

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3730740号  
(P3730740)

(45) 発行日 平成18年1月5日(2006.1.5)

(24) 登録日 平成17年10月14日(2005.10.14)

(51) Int. Cl.

F I

**G06F 9/50 (2006.01)**

G06F 9/46 4 6 5 Z

**G06F 9/48 (2006.01)**

G06F 9/46 4 5 2 Z

請求項の数 10 (全 42 頁)

(21) 出願番号 特願平9-38903  
 (22) 出願日 平成9年2月24日(1997.2.24)  
 (65) 公開番号 特開平10-240549  
 (43) 公開日 平成10年9月11日(1998.9.11)  
 審査請求日 平成13年1月26日(2001.1.26)

(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 100077816  
 弁理士 春日 譲  
 (72) 発明者 中谷 明弘  
 神奈川県川崎市麻生区王禅寺1099番地  
 株式会社 日立製作所 シス  
 テム開発研究所内  
 (72) 発明者 西門 隆  
 神奈川県川崎市麻生区王禅寺1099番地  
 株式会社 日立製作所 シス  
 テム開発研究所内

最終頁に続く

(54) 【発明の名称】 並列ジョブ多重スケジューリング方法

(57) 【特許請求の範囲】

【請求項1】

同等の機能を有する複数の演算プロセッサを有する計算機において、逐次処理部分と並列処理部分からなる並列ジョブを、複数、多重実行すべく、スケジューリングを行う、並列ジョブ多重スケジューリング方法において、

上記複数の演算プロセッサを、一の並列ジョブの逐次処理部分あるいは並列処理部分を実行する逐次演算プロセッサと、当該並列ジョブの並列処理部分を並列的に実行する複数個の演算プロセッサからなる並列演算プロセッサ群とに、論理的に分割する第1のステップと、

一の上記逐次演算プロセッサが、一の上記並列ジョブの逐次処理部分を実行する第2のステップと、

別の上記逐次演算プロセッサが、別の上記並列ジョブの逐次処理部分を実行する第3のステップと、

上記一の逐次演算プロセッサが、自己を含む上記並列演算プロセッサ群に対し、上記別の逐次演算プロセッサを含む他の逐次演算プロセッサを排して、同期処理を行う第4のステップと、

上記並列ジョブの並列処理部分が実行される第5のステップと、

上記並列演算プロセッサ群が終了通知を発行し、上記一の逐次演算プロセッサが、上記一の並列ジョブの逐次処理部分を実行する第6のステップと、  
 を有することを特徴とする並列ジョブ多重スケジューリング方法。

10

20

**【請求項 2】**

請求項 1 記載の並列ジョブ多重スケジューリング方法において、

上記第 5 のステップの並列処理部分の実行中に、上記別の逐次演算プロセッサが、上記並列処理部分を実行する上記並列演算プロセッサ群に対し、同期処理を要求する第 7 のステップと、

上記逐次演算プロセッサが、複数の並列ジョブを多重に実行するためのスケジューラを起動する第 8 のステップと、

上記同期処理の要求による割り込みを契機に、上記別の逐次演算プロセッサが、上記並列ジョブの逐次処理部分を実行待ち状態とする第 9 のステップと、

上記スケジューラが、あらかじめ定められた規則に基づく次のジョブスケジューリングイベントにて、別の実行待ち状態の並列ジョブを取り出す第 10 のステップと、

上記スケジューラが、上記一の逐次演算プロセッサが上記第 5 のステップの並列処理部分を実行する並列演算プロセッサ群を使用することを不可とする第 11 のステップと、

上記スケジューラが、上記別の逐次演算プロセッサが上記第 10 のステップで取り出した並列ジョブの並列処理部分を実行する並列演算プロセッサ群を使用することを許可する第 12 のステップとを有し、

上記取り出した並列ジョブを実行することを特徴とする並列ジョブ多重スケジューリング方法。

**【請求項 3】**

請求項 2 記載の並列ジョブ多重スケジューリング方法において、

上記ジョブスケジューリングイベントは、上記第 5 のステップの並列処理部分を実行する上記並列演算プロセッサ群が発行する終了通知であることを特徴とする並列ジョブ多重スケジューリング方法。

**【請求項 4】**

請求項 2 記載の並列ジョブ多重スケジューリング方法において、

上記ジョブスケジューリングイベントは、1つの並列処理部分に引き続く逐次処理部分の処理時間が一定時間より長いと、コンパイラが設定できたとき、又はユーザがプログラム中で指定したとき、に発行される当該並列処理部分の終了通知であることを特徴とする並列ジョブ多重スケジューリング方法。

**【請求項 5】**

請求項 2 記載の並列ジョブ多重スケジューリング方法において、

上記ジョブスケジューリングイベントは、同期処理の対象となっている上記逐次演算プロセッサの全てが、シテムコール処理時の待ち状態と設定されていた場合であることを特徴とする並列ジョブ多重スケジューリング方法。

