

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4185103号  
(P4185103)

(45) 発行日 平成20年11月26日(2008.11.26)

(24) 登録日 平成20年9月12日(2008.9.12)

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

請求項の数 8 (全 12 頁)

<p>(21) 出願番号 特願2006-39467 (P2006-39467)                  (22) 出願日 平成18年2月16日(2006.2.16)                  (65) 公開番号 特開2006-244479 (P2006-244479A)                  (43) 公開日 平成18年9月14日(2006.9.14)                  審査請求日 平成18年2月16日(2006.2.16)                  (31) 優先権主張番号 11/067852                  (32) 優先日 平成17年2月28日(2005.2.28)                  (33) 優先権主張国 米国 (US)</p>	<p>(73) 特許権者 503003854                  ヒューレット・パッカード デベロップメント カンパニー エル. ピー.                  アメリカ合衆国 テキサス州 77070                  ヒューストン 20555 ステイト                  ハイウェイ 249                  (74) 代理人 110000039                  特許業務法人アイ・ピー・エス                  (72) 発明者 スコット・エイ・ライン                  アメリカ合衆国カリフォルニア州 パロアルト                  ハノーバー・ストリート 3000                  ヒューレット・パッカード・カンパニー                  内                  審査官 鈴木 修治</p>
--	---

最終頁に続く

(54) 【発明の名称】 実行可能プログラムをスケジューリングするためのシステム及び方法

(57) 【特許請求の範囲】

【請求項 1】

複数のプロセッサ(131)と、

実行可能プログラムの複数のグループ(143)であって、各グループ(143)について、前記グループ(143)の実行可能プログラムをサポートするプロセッサ資源の量を表す各シェアパラメータ(122)が定義される、実行可能プログラムの複数のグループ(143)と、

前記シェアパラメータ(122)を使用して複数の重みを配分し、前記複数のプロセッサ(131)間における前記重みの配分を生成するソフトウェアルーチン(124)であって、前記配分が、各グループ(143)についてのプロセッサのサブセットと、前記グループ(143)の実行可能プログラムとをスケジューリングするための、前記サブセット内の各プロセッサの比率を定義するソフトウェアルーチン(124)と、

前記配分に従ったスケジューリング間隔の期間中に、前記複数のグループ(143)の各実行可能プログラムを前記複数のプロセッサ(131)の特定のプロセッサにスケジューリングするためのスケジューリングソフトウェアルーチン(125)と

を備え、

前記ソフトウェアルーチン(124)は、

複数の整数区画問題(IPP; integer partition problem)アルゴリズムを使用して、前記複数のグループ(143)の少なくとも1つに割り当てられるシェアパラメータ(122)から生成された重みの個数を変化させることによって複数の配分を生成し、生成

された前記複数の配分をスケジューリングに使用する

コンピュータシステム。

【請求項 2】

前記可変の個数の重みは、

デフォルトのグループに割り当てられるシェアパラメータ ( 1 2 2 ) から生成される請求項 1 に記載のコンピュータシステム。

【請求項 3】

前記スケジューリングソフトウェアルーチン ( 1 2 5 ) は、

前記複数のグループ ( 1 4 3 ) の間のスケジューリングの相違を補償するように前記複数の配分を交互に行う

請求項 1 に記載のコンピュータシステム。

【請求項 4】

前記複数のグループ ( 1 4 3 ) の実行可能プログラムのスケジューリングパラメータを保持するためのソフトウェアルーチンであって、各スケジューリングパラメータが、グループの平均に関連した ( relate )、各実行可能プログラムが受け取るプロセッサ資源の量を示すソフトウェアルーチン

をさらに備える

請求項 1 に記載のコンピュータシステム。

【請求項 5】

前記スケジューリングソフトウェアルーチン ( 1 2 5 ) は、

前記スケジューリングパラメータの値に従って、前記複数のグループ ( 1 4 3 ) の実行可能プログラムのサブセットを、割り当て期間内にプロセッサ資源を受け取る付加的な機会を実行可能プログラムの前記サブセットに提供する 1 つ又は数個のプロセッサに割り当てる

請求項 4 に記載のコンピュータシステム。

【請求項 6】

複数のグループの実行可能プログラムをスケジューリングするためのプロセッサ資源の量を表す複数のシェアパラメータを定義すること ( 2 0 2 ) と、

前記複数のシェアパラメータを使用して、整数区画問題 ( I P P ; integer partition problem ) に従い複数の重みを生成すること ( 2 0 4 ) と、

