

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4963018号
(P4963018)

(45) 発行日 平成24年6月27日(2012.6.27)

(24) 登録日 平成24年4月6日(2012.4.6)

(51) Int.Cl. F 1
G 0 6 F 9/48 (2006.01) G 0 6 F 9/46 4 5 2 C

請求項の数 10 (全 15 頁)

(21) 出願番号	特願2005-235582 (P2005-235582)	(73) 特許権者	310021766
(22) 出願日	平成17年8月15日(2005.8.15)		株式会社ソニー・コンピュータエンタテインメント
(65) 公開番号	特開2007-52511 (P2007-52511A)		東京都港区港南1丁目7番1号
(43) 公開日	平成19年3月1日(2007.3.1)	(74) 代理人	100105924
審査請求日	平成20年7月22日(2008.7.22)		弁理士 森下 賢樹
		(74) 代理人	100109047
			弁理士 村田 雄祐
		(74) 代理人	100109081
			弁理士 三木 友由
		(74) 代理人	100134256
			弁理士 青木 武司

最終頁に続く

(54) 【発明の名称】 スケジューリング方法およびスケジューリング装置

(57) 【特許請求の範囲】

【請求項 1】

マルチプロセッサシステムに含まれる複数のプロセッサのいずれにおいても実行可能な実行対象となる各実行単位に、当該実行単位とは分離して保存される識別子を付与し、

各実行単位のそれぞれが実行可能な状態にあるか、あるいはあるプロセッサで実行中であって別のプロセッサからは実行不可能な状態であるかを前記実行単位の識別子に対応づけて示す実行可否情報と、各実行単位のうちの、直近に実行された実行単位の識別子を示す直近実行情報とを含む実行単位情報を保持し、

実行単位情報に基づいて、直近に実行された実行単位の識別子以外の識別子を優先的に選出する制約の下でいずれかの実行可能な実行単位の識別子を、プロセッサによって実行される実行単位の識別子として選出するとともに、実行単位情報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 2】

前記実行単位情報が、各プロセッサがアクセス可能なメモリに保持され、

前記選出および更新は、選出された識別子に対応する実行単位を実行するプロセッサ自身によって行われることを特徴とする請求項 1 記載のスケジューリング方法。

【請求項 3】

実行可否情報は、各実行単位について識別子として1ビットを割り当てたビット列として保持され、

前記選出および更新は、不可分操作によって行われることを特徴とする請求項 1 または

2 記載のスケジューリング方法。

【請求項 4】

直近に実行された実行単位に対応するビットが末尾になるようにビット列をローテートし、

ローテートされたビット列の先頭から順に実行可能な実行単位のビットを検索することによって前記選出を行うことを特徴とする請求項 3 記載のスケジューリング方法。

【請求項 5】

マルチプロセッサシステムに含まれる複数のプロセッサのいずれにおいても実行可能な実行対象となる各実行単位のそれぞれが実行可能な状態にあるか、あるいはあるプロセッサで実行中であって別のプロセッサからは実行不可能な状態であるかを、それぞれの実行 10
単位に対して付与された識別子に対応づけて示す実行可否情報と、各実行単位のうちの、直近に実行された実行単位の識別子を
示す直近実行情報とを含む実行単位情報を、前記各
実行単位とは分離して保持する実行単位情報保持部と、

実行単位情報に基づいて、直近に実行された実行単位の識別子以外の識別子を優先的に選出する制約の下でいずれかの実行可能な実行単位の識別子を、プロセッサによって実行される実行単位の識別子として選出する実行単位選出部と、

前記選出に伴い、実行単位情報を更新する実行単位情報更新部とを備えることを特徴とするスケジューリング装置。

【請求項 6】

実行単位情報保持部は、各プロセッサがアクセス可能なメモリに実行単位情報を保持し 20

、
 実行単位選出部および実行単位情報更新部は、選出された識別子に対応する実行単位を実行するプロセッサ自身により構成されることを特徴とする請求項 5 記載のスケジューリング装置。

【請求項 7】

実行単位情報保持部は、実行可否情報を、各実行単位について識別子として 1 ビットを割り当てたビット列として保持し、

実行単位選出部および実行単位情報更新部は、不可分操作によって前記選出および更新を行うことを特徴とする請求項 5 または 6 記載のスケジューリング装置。

【請求項 8】

実行単位選出部は、直近に実行された実行単位に対応するビットが末尾になるようにビット列をローテートし、

ローテートされたビット列の先頭から順に実行可能な実行単位のビットを検索することによって前記選出を行うことを特徴とする請求項 7 記載のスケジューリング装置。

【請求項 9】

マルチプロセッサシステムに含まれる複数のプロセッサのいずれにおいても実行可能な実行対象となる各実行単位のそれぞれが実行可能な状態にあるか、あるいはあるプロセッサで実行中であって別のプロセッサからは実行不可能な状態であるかを、それぞれの実行 40
単位に対して付与された識別子に対応づけて示す実行可否情報と、各実行単位のうちの、直近に実行された実行単位の識別子を
示す直近実行情報とを含む実行単位情報を、前記各
実行単位とは分離して保持する機能と、

実行単位情報に基づいて、直近に実行された実行単位の識別子以外の識別子を優先的に選出する制約の下でいずれかの実行可能な実行単位の識別子を、プロセッサによって実行される実行単位の識別子として選出するとともに、実行単位情報を更新する機能とをコンピュータに実行せしめることを特徴とするプログラム。

【請求項 10】

プログラムを記憶した記憶媒体であって、

前記プログラムは、

マルチプロセッサシステムに含まれる複数のプロセッサのいずれにおいても実行可能な 50
実行対象となる各実行単位のそれぞれが実行可能な状態にあるか、あるいはあるプロセッ

サで実行中であって別のプロセッサからは実行不可能な状態であるかを、それぞれの実行単位に対して付与された識別子に対応づけて示す実行可否情報と、各実行単位のうちの、直近に実行された実行単位の識別子を示す直近実行情報とを含む実行単位情報を、前記各実行単位とは分離して保持する機能と、

