

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5469940号
(P5469940)

(45) 発行日 平成26年4月16日(2014.4.16)

(24) 登録日 平成26年2月7日(2014.2.7)

(51) Int. Cl. F I
G 0 6 F 9/48 (2006.01) G O 6 F 9/46 4 5 2 Z
G 0 6 F 9/46 (2006.01) G O 6 F 9/46 3 5 0

請求項の数 10 (全 29 頁)

(21) 出願番号	特願2009-164360 (P2009-164360)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成21年7月13日(2009.7.13)	(74) 代理人	100100310 弁理士 井上 学
(65) 公開番号	特開2011-22627 (P2011-22627A)	(74) 代理人	100098660 弁理士 戸田 裕二
(43) 公開日	平成23年2月3日(2011.2.3)	(72) 発明者	松本 周平 神奈川県秦野市堀山下1番地 株式会社日立製作所 エンタープライズサーバ事業部内
審査請求日	平成24年2月23日(2012.2.23)	(72) 発明者	井上 裕功 神奈川県秦野市堀山下1番地 株式会社日立製作所 エンタープライズサーバ事業部内

最終頁に続く

(54) 【発明の名称】 計算機システム、仮想計算機モニタ及び仮想計算機モニタのスケジューリング方法

(57) 【特許請求の範囲】

【請求項1】

1つ以上の物理プロセッサを有する物理計算機と、仮想計算機モニタと、前記仮想計算機モニタにより前記1つ以上の物理プロセッサを時分割して構築された複数の仮想プロセッサを有する仮想計算機とを備える計算機システムにおいて、

前記仮想計算機モニタは、

前記複数の仮想プロセッサを前記物理プロセッサに割り当て、前記仮想プロセッサごとの仮想的な通常動作状態と仮想的な休止状態との切り替えを制御し、

第一の物理プロセッサに割り当てられ実行される第一の仮想プロセッサを、前記仮想的な通常動作状態から前記仮想的な休止状態に切り替えるとき、前記仮想的な休止状態として、

前記第一の仮想プロセッサを前記第一の物理プロセッサに割り当てたまま、前記第一の仮想プロセッサを実行不可とする第一の仮想的な休止状態と、

前記第一の仮想プロセッサを前記第一の物理プロセッサへの割り当てから外し、前記第一の仮想プロセッサを実行不可とする第二の仮想的な休止状態と、のいずれか一方を選択し、

前記選択するとき、前記第一の仮想プロセッサにおける前記仮想的な休止状態の解除が予測される時刻が、前記複数の仮想プロセッサのうちの他の仮想プロセッサにおける前記仮想的な休止状態の解除が予測される時刻より遅い場合は前記第二の仮想的な休止状態を選択する

ことを特徴とする計算機システム。

【請求項 2】

前記仮想計算機モニタは、

前記仮想プロセッサが前記仮想的な通常動作状態にあるときは、該仮想プロセッサ上のプログラムを実行可とする

ことを特徴とする請求項 1 記載の計算機システム。

【請求項 3】

前記仮想計算機モニタは、

前記第一の仮想プロセッサ上のプログラムが休止状態への切り替え要求を発行すると、前記切り替え要求を検出し、前記第一の仮想プロセッサを前記仮想的な通常動作状態から前記仮想的な休止状態に切り替え、

前記複数の仮想プロセッサのうちの他の仮想プロセッサが前記第一の仮想プロセッサへ割り込みを発行すると、前記割り込みを検出し、前記第一の仮想プロセッサを前記仮想的な休止状態から前記仮想的な通常動作状態に切り替える

ことを特徴とする請求項 2 記載の計算機システム。

【請求項 4】

前記仮想計算機モニタは、

前記仮想プロセッサごとに仮想プロセッサコンテキスト領域を有し、

前記第一の仮想プロセッサを前記第一の物理プロセッサに割り当てるとき、前記第一の仮想プロセッサの仮想プロセッサコンテキストを前記第一の物理プロセッサのレジスタに保持し、

前記第一の仮想プロセッサを前記第一の物理プロセッサに割り当てないとき、前記第一の仮想プロセッサの仮想プロセッサコンテキストを前記第一の仮想プロセッサの仮想プロセッサコンテキスト領域に保持する

ことを特徴とする請求項 3 記載の計算機システム。

【請求項 5】

前記仮想計算機モニタは、

前記仮想プロセッサごとに前記仮想プロセッサの直近の N 回の仮想プロセッサ休止時間を保持し、

前記仮想プロセッサの直近の N 回の前記仮想プロセッサ休止時間の平均値を該仮想プロセッサの仮想プロセッサ休止時間予測値とし、

前記仮想プロセッサが前記仮想的な休止状態の解除を待つときに、現在時刻に該仮想プロセッサの前記仮想プロセッサ休止時間予測値を加算して、該仮想プロセッサの前記仮想的な休止状態の解除が予測される時刻を算出する

ことを特徴とする請求項 4 記載の計算機システム。

【請求項 6】

前記仮想計算機モニタは、

前記第一の物理プロセッサに割り当てられ実行される前記第一の仮想プロセッサを、前記仮想的な通常動作状態から前記仮想的な休止状態に切り替えるとき、

前記第一の仮想プロセッサの前記仮想プロセッサ休止時間予測値が、所定の処理時間コスト未満であるならば、前記第一の仮想的な休止状態を選択する

ことを特徴とする請求項 5 記載の計算機システム。

【請求項 7】

前記仮想計算機モニタは、

前記第一の物理プロセッサに割り当てられ実行される前記第一の仮想プロセッサを、前記仮想的な通常動作状態から前記仮想的な休止状態に切り替えるとき、

前記第一の仮想プロセッサの前記仮想プロセッサ休止時間予測値が、前記所定の処理時間コスト以上であり、かつ第一の物理プロセッサへの割り当て対象の仮想プロセッサが存在するならば、前記第二の仮想的な休止状態を選択する

ことを特徴とする請求項 6 記載の計算機システム。

10

20

30

40

50

【請求項 8】

前記仮想計算機モニタは、

前記第一の物理プロセッサに割り当てられ実行される前記第一の仮想プロセッサを、前記仮想的な通常動作状態から前記仮想的な休止状態に切り替えるとき、

前記第一の仮想プロセッサの前記仮想プロセッサ休止時間予測値が、所定のしきい値以上であるならば、前記第二の仮想的な休止状態を選択する

ことを特徴とする請求項 7 記載の計算機システム。

【請求項 9】

物理計算機が有する 1 つ以上の物理プロセッサを時分割して複数の仮想プロセッサを有する仮想計算機を構築し、前記複数の仮想プロセッサを前記物理プロセッサに割り当て、前記仮想プロセッサの仮想的な通常動作状態と仮想的な休止状態との切り替えを制御する仮想計算機モニタであって、

第一の物理プロセッサに割り当てられ実行される第一の仮想プロセッサを、前記仮想的な通常動作状態から前記仮想的な休止状態に切り替えるとき、前記仮想的な休止状態として、

前記第一の仮想プロセッサを前記第一の物理プロセッサに割り当てたまま、前記第一の仮想プロセッサを実行不可とする第一の仮想的な休止状態と、

前記第一の仮想プロセッサを前記第一の物理プロセッサへの割り当てから外し、前記第一の仮想プロセッサを実行不可とする第二の仮想的な休止状態と、のいずれか一方を選択し、

前記選択するとき、前記第一の仮想プロセッサにおける前記仮想的な休止状態の解除が予測される時刻が、前記複数の仮想プロセッサのうちの他の仮想プロセッサにおける前記仮想的な休止状態の解除が予測される時刻より遅い場合は前記第二の仮想的な休止状態を選択する

ことを特徴とする仮想計算機モニタ。

【請求項 10】

物理計算機が有する 1 つ以上の物理プロセッサを時分割して複数の仮想プロセッサを有する仮想計算機を構築し、前記複数の仮想プロセッサを前記物理プロセッサに割り当て、前記仮想プロセッサの仮想的な通常動作状態と仮想的な休止状態との切り替えを制御する仮想計算機モニタのスケジューリング方法であって、

第一の物理プロセッサに割り当てられ実行される第一の仮想プロセッサを、前記仮想的な通常動作状態から前記仮想的な休止状態に切り替えるとき、前記仮想的な休止状態として、

前記第一の仮想プロセッサを前記第一の物理プロセッサに割り当てたまま、前記第一の仮想プロセッサを実行不可とする第一の仮想的な休止状態と、

前記第一の仮想プロセッサを前記第一の物理プロセッサへの割り当てから外し、前記第一の仮想プロセッサを実行不可とする第二の仮想的な休止状態と、のいずれか一方を選択し、

前記選択するとき、前記第一の仮想プロセッサにおける前記仮想的な休止状態の解除が予測される時刻が、前記複数の仮想プロセッサのうちの他の仮想プロセッサにおける前記仮想的な休止状態の解除が予測される時刻より遅い場合は前記第二の仮想的な休止状態を選択する

ことを特徴とするスケジューリング方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、仮想計算機システムに係わり、特に仮想計算機モニタが、仮想プロセッサのスケジューリングを制御する方法に関する。

【背景技術】

【0002】

10

20

30

40

50

プロセッサの高多重化や、キャッシュの大容量化が進み、プロセッサの性能が上がると共に消費電力も増えてきている。そのため、前記技術の進展と共にプロセッサの電力制御技術も進展している。近年のプロセッサは通常動作状態と休止状態を備え、前記の2つの状態を切り替えるプロセッサ電源制御機能を有する。プロセッサが休止状態にあるときは、該プロセッサは実行を停止し該プロセッサの資源を解放することで、該プロセッサの消費電力を減らす。汎用オペレーティングシステムは、前記プロセッサ電源制御機能をプロセッサ上で実行するプログラムが存在しない場合に使用し、該プロセッサを休止させることで、消費電力を節約する。

【0003】

物理プロセッサ数を超える数の仮想プロセッサをもつ仮想計算機システムにおいては、プロセッサ共有の効率性は非常に重要である。下記非特許文献1及び特許文献1では、仮想計算機システムにおいて仮想プロセッサの休止要求を、仮想計算機モニタが検出し、前記仮想プロセッサをブロックさせ、物理プロセッサ資源を他の仮想プロセッサに与える方法が記載されている。

10

【0004】

ところが、前記ブロック待ち方法は処理時間コストが大きい。該コンテキスト切り替えブロック待ち方法の処理時間コストは、仮想プロセッサの休止状態の解除のレイテンシーに影響する。ワークロードによっては、前記レイテンシーの大きさが性能問題となっている。

【0005】

20

このブロック待ち方法の処理時間コストを削減する方法が下記特許文献2に記載されている。また、ブロック待ちにおいて、プロセッサ利用率に基づいてブロック待ちを行使するのを遅らせる方法が、下記特許文献3に記載されている。

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開2008-186210号公報

【特許文献2】特開平06-044087号公報

【特許文献3】特開2009-009275号公報

【非特許文献】

30

【0007】

【非特許文献1】「How virtualization makes power management different」、[online]、Intel.corp.発行、[平成21年6月29日検索]、インターネット<URL: <http://ols.108.redhat.com/2007/Reprints/ke-Reprint.pdf>>

【発明の概要】

【発明が解決しようとする課題】

【0008】

非特許文献1の従来の仮想計算機システムでは、仮想計算機モニタが、仮想プロセッサの休止要求を検出した場合は、該仮想プロセッサの仮想的な休止状態が解除されるまで、ブロックさせている。前記ブロック待ち方法では、仮想プロセッサの仮想的な休止状態の解除のレイテンシーが大きく、仮想プロセッサの休止が頻発する場合には、該レイテンシーの大きさが性能問題となっていた。

40

【0009】

特許文献2の従来の方法は、あるプロセスがプロセッサを譲り渡すときに、切り替え先プロセスが見つからないならば、コンテキスト切り替えをスキップして前記プロセスのコンテキストのままアイドルプロセスを実行させる方法が提案されている。しかし、この方法では、仮想プロセッサのコンテキスト退避を先送りさせるだけであり、多くの仮想プロセッサが動作する仮想計算機システムでは、物理プロセッサの利用効率を下げる可能性が高いため、仮想計算機システムには適用できない。また、コンテキスト切り替えはスキッ

50

プするが、プロセス制御はスキップしないので、削減できるコストは部分的なものである。

【0010】

特許文献3の従来の方法では、プロセスがブロック待ち方法を実行するときに、プロセッサを譲り渡した後にアイドルのプロセッサが残る場合は、プロセッサ利用率が所定のしきい値よりも大きいならば、ブロック待ち方法の実行を遅らせて、該プロセッサ上に留まることで、ブロック待ち方法を回避する方法が提案されている。しかし、ワークロードには、仮想プロセッサの休止は頻発するが、仮想プロセッサ利用率自体は低いものがある。この場合は、プロセッサ利用率を見ても、仮想プロセッサの仮想的な休止状態から仮想的な通常動作状態への切り替え処理時間の削減が必要であることは分からない。

10

【0011】

一方、非特許文献1及び特許文献1の従来の方法では、仮想プロセッサをブロックさせている間は物理プロセッサを他の仮想プロセッサに割り当てることで高い物理プロセッサ利用効率を実現している。また、仮想プロセッサがブロックした後にアイドルの物理プロセッサが残った場合には、割り当てる仮想プロセッサが見つかるまで該物理プロセッサを休止させることで、消費電力を削減している。この物理プロセッサの高い利用効率と消費電力の削減は、複数の計算機を一台の物理計算機に集約するサーバ統合を主用途とする仮想計算機システムにおいては、非常に重要である。