複数の I P P アルゴリズムを使用して、前記重みの個数を変化させることによって、複数のプロセッサ間における前記重みの配分を複数決定すること ( 2 0 6 ~ 2 1 0 ) と、

前記複数の配分を使用して、前記複数のプロセッサにグループの実行可能プログラムをスケジューリングすること ( 3 0 4 ) と

を含む方法。

【請求項 7】

前記複数のグループの実行可能プログラムのスケジューリングパラメータを保持することであって、各スケジューリングパラメータが、グループの平均に関連した、各実行可能プログラムが受け取るプロセッサ資源の量を示すこと

をさらに含む請求項 6 に記載の方法。

【請求項 8】

前記スケジューリングすることは、

1 つ又は数個のプロセッサについて、前記スケジューリングパラメータの値に従って実行可能プログラムを選択することであって、前記 1 つ又は数個のプロセッサが、前記選択された実行可能プログラムに、スケジューリング間隔内でプロセッサ資源を受け取る付加的な機会を提供すること

を含む請求項 7 に記載の方法。

【発明の詳細な説明】

【技術分野】

【 0 0 0 1 】

10

20

30

40

50

本出願は、包括的には、コンピュータ資源に対するアクセスのスケジューリングに関する。

【背景技術】

【0002】

多くの企業は、その企業の組織内で使用されるコンピュータ及びアプリケーションの数の劇的な増加を経験してきた。企業の事業グループが新たなアプリケーションを配備する時、その新たなアプリケーションをホストする1つ又は複数の専用サーバプラットフォームを追加することが可能である。このタイプの環境は、時に、「one-app-per-box（ボックスごとに1つのアプリケーション）」と呼ばれる。より多くの事業プロセスがデジタル化されるにつれて、「one-app-per-box」環境では、サーバプラットフォームの個数が過度なものになる。その結果、サーバプラットフォームの管理運営コストは大幅に増加する。その上、サーバプラットフォームの資源が実際に使用される時間のパーセンテージ（利用率）は、かなり低くなる可能性がある。これらの問題に対処するために、多くの企業は、複数のアプリケーションを共通のサーバプラットフォームに一元化して、プラットフォームの個数を削減し、システムの利用率を向上させてきた。このような一元化が行われた場合、アプリケーション及び他の実行可能プログラムがいつプロセッサ資源にアクセスするかを決定するための或る機能が通常設けられる。このような機能は、通常、「スケジューリング」と呼ばれる。

10

【0003】

さまざまな複雑さを有するいくつかのスケジューリングアルゴリズムが存在する。おそらく、最も簡単なスケジューリングは、先着順アルゴリズムである。優先順位に基づくアルゴリズムは、優先順位をプロセスに割り当て、最も高い優先順位を有するプロセスが、適切な時刻に選択されて実行される。プリエンティブスケジューリングアルゴリズムは、高い優先順位のプロセスの実行の準備ができると、低い優先順位のプロセスをプロセッサから移す（remove）のに使用することができる。ラウンドロビンスケジューリングアルゴリズムは、プロセスが或る時間間隔の満了まで実行することを可能にし、次いで、別の実行可能プログラムが選択されて、各プロセッサで実行される。これに加えて、フェアシェアスケジューラは、パーセント又はシェアを定義し、定義したシェアに比例して、プロセッサ資源にアクセスする機会をプロセスに提供する。

20

【発明の開示】

30

【課題を解決するための手段】

【0004】

一実施の形態では、コンピュータシステムは、複数のプロセッサと、実行可能プログラムの複数のグループであって、各グループについて、そのグループの実行可能プログラムをサポートするプロセッサ資源の量を表す各シェアパラメータが定義される、実行可能プログラムの複数のグループと、シェアパラメータを使用して複数の重みを生成し、且つ、複数のプロセッサ間における重みの配分を生成する、ソフトウェアルーチンであって、この配分が、各グループについてのプロセッサのサブセット、及び、そのグループの実行可能プログラムをスケジューリングするための、サブセット内の各プロセッサの比率を定義する、ソフトウェアルーチンと、配分に従ったスケジューリング間隔の期間中に、複数のグループの各実行可能プログラムを複数のプロセッサの特定のプロセッサにスケジューリングするためのスケジューリングソフトウェアルーチンと、を備える。

40

【0005】

別の実施の形態では、方法は、複数のグループの実行可能プログラムをスケジューリングするためのプロセッサ資源の量を表す複数のシェアパラメータを定義すること、これらの複数のシェアパラメータを使用して、整数区画問題（IPP）に従い複数の重みを生成すること、IPPアルゴリズムを使用して、複数のプロセッサ間における重みの配分を決定すること、及び、この配分を使用して、複数のプロセッサにグループの実行可能プログラムをスケジューリングすること、を含む。