実行単位情報に基づいて、直近に実行された実行単位の識別子以外の識別子を優先的に選出する制約の下でいずれかの実行可能な実行単位の識別子を、プロセッサによって実行される実行単位の識別子として選出するとともに、実行単位情報を更新する機能とをコンピュータに実行せしめることを特徴とする記憶媒体。

【発明の詳細な説明】

【技術分野】

10

【0001】

本発明は、マルチプロセッサシステムにおける並列処理の実行単位のスケジューリング方法および装置に関する。

【背景技術】

【0002】

最近のマルチタスクをサポートするオペレーティングシステムは、複数のプロセスを同時に実行することができるマルチタスク環境を実現するとともに、これらのプロセスがプロセス内部で複数のスレッドを生成して並行処理を行うことができるマルチスレッド技術を搭載している。プロセスは実行時に固有のリソースやアドレス空間が割り当てられ、他のプロセスの領域にアクセスすることができない。これに対してスレッドは、プロセス内部で生成される実行単位であり、各スレッドはプロセス内の領域に互いに自由にアクセスすることができる。スレッドは、オペレーティングシステムがCPU(Central Processing Unit)の実行時間を割り当てる基本的な実行単位となる。本明細書において、スレッドの割り当てをスケジューリングという。

20

【0003】

1つのシステム内に複数のプロセッサを搭載したマルチプロセッサシステムでは、並列に、あるいは協調して処理を実行して処理全体の高速化を図ることができる。通常、マルチプロセッサシステムにおいて、共有メモリにタスクキューと呼ばれる待ち行列に実行可能なスレッドが保持される。これらのスレッドはいずれかのプロセッサに割り当てられて、実行される。スレッドを割り当てる方法によって、プロセスの実行速度やメモリ消費量などの性能が変わってくるため、マルチプロセッサシステムでは、シングルプロセッサシステムにおけるスケジューリングと異なった工夫が必要となる。

30

【0004】

マルチプロセッサシステムにおけるスケジューリングについて、たとえば、スレッドを管理する役割を担う管理ユニットが各プロセッサと通信し、それぞれのプロセッサに割り当てるスレッドをスケジューリングする方法が考えられる。

管理ユニットによってスケジューリングを行うこの方法では、管理ユニットとプロセッサ間の通信において、メッセージ遅延が起きることが多く、スレッドの実行を遅延させてしまうという問題がある。この問題を解決するために、各プロセッサによって自己支配的にスケジューリングする方法が考えられる。この方法では各プロセッサ上でおのこのスケジューラを実行し、共有メモリにあるタスクキューにアクセスして、実行するスレッドを選択する。

40

【0005】

この方法では、各プロセッサがスケジューラを実行している間、すなわちスケジューリング中において、共有メモリ上のタスクキューが他のプロセッサによって修正されることを防ぐために、スケジューラがタスクキューをロックする。この間、割り込みを禁止するか否かはシステムの設計者にとってジレンマである。

【0006】

スケジューリング中に割り込みを禁止しない場合には、スケジューラがタスクキューをロックしている最中に、割り込みを受け付けてしまう。この場合、割り込み処理が終了し

50

、ロックが開放されるまで、スケジューリングが実行されない結果になってしまい、システムの効率低下につながる。

【0007】

一方、スケジューリング中に割り込みを禁止すると、システムの割り込み応答性能が低下してしまう問題がある。

【0008】

これらの問題はこれまで述べたマルチプロセッサ上で動作するカーネルスケジューラだけの問題でなく、マルチプロセッサ・マルチスレッド環境でのユーザレベルスケジューラ（すなわち、各プロセッサで動作するスレッド上でスケジューラを実行することで、ユーザレベルにおいてマルチスレッドを実現する方式）の場合も同様である。

特に、ユーザレベルスケジューラを実行しているカーネルスレッドは各プロセッサ上のカーネルレベルスケジューラによってスケジューリングされ、他のカーネルスレッドによってプリエンプト（横取り）される可能性があるため、タスクキューのロックはより深刻な問題を引き起しかねない。この場合、横取りされた際にそのカーネルスレッドで動作するユーザレベルスケジューラがロックを取得していた場合には、他のプロセッサのスレッド上で動作するユーザレベルスケジューラがロックが解除されるまでスケジューリングすることができなくなる。

【発明の開示】

【発明が解決しようとする課題】

【0009】

本発明は上記事情に鑑みてなされたものであり、その目的は、マルチプロセッサシステムにおける並列処理の実行単位をプロセッサに割り当てる順序を制御して処理効率を向上させることができるスケジューリング技術を提供することにある。

【課題を解決するための手段】

【0010】

本発明にかかる態様は、マルチプロセッサシステムにおいて、プロセッサの実行対象となる実行単位をスケジューリングする方法に関する。このスケジューリング方法は、マルチプロセッサシステムに含まれる複数のプロセッサの実行対象となる各実行単位に識別子を付与し、各実行単位のそれぞれが実行可能な状態にあるか否かを実行単位の識別子に対応づけて示す実行可否情報と、各実行単位のうちの、直近に実行された実行単位の識別子を示す直近実行情報とを含む実行単位情報を保持する。そして、実行単位情報に基づいて、直近に実行された実行単位の識別子以外の識別子を優先的に選出する制約の下でいずれかの実行可能な実行単位の識別子を、プロセッサによって実行される実行単位の識別子として選出するとともに、実行単位情報を更新する。

【0011】

この態様では、実行可否情報と直近実行情報を含む実行単位情報が保持されるとともに、次に実行する実行単位の選出に伴って更新される。実行単位の選出は、実行単位情報を参照してなされる。そのため、実行単位の選出にあたって、実行単位情報を参照すればよく、実行単位の実体を記憶したたとえばタスクキューをロックする必要がない。そのため、スケジューリング中に割り込み処理を受け付けるか否かのジレンマが解消される。

【0012】

また、直近実行情報が保持され、実行単位の識別子を選出する際に、直近実行情報により示される識別子以外の識別子を優先的に選出するとともに、この選出に伴って直近実行情報を更新するようにしているので、直近に実行された実行単位がまた実行可能になっても、他の実行可能な実行単位が優先的に実行されるので、実行単位のスケジューリングにおいて重要な公平性を保つことができる。