【0012】

本発明の目的は、仮想計算機システムにおいて、従来方法と同等の物理プロセッサの利用効率及び消費電力の削減を保証する限りにおいて、仮想プロセッサの休止が頻発する場合に、該仮想プロセッサの仮想的な休止状態から仮想的な通常動作状態への切り替えの処理時間を小さくして、該仮想プロセッサ上のプログラムの起動を速くすることである。

20

【課題を解決するための手段】

【0013】

本発明は、仮想計算機モニタが、仮想プロセッサの仮想的な通常動作状態と仮想的な休止状態の切り替え及び仮想プロセッサのスケジューリングを制御する方法である。

【0014】

本発明では、仮想プロセッサが仮想的な休止状態が解除されるのを待ち始めてから該仮想的な休止状態が解除されるまでの間の時間を仮想プロセッサ休止時間とし、仮想プロセッサの直近のN回の前記仮想プロセッサ休止時間の平均値を該仮想プロセッサの仮想プロセッサ休止時間の予測値とし、次回仮想プロセッサが休止するときの該仮想プロセッサの仮想プロセッサ休止時間予測値を、1つ前の仮想プロセッサ休止時間予測値に、新たな仮想プロセッサ休止時間をNで割った値を加算し、最も古い仮想プロセッサ休止時間をNで割った値を除くことで、新た直近の複数回の仮想プロセッサ休止時間の総和を求めなおす必要なく、該仮想プロセッサの情報のみを用いて、一定の処理時間で算出する。

30

【0015】

仮想プロセッサが仮想的な休止状態の解除を待つときに、該仮想プロセッサの仮想プロセッサ休止解除時刻予測値を、現在時刻に前記プロセッサ休止時間予測値を加算することで算出する。

40

【0016】

仮想プロセッサの仮想的な休止状態の解除を待つ方法としてブロック待ち方法を選択したときに、前記算出した仮想プロセッサ休止解除時刻予測値を物理プロセッサごとの仮想プロセッサ休止解除時刻管理キューに期限が近い順に昇順でソートして挿入し、該仮想プロセッサの仮想的な休止状態が解除されるときに、前記管理キューから前記仮想プロセッサ休止解除時刻予測値を取り出す。

【0017】

仮想計算機モニタは、仮想的な休止状態にある仮想プロセッサが、仮想的な休止状態が解除されるのを待つ方法として、仮想的な休止状態が解除されるまで繰り返し調べるビジー待ち方法と、該仮想プロセッサの実行を停止させて動作していた物理プロセッサを譲り

50

渡し、仮想的な休止状態が解除されるまで該仮想プロセッサを物理プロセッサへの割り当て対象から外し、その間物理プロセッサ上では他の仮想プロセッサを動作させるブロック待ち方法の2つの方法を有し、前記算出の該仮想プロセッサの仮想プロセッサ休止時間予測値を用いてビジー待ち方法とブロック待ち方法の2つの方法から1つの方法を動的に選択して実行する。

【0018】

仮想計算機モニタは、前記選択処理について、該仮想プロセッサの仮想プロセッサ休止時間予測値が、所定の前記ブロック待ち方法の処理時間コスト未満であるならば、現在動作中の物理プロセッサへの割り当て対象の仮想プロセッサが存在するか否かを問わず、ビジー待ち方法を選択し、前記ブロック待ち方法処理時間コスト以上であり、かつ現在動作中の物理プロセッサへの割り当て対象の仮想プロセッサが存在する場合には、ブロック待ち方法を選択し、該仮想プロセッサの仮想プロセッサ休止時間予測値が、該仮想プロセッサが現在動作中の物理プロセッサ上での動作を継続するかを決める所定の仮想プロセッサ休止時間しきい値以上である場合には、ブロック待ち方法を選択し、該仮想プロセッサの仮想プロセッサ休止時間予測値が、前記ブロック待ち方法処理時間コスト以上でかつ前記仮想プロセッサ休止時間しきい値未満であり、現在動作中の物理プロセッサへの割り当て対象の仮想プロセッサが存在せず、該仮想プロセッサの仮想プロセッサ休止解除時刻予測値が、該仮想プロセッサが動作する物理プロセッサの仮想プロセッサ休止解除時刻管理キューの先頭の時刻以上であるならば、ブロック待ち方法を選択し、該仮想プロセッサの仮想プロセッサ休止時間予測値が、前記ブロック待ち方法処理時間コスト以上でかつ前記仮想プロセッサ休止時間しきい値未満であり、現在動作中の物理プロセッサへの割り当て対象の仮想プロセッサが存在せず、該仮想プロセッサの仮想プロセッサ休止解除時刻予測値が、該仮想プロセッサが動作する物理プロセッサの仮想プロセッサ休止解除時刻管理キューの先頭の時刻未満である場合には、ビジー待ち方法を選択する。

【発明の効果】

【0019】

本発明では、仮想的な休止状態にある仮想プロセッサが、前記仮想的な休止状態が解除されるのを待つ方法として、ビジー待ち方法を選択して実行した場合には、仮想プロセッサの仮想的な休止状態から仮想的な通常動作状態への切り替えの処理時間コストが削減され、該仮想プロセッサ上のプログラムの実行再開が高速化される。

【0020】

ただし、仮想的な休止状態の解除を待つ仮想プロセッサの他に物理プロセッサ上で実行させる仮想プロセッサがある場合は、ブロック待ち方法を選択して実行することで、物理プロセッサの利用効率に関して従来方法と同等を保証する。

【0021】

また、物理プロセッサ上でアイドルプロセスが実行している場合、あるいは仮想プロセッサが仮想的な休止状態の解除をビジー待ち方法で待つ場合は、物理プロセッサを休止させることによって、物理プロセッサの消費電力の削減に関して従来方法と同等を保証する。

【図面の簡単な説明】

【0022】

【図1】本発明を適用した仮想計算機システムの一実施形態を示す構成図。

【図2】汎用オペレーティングシステムのプロセッサ電源制御機能の用途を時系列で示す概念図。

【図3】本発明を適用した一実施形態における、仮想計算機モニタの構成図。

【図4】仮想プロセッサの仮想プロセッサ状態SS、スケジュール状態VS、及び仮想プロセッサ休止解除通知フラグNFに関する状態遷移図。

【図5】本発明を適用した一実施形態における、仮想プロセッサの休止処理の1つの例を示す概念図。

【図6】仮想プロセッサ休止解除待ち処理のフローチャート。

【図 7】仮想プロセッサ休止解除通知処理のフローチャート。

【図 8】仮想プロセッサ休止解除待ち処理と仮想プロセッサ休止解除通知処理の流れを時系列で示した概念図。

【図 9】仮想プロセッサ休止時間と該時間に関する予測値の算出処理と管理処理を示すフローチャート。

【図 10】仮想プロセッサ休止時間予測値と仮想プロセッサ休止時間の関係を示す概念図。

【図 11】仮想プロセッサ休止解除待ち方法選択処理のフローチャート。

【図 12】仮想プロセッサビジー待ち方法継続チェック処理のフローチャート。

【図 13】仮想プロセッサ休止解除待ち方法選択処理の各条件判定における変数の関係を示す概念図。

10

【図 14】物理プロセッサの電源制御を組み込んだアイドルループ処理のフローチャート。

【図 15】物理プロセッサの電源制御を組み込んだ仮想プロセッサ休止解除待ち処理におけるビジー待ち処理と該処理に対する仮想プロセッサ休止解除通知処理のフローチャート。

【発明を実施するための形態】

【0023】

本発明を適用した一実現形態を、図面を用いて詳細に説明する。

【0024】

20

図 1 は本発明を適用した仮想計算機システムの概要を示す構成図である。仮想計算機システムは、大きくは、物理計算機 001、仮想計算機モニタ 100、及び第 1 の仮想計算機 200 と第 2 の仮想計算機 300 に分類される。

【0025】

物理計算機は、物理プロセッサ 0 (002a) と物理プロセッサ 1 (002b)、メモリ 004、1 つ以上の I/O デバイス 005、外部割り込み機構 006、システム時刻 007、及びタイマ 008 を含む。

【0026】

物理プロセッサ 0 (002a) と物理プロセッサ 1 (002b) は、汎用オペレーティングシステムが動作可能なものであり、電源制御機能 003a と電源制御機能 003b をそれぞれ含む。図 1 においては、物理プロセッサの数を 2 つとしているが、本発明は物理プロセッサの数を限定するものではなく、1 以上の任意の数の物理プロセッサを有する物理計算機に適用可能である。

30

【0027】

ここで、近年の物理プロセッサは、標準でプロセッサ仮想化アシスト機能を有するが、本発明は、物理プロセッサがプロセッサ仮想化アシスト機能を有するかどうかは問わない。説明を単純にする為に、図 1 はプロセッサ仮想化機能を省略する。

【0028】

メモリ 004 は、プログラムとデータの為の揮発性記憶装置として使用する。

【0029】

40

I/O デバイス 005 は、ディスクドライブを動かすディスクコントローラや、ネットワークを通して外部の計算機システムと通信する為のネットワークインタフェースを含む。本実施形態においては、I/O デバイス 005 の詳細は問わない。

【0030】

外部割り込み機構 006 は、物理プロセッサ 0 (002a) と物理プロセッサ 1 (002b) に非同期でイベント発生を通知する機能を有する。外部割り込み機構 006 は、I/O デバイス 005 から物理プロセッサ 0 (002a) あるいは物理プロセッサ 1 (002b) へのコマンド完了通知、タイマ 008 の期限が来たことの物理プロセッサ 0 (002a) あるいは物理プロセッサ 1 (002b) への通知、及び物理プロセッサ 0 (002a) あるいは物理プロセッサ 1 (002b) から前記いずれかの物理プロセッサへのプロ

50

セッサ間割り込みなどで使用されることを想定するが、本実施形態では詳細の用途を問わない。

【0031】

システム時刻007は、仮想計算機モニタ001が現在時刻と経過時間を参照する為に使用する全ての物理プロセッサの間の共通時刻であり、物理計算機単独で、あるいは物理計算機001と仮想計算機モニタ100で連携して常に複数の物理プロセッサ間で同期することを保証する。本実施形態においては、システム時刻の実現方式は問わない。

【0032】

タイマ008は、指定した時刻に外部割り込み機構006を介して物理プロセッサ0(002a)及び物理プロセッサ1(002b)に外部割り込みを発生させる機能を有し、仮想計算機モニタ100がタイムスライスによって、物理計算機001の物理プロセッサ0(002a)と物理プロセッサ1(002b)を時分割共有する為に使用する。

【0033】

上述のとおり、物理計算機の物理プロセッサ0(002a)と物理プロセッサ1(002b)は、それぞれ電源制御機能003aと003bを有する。電源制御機能は、通常の動作状態と休止状態、及び前記の2つの状態を切り替える機能を有する。物理プロセッサが休止状態にあるときは、該物理プロセッサは実行を停止し該物理プロセッサの資源を解放することで、該物理プロセッサの消費電力を減らす。汎用オペレーティングシステムは、電源制御機能を使用し、プロセッサ上で実行するプログラムが存在しない場合に該プロセッサを休止させることで消費電力を節約する。プロセッサの通常動作状態から休止状態への切り替えは、プロセッサが標準で備えるプロセッサ休止命令の発行を介して、プロセッサ上のプログラムから該プロセッサに要求を発行することで行われる。プロセッサの休止状態から通常動作状態への切り替えは、外部割り込みの発行を介して、外部割り込み機構からプロセッサに要求を発行することで行われる。

【0034】

ここで、汎用オペレーティングシステムのプロセッサ電源制御機能の基本的な用途を説明する。本発明では、後述の通り、仮想計算機の仮想プロセッサも仮想的な電源制御機能を有する。以下の説明は、物理計算機と仮想計算機の両方に当てはまるように、物理あるいは仮想を用語の頭に付けない。

【0035】

汎用オペレーティングシステムが電源制御機能を使用してプロセッサを休止させるのは、大きくは、プロセッサがアイドルである第一の場合と、ワークロード実行中のI/Oデバイスからのコマンド完了通知待ちあるいは並列処理でのロック解放待ちである第二の場合に分けられる。以降、第一の場合をプロセッサがアイドルの場合、第二の場合をプロセッサが同期化イベント待ちの場合と略記する。

【0036】

図2(a)は、上記第一の場合であるプロセッサがアイドルの場合の通常動作状態と休止状態の切り替えを時系列で示した概念図である。汎用オペレーティングシステムは、プロセッサを休止状態に置き(ステップ401)、十分に長い周期のタイマ割り込みが該プロセッサに伝達されると休止状態が解除され通常動作状態に戻り(ステップ402)、汎用オペレーティングシステムはタイマ割り込み処理を行い(ステップ403)、タイマ割り込み処理を完了すると再びプロセッサを休止状態に切り替える(ステップ404)という動作を繰り返す。

【0037】

図2(b)は、上記第二の場合であるプロセッサが同期化イベント待ちの場合の通常動作状態と休止状態の切り替えの動作を時系列で示した概念図である。汎用オペレーティングシステムは、プロセッサを休止状態に置き(ステップ407)、非常に短いほぼ一定周期のI/O割り込みあるいはプロセッサ間割り込みが該プロセッサに伝達されると休止状態が解除されて通常動作状態に戻り(ステップ408)、該外部割り込みで要求された処理を行い(ステップ409)、要求された処理を完了すると再びプロセッサを休止状態に