【0006】

50

別の実施の形態では、コンピュータシステムは、複数の資源デバイスと、複数の実行可能プログラムの複数のグループであって、各グループについて、そのグループの実行可能プログラムをサポートする複数の資源デバイスへのアクセス量を表す各シェアパラメータが定義される、複数の実行可能プログラムの複数のグループと、シェアパラメータを使用して複数の重みを生成し、且つ、複数の資源デバイス間における重みの配分を生成する、ソフトウェアルーチンであって、この配分が、各グループについての資源デバイスのサブセット、及び、そのグループの実行可能プログラムをスケジューリングするための、サブセット内の各資源デバイスの比率を定義する、ソフトウェアルーチンと、配分に従って、複数の資源デバイスの特定の資源デバイスに複数のグループの各実行可能プログラムをスケジューリングするためのスケジューリングソフトウェアルーチンと、を備える。

10

【0007】

別の実施の形態では、コンピュータシステムは、整数区画問題（IPP）アルゴリズムを使用して、コンピュータシステムの複数の資源デバイス間における重みの配分を生成するための手段であって、これらの重みが、それぞれが実行可能プログラムの各グループに提供される複数の資源デバイスへのアクセス量を表す複数のシェアパラメータから生成され、この配分が、各グループについての資源デバイスのサブセット、及び、そのグループの実行可能プログラムをスケジューリングするための、サブセット内の各資源デバイスの比率を定義する、生成するための手段と、この配分に従って、グループの各実行可能プログラムを資源デバイスにスケジューリングするための手段と、を備える。

20

【発明を実施するための最良の形態】

【0008】

いくつかの代表的な実施の形態は、整数区画問題（IPP）アルゴリズムを使用して、シェアに基づくワークロードグループのオペレーションのスケジューリングを実行する。各グループには、そのグループに割り当てられたシステム資源の「シェア」を表すパラメータの値が与えられる。ソフトウェアモジュールは、IPPアルゴリズムを使用して各グループを1つ又は数個のプロセッサにマッピングする。具体的には、グループのシェアは「重み」に区分（separate）され、これらの重みは、各プロセッサに関連した重みがほぼ等しくなるようにプロセッサ（「ビン（bin）」）に配分される。

【0009】

シェアの重みへの区分は、ワークロードのいくつかをサポートするのに使用される複数の「仮想プロセッサ」を考慮することができる。たとえば、或るグループに4つの仮想CPUが割り当てられ、各仮想CPUが物理CPUの約75パーセントの能力を有する場合、そのグループは、それぞれが75の4つの個々の重みを生成する。各CPUを100よりも多くのシェア又は少ないシェアでスケジューリングすることができるので、重みは、資源のパーセンテージに正確に対応するものではない。特定のCPUの実際のスケジューリングパーセンテージは、そのCPUで現在実行されているすべてのジョブの全重みを使用して求められる。

30

【0010】

また、デフォルトのグループのシェアパラメータ又は最も低い優先順位のグループのシェアパラメータを複数の重みに区分することを可変に行って、プロセッサ間における重みの最適な配分を達成する確率を改善することができる。このデフォルトのグループは、具体的な重み又は優先順位を有しないすべての資源要求を保持するのに使用することができる。1つのインプリメンテーションでは、デフォルトのグループのすべてのメンバーは、他のグループにまだ割り当てられていない資源を均等に分け合う。

40

【0011】

IPPアルゴリズムによって生成された配分は、各グループの物理CPUと、各グループがスケジューリング間隔内で受け取るそれらCPUの比率とのリストを提供する。これに加えて、各ジョブが受け取るプロセッサ時間の量が、ジョブスケジューリングパラメータを使用して追跡される。時間サンプリング間隔でより多くのプロセッサチェックを累積したジョブは、自身のパラメータを減らす。平均のプロセッサチェックよりも少ないプロセッ

50

サチックを累積したジョブは、自身のパラメータをインクリメントする。新たな各スケジューリング間隔時には、最も高いパラメータの値を有するジョブが、より多くのプロセッサチックをこれらのジョブに提供する利用可能な物理CPU（すなわち、全スケジューリング重みが最も低いCPU（複数可））に対して選択される。また、2つのCPUのスケジューリング重みが等しい場合、履歴に基づく最も低い使用量が、より良いCPUを選択するのに使用される。