【0013】

本発明のこの態様は、管理ユニットによってスケジューリングするシステムに適用してもよく、各プロセッサ自身によってスケジューリングを行うシステムに適用してもよい。

また、本発明のこの態様は、各プロセッサ上で直接動作するスケジューラだけではなく

10

20

30

40

50

、各プロセッサ上のスケジューラが提供するスレッド上で動作するユーザレベルスケジューリングを用いるシステムに適用してもよい。

【0014】

ここで、実行可否情報は、各実行単位について識別子として1ビットを割り当てたビット列として保持され、ビットの選出および実行単位情報の更新を不可分操作、すなわちアトミック操作によって行うようにしてもよい。

【0015】

「不可分操作」は、これ以上分割できない最小単位の操作を意味し、マルチプロセッサシステムにおいて、他のプロセッサが行う操作との相互作用なしに実行されることが保証される操作である。

10

【0016】

ビットの選出に当たっては、直近に実行された実行単位に対応するビットが末尾になるようにビット列をローテートし、ローテートされたビット列の先頭から順に実行可能な実行単位のビットを検索することによって行うようにしてもよい。

【0017】

なお、以上の構成要素の任意の組合せ、本発明をシステム、プログラム、プログラムを記憶した記憶媒体として表現したのも、本発明の態様としては有効である。

【発明の効果】

【0018】

本発明は、マルチプロセッサシステムにおける並列処理の実行単位をスケジューリングすることにおいて有利である。

20

【発明を実施するための最良の形態】

【0019】

マルチプロセッサシステムにおいて、各プロセッサの実行単位たとえばスレッドのスケジューリング中に、スレッドを記憶したタスクキューがほかのプロセッサにより修正されることを防ぐために、タスクキューをロックする必要があった。

【0020】

このロックによって、スケジューリング中に、割り込みを禁止するとシステムの割り込み応答性能が低下し、割り込みを禁止しないと、割り込み処理が終わり、ロックが解除されるまでスケジューリングができないというジレンマがあった。

30

【0021】

上述した問題を解決するために、本発明者は、下記の技術を提案する。

【0022】

各スレッドに識別子を付与し、これらのスレッドのそれぞれが実行可能な状態にあるか否かをスレッドの識別子に対応づけて示す実行可否情報と、各スレッドのうちの、直近に実行されたスレッドの識別子を示す直近実行情報とを含むスレッド情報を保持する。そして、スレッド情報に基づいて、直近に実行されたスレッドの識別子以外の識別子を優先的に選出するという制約の下で、いずれかの実行可能なスレッドの識別子を、プロセッサに割り当てるスレッドの識別子として選出する。

【0023】

40

ここで、スレッドの実体（以下スレッド実体という）を、プロセッサが選出された識別子に基づいて取得することができるいかなる方法で保持してもよい。たとえばスレッド実体を識別子と対応づけて保存するようにしてもよいし、各スレッド実体をメモリのそれぞれの所定の領域に保存するとともに、スレッド実体が保存された領域の開始アドレスとその識別子とを対応づけてメモリに保存するようにしてもよい。

【0024】

この技術は、スレッド実体と、いずれのスレッドを実行するかを選出を行うために必要なスレッド情報とに分けてメモリに保存する。そのために、スケジューリングする際に当たって、スレッド実体を記憶した領域をロックする必要がなく、スケジューリング中に割り込みを禁止するか否かのジレンマを解消することができる。

50

【 0 0 2 5 】

この技術は、ユーザレベルスケジューラを実行するシステムにおいても、上記ジレンマを解消することができるとともに、ユーザレベルスケジューラゆえの問題も解決することができる。

【 0 0 2 6 】

たとえば、マルチプロセッサシステム、特にOS（オペレーティングシステム）機能を実行できるプロセッサに限られるような非対称のマルチプロセッサシステムにおいては、各プロセッサで動作するOSが提供するカーネルスレッド上にユーザレベルスレッドを作成して、自己支配的にスケジューリングする方法、すなわちユーザレベルスケジューラを用いる方法は、マルチプロセッサシステムの処理効率を向上させる有効な方法であると考えられる。この方法では、ユーザレベルスケジューラを用いることに起因するスケジューリングの破綻という問題がある。

10

【 0 0 2 7 】

たとえば、ユーザレベルスケジューラを実行しているスレッドは、各プロセッサの上のカーネルスケジューラによってスケジューリングされ、ほかのカーネルスレッドによってプリエンプトされる可能性がある。そのため、プリエンプトされた際にそのスレッドで動作するユーザレベルスケジューラがロックを取得していた場合には、ほかのプロセッサのスレッド上で動作するユーザレベルスケジューラがロックが解除されるまで、スケジューリングができなくなり、プロセッサの処理効率を下げってしまう。さらに、ロックが解除されるまでの時間は、状況に依存してしまい、見積もることが困難になるため、システムの不安定を引き起こす原因にもなりえる。

20

【 0 0 2 8 】

ここで、マルチプロセッサシステムにおいて各プロセッサが自律的に動作するカーネルスケジューラを実行する従来のシステムについて考える。タスクキューにスレッド1、スレッド2、スレッド3が入っており、複数のプロセッサのうちのプロセッサAは、タスクキューにあるスレッドを実行することができる状態になった場合を想定する。

【 0 0 2 9 】

プロセッサAは、実行するスレッドを選出するために、タスクキューを一旦自分のローカルメモリにコピーし、タスクキューからたとばスレッド1を選出した後に、タスクキューからスレッド1を削除する更新処理を行って、更新されたタスクキューをメインメモリに書き戻す。タスクキューのコピーが開始するときから、タスクキューの書き戻しが終了するまでの間、共有メモリ上のタスクキューが他のプロセッサによって修正されることを防ぐために、タスクキューをロックする。タスクキューがロックされている間、他のプロセッサは、タスクキューに入ったスレッドを実行できる状態になったとしても、タスクキューを利用することができないため、ロックの期間が長いほどシステムの処理効率が低下する。

