

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4705051号  
(P4705051)

(45) 発行日 平成23年6月22日(2011.6.22)

(24) 登録日 平成23年3月18日(2011.3.18)

(51) Int.Cl. F I  
**G 0 6 F 9/46 (2006.01)** G O 6 F 9/46 3 5 0  
**G 0 6 F 9/50 (2006.01)** G O 6 F 9/46 4 6 5 Z

請求項の数 7 (全 11 頁)

<p>(21) 出願番号 特願2007-17659 (P2007-17659)                  (22) 出願日 平成19年1月29日 (2007. 1. 29)                  (65) 公開番号 特開2008-186136 (P2008-186136A)                  (43) 公開日 平成20年8月14日 (2008. 8. 14)                  審査請求日 平成21年7月27日 (2009. 7. 27)</p>	<p>(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                  (74) 代理人 110000350                  ポレール特許業務法人                  (72) 発明者 早川 典充                  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                  (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                  審査官 北元 健太</p>
--	---

最終頁に続く

(54) 【発明の名称】 計算機システム

(57) 【特許請求の範囲】

【請求項1】

複数の物理CPUと、該複数の物理CPUのいずれかが割り当てられ、それぞれゲストOS上でプログラムを実行する複数のLPARと、該複数のLPARを管理する管理プログラムとを有し、全ての該LPARから必要とされる共有処理を、該管理プログラムによって該複数の物理CPUに行なわせることができる計算機システムにおいて、

該LPAR上のゲストOSの状態を把握することができるゲストOS監視部と、前記LPARに該物理CPUを割り当てるディスパッチ処理部と、を有し、

前記ディスパッチ処理部は、

各物理CPUが該共有処理に割り当てられた回数を管理するカウンタを有し、

該共有処理が発生した場合、前記カウンタを参照し、前記カウンタの回数が一番少ない物理CPUに対して該共有処理をディスパッチし、

該共有処理をディスパッチされた該物理CPUが処理したLPARの状態を前記ゲストOS監視部から取得し、該LPARがアイドル状態であった場合は該カウンタを更新せず、該LPARがアイドル状態でなかった場合は該カウンタを更新することを特徴とする計算機システム。

【請求項2】

各物理CPUに対応して、該管理プログラムによる共有処理及び共有処理以外の処理の優先度を管理するディスパッチ優先度テーブルを有し、該共有処理を実行した後、関連する該ディスパッチ優先度テーブルを更新することを特徴とする請求項1の計算機システム。

## 【請求項 3】

該ディスパッチ優先度テーブルの内容を更新するテーブル更新部を有し、該テーブル更新部は、共有処理の優先度をアイドル状態の L P A R よりも高い優先度に設定し、該共有処理の実行が終わった後、該優先度を低い状態に設定することにより優先度の高い共有処理に物理 C P U を占有させず、該物理 C P U が処理している L P A R を沈み込ませないことを特徴とする請求項 2 の計算機システム。

## 【請求項 4】

物理 C P U の数は L P A R の数以下であることを特徴とする請求項 1 乃至 3 のいずれかの計算機システム。

## 【請求項 5】

複数の物理 C P U と、該複数の物理 C P U のいずれかが割り当てられ、それぞれゲスト O S 上でプログラムを実行する複数の L P A R と、該複数の L P A R を管理する管理プログラムと、該 L P A R 上のゲスト O S の状態を把握することができるゲスト O S 監視部及び前記 L P A R に該物理 C P U を割り当てるディスパッチ処理部を含むスケジューリング処理部とを有し、全ての該 L P A R から必要とされる共有処理を、該管理プログラムによって該複数の物理 C P U に行なわせることができる計算機システムにおいて、全ての該 L P A R から必要とされる共有処理のスケジューリングを行うスケジューリング方法であって、

該スケジューリング処理部は；該ディスパッチ処理部が有するカウンタを用いて、各物理 C P U が該共有処理に割り当てられた回数を管理し、

該スケジューリング処理部は；該共有処理が発生した場合、該カウンタを参照して、該カウンタの回数が一番少ない物理 C P U に対して該共有処理をディスパッチし、

該スケジューリング処理部は；該共有処理をディスパッチされた該物理 C P U が処理した L P A R の状態を前記ゲスト O S 監視部から取得し、該 L P A R がアイドル状態であった場合は該カウンタを更新せず、該 L P A R がアイドル状態でなかった場合は該カウンタを更新する

ことを特徴とするスケジューリング方法。

## 【請求項 6】

該スケジューリング処理部は、各物理 C P U に対応して備えられたディスパッチ優先度テーブルを用いて、該管理プログラムによる共有処理及び共有処理以外の処理の優先度を管理し、該共有処理を実行した後、関連する該ディスパッチ優先度テーブルを更新することを特徴とする請求項 5 のスケジューリング方法。