10

20

30

40

50

切り替える（ステップ410）という動作を繰り返す。

【0038】

プロセッサがアイドルの場合とプロセッサが同期化イベント待ちの場合では、プロセッサ休止の周期とプロセッサ休止時間の長さの違いがある。プロセッサが同期化イベント待ちの場合のプロセッサ休止の周期412とプロセッサ休止時間411は、プロセッサがアイドルの場合のプロセッサ休止の場合と比べて共に非常に短い。共通点としては、プロセッサが同期化イベント待ちの場合とプロセッサがアイドルの場合のプロセッサ休止時間は共に定常的にほぼ一定の値をとる。

【0039】

仮想計算機モニタ100は、物理計算機001が有する資源を制御して、汎用オペレーティングシステムに代表されるプログラムが1つの計算機と認識するプログラム実行環境である仮想計算機を構築し、仮想計算機上でプログラムを、物理計算機001と比較して、少しの性能低下とタイミングの違いを除けば物理計算機001上と同等に動作させる役割を担う制御プログラムである。

10

【0040】

仮想計算機モニタ100は、物理計算機001が有する資源を制御して、仮想プロセッサ0(202a)と仮想プロセッサ1(202b)、仮想メモリ204、仮想I/Oデバイス205、仮想外部割り込み機構206を有する第1の仮想計算機200と、仮想プロセッサ0(302a)と仮想プロセッサ1(302b)、仮想メモリ304、仮想I/Oデバイス305、仮想外部割り込み機構306を有する第2の仮想計算機300を構築する。

20

【0041】

第1の仮想計算機200と第2の仮想計算機300上ではそれぞれ第1の汎用オペレーティングシステム201と第2の汎用オペレーティングシステム301が動作する。第1の汎用オペレーティングシステム201と第2の汎用オペレーティングシステム301は上述の用途で後述の仮想プロセッサの仮想電源制御機能を使うことを想定する。

【0042】

第1の仮想計算機200の仮想メモリ204と仮想I/Oデバイス205、第2の仮想計算機300の仮想メモリ304と仮想I/Oデバイス305は、仮想計算機モニタ100が対応する物理計算機001の有する資源を論理分割して実現する場合もあれば、物理計算機001の資源を共有して実現する場合もあれば、物理計算機001に存在しない資源を仮想計算機モニタ100が仮想的につくり出して実現する場合もある。本実施形態においては、仮想メモリ204と304、仮想I/Oデバイス205と305は、第1の仮想計算機200と第2の仮想計算機300の上のプログラムから見て物理計算機001と同等の機能を有するとし、その詳細は問わない。

30

【0043】

第1の仮想計算機200の仮想外部割り込み機構206と第2の仮想計算機300の仮想外部割り込み機構306は、それぞれ第1の仮想計算機200と第2の仮想計算機300の上のプログラムから見て、共に物理計算機001の物理外部割り込み機構006と同等の機能を有し、該物理機構と独立に動作するように仮想計算機モニタ100が仮想的に実現する。仮想外部割り込み機構は仮想プロセッサに非同期でイベント発生を通知する機能を有する。また、仮想外部割り込み機構は、仮想I/Oデバイスから仮想プロセッサへのコマンド完了通知、及びある仮想プロセッサから自身を含め指定した仮想プロセッサへの仮想プロセッサ間割り込みなどを実現するに足る機能を有する。本実施形態においては、仮想外部割り込み機構の実現方式の詳細は問わない。

40

【0044】

第1の仮想計算機200の仮想プロセッサ0(202a)と仮想プロセッサ1(202b)、第2の仮想計算機300の仮想プロセッサ0(302a)と仮想プロセッサ1(302b)は、それぞれ仮想電源制御機能203a、仮想電源制御機能203b、仮想電源制御機能303a、及び仮想電源制御機能303bを有する。前記仮想電源制御機能は、物

50

理プロセッサの電源制御機能と同様に、仮想的な通常動作状態と仮想的な休止状態の2つの状態と前記2つの状態の切り替え機能を有する。

【0045】

本発明の仮想計算機システムにおいては、仮想計算機モニタ100は、仮想プロセッサ上のプログラムが発行する上述のプロセッサ休止命令をハードウェアに直接実行させずに、検出することで、仮想プロセッサの通常動作状態から仮想的な休止状態への切り替えを物理プロセッサの電源制御機能と独立に仮想的に実現する。

【0046】

また、仮想計算機モニタ100は、前記仮想外部割り込み機構を使用することによって、仮想プロセッサの仮想的な休止状態から仮想的な通常動作状態への切り替えを、物理プロセッサの電源制御機能と独立に仮想的に実現する。

10

【0047】

仮想計算機モニタ100は、第1の仮想計算機200の仮想プロセッサ0(202a)と仮想プロセッサ1(202b)及び第2の仮想計算機300の仮想プロセッサ0(302a)と仮想プロセッサ1(302b)の間で、物理計算機001の物理プロセッサ0(002a)と物理プロセッサ1(002b)を時分割共有する。それぞれの仮想プロセッサ上のプログラムは、仮想計算機モニタ100が該仮想プロセッサを、物理計算機001の物理プロセッサ0(002a)と物理プロセッサ1(002b)のいずれかに割り当てることで実行される。

【0048】

20

仮想計算機モニタ100は、物理プロセッサ0(002a)と物理プロセッサ1(002b)にそれぞれアイドルプロセス020aとアイドルプロセス020bを対応付ける。仮想計算機モニタ100は、物理プロセッサに割り当てる仮想プロセッサが見つからなかったときに、該物理プロセッサに割り当てる仮想プロセッサが見つかるまで、該物理プロセッサ上で対応するアイドルプロセスを実行させる。

【0049】

仮想計算機モニタ100は、物理プロセッサに仮想プロセッサを割り当てるときに、該仮想プロセッサにタイムスライスと呼ぶ短いプロセッサ時間を与える。仮想プロセッサ上で動作する汎用オペレーティングシステムは、与えられたタイムスライスを使い切る前に、該仮想プロセッサの休止状態への切り替え要求を発行して、該仮想プロセッサを休止させることがある。この要因には、上述のプロセッサがアイドルの場合とワークロード実行中の同期化イベント待ちの場合がある。

30

【0050】

ワークロード実行中の同期化イベント待ちで仮想プロセッサを休止させる場合は、仮想プロセッサ休止の周期と仮想プロセッサ休止時間が共に仮想プロセッサがアイドルの場合のものに比べて非常に小さく、前記ワークロードは性能上、仮想プロセッサ休止解除から該仮想プロセッサ上のプログラムの実行再開までのレイテンシーの小ささを要求する。

【0051】

一方、仮想プロセッサを休止させるときに、該仮想プロセッサに物理プロセッサを直ちに譲り渡させ、該仮想プロセッサの休止が解除されるまで、前記物理プロセッサを他の仮想プロセッサに割り当てることで、物理プロセッサを効率よく使用することができる。また、前記譲り渡し処理でアイドルの物理プロセッサが残り、アイドルプロセスを実行することがあるが、この場合はアイドルプロセスが物理プロセッサを休止状態におくことで消費電力を削減できる。複数の計算機を1台の物理計算機に集約するサーバ統合の用途で仮想計算機システムを使用する場合は、この物理プロセッサの省電力を含めた利用効率が特に重要である。

40

【0052】

このように仮想計算機システムにおいては、仮想プロセッサの休止に関して仮想プロセッサ単体の性能と、物理プロセッサの省電力を含めた利用効率が共に重要である。

【0053】

50

本発明の仮想プロセッサスケジューリング方式は、物理プロセッサの利用効率を保證できる限りにおいて、仮想プロセッサがワークロード実行中のイベント待ちで該仮想プロセッサを休止させるときに、該仮想プロセッサの休止解除から該仮想プロセッサ上のプログラムの実行再開までのレイテンシーを小さくすることを目的とする。本発明の仮想プロセッサスケジューリング方式は、仮想プロセッサ休止解除待ちの方法として、プロセッサ利用効率を保證するブロック待ち方法と仮想プロセッサ休止解除のレイテンシーを小さくするビジー待ち方法の2つの方法を有し、仮想プロセッサ休止時間の情報に基づいて、仮想プロセッサを休止させるたびに前記2つの方法から1つを動的に選択して実行するものである。

【0054】

10

図3は、本実施形態における、仮想計算機モニタ100の概要を示すブロック図である。仮想計算機モニタ100が有するブロックは、仮想計算機システム、物理プロセッサ、仮想計算機、及び仮想プロセッサの区分に分類される。

【0055】

仮想計算機モニタ001は、仮想計算機システムに対して、ブロック待ち処理時間コストTB(101)、ビジー待ち実行時間しきい値TH(102)を有する。

【0056】

仮想計算機モニタ100は、物理プロセッサに対して、物理プロセッサごとの、物理プロセッサ構造体0(110)と物理構造体1(120)を有する。物理プロセッサ0(002a)の物理プロセッサ構造体0(110)は、物理プロセッサ状態PS(111)、後述の仮想プロセッサ休止解除時刻予測値を物理プロセッサごとに管理する仮想プロセッサ休止解除時刻管理キュー112、当該物理プロセッサ上で走行する仮想プロセッサへのポインタ、ラン状態仮想プロセッサポインタ113、当該物理プロセッサに対応付けられたアイドルプロセスへのポインタ114001を有する。物理プロセッサ1(002b)の物理プロセッサ構造体1(120)も同様である。

20

【0057】

仮想計算機モニタ100は、仮想計算機に対して、仮想計算機ごとに、仮想プロセッサレディキュー130と140、仮想プロセッサブロックキュー131と141を有する。

【0058】

仮想計算機モニタ100は、仮想プロセッサごとに、仮想プロセッサ構造体150a、150b、150c、150dを有する。仮想プロセッサ構造体150aは、仮想プロセッサ状態VS(151a)、スケジュール状態SS(152a)、仮想プロセッサコンテキスト領域153a、仮想プロセッサ休止解除通知フラグNF(154a)、仮想プロセッサ休止時間の予測値E(155a)、仮想プロセッサ休止解除時刻予測値F(156a)、仮想プロセッサ休止開始時刻TS(157a)、仮想プロセッサ休止時間X(158a)、最近N回の仮想プロセッサ休止時間を格納する配列、仮想プロセッサ休止時間履歴テーブルR[i]、 $i = 0, 1, \dots, N - 1$ 159a、前記仮想プロセッサ休止時間履歴テーブルの中で最も古い仮想プロセッサ休止時間が格納されたエントリのインデックス、最古仮想プロセッサ休止インデックスK(160a)、最近N回の仮想プロセッサ休止時間の総和S(161a)を有する。仮想プロセッサ構造体150b、150c、及び150dも同様である。

30

40

【0059】

以上の図1と図3記載の情報をを用いて、本実施形態における仮想プロセッサスケジューリング方式を説明する。

【0060】

最初に、仮想プロセッサのコンテキスト切り替え方式を説明する。第1の仮想計算機200の仮想プロセッサ0(202a)と第2の仮想計算機300の仮想プロセッサ0(302a)、及び物理プロセッサ0(002a)を例とする。仮想プロセッサのコンテキスト切り替えは、仮想プロセッサ間の連続した切り替えとは限らず、仮想プロセッサからアイドルプロセスに切り替え、アイドルプロセスから仮想プロセッサに切り替え、間にアイ

50

ドルプロセスを挟む場合があるが、仮想プロセッサに関して言えば処理に違いがないため、この場合の説明は省略する。

【 0 0 6 1 】

仮想プロセッサコンテキストは、仮想プロセッサのプロセッサレジスタ値と仮想プロセッサ制御情報のスナップショットである。第1の仮想計算機200の仮想プロセッサ0(202a)の仮想プロセッサ構造体150aの仮想プロセッサコンテキスト領域153aは、該仮想プロセッサの仮想プロセッサコンテキストを保持する。仮想計算機システムには複数のアクティブな仮想プロセッサが存在するが、1つの物理プロセッサのハードウェアレジスタは一面である。仮想計算機モニタ100は、仮想プロセッサ0(202a)が物理プロセッサ0(002a)上で走行中の場合は、仮想プロセッサ0(202a)の仮想プロセッサコンテキストを物理プロセッサ0(002a)のハードウェアレジスタに保持する。仮想プロセッサ0(202a)が走行中でない場合は、仮想プロセッサコンテキスト領域153aに仮想プロセッサ0(202a)のコンテキストを保持する。

10

【 0 0 6 2 】

仮想計算機モニタ100は、物理プロセッサ0(002a)に割り当てる仮想プロセッサを第1の仮想計算機200の仮想プロセッサ0(202a)から第2の仮想計算機300の仮想プロセッサ0(302a)に切り替えるときは、まず第1の仮想計算機200の仮想プロセッサ0(202a)の仮想プロセッサコンテキストを仮想プロセッサコンテキスト領域153aに保存し、それから第2の仮想計算機300の仮想プロセッサ0(302a)の仮想プロセッサコンテキストを仮想プロセッサコンテキスト領域153bから物理プロセッサ0(002a)にロードする。その後、物理プロセッサ0(002a)はロードされた仮想プロセッサコンテキストから第2の仮想計算機300の仮想プロセッサ0(302a)の実行を開始する。

20

【 0 0 6 3 】

以上のハードウェアレベルでの仮想計算機モニタ100による仮想プロセッサの切り替えをコンテキスト切り替えと呼ぶ。仮想プロセッサコンテキストは、プロセッサ全体を表すので数も種類も多く、さらには物理プロセッサがプロセッサ仮想化アシスト機能を有する場合には、該機能の制御データの退避回復も必要であるため、コンテキスト切り替えは処理時間の大きな処理である。

