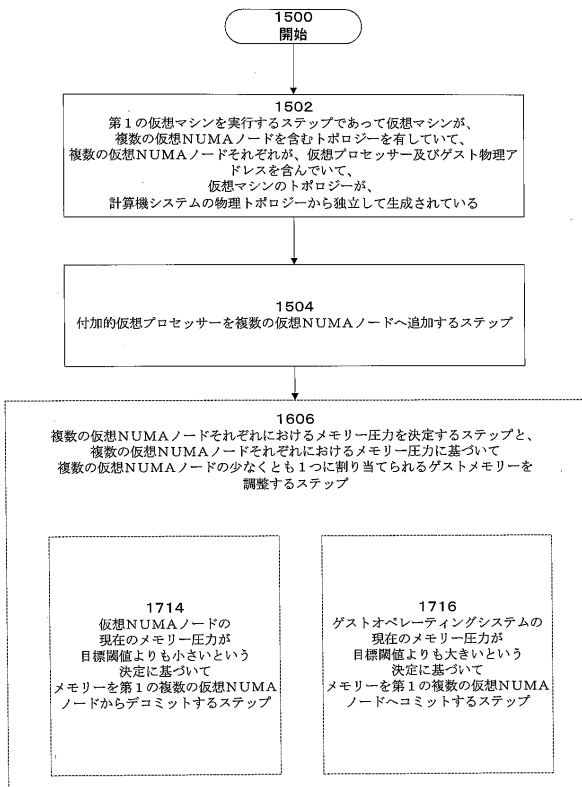
【図 15】



【図 16】



【図 17】



フロントページの続き

(72)発明者 オシンズ, ジェイコブ

アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9 , 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, エルシーエイ - インターナショナル・パテント

審査官 篠塚 隆

(56)参考文献 特開 2 0 0 6 - 1 7 8 9 3 3 (J P , A)

特開 2 0 0 7 - 2 5 7 0 9 7 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 9 / 4 6

9 / 4 8

9 / 5 0 - 9 / 5 2

9 / 5 4