## 【請求項 7】

該スケジューリング処理部は、各物理 C P U が共有処理に割り当てられた回数を管理するカウンタを参照して、カウンタ値の小さいカウンタに対応する物理 C P U から順に共有処理をディスパッチし、かつ、該カウンタの値を更新する時、ディスパッチされていた該物理 C P U が処理している L P A R の状態を前記ゲスト O S 監視部から取得し、該 L P A R がアイドル状態である時は該カウンタを更新せず、該 L P A R がアイドル状態でない時に該カウンタを更新することを特徴とする請求項 5 又は 6 の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、計算機システムに係り、特に計算機の物理リソースを分割された複数の論理的なパーティション（以下、L P A R という）で利用するシステムにおける共有処理のスケジューリングに関するものである。

## 【背景技術】

## 【0002】

高性能な計算機のハードウェアリソースを分割して L P A R と称する単位に分け、この

10

20

30

40

50

パーティションをそれぞれ1台のサーバとして運用する「サーバ分割」が実用化されている。サーバ分割は、L P A Rの動作環境によって、ハイパーバイザ型とホストOS型に分類される。前者は、複数のL P A Rをベアマシン上に存在する特殊な管理プログラム（ハイパーバイザ）上で動作させる方法である。後者は一般的なOS上（ホストOS）で管理プログラムをアプリケーションとして稼働させ、その上でL P A Rを動作させる方法である。何れの方式に関わらずサーバ分割には、V M M (Virtual Machine Monitor)と呼ばれる制御機能と、その上で動作するL P A Rが存在する。L P A Rは「仮想計算機」と呼ばれ、その上で動作するOSは「ゲストOS」と呼ばれる。ゲストOSの動作はサーバ分割していないサーバ上での動作と機能的に同じものとなる。

【0003】

10

ハイパーバイザは、主として「命令エミュレーション」、「メモリ管理」、「I/O処理」、「スケジューリング」の処理を司る。L P A Rではホストのリソースを共有して使用することもあり、ハイパーバイザには全てのL P A Rから必要とされる共有処理が存在する。一般にこのような処理は、各L P A Rの処理よりも優先して処理しなければならない。この共有処理用にC P Uが割り当てられていないと、L P A R処理用に割り当てられているC P Uリソースを独占してしまうことがあり、これによってL P A R側の処理が続行できなくなってしまふという問題が生じる。

【0004】

上記のように、ある特定の処理のため、他の処理に物理C P Uリソースが配分されない問題を解決する方法としては、処理の公平性を保証するスケジューリングアルゴリズムとしてタイムシェアリングが知られている。この技術を用いることにより共有処理を沈み込ませること無く、全ての処理に対して物理C P Uを割り当てることが可能となる。しかし、このタイムシェアリングは、リソースを奪取した物理C P Uが、ビジー状態のL P A Rを処理している物理C P Uであるか、内部的にはアイドル状態であるL P A Rを処理している物理C P Uであるかを考慮することができない。

20

【0005】

一方、優先度に応じてリソース配分を変動させ、特定のジョブの沈み込みを回避する手段として、「ジョブクラス」を定義する方法が知られている。この方法は、システムには予め定義された、複数の階層的な優先度があり、プロセス実行時に当該ジョブクラスを指定して実行させる方法である。これにより、優先度を考慮したC P Uリソースの配分が可能となる。

30

【0006】

しかし、この方法は一つの計算機（OS）内の閉じた環境では有効であるが、L P A Rとハイパーバイザを跨ぐような環境への適応は困難である。それは、ハイパーバイザからは、全てのL P A Rの優先度は静的には同一のものでありながら、L P A Rの内部的な状態はビジー状態からアイドル状態へまたはその逆へ動的に遷移するためである。

【0007】

また、特許文献1（特開2006-244483公報）には、C P Uリソースを大量に消費するプロセスを検出し、リソース配分が他のL P A Rよりも少なく定義してある「隔離L P A R」に割り当ててすることで、他のL P A Rへの影響を防ぐ方法が開示されている。しかし、この方法ではL P A Rの処理よりも優先度を上げて実行しなくてはならないハイパーバイザの共有処理には適用できない。

40

【0008】

【特許文献1】特開2006-244483号公報

【発明の開示】

【発明が解決しようとする課題】

【0009】

図1～3を参照して、発明が解決しようとしている課題について説明する。

図1は、タイムシェアリング方式による共有処理のディスパッチについて示す。

3つの物理C P Uを占有的に3つのL P A Rに割り当てた状態を示している。物理C P U

50

0は占有的にLPAR0に割り当てられ、同様に物理CPU1はLPAR1に割り当てられ、物理CPU2はLPAR2に割り当てられている。LPAR0、LPAR2上で動作しているゲストOSは終始ビジー状態であり、LPAR1のそれはアイドル状態である。