30

【 0 0 6 4 】

仮想プロセッサのスケジューリングを行うのに用いる、仮想プロセッサごとに有する状態情報、仮想プロセッサ状態VS、仮想プロセッサ休止解除通知フラグNF、スケジュール状態SSを説明する。第1の仮想計算機001の仮想プロセッサ0001を例とする。

【 0 0 6 5 】

仮想プロセッサ状態VS(151a)は、仮想通常動作状態VEXEと仮想休止状態HHLTの2値をとる。仮想計算機モニタ100は、仮想プロセッサ0(202a)が、仮想通常動作状態VEXEにあるときは、該仮想プロセッサ上のプログラムを実行可とし、仮想プロセッサ0(202a)が仮想休止状態HHLTにあるときは、該仮想プロセッサ上のプログラムを実行不可とする。仮想プロセッサ状態VS(152a)の仮想通常動作状態VEXEから仮想休止状態HHLTへの変更と、仮想休止状態HHLTから仮想通常動作状態VEXEへの変更は共に、仮想プロセッサ0(202a)が物理プロセッサに割り当てられたときに該物理プロセッサ上でのみ行う。

40

【 0 0 6 6 】

仮想プロセッサ休止解除通知フラグNF(154a)は、仮想プロセッサ状態VS(151a)の制御を同期化するためのものである。仮想プロセッサ休止解除通知フラグNF(154a)は、真と偽の2値をとる。仮想プロセッサ休止解除通知フラグNF(154a)の偽から真への変更は仮想プロセッサ0(202a)の物理プロセッサ上での動作とは独立に非同期に任意の時点で行い、仮想プロセッサ休止解除通知フラグNF(154a)の真から偽への変更は、仮想プロセッサ0(202a)が物理プロセッサに割り当てら

50

れたときに、該物理プロセッサ上でのみ行う。

【 0 0 6 7 】

スケジュール状態 $SS(152a)$ は、ラン状態 RUN 、レディ状態 RDY 、及びブロック状態 BLK の 3 値をとる。ラン状態 RUN は、仮想プロセッサ $0(202a)$ が物理プロセッサに割り当てられて、該物理プロセッサを使用している状態である。レディ状態 RDY は、仮想プロセッサ $0(202a)$ が物理プロセッサに割り当て可能だが、別の仮想プロセッサが物理プロセッサ上で走行中の為物理プロセッサが空いていないので停止している状態である。ブロック状態 BLK は、特定のイベントが発生するまではたとえ物理プロセッサが空いていても仮想プロセッサ $0(202a)$ が走行不可の状態である。

【 0 0 6 8 】

仮想計算機モニタ 100 は、仮想プロセッサ $0(202a)$ のスケジュール状態 $SS(152a)$ が物理プロセッサ $0(002a)$ 上でラン状態 RUN にあるときは、該仮想プロセッサを、物理プロセッサ $0(002a)$ の物理プロセッサ構造体 $0(110)$ のラン状態仮想プロセッサポインタ 113 で管理し、仮想プロセッサ $0(202a)$ のスケジュール状態 $SS(152a)$ がレディ状態 RDY 及びブロック状態 BLK にあるときは、それぞれ該仮想プロセッサが属する第 1 の仮想計算機 200 の仮想プロセッサレディキュー 130 と仮想プロセッサブロックキュー 140 で管理する。

【 0 0 6 9 】

仮想計算機モニタ 100 は、物理プロセッサ $0(002a)$ に仮想プロセッサを割り当てるときは、所定のアルゴリズムに従い、いずれかの仮想計算機の仮想プロセッサレディキューから 1 つ仮想プロセッサを取り出して、物理プロセッサ $0(002a)$ に割り当てる。その際、前記仮想プロセッサを物理プロセッサ $0(002a)$ の物理プロセッサ構造体 110 のラン状態仮想プロセッサポインタ 113 にポイントさせる。

【 0 0 7 0 】

仮想プロセッサの仮想プロセッサ状態 SS 、スケジュール状態 VS 、及び仮想プロセッサ休止解除通知フラグ NF を連動させた、仮想プロセッサの状態遷移を説明する。図 4 に仮想プロセッサの状態遷移図を示す。

【 0 0 7 1 】

仮想計算機モニタは仮想プロセッサを物理プロセッサに割り当て、該仮想プロセッサ上のプログラムを実行させるときは、該仮想プロセッサの仮想プロセッサ状態 SS をラン状態 RUN とし、スケジュール状態 VS を仮想通常動作状態 $VEXE$ とする。仮想プロセッサ休止解除通知フラグ NF の値は問わない(状態 1501)。

【 0 0 7 2 】

ステップ 1511 は、前記状態 1501 の仮想プロセッサがタイムスライスを使い切ったかあるいは他に優先度の高い仮想プロセッサが現れた場合に、仮想プロセッサの仮想プロセッサ状態 VS を $VEXE$ にしたままスケジュール状態 SS をレディ状態 RDY に変更し、該仮想プロセッサを属する仮想計算機の仮想プロセッサレディキューに挿入する(状態 2502) 処理を示す。

【 0 0 7 3 】

ステップ 2512 は、前記状態 2502 の仮想プロセッサを物理プロセッサへの割り当て対象に選択したときに、属する仮想計算機の仮想プロセッサレディキューから取り出し、割り当て先の物理プロセッサのラン状態仮想プロセッサポインタに登録して、スケジュール状態 SS をラン状態 RUN に変更して前記状態 1501 に戻る処理を示す。

【 0 0 7 4 】

ステップ 3513 は、前記状態 1501 の仮想プロセッサが仮想休止状態への切り替え要求を受けると、まずスケジュール状態は変更せずに、仮想プロセッサ状態 VS を仮想休止状態 $VHLT$ に変更する(状態 3503) 処理を示す。

【 0 0 7 5 】

ステップ 4514 は、後述のビジー待ち方法によって仮想休止状態の解除を待つ場合に、状態 3503 のまま仮想プロセッサ休止解除通知フラグ NF が真にされるのを繰り返し

10

20

30

40

50

返し調べる処理を示す。

【 0 0 7 6 】

ステップ 5 5 1 5 は、前記のステップ 4 5 1 4 を途中で中断するかあるいは初めから後述のブロック待ち方法によって前記仮想休止状態の解除を待つ場合に、スケジュール状態 S S をブロック状態 B L K に変更して、該仮想プロセッサを属する仮想計算機の仮想プロセッサブロックキューに挿入する（状態 4 5 0 4 ）処理を示す。

【 0 0 7 7 】

ステップ 6 5 1 6 は、前記状態 3 5 0 3 の仮想プロセッサに対して仮想休止状態の解除を行うときに、仮想プロセッサの仮想プロセッサ休止解除通知フラグ N F を真に設定する（状態 5 5 0 5 ）処理を示す。

10

【 0 0 7 8 】

ステップ 6 5 1 6 は、前記状態 4 0 0 1 の仮想プロセッサに対して仮想休止状態の解除を行うときに、仮想プロセッサの仮想プロセッサ休止解除通知フラグ N F を真に設定し、属する仮想計算機の仮想プロセッサブロックキューから取り出し、スケジュール状態 S S をレディ状態 R D Y に変更して、属する仮想計算機の仮想プロセッサレディキューに挿入する（状態 6 5 0 6 ）処理を示す。

【 0 0 7 9 】

ステップ 7 5 1 7 は、前記状態 6 5 0 6 の仮想プロセッサを物理プロセッサへの割り当て対象に選択すると、該仮想プロセッサを属する仮想計算機の仮想プロセッサレディキューから取り出し、割り当て先の物理プロセッサのラン状態仮想プロセッサポイントに登録して、状態 5 5 0 5 に遷移する処理を示す。

20

【 0 0 8 0 】

ステップ 8 5 1 8 は、前記状態 5 5 0 5 にある仮想プロセッサは、仮想プロセッサ休止解除通知フラグ N F が真にされていることを検出して、前記フラグを偽に変更し、仮想プロセッサ状態 V S を V E X E に変更して、前記状態 1 0 0 1 に戻す処理を示す。

【 0 0 8 1 】

ステップ 1 5 1 1 からステップ 9 5 1 9 までの処理のうち、ステップ 1 5 1 1、ステップ 3 5 1 3、ステップ 4 5 1 4、ステップ 5 5 1 5、ステップ 8 5 1 8、及びステップ 9 5 1 9 は、仮想プロセッサが物理プロセッサに割り当てられたときに該物理プロセッサ上で、前記仮想プロセッサと同期して行う。ステップ 2 5 1 2、ステップ 6 5 1 6、及びステップ 7 5 1 7 は、仮想プロセッサの物理プロセッサ上の動作とは独立に非同期に任意の時点で行う。

30

【 0 0 8 2 】

以上の説明を基に、仮想プロセッサのスケジュールリングと同期化を含んだ仮想プロセッサの仮想休止状態と仮想通常動作状態の間の切り替え方を説明する。仮想プロセッサの仮想休止状態と仮想通常動作状態の間の切り替えは、仮想プロセッサ休止解除待ち処理と仮想プロセッサ休止解除通知処理の 2 つの処理から成る。

【 0 0 8 3 】

図 5 に時系列で示した、第 1 の仮想計算機 2 0 0 の仮想プロセッサ 0 (2 0 2 a) の仮想休止状態と仮想通常動作状態の間の切り替えを例に取り上げ、仮想プロセッサ休止解除待ち処理と仮想プロセッサ休止解除通知処理を説明する。

40

【 0 0 8 4 】

まず、図 5 の内容を簡単に説明する。第 1 の仮想計算機 2 0 0 の仮想プロセッサ 0 (2 0 2 a) と仮想プロセッサ 1 (2 0 2 b) がそれぞれ物理プロセッサ 0 (0 0 2 a) と物理プロセッサ 1 (0 0 2 b) 上で走行中である。第 2 の仮想計算機 3 0 0 の仮想プロセッサ 0 3 0 2 と仮想プロセッサ 1 (3 0 2 b) は共に仮想休止状態にあり、物理プロセッサに割り当てられていないものとする。説明を単純にするために図 5 では省略する。

【 0 0 8 5 】

仮想プロセッサ 0 (2 0 2 a) 上のプログラムが休止状態への切り替え要求を発行すると、仮想計算機モニタ 1 0 0 は該要求をハードウェアに直接実行させずに検出し、仮想プ

50

ロセッサ状態を仮想休止状態に変更し、仮想プロセッサ0(202a)上のプログラムの実行を停止させる。仮想プロセッサ1(202b)上のプログラムが仮想プロセッサ0(202b)に仮想プロセッサ間割り込みを発行すると、仮想計算機モニタ100は、該割り込みを検出して、該割り込みの動作を仮想的に実現すると共に、仮想プロセッサ0(202a)の仮想プロセッサ状態を仮想通常動作状態に切り替え、仮想プロセッサ0(202a)上のプログラムの実行を再開する。

【0086】

図6に仮想プロセッサ休止除待ち処理のフローチャートを示す。図5の例に沿って、図6のフローチャートに従い、仮想プロセッサ休止解除待ち処理を説明する。

【0087】

仮想プロセッサ0(202a)と仮想プロセッサ1(202b)が共に状態1501にあるとする。仮想プロセッサ0(202a)上のプログラムが休止状態への切り替え要求を発行すると、仮想計算機モニタ100は、該要求をハードウェアに直接実行させずに検出して、仮想プロセッサ状態VS(151a)を仮想休止状態VHLTに変更する(ステップA600)。

【0088】

仮想計算機モニタ100は、後述する仮想プロセッサ休止解除待ち方法選択処理にて仮想プロセッサ0の仮想プロセッサ休止解除待ち処理としてブロック待ち方法を選択したならば(ステップB602)、スケジュール状態SS(152a)をBLKに変更して、仮想プロセッサ0(202a)を第1の仮想計算機200の仮想プロセッサブロックキュー131に挿入して(ステップC602)、コンテキスト切り替えを実行して(ステップD603)物理プロセッサ0(002a)を譲り渡す。

【0089】

仮想計算機モニタ100は、物理プロセッサ0(002a)に割り当てる仮想プロセッサが見つからないので、物理プロセッサ0(002a)の物理プロセッサ構造体0(110)のアイドルプロセスポインタ114がポイントするアイドルプロセスを実行する。

【0090】

一方、仮想計算機モニタ100は、後述する仮想プロセッサ休止解除待ち方法選択処理にてビジー待ち方法を選択したならば(ステップB601)、仮想プロセッサ休止解除通知フラグNF(154a)が真になるのを繰り返し調べる(ステップE604)。仮想プロセッサ休止解除通知フラグNF(154a)が真に設定されるのを確認する前に後述のビジー待ち方法継続チェック処理にてビジー待ち方法を途中で中断した場合は(ステップF605)、初めからブロック待ち方法を選択したのと同様に、ステップB602とステップD603を実行して物理プロセッサ0(002a)を譲り渡す。

【0091】

ビジー待ち方法を選択しステップE604にて後述の仮想プロセッサ休止解除通知処理によって、仮想プロセッサ休止解除通知フラグNF(154a)が真にされたのを確認すると、仮想プロセッサ休止解除通知フラグNF(154a)を偽にして(ステップH607)、仮想プロセッサ状態VS(151a)を仮想実行状態VEXEに変更して(ステップI608)、仮想プロセッサ0(202a)上のプログラムの実行を再開する。

【0092】

ブロック待ち方法を選択したときは、後述の仮想プロセッサ休止解除通知処理によって起こされ、さらに物理プロセッサ0(002a)への割り当て対象として選択され、コンテキスト切り替えにされて物理プロセッサ0(002a)上で実行を再開すると(ステップG606)、仮想プロセッサ休止解除通知フラグNF(154a)を偽にして(ステップH607)、仮想プロセッサ状態VS(151a)を仮想実行状態VEXEに変更して(ステップI608)、仮想プロセッサ0(202a)上のプログラムの実行を再開する。