#### 【0012】

次に図面を参照して、図1は、代表的な一実施の形態によるシステム100を示している。このシステム100は、プラットフォームのハードウェアレイヤ130に対する低レベルのアクセスを制御するホストオペレーティングシステム120を含む。一実施の形態では、ホストオペレーティングシステム120は、一例としてそのカーネル内に仮想化レイヤ121を含む。仮想化レイヤ121は、ハードウェアレイヤ130の物理資源に対応するソフトウェア構成（software constructs）（論理デバイス）を作成する。ハードウェアレイヤ130は、CPU131-1から131-N、メモリ132、ネットワークインターフェース133、入出力（I/O）インターフェース134等の任意の個数の物理資源を含むことができる。

10

#### 【0013】

一実施の形態では、仮想資源（たとえば、1又は数個の仮想CPU、仮想メモリ、仮想ネットワークインターフェースカード、仮想I/Oインターフェース等）が、各仮想マシン141に割り当てられる。仮想CPUの個数は、物理CPU131の個数を超えてもよい。各仮想マシン141は、自身に割り当てられた仮想資源に従ってオペレーティングシステム120上でプロセスとして実行される。CPUの仮想化は、各仮想マシン141がそれ自身のCPU又はそれ自身の1組のCPUで実行するように見えるよう行うことができる。CPUの仮想化は、各仮想CPUについて、1組のレジスタ、変換索引バッファ、及び他の制御構造体を設けることによって実施することができる。したがって、各仮想マシン141は、他の仮想マシン141から隔離される。これに加えて、各仮想マシン141は、各ゲストオペレーティングシステム142を実行するのに使用される。仮想マシン141に割り当てられた仮想資源は、ゲストオペレーティングシステム142には、物理サーバのハードウェア資源として見える。次に、ゲストオペレーティングシステム142は、1つ又は数個のアプリケーション143を実行するのに使用することができる。

20

30

#### 【0014】

スケジューリングルーチン125は、仮想マシン141に関連したどの実行可能なスレッドが各プロセッサ131で実行されるかを決定する。実行可能なスレッドには、保証された時間の比率で、各プロセッサ131で実行される機会が与えられる。この比率は、所与のスケジューリング間隔については、部分的に、実行可能なスレッドのグループに従って定義される。たとえば、各仮想マシン141をグループに割り当てることができ、シェア122は、さまざまなグループについて定義される。各シェアパラメータは、各グループの仮想マシン141が平均して受け取るべきプロセッサ「チック」の最小量を表す。

#### 【0015】

シェアは、仮想マシングループの現在の要求と組み合わせられて、重み付き資源要求に変換される。IPPアルゴリズム124は、これらの重みを使用して、各グループを1組の物理CPUにマッピングする。このマッピングは、（要素123に記憶された）配分と呼ばれる。各グループについて、このマッピングは、CPU、いくつのスレッドが各CPUで実行されるか、及び、時間の比率は何であるかのリストを含む。IPPアルゴリズム124によって生成された配分は、各CPUがサービス提供する全重みが可能な限り均一となるようにする。

40

#### 【0016】

所与のスケジューリング間隔内において、スケジューリングルーチン125は、各配分123及びスケジューリングパラメータ126を使用して、各グループ内のどの実行可能プログラムが各プロセッサ131で実行されるかを決定する。前述したように、選択され

50

た配分 1 2 3 は、各グループに利用可能な物理 CPU を定義する。スケジューリングルーチン 1 2 5 は、スケジューリングパラメータ 1 2 6 を使用して、各グループからのどの特定のスレッドが、この間隔の間、そのリストのどの CPU で実行されるかを決定する。スケジューリングパラメータ 1 2 6 は、さまざまな実行可能プログラムが受け取ったプロセッサチックの受け取り履歴を示す。最も高いパラメータの値を有する実行可能プログラムが、最良の利用可能な物理 CPU に対して選択される。スケジューリング間隔の完了時に、累積したプロセッサチックが平均プロセッサチックよりも少ない実行可能プログラムは、自身のパラメータをインクリメントし、平均よりも多くのプロセッサチックを累積した実行可能プログラムは、自身のパラメータを減らす。

【 0 0 1 7 】

仮想プロセッサに関連したマッピング及びスケジューリングが考察されたが、他の代表的な実施の形態は、あらゆる適切なマルチプロセッサコンピュータシステムのあらゆるタイプの実行可能プログラムをスケジューリングするのに使用することができる。これに加えて、マッピング及びスケジューリングは、コンピュータのタイムスライスされたあらゆるタイプの資源（たとえば、ネットワーク接続カード、ディスク I/O チャンネル、暗号化デバイス等）について行うことができる。