#### 【0010】

ある時点100でハイパーバイザの共有処理が発生した時、タイムシェアリング方式を用いて共有処理をディスパッチさせるとすると(101~103)、各々の物理CPUを用いて時分割で共有処理を処理するので特定のLPAR処理の沈み込みは回避できる。しかし、全てのLPARの優先度が等しいとすると、内部的にはアイドル状態のLPARが存在するにも拘わらず、104に示すように、ゲストOSがビジー状態で共有処理の合間にLPARの処理を実行しているような物理CPUからもリソースを均等に奪取してしま

10

#### 【0011】

そこで、図2に示すように、タイムシェアリング方式を用いず、アイドル状態のLPARを検知できる機能と、時分割ではなくディスパッチされた回数を基に、各物理CPUに分散させてディスパッチする方式がある。ある時点200でハイパーバイザの共有処理が発生したとすると、アイドル状態のLPARに共有処理を割り当てる。しかしながら、その後、201においてLPAR側のOSやアプリケーションでプロセス実行要求が発生しても、共有処理の方がLPARの処理よりも優先度が高いため共有処理が終了する202までCPUが明け渡されない。そのため、203に示すように、LPAR上のOSが沈み込む状態が発生する。

20

#### 【0012】

ここで、図3に示すように、ゲストOSがアイドル状態のLPAR1に割り当てられている物理CPU1が、ある時点300において発生した共有処理を複数回処理した後に、301においてアイドル状態からビジー状態へ遷移すれば、上記問題は解決され、他のビジー状態のLPARを処理している物理CPUに共有処理を迅速に引き継げたと考えられる。

しかし、カウント方式にてスケジューリングしている場合、大域的には平準化されるわけだが、局所的には一方のLPARを処理している物理CPUに共有処理が集中してディスパッチされてしまうことが起こる可能性がある。これでは、302で平準化されるまで、一方の物理CPUに集中的に割り当てられてしまい、その物理CPUが処理しているLPARが沈み込んでしまう。これは、アイドル状態のLPARに割り当てられている物理CPUで処理された方では、カウンタ値が積算されているため、ビジー状態に遷移した後は平準化のため処理が割り当てられず発生する問題である。

30

#### 【0013】

本発明の目的は、共有処理のためにLPARの沈み込みが生じることを防止した計算機システムを提供することにある。より具体的には、本発明は、アイドル状態のLPARに割り当てられているCPUを有効利用して、ビジー状態のLPAR処理を沈み込ませることを回避して、LPAR処理よりも優先度の高い共有処理のスケジューリングを行うことにある。

40

#### 【課題を解決するための手段】

#### 【0014】

本発明に係る計算機システムは、好ましくは、複数の物理CPUと、複数の物理CPUのいずれかが割り当てられ、それぞれゲストOS上でプログラムを実行する複数のLPARと、複数のLPARを管理する管理プログラムとを有する計算機システムにおいて、LPAR上のゲストOSの状態を把握することができるゲストOS監視部と、LPARの処理よりも優先度が高い管理プログラムによる共有処理を、いずれかの物理CPUに割り当てるディスパッチ処理部と、を有し、ゲストOS監視部による監視に基づいて、ディスパッチ処理部は、アイドル状態のLPARを処理している物理CPUに対して共有処理を優先的にディスパッチする計算機システムとして構成される。

50

## 【 0 0 1 5 】

好ましい例では、各物理CPUに対応して、管理プログラムによる共有処理及び共有処理以外の処理の優先度を管理するディスパッチ優先度テーブルを有し、共有処理を実行した後、関連するディスパッチ優先度テーブルを更新する計算機システムとして構成される。

また、好ましくは、ディスパッチ優先度テーブルの内容を更新するテーブル更新部を有し、テーブル更新部は、共有処理の優先度をアイドル状態のLPARよりも高い優先度に設定し、共有処理の実行が終わった後、優先度を低い状態に設定することにより優先度の高い共有処理に物理CPUを占有させず、物理CPUが処理しているLPARを沈み込ませない計算機システムとして構成される。

また、好ましくは、前記ディスパッチ処理部は、各物理CPUが共有処理に割り当てられた回数を管理するカウンタを有し、カウンタの値を更新する時、ディスパッチされていた物理CPUが処理しているLPARの状態を前記ゲストOS監視部から取得し、LPARがアイドル状態である時は、カウンタを更新しない計算機システムとして構成される。

また、好ましくは、前記複数のLPARの内少なくとも1つの所定のLPARは主として共有処理を行うように割り当てられており、共有処理が発生した場合、前記ディスパッチ処理部は、所定のLPAR又はアイドル状態のLPARを処理している物理CPUに対して共有処理を優先的にディスパッチする計算機システムとして構成される。

## 【 0 0 1 6 】

また、本発明に係るスケジューリング方法は、好ましくは、ハイパーバイザによって複数の物理CPUを複数のLPARに割り当てると共に、全てのLPARに関係する共有処理のスケジューリングを行うスケジューリング方法であって、ゲストOS監視部によってLPAR上のゲストOSの状態を監視し、共有処理が発生した場合、ゲストOS監視部による監視に基づいて、アイドル状態のLPARを処理している物理CPUに対して優先的に共有処理をディスパッチすることを特徴とするスケジューリング方法として構成される。