【0093】

図7に仮想プロセッサ休止解除通知処理のフローチャートを示す。続けて仮想プロセッサ

10

20

30

40

50

サ休止解除通知処理を説明する。

【 0 0 9 4 】

仮想プロセッサ 1 (2 0 2 b) が仮想プロセッサ 0 (2 0 2 a) に仮想プロセッサ間割り込みを発行すると、仮想計算機モニタ 1 0 0 が該割り込み発行を検出し、該割り込みの動作を仮想的に実現する。それから仮想プロセッサ休止解除通知処理を開始する。仮想プロセッサの状態によらずに仮想プロセッサ休止解除通知フラグ NF (1 5 4 a) を真にする (ステップ J 6 0 9)。仮想プロセッサ 0 (2 0 2 a) のスケジュール状態 SS (1 5 2 a) が B L K であるならば (ステップ K 6 1 0)、仮想プロセッサ 0 (2 0 2 a) を第 1 の仮想計算機 2 0 0 の仮想プロセッサブロックキュー 1 3 1 から取り出して、スケジュール状態 SS (1 5 2 a) をレディ状態 R D Y に変更して第 1 の仮想計算機 2 0 0 の仮想プロセッサレディキュー 1 3 0 に挿入する (ステップ K 6 1 1)。

10

【 0 0 9 5 】

仮想プロセッサ 0 (2 0 2 a) のスケジュール状態 SS (1 5 2 a) が B L K 以外、すなわちラン状態 R U N あるいはレディ状態 R D Y の場合は以降何もしない。仮想計算機モニタ 1 0 0 は物理プロセッサ 0 (0 0 2 a) への割り当て対象として仮想プロセッサ 0 (2 0 2 a) を選択すると、スケジュール状態 SS (1 5 2 a) をラン状態 R U N に変更して、仮想プロセッサ 0 (2 0 2 a) の仮想プロセッサコンテキストを物理プロセッサ 0 (0 0 2 a) に回復して、仮想プロセッサ 0 (2 0 2 a) を起動する (ステップ G 6 0 6)。

【 0 0 9 6 】

図 5 の例を基にブロック待ち方法とビジー待ち方法の性能と処理コストを考察する。性能は、仮想プロセッサ休止解除レイテンシーの大きさで見る。処理コストは、仮想プロセッサが休止解除待ちを開始してから該仮想プロセッサ上のプログラムの実行が再開されるまでの間に、物理プロセッサを該仮想プロセッサの休止を制御すること以外に使えるプロセッサ時間の割合、空きプロセッサ利用率で見る。

20

【 0 0 9 7 】

図 6 と図 7 のステップ A 6 0 0 からステップ L 6 1 1 までのそれぞれの処理時間コストの大小関係は次のとおりである。

【 0 0 9 8 】

ステップ A コスト、ステップ B コスト、ステップ F コスト、ステップ G コスト、ステップ H コスト、ステップ I コスト < < ステップ C コスト、ステップ D コスト、ステップ E コスト、ステップ J コスト、ステップ K コスト

30

ステップ D 6 0 3 とステップ K 6 1 0 はコンテキスト切り替え処理のため処理時間が非常に大きい。ステップ C 6 0 2 とステップ J 6 0 9 はキュー操作処理であり、コンテキスト切り替え処理に比べれば小さいが処理時間は大きい。ステップ E 6 0 4 はビジー待ち処理のため処理時間は最大になりえるが、このステップ E 6 0 4 がコストとなるかどうかは、その間他に行う処理があるかどうかによって決まる。残りのステップ A 6 0 0、ステップ B 6 0 1、ステップ F 6 0 5、ステップ G 6 0 6、ステップ H 6 0 7、及びステップ I 6 0 8 はビット参照・操作処理のため、処理時間は全体からすれば非常に小さい。

40

【 0 0 9 9 】

図 8 (a) はブロック待ち方法を選択した場合、図 8 (b) はビジー待ち方法を選択したが途中でブロック待ちに変更した場合、図 8 (c) はビジー待ち方法を選択した場合の、仮想プロセッサ休止解除待ち処理と仮想プロセッサ休止解除通知処理を時系列で示したものである。それぞれの場合の仮想プロセッサ休止解除レイテンシーと物理プロセッサ 0 の空きプロセッサ利用率を示す。以下では図と同様に例えばステップ A コストを A と略記することで記述を簡略化する。

【 0 1 0 0 】

図 8 (a) に示すブロック待ち方法を選択した場合の仮想プロセッサ休止解除レイテンシーと空きプロセッサ利用率はそれぞれ次のとおりである。

50

ブロック待ちの仮想プロセッサ休止解除レイテンシー = $(H + I + J) + (K + G)$
 ブロック待ちの物理プロセッサ 0 の空きプロセッサ利用率 = $(VP \text{ 休止時間} - (A + B + C + D) + (H + I + J)) / (VP \text{ 休止時間} + (H + I + J) + (K + G))$

【0101】

図8(b)に示すビジー待ち方法を選択したが途中でブロック待ち方法に変更した場合の仮想プロセッサ休止解除レイテンシーと空きプロセッサ利用率はそれぞれ次のとおりである。

ビジー待ち ブロック待ちの仮想プロセッサ休止解除レイテンシー = $(H + I + J) + (K + G)$

ビジー待ち ブロック待ちの物理プロセッサ 0 の空きプロセッサ利用率 = $(VP \text{ 休止時間} - (A + B + E + F + C + D) + (H + I + J)) / (VP \text{ 休止時間} + (H + I + J) + (K + G))$

10

【0102】

図8(c)に示すビジー待ち方法を選択した場合の仮想プロセッサ休止解除レイテンシーと空きプロセッサ利用率はそれぞれ次のとおりである。

ビジー待ちの仮想プロセッサ休止解除レイテンシー = $H + (K + G)$

ビジー待ちの物理プロセッサ 0 の空きプロセッサ利用率 = 0

【0103】

仮想プロセッサ休止解除レイテンシーを比較すると、以下のとおりである。

ビジー待ち << ブロック待ち = (ビジー待ち ブロック待ち)

20

【0104】

同様に空きプロセッサ利用率を比較すると、以下のとおりである。

0 = ビジー待ち << (ビジー待ち ブロック待ち) < ブロック待ち

【0105】

ブロック待ち方法は、空きプロセッサ利用率が最も高い、つまり物理プロセッサ利用効率の高い方法である。ビジー待ち方法は、仮想プロセッサ休止解除レイテンシーが最も小さく性能面で優れている。ビジー待ち方法は空きプロセッサ利用率が0であり、物理プロセッサを他の用途には使うことはできないが、後述のように、ビジー待ち実行中は物理プロセッサを休止させることで消費電力を減らすことができるので、他に動作させる仮想プロセッサが存在しない場合は問題とならない。ビジー待ち方法を途中でやめてブロック待ち方法に切り替えた場合は、最初からブロック待ち方法を選択した場合に比べて空きプロセッサ利用率は低いが、仮想プロセッサ休止解除レイテンシーは最初からブロック待ち方法を選択した場合と同じである。ビジー待ち方法から途中でブロック待ち方法に切り替えることは避けなければならない。

30

【0106】

以上から、前述の仮想プロセッサ休止解除待ち方法選択処理の方針を、プロセッサ利用効率が要求される場合はブロック待ち方法を選択し、他に動作させる仮想プロセッサが存在せずかつ仮想プロセッサ休止解除レイテンシーの小ささが要求される場合にビジー待ち方法を選択し、ビジー待ち方法から途中でブロック待ち方法に切り替えることを避けることと定める。

40

【0107】

次に、この方針を実現するために使用する、仮想プロセッサ休止時間、仮想プロセッサ休止時間予測値、及び仮想プロセッサ休止解除時刻予測値の算出方法と管理方法を説明する。第1の仮想計算機200の仮想プロセッサ0(202a)を例とする。

【0108】

仮想プロセッサ0(202a)の仮想プロセッサ休止時間を、仮想プロセッサ0(202a)の仮想プロセッサ休止解除通知フラグNF(154a)が真に設定されるのを待ち始めてから、該仮想プロセッサ休止解除通知フラグNFが真に設定されるまでの時間と定義する。直近のN回の仮想プロセッサ休止時間の平均値を仮想プロセッサ休止時間予測値Eとする。

50

【0109】

仮想プロセッサ0の仮想プロセッサ休止時間Xの算出方法を、図9(a)のフローチャートに基づき説明する。

【0110】

上述の仮想プロセッサ休止解除待ち処理において、仮想プロセッサ0(202a)の仮想プロセッサ状態VS(151a)を仮想休止状態VHLTに変更した後で仮想プロセッサ0(202a)の仮想プロセッサ休止解除通知フラグNF(154a)が真に設定されるのを待ち始める直前に、その時点のシステム時刻を仮想プロセッサ0(202a)の仮想プロセッサ休止開始時刻TS(157a)に保持する(ステップ620)。

【0111】

上述の仮想プロセッサ休止解除通知処理において、仮想プロセッサ0(202a)の仮想プロセッサ休止解除通知フラグNF(154a)を真に設定した直後に、その時点のシステム時刻t1から、保存した仮想プロセッサ休止開始時刻TS(157a)を引くことで、仮想プロセッサ0(202a)のプロセッサ休止時間t1-TSを求め、仮想プロセッサ休止時間X(158a)に保存する(ステップ621)。

【0112】

仮想プロセッサ0(202a)の仮想プロセッサ休止時間予測値E(155a)の算出方法を、図9(b)のフローチャートに基づき説明する。図10にはプロセッサ休止時間予測値E(155a)の算出方法の概念図を示す。

【0113】

仮想プロセッサ休止時間予測値E(155a)の算出は、前記の仮想プロセッサ休止時間X(158a)の算出に続けて、仮想プロセッサ休止解除通知処理の中で行う。

【0114】

上述の仮想プロセッサ休止解除通知処理において、算出した仮想プロセッサ休止時間X(158a)によって、直近N回の仮想プロセッサ休止時間の総和S(161a)を更新する。最古仮想プロセッサ休止インデックスK(160a)は仮想プロセッサ休止時間履歴テーブルR[i]、i=0、1、・・・、N-1 159aの中で最古の仮想プロセッサ休止時間をさしているの、直近N回の仮想プロセッサ休止時間の総和S(161a)から前記最古の仮想プロセッサ休止時間R[K]を引き、今回の仮想プロセッサ休止時間X(158a)を足すことで、直近N回の仮想プロセッサ休止時間の総和S(161a)を更新する、すなわち(式1)と求める(ステップ622)。

$$S_{new} = S_{old} - R[K] + X \quad (\text{式1})$$

【0115】

この更新した直近N回の仮想プロセッサ休止時間の総和SをNで割ることで、最近N回の仮想プロセッサ休止時間の平均値、次の仮想プロセッサ休止時間の予測値E(155a)を求める、すなわち(式2)と求める(ステップ623)。

$$E = S / N \quad (\text{式2})$$

【0116】

最後に、仮想プロセッサ休止時間履歴テーブルR[i]、i=0、1、・・・、N-1 159aを更新する。最古仮想プロセッサ休止インデックスK(160a)が指す仮想プロセッサ休止時間履歴テーブルのエントリR[K]に直近仮想プロセッサ休止時間X(157a)を保存し、すなわち(式3)とし(ステップ624)、

$$R[K] = X \quad (\text{式3})$$

最古仮想プロセッサ休止インデックスKに+1を足し、KがNと等しくなったならば、Kを0とする、すなわち(式4)と求める(ステップ625)。

$$K = (+ + K) \bmod N \quad (\text{式4})$$

【0117】

(式1)、(式2)、(式3)、及び(式4)を見ると分かるように、仮想プロセッサ休止時間予測値E(155a)を、直近N回の仮想プロセッサ休止時間の総和を求めなおす必要なく、仮想プロセッサ0の情報のみを用いて一定の処理時間で算出することができ

10

20

30

40

50

る。

【0118】

仮想プロセッサ0(202a)の仮想プロセッサ休止解除時刻予測値F(156a)の算出方法と管理方法を、図9(c)のフローチャートに基づき説明する。第1の仮想計算機200の仮想プロセッサ0(202a)が物理プロセッサ0(002a)上で動作すると仮定する。

【0119】

上述の仮想プロセッサ休止解除待ち処理において、仮想プロセッサ0(202a)の仮想プロセッサ状態VS(152a)を仮想休止状態VHLTに変更した後で仮想プロセッサ0(202a)の仮想プロセッサ休止解除通知フラグNF(154a)が真に設定されるのを待ち始める直前に、その時点のシステム時刻t0を仮想プロセッサ0の仮想プロセッサ休止開始時刻TS(157a)に保持した後で、前記時刻T0に仮想プロセッサ休止時間予測値E(155a)を加算し、(式5)の通り仮想プロセッサ休止解除時刻予測値F(156a)を求める(ステップ626)。

$$F = TS + E \quad (\text{式5})$$

【0120】

求めた仮想プロセッサ休止解除時刻予測値F(156a)を現在動作中の物理プロセッサ0(002a)の仮想プロセッサ休止解除時刻管理キュー112の適切な位置に挿入する(ステップ627)。仮想プロセッサ0(202a)は仮想プロセッサ解除通知フラグNF(154a)が真に設定されたことを検出すると、前記の物理プロセッサ0(002a)の仮想プロセッサ休止解除時刻管理キュー112から仮想プロセッサ休止解除時刻予測値F(156a)を取り出す(ステップ628)。