30

【 0 0 3 0 】

また、プロセッサAで動作するスケジューラは前述したように割り込みによってプリエンプトされる可能性がある。プロセッサAはタスクキューがロックされた状態でプリエンプトされると、ほかのプロセッサは、プロセッサAが割り込み処理から復帰し、ロックが解除されるまでスケジューリングを行うことができず、スレッド2、3を実行することが不可能となる。これでは、システムは処理効率が低下するとともに、不安定にもなりかねない。

40

なお、上記説明は、プロセッサA、・・・を、プロセッサA、・・・上で動作するユーザレベルスケジューラを実行するスレッド（たとえばスレッドa、・・・）に置き換えれば、ユーザレベルスケジューラを実行するシステムにも適用することができる。具体的には、ユーザレベルスケジューラを実行するシステムにおいて、たとえばプロセッサA上で動作するユーザレベルスケジューラを実行するスレッドaも、割り込みによってプリエンプトされる可能性がある。スレッドaは、タスクキューがロックされた状態でプリエンプトされると、ほかのスレッドは、スレッドaが割り込み処理から復帰し、ロックが解除さ

50

れるまでスケジューリングを行うことができない。すなわち、ユーザレベルスケジューラを用いるシステムにおいても同じように、タスクキューのロック起因して、システムは処理効率が低下し、不安定になりかねない問題がある。

【0031】

また、マルチプロセッサシステムにおいて、各々のプロセッサがそれぞれの処理ユニットに含まれる形で存在する。これらの処理ユニットはPPU(Power Processing Unit)とSPU(Synergistic Processing Unit)とに分けることができる。SPUのすべてが同一のアーキテクチャを用いて実現されてもよく、それぞれ異なる構成を有してもよい。PPUは、SPUに対してローカルに、たとえばSPUと同一のチップ、同一のパッケージ、同一の回路基板、同一の製品に位置してもよいし、SPUに対してリモートに、たとえばバスやインターネットなどの通信ネットワークを介して接続可能な異なる製品に位置してもよい。同様に、SPUは、互いにローカルにまたはリモートに位置してもよい。

10

【0032】

スケジューリング禁止区間の存在と禁止時間の見積もりの困難さはマルチプロセッサシステムの処理効率を下げる要因になる。ユーザレベルスケジューラを実行するスレッドのすべてがSPUのスレッド(以下SPUスレッドという)であれば、この問題を解決する方法としては、スレッドをグループ化してグループ単位でスケジューリングする方法が考えられる。ここで、図9に示すマルチプロセッサシステムを例にして説明する。

【0033】

図9に示すマルチプロセッサシステムは、複数の処理ユニット110とメインメモリ130と有し、それらはメインバス120に接続されている。各処理ユニット110は、プロセッサ112、ローカルメモリ114、メモリ制御部116を有する。プロセッサ112は、ローカルメモリ114に対してデータを読み書きすることができる。メモリ制御部116は、他の処理ユニット110のプロセッサ112からローカルメモリ114のデータを参照するときのインターフェースを与えるとともに、メモリの同期、排他制御の機能を提供する。

20

【0034】

ここで、メインメモリ130に設けられたタスクキューにSPUスレッドのみが記憶されている場面について考える。この場合、タスクキュー内のスレッドが図10に示すようにグルーピングされる。図10(a)は、3つのスレッドth1a、th1b、th1cを含む第1スレッドグループを示す。図10(b)は、1つのスレッドth2aを含む第2スレッドグループを示す。このようなスレッドが1つだけの場合もスレッドグループとして扱う。同様に、図10(c)は、2つのスレッドth3a、th3bを含む第3スレッドグループを示す。図10(d)は、1つのスレッドth4aを含む第4スレッドグループを示す。

30

【0035】

これらのスレッドのスケジューリングは、同一スレッドグループに属するすべてのスレッドを同時にいずれかのプロセッサ112に割り当ててを条件として行われる。第1スレッドグループがプロセッサ112に割り当てられるときは、第1スレッドグループに属する3つのスレッドth1a、th1b、th1cが同時にいずれかのプロセッサ112に割り当てることができる場合に限られる。3つのスレッドth1a、th1b、th1cの1つまたは2つがプロセッサ112に割り当てられ、残りがメインメモリ130に退避しているという状況は作らない。

40

【0036】

図11は、スレッドグループ単位でスレッドがプロセッサ112に割り当てられる様子を説明する図である。同図は、プロセッサ総数4のマルチプロセッサシステムにおいて、図10に示した4つのスレッドグループに属するスレッドのプロセッサ112への割り当て状態を示している。ある時刻において、第1スレッドグループに属する3つのスレッドth1a、th1b、th1cは、それぞれ第1プロセッサ、第2プロセッサ、第3プロ

50

セッサに割り当てられ、第2スレッドグループに属する1つのスレッドt h 2 aは、第4プロセッサに割り当てられている。それ以外の第3スレッドグループに属する2つのスレッドt h 3 a、t h 3 b、および第4スレッドグループに属する1つのスレッドt h 4 aはメインメモリ112に退避されている。

【0037】

このようなシステムによれば、1つのスレッドグループ内に所属する複数のスレッドを必ず同時にいずれかのプロセッサに割り当てるようにする。スレッドグループ内において、ロックを取ったスレッドだけがプリエンプトされることがないため、スケジューリングの禁止区間を限定することができる。

【0038】

しかし、PPUのスレッド(以下PPUスレッドという)とSPUスレッドが非同期にスケジューリングされる環境では、PPUスレッドとSPUスレッドで同じタスクキューを共有した場合には下記の問題が生じうる。あるPPUスレッドがタスクキューのロックを取ったまま他のPPUスレッドにプリエンプトされると、SPUスレッドはこのPPUスレッドが再実行されるまで待たざるを得ないため、スケジューリングの禁止区間を限定することができなくなってしまう。