## 【 発明の効果 】

## 【 0 0 1 7 】

本発明によれば、複数の物理CPUがLPARに割り当てられている場合、アイドル状態のLPARに割り当てられた物理CPUを有効に利用して共有処理を行うことができる。また、アイドル状態のLPARが存在しない場合は、全てのLPARに割り当てられた物理CPUを均等に使用して特定のLPARを沈み込ませることなく、LPAR処理よりも優先度の高い共有処理を優先したスケジューリングを実現することができる。

## 【 発明を実施するための最良の形態 】

## 【 0 0 1 8 】

本発明の実施例を図面に基づいて詳細に説明する。

## [実施例 1]

図4は一実施例による計算機システムを示す。

この計算機システムは、4つの物理CPU 411～414（総称して41ということがある）をサーバ分割して、同数の4つのLPAR 430～433（総称して43ということがある）が割り当てられる。また、LPARの制御を司るハイパーバイザを一つの物理CPUが専用で処理するのではなく（隠れCPUは存在しない）、複数の物理CPUが協調して処理する。また、共有処理をシステムで共有して使用されるデバイスのI/O制御処理とする。

## 【 0 0 1 9 】

計算機システムの物理的なリソース群（ハードウェアリソース）410は、複数の物理CPU 411～414、主記憶装置415、外部記憶装置416、およびネットワークインタフェースカード417等の構成要素を有する。

ハイパーバイザ420は、ハードウェアリソース410を用いて動作し、その上で動作する複数のLPAR 430～433の管理及び制御を司る。ハイパーバイザ420は、主に

10

20

30

40

50

、命令エミュレーション部 4 2 1、メモリ管理部 4 2 2、スケジューリング処理部 4 2 3、I/O 処理部 4 2 4 の各処理機能を有する。

L P A R 4 3 はそれぞれゲスト OS 4 4 0 ~ 4 4 3 (総称して 4 4 ということがある)と、各ゲスト OS 上で動作するアプリケーションプログラム 4 5 0 ~ 4 5 3 (総称して 4 5 ということがある)を有している。アプリケーションの動作はサーバ分割されていないハードウェア上の動作と同様である。

#### 【 0 0 2 0 】

図 5 はスケジューリング処理部 4 2 3 の構成を示す。

スケジューリング処理部 4 2 3 は、ゲスト OS 監視部 5 1、ディスパッチ処理部 5 2、ディスパッチ優先度テーブル更新部 5 3 を有して構成される。ゲスト OS 監視部 5 1 は共有処理に限らず、各 L P A R 4 3 上で動作している OS のリアルタイムな状態や、共有処理を実施した L P A R の状態を監視して把握する。

10

#### 【 0 0 2 1 】

ディスパッチ処理部 5 2 はスケジューリングアルゴリズムを用いて、ハイパーバイザ 4 2 0 および各 L P A R 4 3 に CPU リソースを割り当てる処理を行う。ディスパッチ処理部 5 2 は、ディスパッチカウンタテーブル 5 2 1 を備え、物理 CPU 4 1 を各 L P A R 4 3 に偏り無く平均的にディスパッチさせるために、各物理 CPU 4 1 が共有処理に割り当てられた回数を計数して、ディスパッチカウンタテーブル 5 2 1 に保持する。例えば、ある時点における、物理 CPU 4 1 1 ~ 4 1 4 に対応するディスパッチカウンタテーブル 5 2 1 の値が、「 2 」 「 3 」 「 4 」 「 1 」 の場合、次回の共有処理は、カウンタ値が最小「 1 」の物理 CPU 4 1 4 が割り当てられるが如きである。

20

#### 【 0 0 2 2 】

ディスパッチ優先度テーブル更新部 5 3 は、ディスパッチ優先度テーブル 5 4 1 ~ 5 4 4 (総称して 5 4 ということがある)を有し、各物理 CPU 0 ~ 3 に対応して現時点でディスパッチすべき処理の優先度を管理する。

より具体的には、図 6 を参照するに、ディスパッチ優先度テーブル 5 4 は、優先度として、ハイパーバイザの処理 (例えば、共有処理)、ビジー状態の L P A R 処理、アイドル状態 L P A R 処理等の優先度を設定して管理する。ディスパッチ優先度テーブル 5 4 が状態 6 1 の時は、初期状態であり、共有処理を含んだハイパーバイザ処理の優先度が最も高いものとなっている。この状態は定常状態 6 1 であり、L P A R 処理よりも優先度の高い共有処理を最優先で実行する。