【 0 0 1 8 】

図 2 は、代表的な一実施の形態による、ソフトウェアジョブのグループのプロセッサへのマッピングを生成するためのフローチャートを示している。図 2 は、適切なコンピュータ可読媒体から取り出されたソフトウェアコード又はソフトウェア命令を使用して実施される。ステップ 2 0 1 では、複数のジョブをサポートする複数のグループが定義される。一実施の形態では、各ジョブは、各仮想マシンによってサポートされる。各仮想マシンは、1 つ又は数個の仮想プロセッサを備える。ステップ 2 0 2 では、それらグループのシェアが定義される。これらのシェアは、各グループが平均して受け取る機会を有するプロセッサ資源の量を定義する。一実施の形態では、シェアは、1 0 0 チックが 1 秒の時間内における単一の物理プロセッサの全能力を表す、プロセッサ「チック」を符号化する。いくつかの実施の形態では、他のグループに明示的に割り当てられていないチックのすべてを受け取る最も優先順位の低いグループ又はデフォルトのグループが定義される。たとえば、或るコンピュータシステムが 6 つのジョブ（A、B、C、D、及び E）をサポートし、2 つのプロセッサ（2 0 0 の全シェアが利用可能である）を有するものと仮定する。ジョブ A は、「高い」優先順位のグループに割り当てられ、8 0 のシェアを受け取る。ジョブ B は、「中程度」の優先順位のグループに割り当てられ、4 5 のシェアを受け取る。ジョブ C、D、及び E は、デフォルトのグループに割り当てられ、このデフォルトのグループには、残りの 7 5 のシェアが割り当てられ、各グループは、平均して約 2 5 のシェアを受け取る。実行可能プログラムをグループに割り当てることによって、整数区画問題の組み合わせ複雑度が明らかに削減される。

【 0 0 1 9 】

ステップ 2 0 3 では、変数（N）が、（i）スケジューリングの目的のためにコンピュータシステムで利用可能な物理プロセッサの個数、及び、（i i）デフォルトのグループのアクティブなジョブの個数、のうちの最小値に等しくなるように設定される。

【 0 0 2 0 】

ステップ 2 0 4 では、グループのシェアパラメータが、異なる重みに区分される。いくつかの実施の形態では、デフォルトのグループのシェアパラメータは、N 個の異なる等しい重み（又は、丸め誤差を考慮してほぼ等しい重み）に分割される。前の 2 つのプロセッサの例を使用すると、最初の繰り返しにおいて、デフォルトのグループのシェアパラメータ（7 5）は、3 7 の第 1 の重み及び 3 8 の第 2 の重みに分割することができる。いくつかの実施の形態では、グループのシェアは、これに加えて、マルチスレッド化されたジョブを考慮した異なる重みに区分される。たとえば、ジョブ A が、2 つの仮想プロセッサを有する仮想マシンを使用して実施されるものと仮定する。高い優先度のグループの 8 0 のシェアは、2 つの仮想プロセッサのスレッドをサポートするために、4 0 の 2 つの重みに

10

20

30

40

50

分割することができる。或るグループ（デフォルトのグループ以外）がマルチスレッド化されたジョブを含まない場合、そのグループには、そのシェアパラメータと等しい単一の重みが生成される。

#### 【 0 0 2 1 】

ステップ 2 0 5 では、プロセッサ間での重みの配分を制限するための制約（constraints）が定義される。これらの制約は、あらかじめ定義された 1 組のルール又は条件に従って自動的に生成することができる。たとえば、マルチスレッド化されたジョブについて複数の資源要求の重みが生成される場合、これらの重みが同じプロセッサに割り当てられることを防止するための制約が定義される。また、特定のシステムについて、たとえば、高可用性アプリケーションに使用される冗長なソフトウェアモジュールを分離する（separate）ための制約を手動で定義することもできる。