【0121】

本発明の仮想プロセッサスケジューリング方式で仮想プロセッサ休止時間予測値を使う理由を説明する。

【0122】

上述のように、汎用オペレーティングシステムがプロセッサを休止させるのは、大きくはプロセッサがアイドルの場合とワークロード実行中のイベント待ちの場合がある。プロセッサがワークロード実行中の同期化イベント待ちでのプロセッサ休止の周期とプロセッサ休止時間は、プロセッサがアイドルの場合と比べて非常に短い。共通点としては、プロセッサがワークロード実行中の同期化イベント待ちの場合とプロセッサがアイドルの場合のプロセッサ休止時間は共に定常的にほぼ一定の値をとる。

【0123】

仮想計算機モニタ100が行いたいのは、仮想プロセッサが休止する要因がワークロード実行中の同期化イベント待ちによるものであるかそうでないかの判定、及び仮想プロセッサの休止が解除される時刻の推定である。

【0124】

仮想プロセッサの動作を判定する情報には、プロセッサ利用率がある。同期化イベント待ちを頻発させるワークロードには、仮想プロセッサの仮想的な休止状態と仮想的な通常状態の切り替えは頻発するが、仮想プロセッサの利用率自体は低いものがある。この場合は、プロセッサ利用率を見ても、仮想プロセッサが休止する要因がワークロード実行中のイベント待ちによるものであることを判定することができない。また、プロセッサ利用率の情報では、仮想プロセッサの休止が解除される時刻を推定することができない。

【0125】

これに対して同期化イベント待ちを頻発するワークロード実行中のプロセッサ休止時間には、上述の通り定常的にほぼ一定でかつ非常に小さな値をとる、という特徴がある。プロセッサ休止時間がこの特徴をもつならば、直近の複数回のプロセッサ休止時間の平均値と、次回のプロセッサ休止時間がほぼ等しくなる。つまり、直近の複数回のプロセッサ休止時間の平均値を求めることで、該平均値を次回のプロセッサ休止時間の推定値として使うことができる。

10

20

30

40

50

【 0 1 2 6 】

また、直近の複数回のプロセッサ休止時間の平均値は、最古のプロセッサ休止時間を引いて最新のプロセッサ休止時間を足すことで、総和を求めなおす必要なく算出することができる。プロセッサ利用率のような課題は生じない。仮に予測が外れたとしても、高々仮想プロセッサのコンテキスト切り替えの処理時間が2倍になる程度であり、ワークロード実行時間から見れば無視できるくらいに小さい。また、仮想プロセッサ休止解除時刻の予測値は、仮想プロセッサが休止を開始するときに現在システム時刻に該仮想プロセッサの仮想プロセッサ休止時間予測値を足すことによって、容易に算出することができる。それゆえ、本発明の仮想プロセッサスケジューリング方法では、仮想プロセッサ休止時間予測値を使用する。

10

【 0 1 2 7 】

前述の方針を仮想プロセッサ休止時間の情報を用いて実現する、前述の仮想プロセッサ休止解除待ち方法選択処理と、前述のビジー待ち方法継続チェック処理を説明する。第1の仮想計算機の仮想プロセッサ0が物理プロセッサ0上で動作する場合を例とする。

【 0 1 2 8 】

ここで、以下ではしきい値情報として、ブロック待ち処理時間コストTB(101)とビジー待ち実行時間しきい値TH(102)を用いる。ブロック待ち処理時間コストTB(101)には、ブロック待ち方法を適用したが待ち時間なしで直ちに起こされた場合の処理時間を予め求めて、該値を設定する。つまり、前述のステップA、B、C、D、H、I、J、K、及びGの合計値を予め求めて設定する。ビジー待ち実行時間しきい値TH(102)は、仮想プロセッサ休止解除レイテンシーの小ささを要求する場合と要求しない場合を分けられる値とする。ある1つの実現形態においては、ビジー待ち実行時間しきい値TH(102)は200マイクロ秒としている。

20

【 0 1 2 9 】

まず、仮想プロセッサ休止解除待ち方法選択処理の条件判定を1つずつ順に説明する。それぞれの条件判定の説明は、条件の意味を例示する図13を用いて行う。これまでと同様、第1の仮想計算機200の仮想プロセッサ0(202a)が物理プロセッサ0(002a)上で動作している場合を例とする。

【 0 1 3 0 】

仮想プロセッサ0(202a)の仮想プロセッサ休止時間予測値E(155a)が、所定のブロック待ち処理時間コストTB(101)未満であるならば、現在動作中の物理プロセッサ0(002a)への割り当て対象の仮想プロセッサが存在するか否かを問わず、ビジー待ち方法を選択する。この場合は、仮想プロセッサ0 002aの休止時間がブロック待ち処理時間コストTB(101)未満のため、プロセッサ利用率と性能の両方の観点でビジー待ち方法の適用が望ましい。それゆえビジー待ち方法を選択する(図13(a)を参照)。

30

【 0 1 3 1 】

次に、仮想プロセッサ0(202a)の仮想プロセッサ休止時間予測値E(155a)が、前記ブロック待ち処理時間コストTB(101)以上であり、物理プロセッサ0(002a)への割り当て対象の仮想プロセッサが存在する場合には、ブロック待ち方法を選択する。この場合は、仮想プロセッサ0 002aをビジー待ちさせるとプロセッサ利用率を下げるので、ブロック待ち方法を選択する。

40

【 0 1 3 2 】

また、仮想プロセッサ0(202a)の仮想プロセッサ休止時間予測値E(155a)が、該仮想プロセッサが現在所定のビジー待ち実行時間しきい値TH(102)以上である場合には、ブロック待ち方法を選択する。この場合は、仮想プロセッサ0(202a)上の第1の汎用オペレーティングシステム201が仮想プロセッサ休止解除レイテンシーの小ささを要求しない状況であるので、プロセッサ利用率を下げるリスクをなくすために、ブロック待ち方法を選択する(図13(b)を参照)。

【 0 1 3 3 】

50

次に、仮想プロセッサ0(202a)の仮想プロセッサ休止時間予測値E(155a)が、前記ブロック待ち処理時間コストTB(101)以上でかつ前記ビジー待ち実行時間しきい値TH(102)未満であり、物理プロセッサ0(002a)への割り当て対象の仮想プロセッサが存在せず、該仮想プロセッサの仮想プロセッサ休止解除時刻予測値E(155a)が、物理プロセッサ0(002a)の仮想プロセッサ休止解除時刻管理キュー112の先頭の仮想プロセッサ休止解除時刻F以上であるならば、ブロック待ち方法を選択する。この場合は、仮想プロセッサ0(202a)が物理プロセッサ0(002a)を譲り渡した直後は、他に該物理プロセッサに割り当てられる仮想プロセッサは存在しないが、仮想プロセッサ0(202a)の休止状態が解除されるより前に他の仮想プロセッサの休止状態が解除されるので、仮想プロセッサ0(202a)が物理プロセッサ0(002a)上にとどまると、仮想プロセッサ0(202a)のコンテキスト切り替えを遅らせることになり、結果として仮想プロセッサ0(202a)の代わりに動作させるべき前記仮想プロセッサ上のプログラムの起動を遅らせる。それゆえ、ブロック待ち方法を選択する(図13(c)を参照)。

10

【0134】

また、仮想プロセッサ0(202a)の仮想プロセッサ休止時間予測値E(155a)が、前記ブロック待ち処理時間コストTB(101)以上でかつ前記ビジー待ち実行時間しきい値TH(102)未満であり、現在動作中の物理プロセッサ0(002a)への割り当て対象の仮想プロセッサが存在せず、仮想プロセッサ0(202a)の仮想プロセッサ休止解除時刻予測値F(156a)が、物理プロセッサ0(002a)の仮想プロセッサ休止解除時刻管理キュー113の先頭の仮想プロセッサ休止解除時刻予測値F未満である場合には、ビジー待ち方法を選択する。

20

【0135】

図11に仮想プロセッサの仮想プロセッサ休止解除待ち方法選択処理のフローチャートを示す。第1の仮想計算機200の仮想プロセッサ0(202a)が物理プロセッサ0(002a)上で動作する場合を例として説明する。

【0136】

まず、 $E < TB$ であるならば(ステップ630)ビジー待ち方法を選択し、そうでないならば次の条件判定に進み、現在動作中の物理プロセッサ0(002a)に割り当てられる仮想プロセッサが他に存在するならば(ステップ631)ブロック待ち方法を選択し、そうでないならば次の条件判定に進み、 $E \geq TH$ であるならば(ステップ632)ブロック待ち方法を選択し、そうでないならば次の条件判定に進み、 $F <$ 物理プロセッサ0(002a)の仮想プロセッサ休止解除時刻管理キュー113の先頭の仮想プロセッサ休止解除時刻F未満であるならば(ステップ633)ビジー待ち方法を選択し、そうでないならばブロック待ち方法を選択する。

30

【0137】

図12に仮想プロセッサの仮想プロセッサ休止解除待ち処理におけるビジー待ち方法継続チェック処理のフローチャートを示す。第1の仮想計算機200の仮想プロセッサ0(202a)を例に説明する。

【0138】

仮想プロセッサ休止解除待ち処理の中のビジー待ち処理において、前記算出した仮想プロセッサ休止開始時刻TS(157a)から現在システム時刻を引いて、ビジー待ち実行時間 = $TS -$ 現在システム時刻を求め(ステップ640)、該ビジー待ち実行時間がビジー待ち実行時間しきい値THよりも大きいならば(ステップ641)、ビジー待ち方法をやめてブロック待ち方法に切り替える。本発明では、このチェック処理にかかる確率が低いことを期待するが、万一予測が外れた場合にもその影響を小さく留めるためにこのチェック処理を設ける。

40

【0139】

ここまでで、仮想プロセッサ休止解除の性能と物理プロセッサの利用効率の両方を考慮した仮想プロセッサスケジューリング方式を説明した。仮想計算機システムにおいては、

50

物理プロセッサの省電力化も重要である。以下では本発明の仮想プロセッサスケジューリング方式における物理プロセッサの電力制御方式を説明する。まず、物理プロセッサの通常動作状態と休止状態の切り替え方式、次に、前記の仮想プロセッサスケジューリングと連動した物理プロセッサ電源制御方式を説明する。以下の方式は必ずしも組み込まなくてもよいが、組み込むことによって物理プロセッサの消費電力を抑えることができる。

【0140】

物理プロセッサの通常動作状態と休止状態の切り替え方式を説明する。物理プロセッサ0(002a)とその物理プロセッサ構造体0(110)を例とする。物理プロセッサ1(002b)とその物理プロセッサ構造体1(120)も同様に説明できる。物理プロセッサの物理プロセッサ状態PS 111は、仮想プロセッサの仮想プロセッサ状態VSと同様に通常動作状態PEXEと休止状態PHLTの2値をとる。仮想計算機モニタ100は、物理プロセッサ0(002a)を通常動作状態から休止状態に切り替えるときは、該切り替えの直前に物理プロセッサ構造体0(110)の物理プロセッサ状態PS 111を休止状態PHLTに設定する。仮想計算機モニタ100は、物理プロセッサ0(002a)を休止状態から通常動作状態に切り替えるときは、物理プロセッサ0(002a)が通常動作状態に戻り、仮想計算機モニタ100が物理プロセッサ0(002a)上で実行を開始した直後に、物理プロセッサ構造体0(110)の物理プロセッサ状態PS 111を通常動作状態PEXEに設定する。

10

【0141】

前記の仮想プロセッサスケジューリングと連動した物理プロセッサ電源制御方式を説明する。

20

【0142】

図14に、物理プロセッサの休止処理を組み込んだアイドルプロセスのアイドルループと、前記物理プロセッサの休止解除処理を組み込んだ仮想プロセッサ割り当て要求通知処理をフローチャートで示す。物理プロセッサ0(002a)上でアイドルプロセス020aが動作する場合を例とする。

【0143】

物理プロセッサ0(002a)上で動作させる仮想プロセッサがなくなると、該物理プロセッサに対応付けられたアイドルプロセス020aが、アイドルループの実行を開始する。アイドルループでは、スピンドル実行による物理プロセッサ資源の浪費を防ぐために物理プロセッサ0(002a)を休止させる。仮想計算機モニタ100が物理プロセッサ0(002a)を指定して仮想プロセッサを割り当てようとするときに、物理プロセッサ状態PS 111が休止状態PHLTであるならば、物理プロセッサ0(002a)にプロセッサ間割り込みを発行し、該物理プロセッサの休止状態を解除する。休止していた物理プロセッサ0(002a)は、休止が解除されるとアイドルループの実行を再開する。該物理プロセッサに割り当てる仮想プロセッサを見つけると、該仮想プロセッサへのコンテキスト切り替えを実行し、該仮想プロセッサを起動する。仮想プロセッサが見つからないならば、再び物理プロセッサ0(002a)を休止させる。

30

【0144】

図15に、物理プロセッサの休止処理を組み込んだ仮想プロセッサの仮想休止状態解除待ち処理の中のビジー待ち処理と、物理プロセッサの休止解除処理を組み込んだ仮想プロセッサの仮想休止状態解除通知処理をフローチャートで示す。第1の仮想計算機の仮想プロセッサ0が物理プロセッサ0上で動作する場合を例とする。

40

【0145】

仮想プロセッサ0は、仮想休止状態解除待ち処理でビジー待ち処理を選択すると、物理プロセッサ0を休止させても前記のビジー待ち方法継続チェック処理が動作するように、その時点からビジー待ち実行時間しきい値THを経過時にタイマ割り込みが物理プロセッサ0上で発生し、物理プロセッサ0の休止状態が解除されるように、タイマを設定する。ビジー待ち方法のループ内で物理プロセッサを休止させる。