30

#### 【 0 0 2 3 】

しかし、共有処理による L P A R の沈み込みを避けるために、状態 6 2 に示すように、ある物理 CPU で共有処理を実行した後に、当該物理 CPU に対する共有処理の優先度を意図的に低いものに変更する。この状態は排他状態である。排他状態 6 2 へ遷移した物理 CPU は、共有処理よりも優先度が高くなった L P A R 処理が存在すれば、一度 CPU リソースを開放する。これにより、L P A R 処理に物理 CPU を割り当てることができ、L P A R の沈み込みを回避することができる。ディスパッチ優先度テーブル 5 4 n を更新した後に、L P A R 処理を再び実行した物理 CPU は、ディスパッチ優先度テーブル 5 4 n を元に戻し、定常状態へ復帰する。

40

#### 【 0 0 2 4 】

図 7 は、ディスパッチ優先度テーブル 5 4 の更新処理のシーケンスを示す。

この図では、便宜上 3 つの物理 CPU 0 ~ 2 から構成される計算機の動作例を示しているが、4 つ以上の物理 CPU を有する計算機の場合の動作例も実質的に同様である。

初期状態では全ての CPU (CPU 0 ~ CPU 2) は定常状態 (7 0 1 ~ 7 0 3) である。ハイパーバイザ側より共有処理が発生した場合 (7 0 4)、何れの CPU においても最も優先度高い処理であるため (即ち定常状態)、直ちに CPU が割り当てられる (7 0 5)。タイムスライス分の実行が終了すると、共有処理を実行していた CPU 0 は排他状態に切り替わる (7 0 6)。排他状態になった CPU 0 は共有処理よりも L P A R 0 の処理が優先的に割り当てられる (7 0 7)。

50

## 【 0 0 2 5 】

一方、定常状態にあるCPU1にはLPAR1の処理よりも共有処理が優先的に割り当てられ、共有処理は引継いで処理される(708)。タイムスライス分の実行が終了すると、排他状態にあったCPU0は定常状態へ復帰し(709)、共有処理を実行していたCPU1は排他状態に切り替わる(710)。このように、共有処理はLPAR処理を実行している複数の物理CPU間を順次渡りながら実行されていく。

## 【 0 0 2 6 】

図8はスケジューリング処理部423の処理シーケンスを示す。

共有処理が発生した場合(801)、スケジューラは、LPARに割り当てられていない空き物理CPUがあるか、をチェックする(802)。チェックの結果、空き物理CPUがあればそのCPUをディスパッチする。一方、空き物理CPUが無ければ、ゲストOS監視部52に問い合わせを実施し(803)、アイドル状態のLPARに割り当てられている物理CPUがあるか、をチェックする(804)。チェックの結果、アイドル状態のCPUがあればそれをディスパッチする。一方、アイドル状態のCPUが存在しない場合は、ディスパッチカウンタテーブル521を参照して(805)、ディスパッチ回数が一番少ない物理CPUに共有処理を割り当てる(806)。

10

## 【 0 0 2 7 】

物理CPUを最終的に割り当てる前に、ディスパッチ優先度テーブル54を参照し(807)、現時点の優先度から判断して共有処理を割り当てることができるか、を確認する(808)。当該物理CPUが排他状態であり、割り当てが不可な場合は、その物理CPUを選択の対象から除外して(809)、再度ディスパッチカウンタテーブル521を参照する。一方、割り当てが可能な場合は、その物理CPUで共有処理を実行する(810)。

20

## 【 0 0 2 8 】

共有処理は、決められた期間(タイムスライス分)実施される(811)。共有処理の実行がタイムスライス分の期間を経過した後、共有処理を実行した物理CPUに対して、ディスパッチ優先度テーブル54nを更新する(812)。共有処理を実行した物理CPUが割り当てられていたLPARの状態を調査するため、ゲストOS監視部51に、アイドル状態のLPARの物理CPUを使用していたか、を問い合わせる(813)。問合せの結果、アイドル状態のLPARの物理CPUを使用していた場合(813Y)には、ディスパッチカウンタテーブル521を更新しない。一方、アイドル状態でないLPARの物理CPUを使用していた場合は(813N)、ディスパッチカウンタテーブル521を更新する(814)。

30

## 【 0 0 2 9 】

その後、共有処理が終了したかを判断し(815)、終了した場合は本スケジューリング処理を終了する(816)。一方、共有処理が終了していない場合、次のタイムスライスにおけるCPUディスパッチとして、当該物理CPUに割り当てられているLPARが継続してアイドル状態であるかチェックし(817)、アイドル状態の場合は同じ物理CPUを使用する。一方、当該LPARがアイドル状態からビジー状態へ遷移した場合、ディスパッチされていた物理CPUを開放する(818)。そして、ある一定の期間ごとにディスパッチカウンタテーブル521をリフレッシュし(819, 820)、リフレッシュのための一定時間が経過していなければ、ステップ830に戻って、上記した処理を繰り返す。