10

#### 【 0 0 2 2 】

ステップ 2 0 6 では、I P P アルゴリズムが使用されて、プロセッサ間における重みの最適なバランスを達成するように、プロセッサのリスト全体にわたって重みの配分が生成される。この生成された配分は、後の解析（ステップ 2 1 0 を参照）のために一時的に記憶される。「貪欲（greedy）」方法等の既知の I P P アルゴリズムを使用することができる。この貪欲方法では、すべての重みが割り当てられるまで、前に割り当てられた全重みが最も低い「ピン」に、残っている最も高い重みが割り当てられる。代替的に、「差分」方法を使用することもできる。この「差分」方法では、異なるサブセットの最大数を配置するとともに、それらの差分を新たな数として挿入することによって、割り当てが行われる。これらの数のすべてがこのようにして割り当てられた後、元の重みの配分が、後ろ向き再帰（backward recursion）によって求められる。I P P アルゴリズムのインプリメンテーションに関する詳細は、いくつかの出典から入手可能である。たとえば、I P P アルゴリズムの概要は、2 0 0 2 年 1 2 月 1 4 日の Haikun Zhu 著の論文「On the Integer Partitioning Problem: Examples, Intuition and Beyond」に与えられている。

20

#### 【 0 0 2 3 】

いくつかの実施の形態によれば、解は、まず、高速の貪欲方法を使用して計算される。この解が完全でない場合、N 次元差分方法が使用され、最も高い正確度を有する解が選択される。プロセッサピンへの重みの配分は、各グループに利用可能な C P U の選択肢を決定する。C P U の全重みによって除算された、特定のジョブに関連した重みは、最終的に提供されるその C P U の部分を決定する。特定のスレッドの C P U への明示的なマッピングはこの段階では行われていないことに留意すべきである。その代わりに、スレッドのグループのみが 1 組の C P U にマッピングされている。これに加えて、個々の配分が生成された後、その配分が妥当かどうか（たとえば、制約が満たされるかどうか）を判断するための論理比較（フローチャートには図示せず）が行われる。配分が妥当でない場合、その特定の配分のさらなる使用を省くことができる。代替的に、プロセッサピンへの重みの割り当て中に、I P P アルゴリズムのピン割り当てロジックの変更によって、制約に対処することができる。

30

#### 【 0 0 2 4 】

ステップ 2 0 7 では、生成された配分が完全である（たとえば、各プロセッサに、計画された作業の同じ全重みが割り当てられている）かどうかを判断するための論理比較が行われる。生成された配分が完全である場合、生成された配分は記憶され、その配分は、その後のスケジューリングオペレーションで利用可能になる（ステップ 2 1 1）。また、前に計算された完全でない配分は、完全な配分の生成時に消去することができる。

40

#### 【 0 0 2 5 】

生成された配分が完全でない場合、プロセスフローは、ステップ 2 0 8 に進み、ステップ 2 0 8 において、変数 N が 1 に等しいかどうかを判断するための別の論理比較が行われる。変数 N が 1 に等しくない場合、プロセスフローはステップ 2 0 9 に進み、ステップ 2 0 9 において、N がデクリメントされて、プロセスフローはステップ 2 0 4 に戻る。したがって、デフォルトのグループに関連した重みの個数は変更され、それらの重みの値は変

50

更される。このように整数区画問題を変更して問題を再度解くことにより、配分の正確度を改善することができ、正確な配分を得る確率が增大する。

【 0 0 2 6 】

ステップ 2 0 8 の論理比較において、N が 1 に等しい場合、プロセスフローはステップ 2 1 0 に進む。ステップ 2 1 0 では、記憶された配分が調べられて、M 個の最良な配分（すなわち、各プロセッサに割り当てられた重み間の差分を最小にする配分）が特定される。ステップ 2 1 1 では、特定された配分が記憶されて、それらの配分がその後のスケジューリングオペレーションで利用可能になる。

【 0 0 2 7 】

いくつかの代表的な実施の形態では、図 2 のプロセスフローは、スケジューリングオペレーションの観点から、比較的低い頻度で実行される。具体的には、利用可能なプロセッサの個数が変化しない限り、又は、グループへのシェアの割り当てが変化しない限り、このプロセスフローの結果は変わらない。したがって、図 2 のプロセスフローは、大きなオーバーヘッドを課すことはなく、ワークロードの性能を低減させない。

【 0 0 2 8 】

図 4 は、代表的な一実施の形態による、図 2 のフローチャートに従って生成できる配分 4 0 0 を示している。或るシステムが 8 つのジョブ（A ~ G）をサポートするものと仮定する。このシステムは、3 つのプロセッサを含み、したがって、3 0 0 のシェアが利用可能である（ $3 * 1 0 0$ ）。また、ジョブ A が 8 0 のシェア値を割り当てられ、2 つの仮想プロセッサの仮想マシンに関連付けられるものと仮定する。さらに、ジョブ B、C、及び D がそれぞれ 6 0 のシェア値を割り当てられるものと仮定する。ジョブ A ~ D は、単一のジョブグループ（I ~ IV）に割り当てられる。ジョブ E ~ G は、デフォルトのグループ（グループ V）に割り当てられる。このデフォルトのグループは、4 0 のシェア値、すなわち、他のグループに割り当てられていないシェア量（ $3 0 0 - 8 0 - 6 0 - 6 0 - 6 0$ ）を受け取る。ジョブ A を含むグループ I のシェア値は、2 つの仮想プロセッサをサポートする 2 つの重みに分割される。これらの重みが同じ物理プロセッサに割り当てられることを防止するための制約も定義される。