【0039】

本発明者が提案したスケジューリング技術は、スレッド情報とスレッド実体を分けて保存し、プロセッサは、いずれのスレッドを実行するかを選出を行う際、スレッド情報のみをロードすればよい。選出を終え、スレッド情報の更新(具体的には、直近に実行されたスレッドの識別子を、選出されたスレッドの識別子に変更する処理と、選出されたスレッドの識別子を該識別子の対応するスレッドが実行不可であることを示すように修正する処理)を終了すれば、更新されたスレッド情報をメインメモリにストアする。その後、プロセッサは、選出した識別子に対応するスレッド実体をコピーするが、コピーする間、他のプロセッサは、スレッド情報を利用することができる。こうすることによって、スレッド情報は1つのプロセッサに占有される時間が短いため、スレッド情報をロックしたとしても、システム全体の処理効率の低下を軽減させることができる。

【0040】

さらに、本発明者は、実行可否情報を、各スレッドについて識別として1ビットを割り当てたビット列として保持することを提案する。こうすることによって、スレッドの選出、スレッド情報の更新などを、アトミック操作やアトミック命令を用いて行うことができ、ロック操作を伴わないすなわちロックレスのタスクキューを実現することができる。ロック操作がなければ、上述した、タスクキューのロックに起因するおのおの問題も解消される。

【0041】

図1は、本発明の実施形態となるマルチプロセッサシステム100の構成を示す。マルチプロセッサシステム100は、複数の処理ユニット10とメインメモリ30と有し、それらはメインバス20に接続されている。各処理ユニット10は、プロセッサ12、ローカルメモリ14、メモリ制御部16を有する。プロセッサ12は、ローカルメモリ14に対してデータを読み書きすることができる。メモリ制御部16は、他の処理ユニット10のプロセッサ12からローカルメモリ14のデータを参照するときのインターフェースを与えるとともに、メモリの同期、排他制御の機能を提供する。

【0042】

処理ユニット10のうちのいずれか1つを、スレッドのスケジューリングに関して他の処理ユニットに対するサービスユニットの役割を担う。サービスユニットの役割としては、たとえばメインメモリ30の割当てや、メインメモリ30内のスレッドに関する最初の記憶に関わることなどである。メインメモリ30の割当ては、たとえばスレッド情報に割り当てられる領域、領域の容量や、スレッドの実体に割り当てられるべきメモリ容量を決定することなどとするすることができる。

【0043】

10

20

30

40

50

なお、このサービスユニットは、いずれの処理ユニット10によって担当されてもよい。

【0044】

ある時刻において、各プロセッサ12には1つのスレッドが動作し、マルチプロセッサシステム100全体で並列に複数のスレッドが実行される。各プロセッサ12において動作しているスレッドは、処理ユニット10内のローカルメモリ14やメモリ制御部16内のレジスタなどのすべての資源を占有して使用することができる。

【0045】

この状態において、処理待ちしているスレッドは、そのコンテキストがメインメモリ30に保持される。スレッドのコンテキストは、そのスレッドが実行された処理ユニット10内で占有するすべての資源の状態であり、スレッドがプロセッサ12において動作しているときに各種レジスタに保持されている値の集合、ローカルメモリ14に保持されたデータ、メモリ制御部16の各種レジスタの内部状態などである。スレッドがプロセッサ12上で動作していないときは、そのスレッドのコンテキストをメインメモリ30にコピーしておき、再度プロセッサ12によって処理可能となったときに、そのコンテキストをメインメモリ30から読み込んで、処理を継続することができる。スレッドのコンテキストは、スレッド実体に該当する。

【0046】

図2は、メインメモリ30により記憶された、スレッドに関する情報を示す。これらの情報は、スレッド情報40、スレッドアドレス情報50と、スレッド実体60であり、マルチプロセッサシステム100におけるタスクキューの役割を担う。なお、これらの情報が記憶される領域は、サービスユニットによって割り当てられ、ほかの処理ユニット10に通知される。

【0047】

スレッド情報40は、実行可否情報と直近実行情報とを含む。実行可否情報は、各スレッドに対してそれぞれ付与された識別子を、それぞれのスレッドが実行可能な状態にあるか否かを示す情報に対応づけたものである。図3は、スレッド情報40の詳細を示す。

【0048】

ビット列(isSchedulable)は、実行可否情報であり、それに含まれる1ビットが1つのスレッドに対応し、ビット番号はスレッドの番号に対応する。なお、マルチプロセッサシステム100は、ビット番号そのものをスレッドの番号として用いて処理の簡潔化を図る。

【0049】

ビット列isSchedulableの各ビットの値はそのビットに対応するスレッドが実行可能状態にあるか否かを示している。ここで、ビットの値の「1」は実行可能であることを示し、「0」は実行不可であることを示すようにされている。ビット列isSchedulableに含まれるビットの数は、実行可否情報の保持用に割り当てられた領域の容量に相当し、ここでは128ビットとする。

【0050】

lastScheduledは、ビット列isSchedulableに含まれる各ビットのうち、直近に実行されたビットの番号を示す整数値を取る変数である。

【0051】

isSchedulableとlastScheduledによって、図4の例のように、実行可能なスレッドの番号、直近に実行されたスレッドの番号が示される。図4の例では、ビット列isSchedulableのうち、値が「1」であるビット(図中矢印B、C、Dが示すビット)が3つあり、この3つのビット番号に対応するスレッドは実行可能スレッドである。値が「0」であるほかのビットについて、それらに対応するスレッドは、実行不可である。また、lastScheduledの値に等しいビット番号(図中矢印Aが示すビットの番号)に対応するスレッドは、直近に実行されたスレッドであり、その状態は実行不可である。

10

20

30

40

50

【 0 0 5 2 】

スレッドアドレス情報 5 0 は、各スレッドの番号ここではビット列 `i s S c h e d u l a b l e` に含まれるビットの番号と、その番号に対応するスレッドの実体が保存された領域の開始アドレスとを対応づけたものである。どの番号のスレッド実体を、どのアドレスの領域に保存するかについては、サービスユニットによって決められる。

【 0 0 5 3 】

処理待ちのスレッドがない状態において、スレッド情報 4 0 のビット列 `i s S c h e d u l a b l e` の各ビットの値が「0」である。

【 0 0 5 4 】

処理ユニット 1 0 の処理が進み、スレッドが生成される。生成されたスレッドは、また
10
いずれかの処理ユニット 1 0 によって実行される。処理待ちのスレッドが生じた際に、タスクキューが利用される。ここで、例として図 4 に示すタスクキューの状態を起点にして、処理ユニット 1 0 が次に実行するスレッドを選出する処理、およびこの選出に伴うスレッド情報の更新処理を図 5 のフローチャートを用いて説明する。