40

## 【 0 0 3 0 】

図9は、共有処理のスケジューリングの例を示す。

LPAR1上のゲストOS1は初期段階でアイドル状態である。ある時点901において、ハイパーバイザの共有処理が発生すると、共有処理は、LPAR1を処理している物理CPU1にディスパッチされる。ゲストOS1(物理CPU1)がアイドルである期間は物理CPU1を有効利用し、共有処理を実行する。一方、ゲストOS0、2、3は当該期間において共有処理に何ら阻害されることなく、LPARの処理に専念できる。

50

## 【 0 0 3 1 】

その後、時点 9 0 2 においてゲスト OS 1 のプロセスが発生すると、次のスケジューリングのタイミングにおいて、物理 CPU 1 を開放し、本来の処理である L P A R 1 の処理に移行する。以後、何れかの物理 CPU を独占することなく、各物理 CPU を均等に割り当てることで共有処理を実行して終わることができる ( 9 0 3 )。

## 【実施例 2】

図 1 0 は、他の実施例による計算機システムの構成を示す。

この例は、各ゲスト OS の I / O 処理を各 L P A R に実行させるのではなく、全ての L P A R の I / O 処理を一括して引き受けるドライバ OS を L P A R 1 0 0 1 上に設けたものである。これにより容易に I / O デバイスを複数の L P A R 間で共有することができ、中

10

間バッファ領域等を設けることで効率よく I / O 処理を実行できる。このようなドライバ OS では複数の I / O 処理 1 0 0 2 を実行するため、一般的には命令処理よりも長期間 CPU を割り当てる必要がある。しかしこのような場合でも、本実施例によればドライバ OS ( 特定の L P A R 上で動作 ) を高い優先度を維持したまま、かつ L P A R 4 3 0 ~ 4 3 3 を処理している CPU の内、特定の CPU のリソースを食い潰すことなく実現できる。

## 【図面の簡単な説明】

## 【 0 0 3 2 】

【図 1】タイムシェアリング方式による共有処理のディスパッチについて説明するための図、

20

【図 2】優先度の高い共有処理による L P A R の沈み込みについて説明するための図、

【図 3】割り当て平準化に伴う L P A R の沈み込みについて説明するための図、

【図 4】本発明の一実施例による計算機システムの構成を示す図、

【図 5】計算機システムにおけるスケジューリング処理部 4 2 3 の構成を示す図、

【図 6】ディスパッチ優先度テーブル 5 4 の構成を示す図、

【図 7】ディスパッチ優先度テーブル更新部 5 3 の処理シーケンスを示す図、

【図 8】計算機システムにおけるスケジューリング処理を示すフローチャート、

【図 9】一実施例によるスケジューリングの例を示す図、

【図 1 0】他の実施例による計算機システムの構成を示す図。

30

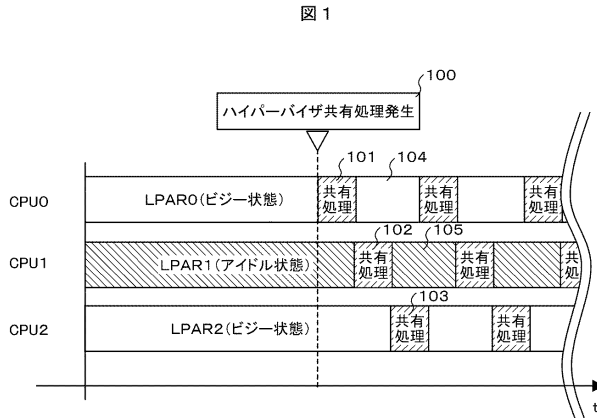
## 【符号の説明】

## 【 0 0 3 3 】

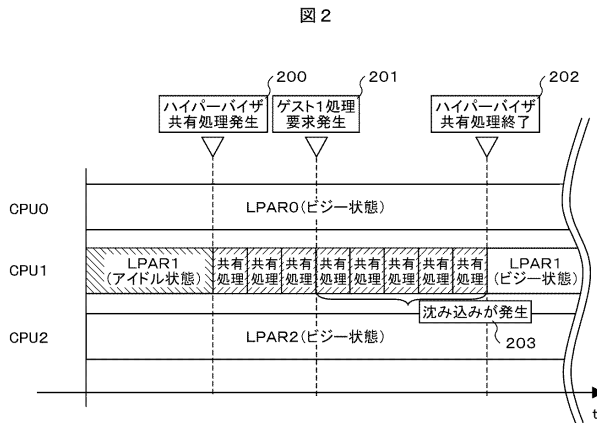
4 1 0 : ハードウェアリソース 4 1 1 ~ 4 1 4 : CPU 4 1 5 : 主記憶装置 4 1 6 : 外部記憶装置 4 1 7 : ネットワークインタフェースカード 4 2 0 : ハイパーバイザ 4 2 1 : 命令エミュレーション部 4 2 2 : メモリ管理部 4 2 3 : スケジューリング処理部 4 2 4 : I / O 処理部 4 3 0 ~ 4 3 3 L P A R 4 4 0 ~ 4 4 3 : ゲスト OS 4 5 0 ~ 4 5 3 : アプリケーションプログラム 5 1 : ゲスト OS 監視部 5 2 : ディスパッチ処理部 5 2 1 : ディスパッチカウンタテーブル 5 3 : ディスパッチ優先度テーブル更新部 5 4 1 ~ 5 4 4 : ディスパッチ優先度テーブル