【0146】

50

前記の仮想プロセッサ0の仮想プロセッサ休止解除通知フラグNFが真であることを検出するか、前記ビジー待ち方法継続チェック処理にかかると、ビジー待ち方法を終了する。このとき、前記設定したタイマがまだ有効であるならば、該タイマをキャンセルする。

【0147】

仮想計算機モニタ100が、仮想プロセッサ0(202a)への仮想的な休止状態の切り替え要求を検出して、仮想プロセッサ0(202a)の仮想プロセッサ休止解除通知フラグNF(154a)を真にした後に仮想プロセッサ0(202a)の状態を調べ、仮想プロセッサ0(202a)が物理プロセッサ0(002a)上でビジー待ち処理を実行中でかつ物理プロセッサ0(002a)が休止状態にあるならば、物理プロセッサ0(002a)にプロセッサ間割り込みを発行する。

10

【0148】

本発明の仮想プロセッサスケジューリング方式によれば、仮想プロセッサ休止時間予測値を用いることにより、ブロック待ち方法のみを有していた従来方法と同等の物理プロセッサ利用効率と消費電力の削減を保証する限りにおいて、仮想プロセッサの休止が頻発する場合には、該仮想プロセッサの仮想的な休止状態から仮想的な通常動作状態への切り替えの処理時間を小さくして、該仮想プロセッサ上のプログラムの起動を速くすることができる。本実施形態においては、図1と図2の構成に基づきこれを示した。

【符号の説明】

【0149】

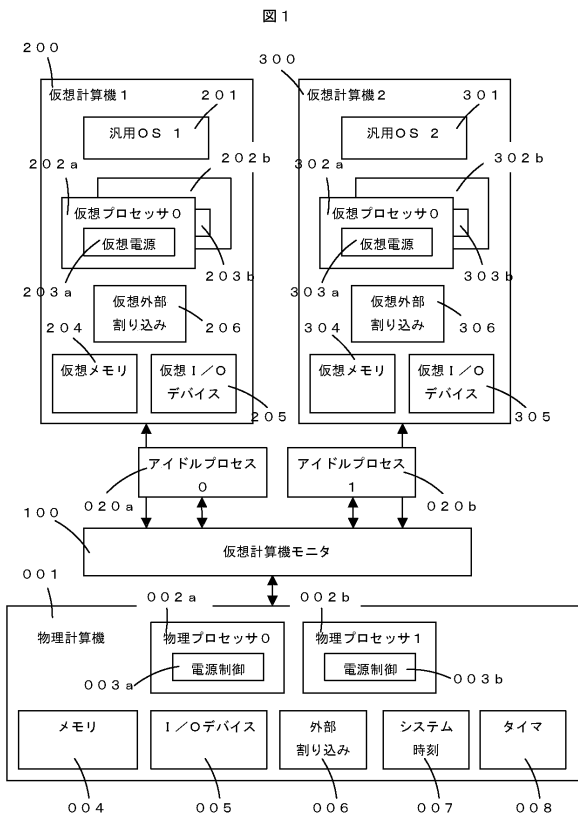
001	物理計算機	20
002a、002b	物理プロセッサ	
003a、003b	電源制御機能	
004	メモリ	
005	I/Oデバイス	
006	外部割り込み機構	
007	システム時刻	
008	タイマ	
100	仮想計算機モニタ	
020a、020b	アイドルプロセス	
200、300	仮想計算機	30
201、301	汎用オペレーティングシステム	
202a、202b、302a、302b	仮想プロセッサ	
203a、203b、303a、303b	仮想電源制御機能	
204、304	仮想メモリ	
205、305	仮想I/Oデバイス	
206、306	仮想外部割り込み機構	
110、120	物理プロセッサ構造体	
111、121	物理プロセッサ状態PS	
112、122	仮想プロセッサ休止解除時刻管理キュー	
113、123	ラン状態仮想プロセッサポインタ	40
114、124	アイドルプロセスポインタ	
130、140	仮想プロセッサレディキュー	
131、141	仮想プロセッサブロックキュー	
150a、150b、150c、150d	仮想プロセッサ構造体	
151a、151b、151c、151d	仮想プロセッサ状態VS	
152a、152b、152c、152d	スケジュール状態VS	
153a、153b、153c、153d	仮想プロセッサコンテキスト領域	
154a、154b、154c、154d	仮想プロセッサ休止解除通知フラグNF	
155a、155b、155c、155d	仮想プロセッサ休止時間予測値E	
156a、156b、156c、156d	仮想プロセッサ休止解除時刻予測値F	50

- 157 a、157 b、157 c、157 d 仮想プロセッサ休止開始時刻 T S
- 158 a、158 b、158 c、158 d 仮想プロセッサ休止時間 X
- 159 a、159 b、159 c、159 d 仮想プロセッサ休止時間履歴テーブル R [i]、i = 0、...、N - 1
- 160 a、160 b、160 c、160 d 最古仮想プロセッサ休止インデックス K
- 161 a、161 b、161 c、161 d 最近 N 回の仮想プロセッサ休止時間の総和 S
- 501、502、503、504、505、506 仮想プロセッサ状態遷移における状態 1 ~ 状態 6
- 511、512、513、514、515、516、517、518、519 仮想プロセッサ状態遷移図における状態遷移のステップ 1 ~ ステップ 9
- 600、601、602、603、604、605、606、607、608、609、610、611 仮想プロセッサ休止解除待ち処理と仮想プロセッサ休止解除通知処理におけるステップ A からステップ L
- 620、621、622、623、624、625、626、627 仮想プロセッサ休止時刻と該時刻の予測値算出処理の処理ステップ
- 630、631、632、633 仮想プロセッサ休止解除待ち方法選択区処理のステップ
- 640、641 仮想プロセッサビジー待ち継続チェック処理のステップ
- 650、651、652、653、654、655、656、657 物理プロセッサ電力制御を組み込んだアイドルループの処理ステップ
- 660、661、662、663、664、665、666、667、668、669、670 物理プロセッサ電力制御を組み込んだ仮想プロセッサ休止解除待ち処理におけるビジー待ち方法の処理ステップ