【 0 0 5 5 】

図 4 に示す状態では、マルチプロセッサシステム 1 0 0 の各処理ユニット 1 0 は、それぞれ処理中のスレッドがあり、処理待ちしている実行可能なスレッドは 3 つあり、この 3 つのスレッドのそれぞれの実体は、メインメモリ 3 0 に記憶されている。

【 0 0 5 6 】

この状態において、ある処理ユニット 1 0 において、処理中のスレッドの処理が終了す
20
ると、この処理ユニット 1 0 のプロセッサ 1 2 は、次に実行するスレッドを選出するために、スレッド情報 4 0 に含まれるビット列 `i s S c h e d u l a b l e` と `l a s t S c h e d u l e d` を、ローカルメモリ 1 4 にロードする (S 1 0)。マルチプロセッサシステム 1 0 0 において、処理ユニット 1 0 は、スレッド情報 4 0 に関わる処理をアトミックコマンドを用いて行い、ここでは、スレッド情報 4 0 をロードするためのコマンドとして、たとえば「`l w a r x`」または「`g e t l l a r`」を用いる。なお、図 6 の A 欄に示すビット列 `i s S c h e d u l a b l e` は、図 4 に示すビット列 `i s S c h e d u l a b l e` である。

【 0 0 5 7 】

プロセッサ 1 2 は、ロードしたビット列 `i s S c h e d u l a b l e` (図 6 の A 欄に示
30
すビット列) を、`l a s t S c h e d u l e d` の値に等しい番号のビット (図 6 の矢印 A 1 が示すビット) が末尾になるように、矢印 L が示す方向、すなわち左方向にローテートする (S 1 4)。これによって、図 6 の A 欄のビット列 `i s S c h e d u l a b l e` は、同図の B 欄のビット列 `i s S c h e d u l a b l e` になる。図示のように、(`l a s t S c h e d u l e d + 1`) の番号のビット (矢印 B 1 が示すビット) が、ビット列 `i s S c h e d u l a b l e` の先頭に位置し、`l a s t S c h e d u l e d` の番号のビット (矢印 B 3 が示すビット) がビット列 `i s S c h e d u l a b l e` の末尾になる。

【 0 0 5 8 】

プロセッサ 1 2 は、続いてローテートされたビット列 `i s S c h e d u l a b l e` に対
40
して先頭から順に、値が「1」であるビットの検索を行い、最も先に検出したビット (矢印 B 2 が示すビット) の番号を、次に実行するスレッドの番号として得る (S 1 8)。ビットの検索に用いるコマンドとして、たとえばビット列 `i s S c h e d u l a b l e` の先頭から連続した、「0」の値を有するビットの数を数える「`C o u n t L e a d i n g Z e r o`」を用いることができる。「`C o u n t L e a d i n g Z e r o`」により得られた値 (図 6 に示す B 欄のビット列 `i s S c h e d u l a b l e` の例では 4) に (`l a s t S c h e d u l e d + 1`) を加算して得た値は、次に実行するスレッドの番号として選出される。

【 0 0 5 9 】

そして、プロセッサ 1 2 は、選出した番号のビットの値を「0」にセットするとともに
50
、`l a s t S c h e d u l e d` をこの番号にセットして、スレッド情報 4 0 の更新を行う

(S20)。更新されたスレッド情報40は、メインメモリ30にストアされる(S24)。ここで、ビットの値の更新は、図6のC欄に示す128ビットのビット列を用いて行う。C欄のビット列は、128ビットを有し、ステップS10において選択されたした番号と同じ番号のビット(矢印C1が示すビット)のみが「1」の値を有する。プロセッサ12は、図6のB欄のビット列*isSchedulable*と、C欄のビット列とを、「AtomicAndc」コマンドで演算することによってB欄のビット列*isSchedulable*を更新する。また、更新されたビット列*isSchedulable*のストアに用いるコマンドは、たとえば「stwcx」または「putllc」とすることができる。

【0060】

10

図6のD欄は、更新されたスレッド情報40を示す。ここで、次に実行するスレッドとして選出されたビット(矢印D1が示すビット)は、直近に実行されたスレッドを示すビットとなり、その値が「1」から「0」になっている。

【0061】

その後、プロセッサ12は、スレッドアドレス情報50を参照して、選出した番号のスレッドに対応するスレッド実体60の開始アドレスを取得するとともに、この開始アドレスにより示される領域からスレッド実体60をローカルメモリ14にロードして処理する。

【0062】

次に図7のフローチャートを用いて、タスクキューに実行可能となったスレッドを追加する処理について説明する。

20

【0063】

処理ユニット10のプロセッサ12は、新たに実行可能となったスレッドをタスクキューに追加するために、スレッド情報40に含まれるビット列*isSchedulable*をローカルメモリ14にロードする(S50)。このビット列*isSchedulable*はたとえば図8のA欄に示すビット列*isSchedulable*である。プロセッサ12は、A欄のビット列*isSchedulable*の各ビットのうちの、値が「0」であるいずれかのビット(たとえば矢印A1が示すビット)の番号を、追加するスレッドのビット番号として選出する。そして、プロセッサ12は、A欄のビット列*isSchedulable*と、B欄に示すビット列とを「AtomicOr」コマンドで演算してC欄に示すビット列*isSchedulable*を得る(S54)。図8のB欄に示すビット列は、128ビットを有し、プロセッサ12により選出された番号と同じ番号のビット(矢印B1が示すビット)のみが「1」の値を有する。

30

【0064】

続いて、プロセッサ12は、ステップS54により得られたビット列*isSchedulable*(図8のC欄のビット列)をメインメモリ30にストアして、スレッドを追加するためのスレッド情報40の更新を終了する(S58)。

【0065】

プロセッサ12は、この後、追加するスレッドの実体を、ステップS54において選出したビット番号に対して割り当てられた領域にコピーして、スレッドを追加する処理を終了する。