【 0 0 2 9 】

これらの重み及び制約についての I P P 解法プロセスは、図 4 に示すような配分 4 0 0 という結果になる場合がある。グループ I の第 1 の重みはプロセッサ 1 に割り当てられ、グループ I の第 2 の重みはプロセッサ 2 に割り当てられる。グループ I I の重み、グループ I I I の重み、及びグループ I V の重みは、それぞれ、プロセッサ 1、2、及び 3 に割り当てられる。この場合、グループ V のシェア値は、複数の重みに分割されず、グループ V の単一の重みがプロセッサ V に割り当てられる。次に、グループ A ~ G の実行可能プログラムのスケジューリングが、配分 4 0 0 で特定されたプロセッサに行われる。

【 0 0 3 0 】

図 3 は、代表的な一実施の形態による、特定の物理 C P U への個々のジョブのスケジューリングを実行するためのフローチャートを示している。図 3 は、適切なコンピュータ可読媒体から取り出されたソフトウェアコード又はソフトウェア命令を使用して実施される。たとえば、システム割り込みに応答して呼び出された、オペレーティングシステムのスケジューリングソフトウェアルーチンを、図 3 のフローチャートを実施するのに使用することができる。

【 0 0 3 1 】

ステップ 3 0 1 では、ジョブスケジューリングパラメータが、ジョブによるプロセッサチェックの受け取りに従って更新される。時間サンプリング間隔中に受け取るジョブがグループの平均よりも少ない場合、それらジョブのパラメータはインクリメントされる。受け取るジョブがグループの平均よりも少ない場合、それらジョブのパラメータはデイクメントされる。アイドル状態にあるか、又は、要求の少ないジョブに関連したパラメータの値は、次第に 0 に減衰させることが可能な場合がある。

【 0 0 3 2 】

10

20

30

40

50

ステップ302では、グループの1つ又は複数の誤差が計算される(もしあれば)。ステップ303では、累積したあらゆるグループの誤差を補正する配分が選択される。具体的には、正確な配分が特定されなかったことから、複数の配分が生成されている場合、プロセスフローのさまざまな繰り返しにおいて、配分を交互に行う(alternation)ことができる。たとえば、配分Aがグループ1に有利であり、配分Bがグループ2に有利である場合、それら2つの配分を交互に行うことによって、より正確な方法でジョブ間のスケジューリングを行うことが可能になる。正確な配分が特定された場合、その正確な配分が使用される。

#### 【0033】

ステップ304では、各グループのジョブが、選択された配分に従い、且つ、各ジョブスケジューリングパラメータを使用してスケジューリングされる。具体的には、各グループについて、そのグループのジョブが、それらジョブの各ジョブスケジューリングパラメータによって順序付けられる。配分によって定義されたそのグループのCPUのリストが「望ましさ」によって順序付けられる。具体的には、全スケジューリング重みの低いCPUほど、より大きな望ましさを所有する。その理由は、このようなCPUの処理能力は、異なるグループの実行可能プログラム用に比較的大きなセグメント又は部分に分割されるからである。複数のCPUの全スケジューリング重みが等しい場合、CPUの履歴に基づく使用量を使用して、相対的な望ましさを求めることができる。具体的には、或るCPUが、履歴に基づく使用量が少ないことを示している場合、或るジョブが、当該ジョブにスケジューリングされた処理能力の部分を使用しない確率は高くなり、このような能力は別のジョブが使用することができる。

#### 【0034】

IPPアルゴリズムを使用して実行可能プログラムのグループをプロセッサへマッピングすること、及び、実行可能プログラムによる処理資源の受け取りを監視することによって、各グループ内の各ジョブは、同じ量のプロセッサ能力を経験することが可能になる。したがって、いくつかの代表的な実施の形態は、他の既知のマルチプロセッサスケジューリングアルゴリズムよりもはるかに「公平な」スケジューリングアルゴリズムを提供する。これに加えて、不完全な配分及び要求が少ないジョブは、限られた個数の間隔の間しかジョブに影響を与えない。具体的には、ジョブスケジューリングパラメータを使用して個々のジョブを特定のプロセッサにマッピングすることによって、このような問題を、ジョブのサブセットの損害につながる、スケジューリングオペレーションの恒久的なゆがみから防止する。グループ間の不完全性は、IPPアルゴリズムによって生成された複数の配分を交互に繰り返すことを使用して対処することができる。これに加えて、グループマッピングを実行可能プログラムの割り当てから分離することによって、いくつかの代表的な実施の形態が課すオーバーヘッドは比較的低くなり、それによって、アプリケーションからスケジューリングオペレーションへのプロセッサ資源の転換が省かれる。