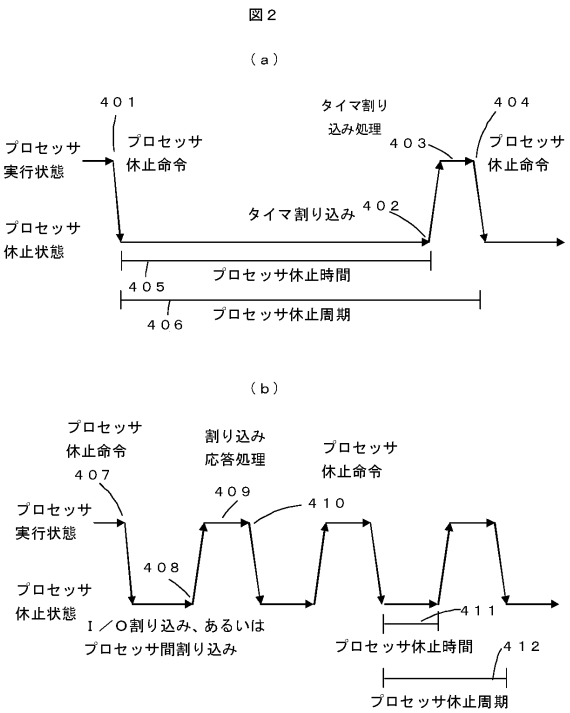
10

20

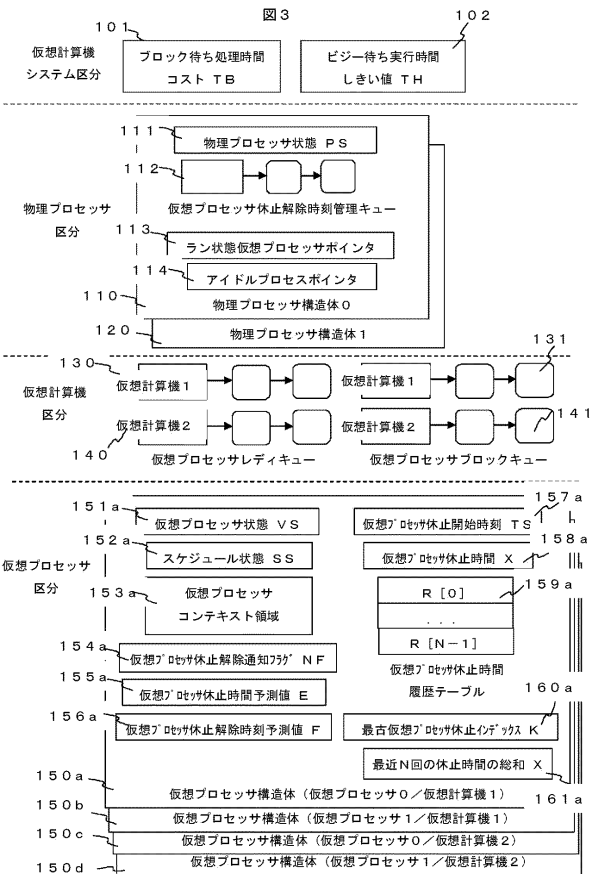
【図 1】



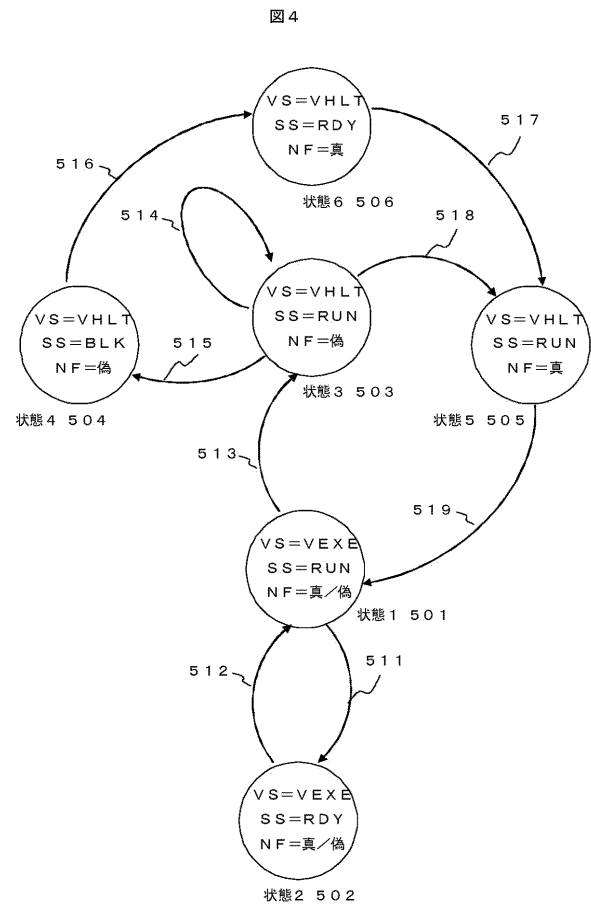
【図 2】



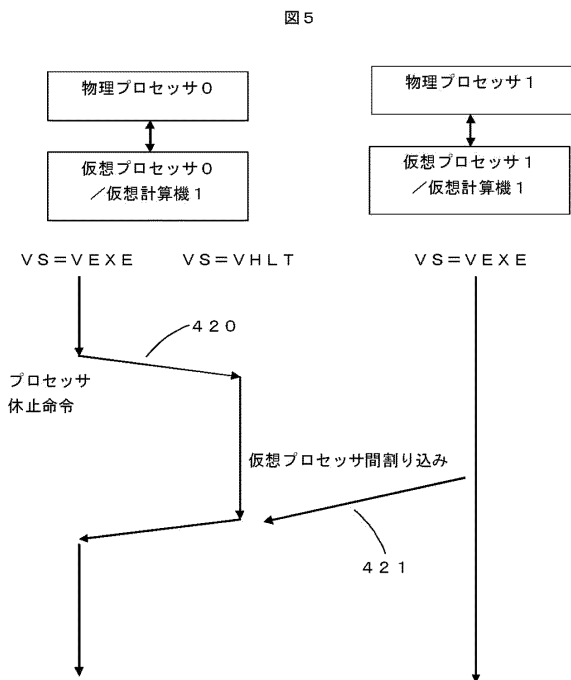
【図3】



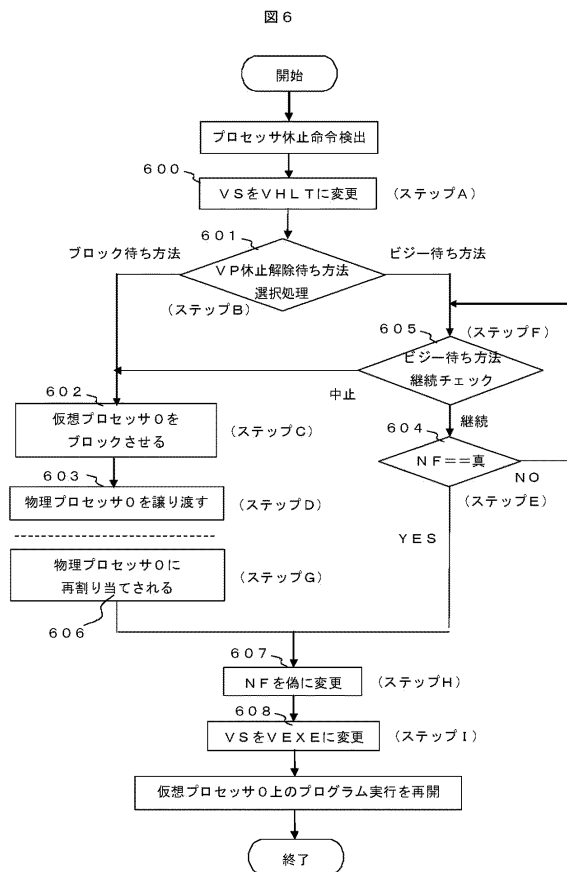
【図4】



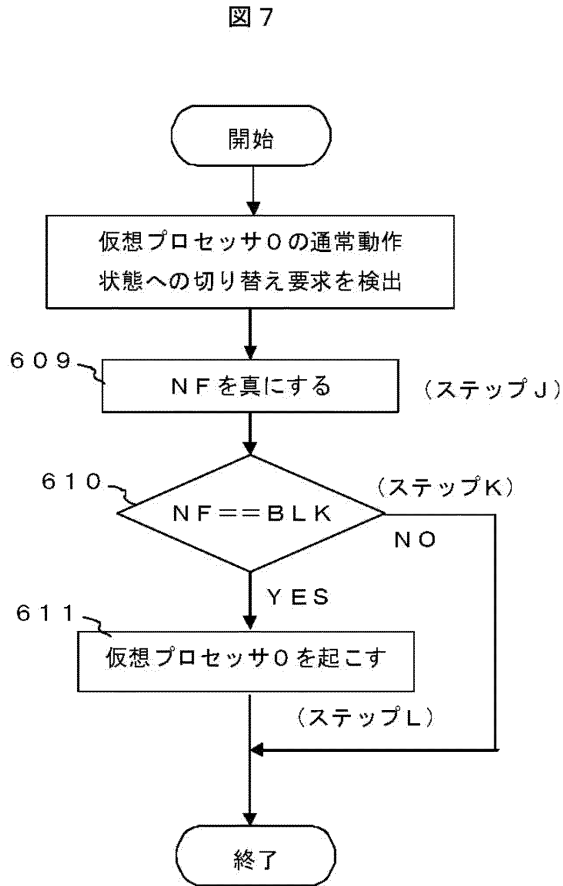
【図5】



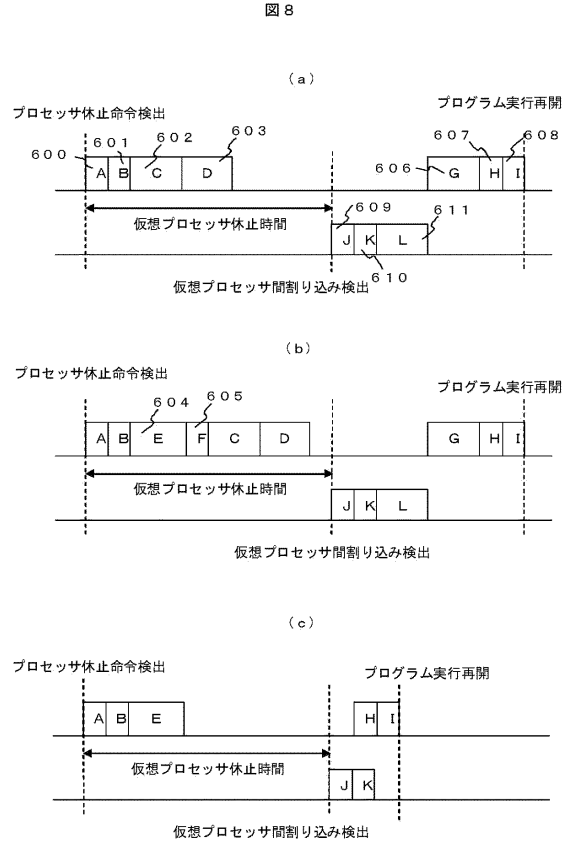
【図6】



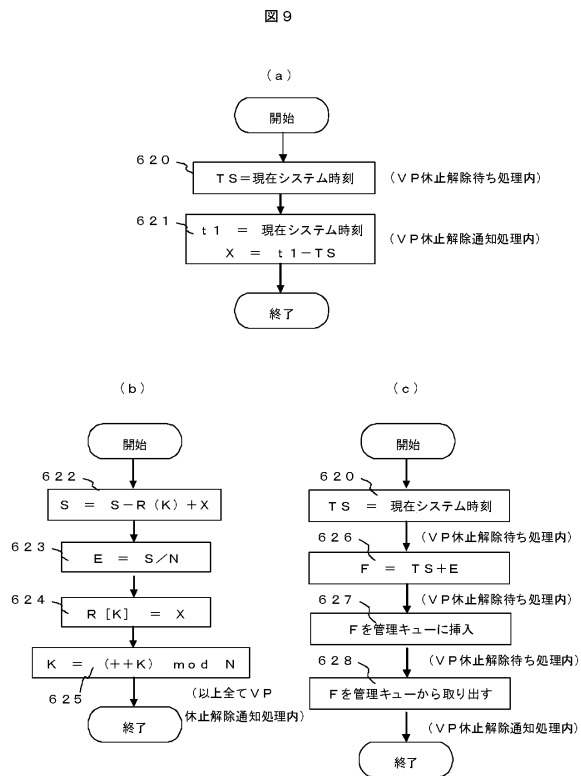
【図7】



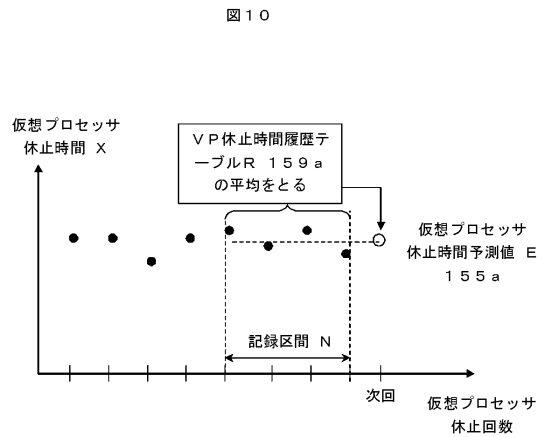
【図8】



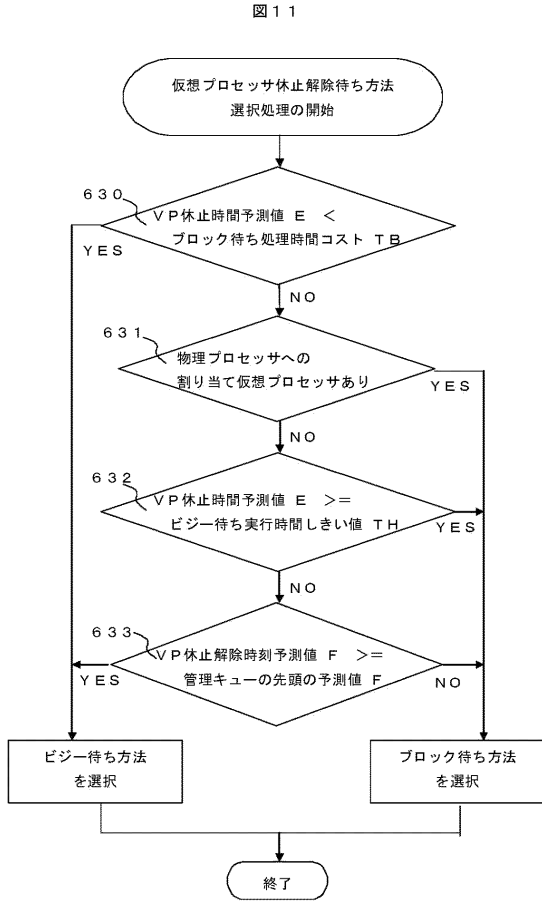
【図9】



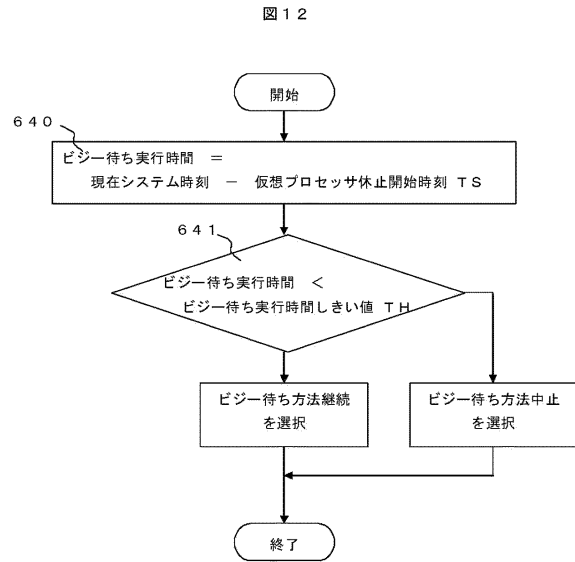
【図10】



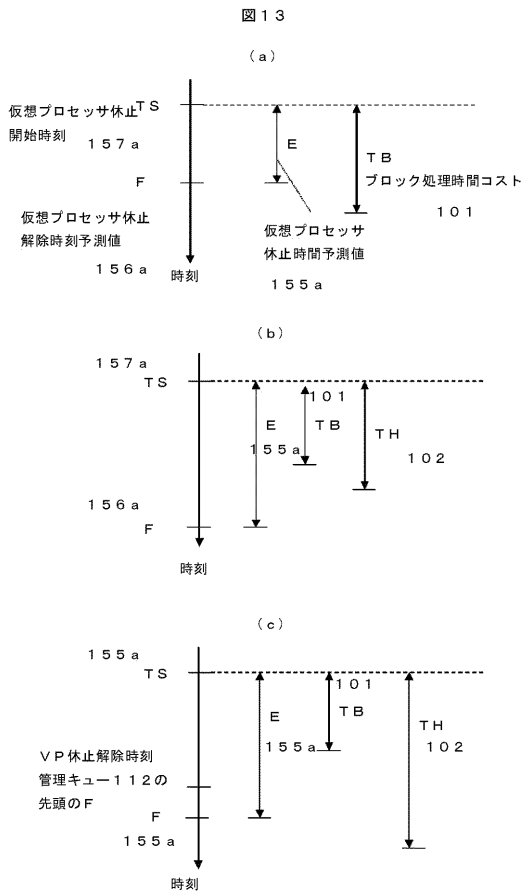
【図11】



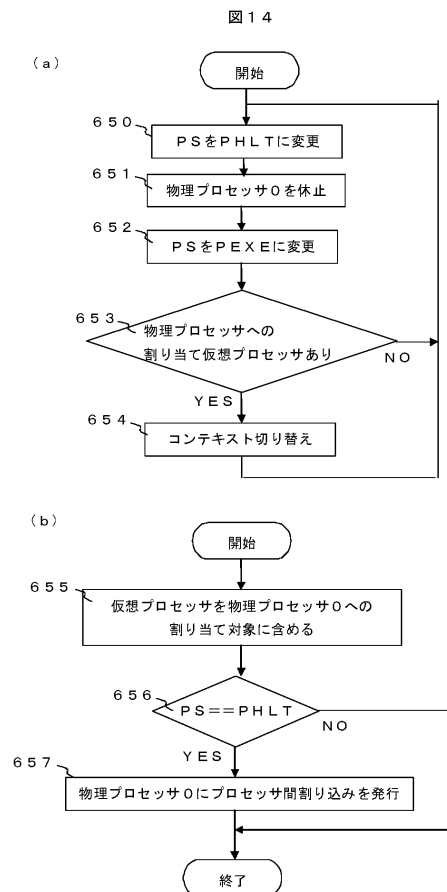
【図12】



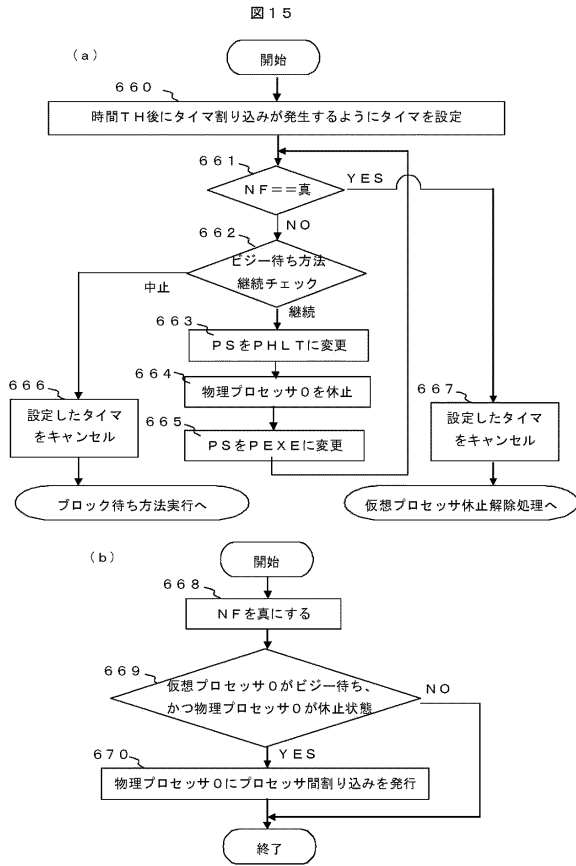
【図13】



【図14】



【図 15】



フロントページの続き

(72)発明者 和田 慎太郎

神奈川県秦野市堀山下1番地 株式会社日立製作所 エンタープライズサーバ事業部内

審査官 篠塚 隆

(56)参考文献 特開2006-18820(JP,A)

特開2008-186210(JP,A)

特開2001-290568(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F9/46-9/54