40

【0066】

このように、図1に示すマルチプロセッサシステム100によれば、各プロセッサ12自身によって、タスクキューからスレッドを選出する処理、タスクキューを更新する処理を行っているので、システム全体の処理効率を向上させることができる。

【0067】

また、タスクキューを構成する際に、スレッド情報とスレッド実体とに分けて保存し、スレッドの選出、更新は、スレッド情報のみを用いて行うことができるようにしたので、より効率の良いマルチプロセッサシステムを実現している。

【0068】

50

さらに、スレッド情報としてビット列 `isSchedulable` と、`lastScheduled` との2つの変数を用いることによって、スレッドの選出と更新を、アトミックコマンドで行うことを可能とした。これによって、ロックレスのタスクキューを実現している。

【0069】

ロックレスのタスクキューの実現により、スケジューリング中に割り込みを受け付けるか否かのジレンマを解消することができる。

ここで、図1に示すマルチプロセッサシステム100は、カーネルスケジューラのみを用いるマルチプロセッサシステムであるが、マルチプロセッサシステム100に用いられたスケジューリング方法は、各プロセッサ上で動作するスレッド上で実現されたユーザレベルスケジューラを用いたシステムにも適用することができる。その場合、スケジューリング中に割り込みを受け付けるか否かのジレンマを解消できるとともに、前述した、ユーザレベルスケジューラを用いることに起因する問題も解消することができる。

【0070】

以上、本発明を実施の形態をもとに説明した。実施の形態は例示であり、それらの各構成要素や各処理プロセスの組合せにいろいろな変形例が可能なこと、またそうした変形例も本発明の範囲にあることは当業者に理解されるところである。

【0071】

また、図1に示す実施形態は、各プロセッサが自律的にスケジューリングを行うシステムであるが、本発明のスケジューリング方法は、このようなシステムに限らず、たとえばひとつの管理ユニットによってスケジューリングするシステムにも適用することができる。

【0072】

また、本発明を適用したデバイスも、本発明の範囲にある。これらのデバイスには、パーソナルコンピュータやサーバなどに限らず、携帯電話、ゲーム機、モバイルコンピュータ、個人携帯情報機器(PDA)、デジタルテレビなどが含まれる。

【図面の簡単な説明】

【0073】

【図1】本発明にかかる実施形態のマルチプロセッサシステムを示す図である。

【図2】図1に示すマルチプロセッサシステムのタスクキューを示す図である。

【図3】図2に示すタスクキューに含まれるスレッド情報の構成を示す図である。

【図4】図3に示すスレッド情報の詳細を説明するための図である。

【図5】プロセッサによりスレッドを選出する処理を示すフローチャートである。

【図6】図5に示す処理に伴うスレッド情報の変化を示す図である。

【図7】プロセッサによりスレッドを追加する処理を示すフローチャートである。

【図8】図7に示す処理に伴うスレッド情報の変化を示す図である。

【図9】ユーザレベルスケジューラを用いるマルチプロセッサシステムの例を示す図である。

【図10】図9に示すマルチプロセッサシステムにおけるスレッドのグルーピングを示す図である。

【図11】図9に示すマルチプロセッサシステムにおけるスレッドのスケジューリングを示す図である。

【符号の説明】

【0074】

10 処理ユニット、 12 プロセッサ、 14 ローカルメモリ、 16 メモリ制御部、 20 メインバス、 30 メインメモリ、 40 スレッド情報、 50 スレッドアドレス情報、 60 スレッド実体、 100 マルチプロセッサシステム、 110 処理ユニット、 112 プロセッサ、 114 メモリ制御部、 120 メインバス、 130 メインメモリ。