#### 【図面の簡単な説明】

#### 【0035】

【図1】代表的な一実施の形態による、複数の物理プロセッサに仮想プロセッサをスケジューリングするシステムを示す図である。

【図2】代表的な一実施の形態による、1組のCPUに実行可能プログラムの各グループをマッピングする1つ又は数個の配分を生成して、スケジューリングオペレーションをサポートするIPPアルゴリズムを含むフローチャートである。

【図3】代表的な一実施の形態による、特定の物理CPUへの個々のジョブのスケジューリングを行うためのフローチャートである。

【図4】実行可能プログラムのグループと複数のプロセッサとの間のマッピングを定義する配分を示す図である。

#### 【符号の説明】

#### 【0036】

120・・・ホストOS

10

20

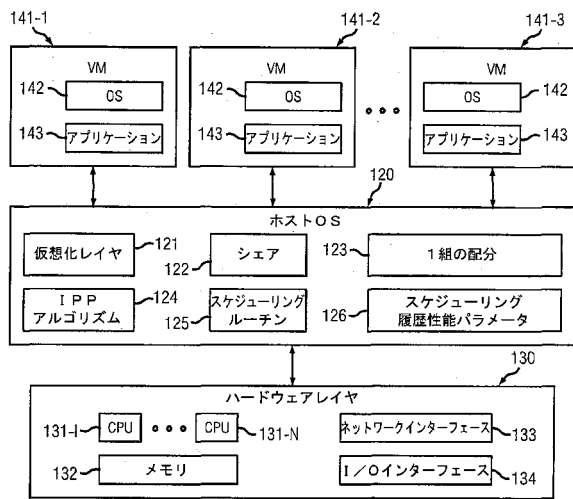
30

40

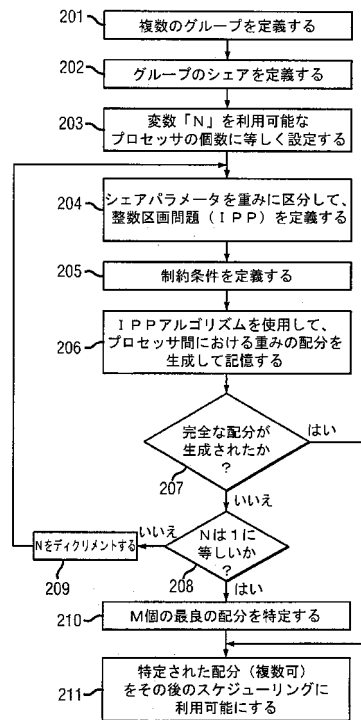
50

- 1 2 1 . . . 仮想化レイヤ
- 1 2 2 . . . シェア
- 1 2 3 . . . 1組の配分
- 1 2 4 . . . IPPアルゴリズム
- 1 2 5 . . . スケジューリングルーチン
- 1 2 6 . . . スケジューリング履歴性能パラメータ
- 1 3 0 . . . ハードウェアレイヤ
- 1 3 2 . . . メモリ
- 1 3 3 . . . ネットワークインターフェース
- 1 3 4 . . . I/Oインターフェース
- 1 4 3 . . . アプリケーション

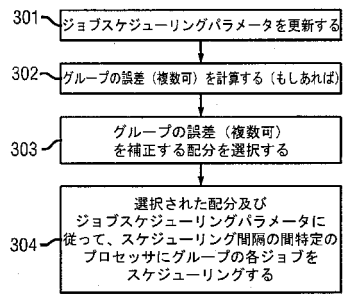
【図1】



【図2】



【 図 3 】



【 図 4 】

400

	プロセッサ		
	1	2	3
グループ及び重み	II (60)	III (60)	IV (60)
	I (40)	I (40)	V (40)

フロントページの続き

(56)参考文献 特開2002-202959(JP,A)  
特開2003-157177(JP,A)

(58)調査した分野(Int.Cl., DB名)  
G06F 9/46 - 9/54