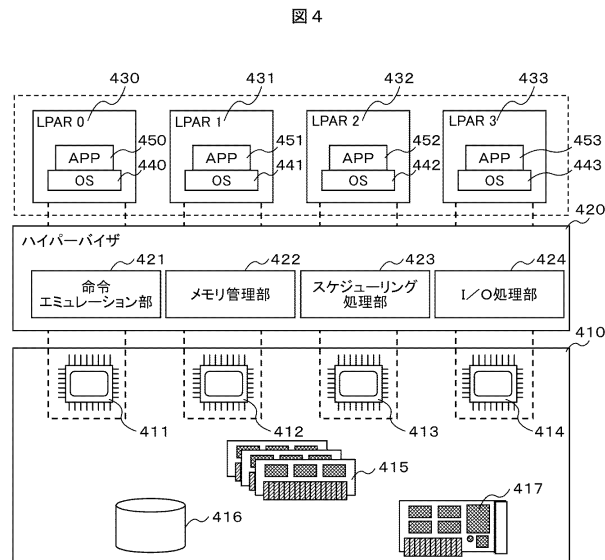
【図1】



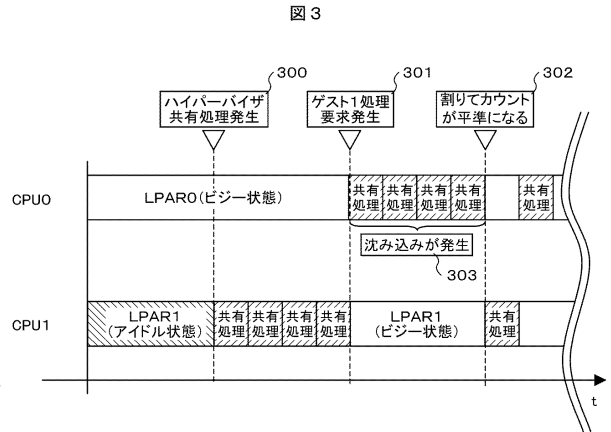
【図2】



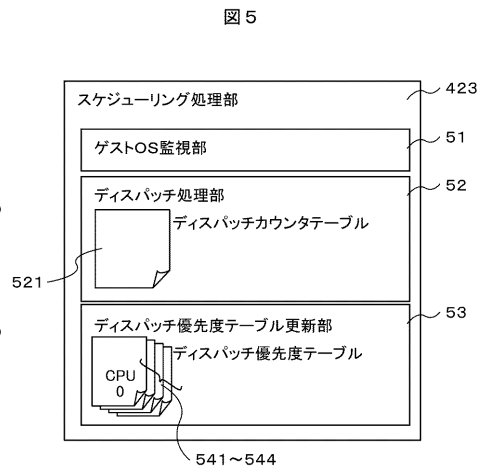
【図4】



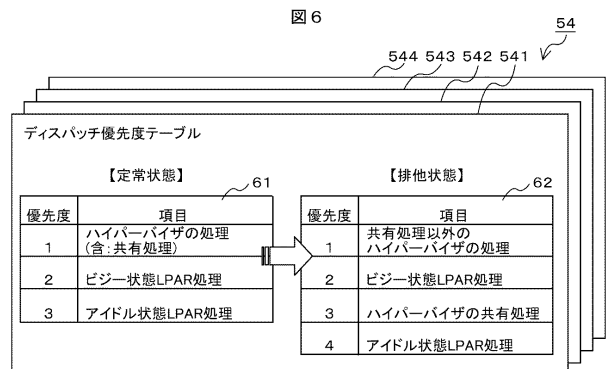
【図3】



【図5】

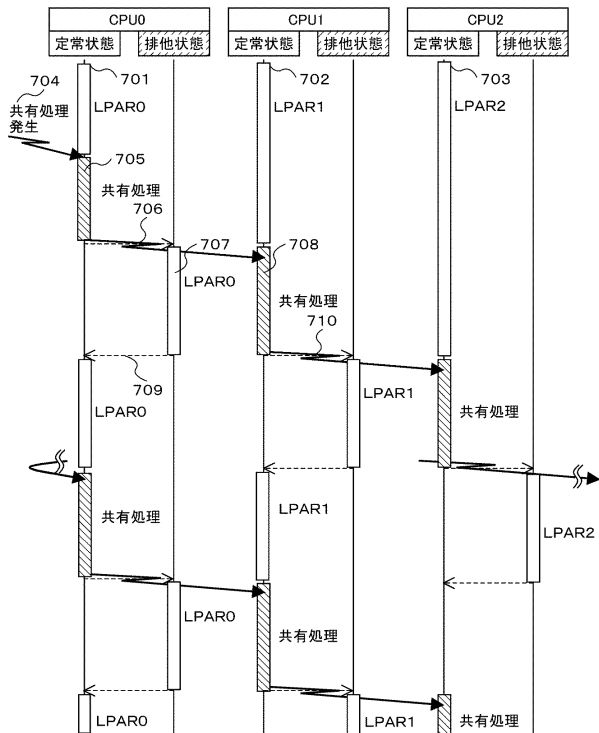


【図6】



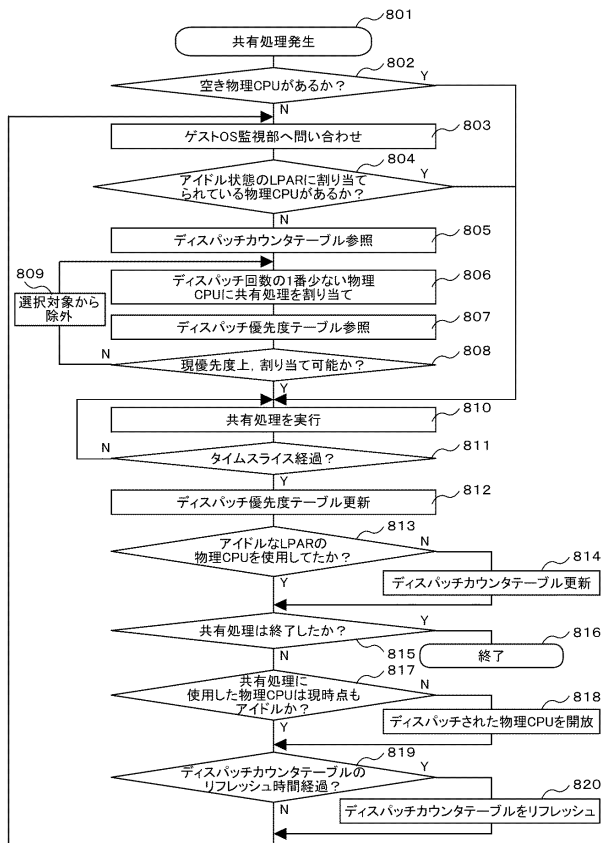
【図7】

図7



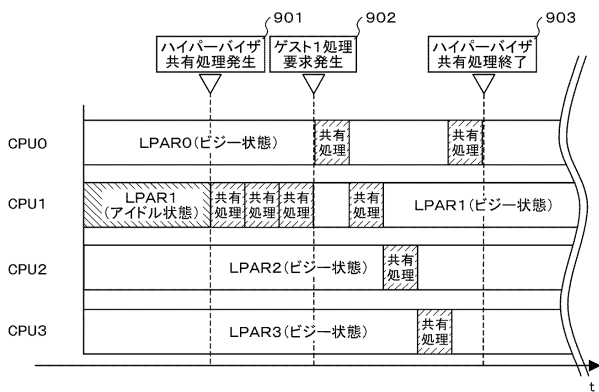
【図8】

図8



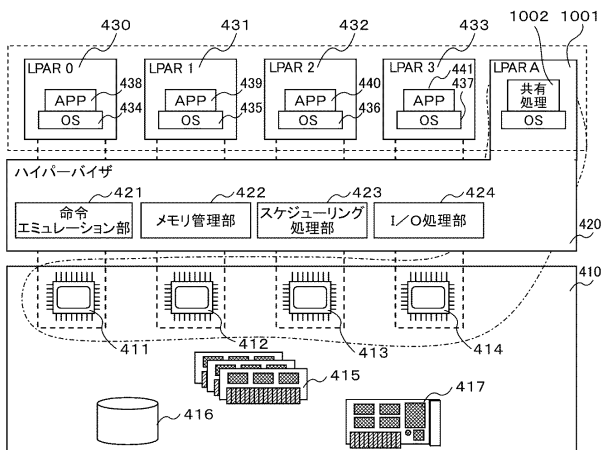
【図9】

図9



【図10】

図10



---

フロントページの続き

(56)参考文献 国際公開第2005/036367(WO, A2)

特表2007-508623(JP, A)

特開2004-252526(JP, A)

特開平8-77025(JP, A)

特開平11-259316(JP, A)

特開2002-318699(JP, A)

特開昭62-221041(JP, A)

特開2006-127524(JP, A)

特開2005-71161(JP, A)

特開平4-148460(JP, A)

Tivoli パフォーマンス・レポーター OS/390 CICSパフォーマンス・フィーチャーガイドおよび  
リファレンス リリース5, 日本アイ・ビー・エム株式会社, 2000年11月, 第1刷, pp.1  
53 - 158

(58)調査した分野(Int.Cl., DB名)

G06F 9/46 - 9/54