10

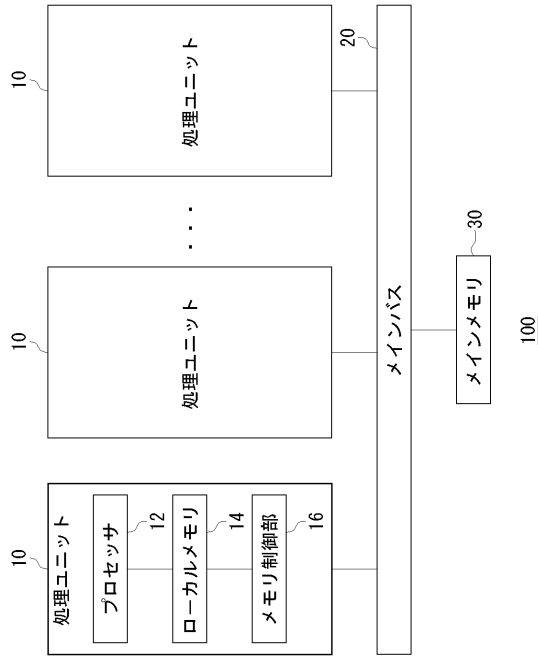
20

30

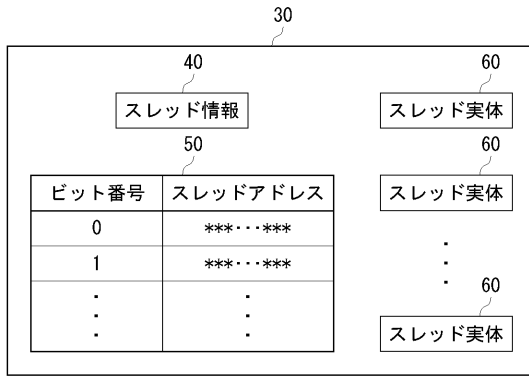
40

50

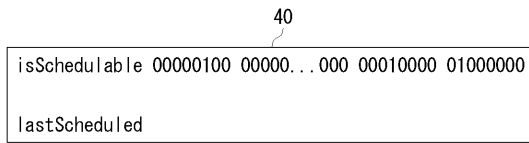
【図1】



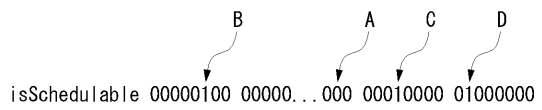
【図2】



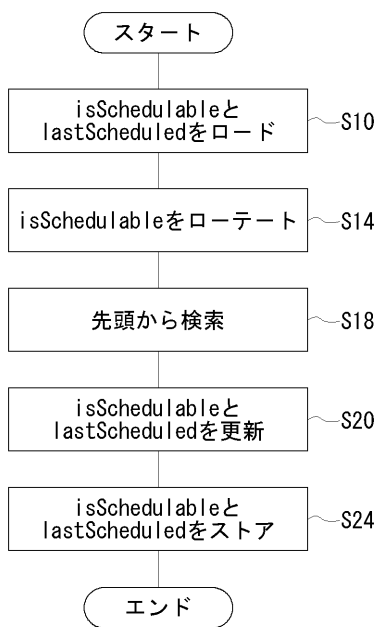
【図3】



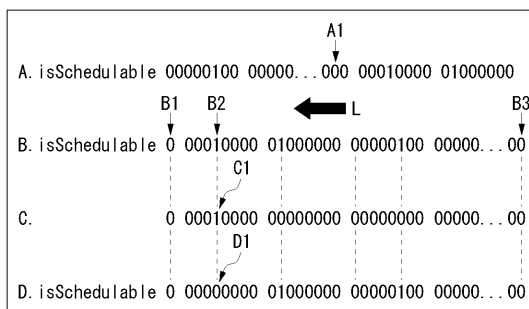
【図4】



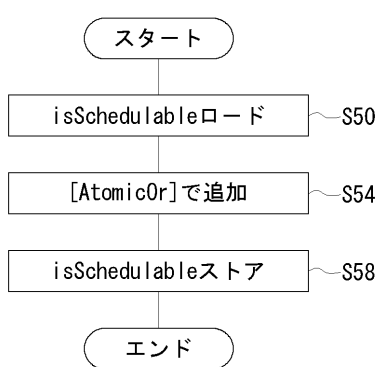
【図5】



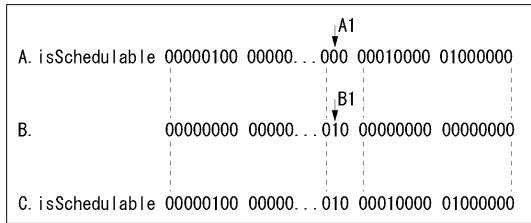
【図6】



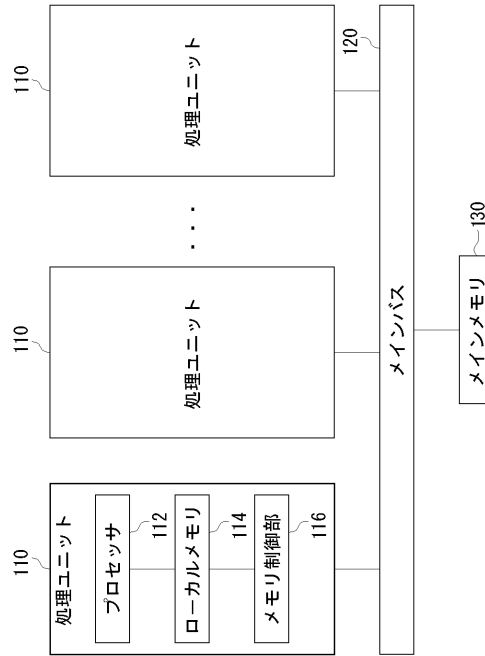
【図7】



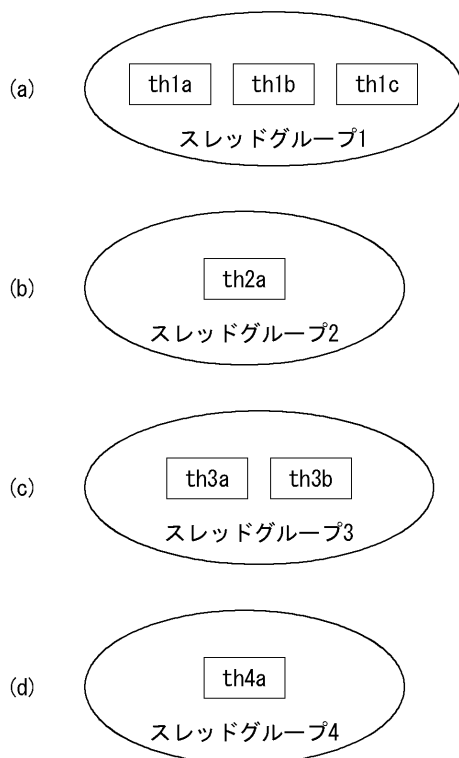
【図8】



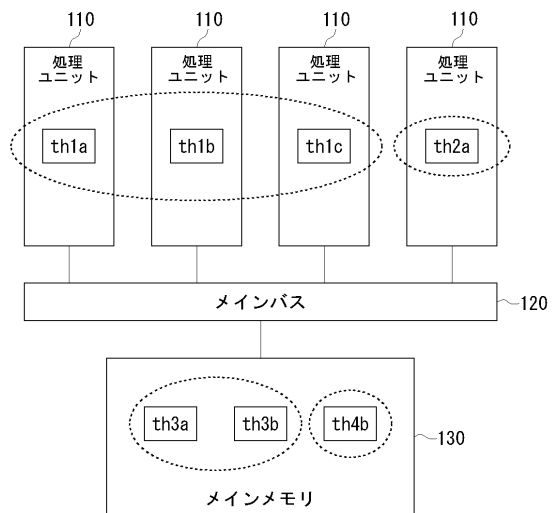
【図9】



【図10】



【図11】



フロントページの続き

- (72)発明者 井上 敬介
東京都港区南青山2丁目6番21号 株式会社ソニー・コンピュータエンタテインメント内
- (72)発明者 村田 誠二
東京都港区南青山2丁目6番21号 株式会社ソニー・コンピュータエンタテインメント内

審査官 鈴木 修治

- (56)参考文献 特開平02-242434(JP,A)
特開平04-024828(JP,A)
特許第2804478(JP,B2)
特許第2505526(JP,B2)

- (58)調査した分野(Int.Cl., DB名)
G06F 9/48