**【請求項 6】**

請求項 2 記載の並列ジョブ多重スケジューリング方法において、

上記ジョブスケジューリングイベントは、ある並列処理部分と引き続く並列処理部分の間の逐次処理部分の実行に要した計算時間を統計的に記録し、これから実行する逐次処理部分が統計的な時間より長い可能性が高いと判断されたとき行う実行中の並列処理部分の終了通知であることを特徴とする並列ジョブ多重スケジューリング方法。

**【請求項 7】**

請求項 2 記載の並列ジョブ多重スケジューリング方法において、

上記並列演算プロセッサ群の各並列演算プロセッサが実行していた上記一の並列ジョブの使用していたレジスタのうち、最小限のレジスタの内容のみを退避し、

さらに、次に上記並列演算プロセッサ群を使用する上記別の並列ジョブに対応して退避されたレジスタの内容を回復し、当該別の並列ジョブの実行を再開することを特徴とする並列ジョブ多重スケジューリング方法。

**【請求項 8】**

請求項 1 記載の並列ジョブ多重スケジューリング方法において、

上記複数の演算プロセッサの、上記逐次演算プロセッサと上記並列演算プロセッサ群へ

10

20

30

40

50

の割り当てのプロセッサ構成の比率を動的に変更することを特徴とする並列ジョブ多重スケジューリング方法。

【請求項 9】

請求項 8 記載の並列ジョブ多重スケジューリング方法において、

並列ジョブは、当該ジョブが使用する上記複数の演算プロセッサの構成に従って、予めユーザによって所属する階級が指定され、上記第 5 のステップにおいて同一の階級の並列ジョブのみが実行されるよう、当該指定に従って、演算プロセッサを切り替えることで上記プロセッサ構成の比率を動的に変更することを特徴とする並列ジョブ多重スケジューリング方法。

【請求項 10】

請求項 1 記載の並列ジョブ多重スケジューリング方法において、

上記第 4 のステップの同期処理を行った結果、予め指定された最大待ち時間以内に、同期処理の対象となった全ての演算プロセッサが同期に到達しなかった場合、第三の並列ジョブのスケジューリングを開始することを特徴とする並列ジョブ多重スケジューリング方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、並列ジョブ多重スケジューリング方法に係り、特に、複数の並列化された情報処理プログラムを並行実行するのに好適な並列ジョブ多重スケジューリング方法に関する。

【0002】

【従来の技術】

従来の並列化された情報処理プログラムを実行する方式としては、例えば、富田眞治著「並列コンピュータ工学」昭晃堂発行（1996），第 92 頁から第 104 頁に記載されているように、複数の情報処理プログラムを並行して処理する場合には、計算機システム全体を時分割によって利用するようにしている。

【0003】

【発明が解決しようとする課題】

しかしながら、従来方式では、並列化された複数のジョブを実行するに当たり、計算機システムに備わる並列演算プロセッサを時分割利用することによって並行処理を行うようにしているため、並列演算プロセッサが 1 並列ジョブに割り当てられた場合、その割り当て時間内に並列ジョブの逐次処理部分が実行されるとき、並列演算プロセッサは遊休状態となり、計算機システム全体の稼働率が低下するという問題がある。

【0004】

なお、デュアルスカラプロセッサ方式と呼ばれる方式においては、2 つの逐次演算機構を用意し、それらの間で単一の並列演算機構を共有するようにしているが、ここでも並列演算機構は時分割利用されるため、逐次演算部分の処理中に並列演算機構が遊休し、計算機システム全体の稼働率が低下するという問題がある。

【0005】

本発明の目的は、複数の並列化されたジョブの同時実行時に、システム内の全プロセッサの稼働率の向上した並列ジョブ多重スケジューリング方法を提供することにある。

【0006】

【課題を解決するための手段】

（1）上記目的を達成するために、本発明は、同等の機能を有する複数の演算プロセッサを有する計算機において、逐次処理部分と並列処理部分からなる並列ジョブを、複数、多重実行すべく、スケジューリングを行う、並列ジョブ多重スケジューリング方法において、上記複数の演算プロセッサを、一の並列ジョブの逐次処理部分あるいは並列処理部分を実行する逐次演算プロセッサと、当該並列ジョブの並列処理部分を並列的に実行する複数の演算プロセッサからなる並列演算プロセッサ群とに、論理的に分割する第 1 のステッ

10

20

30

40

50

ブと、一の上記逐次演算プロセッサが、一の上記並列ジョブの逐次処理部分を実行する第2のステップと、別の上記逐次演算プロセッサが、別の上記並列ジョブの逐次処理部分を実行する第3のステップと、上記一の逐次演算プロセッサが、自己を含む上記並列演算プロセッサ群に対し、上記別の逐次演算プロセッサを含む他の逐次演算プロセッサを排して、同期処理を行う第4のステップと、上記並列ジョブの並列処理部分が実行される第5のステップと、上記並列演算プロセッサ群が終了通知を発行し、上記一の逐次演算プロセッサが、上記一の並列ジョブの逐次処理部分を実行する第6のステップと、を有するものである。

かかる方法によって、複数のジョブを並列的に実行できるため、計算機システム内の全プロセッサの稼働率を向上し得るものとなる。

10

#### 【0007】

(2) 上記(1)において、好ましくは、上記第5のステップの並列処理部分の実行中に、上記別の逐次演算プロセッサが、上記並列処理部分を実行する上記並列演算プロセッサ群に対し、同期処理を要求する第7のステップと、上記逐次演算プロセッサが、複数の並列ジョブを多重に実行するためのスケジューラを起動する第8のステップと、上記同期処理の要求による割り込みを契機に、上記別の逐次演算プロセッサが、上記並列ジョブの逐次処理部分を実行待ち状態とする第9のステップと、上記スケジューラが、あらかじめ定められた規則に基づく次のジョブスケジューリングイベントにて、別の実行待ち状態の並列ジョブを取り出す第10のステップと、上記スケジューラが、上記一の逐次演算プロセッサが上記第5のステップの並列処理部分を実行する並列演算プロセッサ群を使用することを不可とする第11のステップと、上記スケジューラが、上記別の逐次演算プロセッサが上記第10のステップで取り出した並列ジョブの並列処理部分を実行する並列演算プロセッサ群を使用することを許可する第12のステップとを有し、上記取り出した並列ジョブを実行するようにしたものである。

20

かかる方法によって、ジョブ切り替えのオーバーヘッドを低減し得るものとなる。

#### 【0008】

(3) 上記(2)において、好ましくは、上記ジョブスケジューリングイベントは、上記第5のステップの並列処理部分を実行する上記並列演算プロセッサ群が発行する終了通知としたものである。

かかる方法によって、オーバーヘッドの少ない切り替えを行い得るものとなる。

30

#### 【0009】

(4) 上記(2)において、好ましくは、上記ジョブスケジューリングイベントは、1つの並列処理部分に引き続く逐次処理部分の処理時間が一定時間より長いと、コンパイラが設定できたとき、又はユーザがプログラム中で指定したとき、に発行される当該並列処理部分の終了通知としたものである。

かかる方法によって、スケジューリングのオーバーヘッドを少なくして、計算機システム全体の効率を向上し得るものとなる。

#### 【0010】

(5) 上記(2)において、好ましくは、上記ジョブスケジューリングイベントは、同期処理の対象となっている上記逐次演算プロセッサの全てが、シテムコール処理時の待ち状態と設定されていた場合としたものである。

40

かかる方法によって、オーバーヘッドの少ない切り替えを行い得るものとなる。

#### 【0011】

(6) 上記(2)において、好ましくは、上記ジョブスケジューリングイベントは、ある並列処理部分と引き続く並列処理部分の間の逐次処理部分の実行に要した計算時間を統計的に記録し、これから実行する逐次処理部分が統計的な時間より長い可能性が高いと判断されたとき行う実行中の並列処理部分の終了通知としたものである。

かかる方法によって、動的な終了通知の発行による並列演算プロセッサ群の割り当て変更を行い得るものとなる。

#### 【0012】

50

(7) 上記(2)において、好ましくは、上記並列演算プロセッサ群の各並列演算プロセッサが実行していた上記一の並列ジョブの使用していたレジスタのうち、最小限のレジスタの内容のみを退避し、さらに、次に上記並列演算プロセッサ群を使用する上記別の並列ジョブに対応して退避されたレジスタの内容を回復し、当該別の並列ジョブの実行を再開するようにしたものである。

かかる方法によって、退避・回復するレジスタの数を大幅に削減することが可能になり、ジョブ切り替えのオーバーヘッドを低減し得るものとなる。

#### 【0013】

(8) 上記(1)において、好ましくは、上記複数の演算プロセッサの、上記逐次演算プロセッサと上記並列演算プロセッサ群への割り当てのプロセッサ構成の比率を動的に変更

10

するようにしたものである。

#### 【0014】

(9) 上記(8)において、好ましくは、並列ジョブは、当該ジョブが使用する上記複数の演算プロセッサの構成に従って、予めユーザによって所属する階級が指定され、上記第5のステップにおいて同一の階級の並列ジョブのみが実行されるよう、当該指定に従って、演算プロセッサを切り替えることで上記プロセッサ構成の比率を動的に変更するようにしたものである。

かかる方法によって、それぞれの階級に割り付けられた時間内で階級に属するジョブをその階級の実行形式に従って処理し得るものとなる。

20

#### 【0015】

(10) 上記(1)において、好ましくは、上記第4のステップの同期処理を行った結果、予め指定された最大待ち時間以内に、同期処理の対象となった全ての演算プロセッサが同期に到達しなかった場合、第三の並列ジョブのスケジューリングを開始するようにしたものである。

かかる方法によって、その同期処理は失敗したものとして、同期失敗を知らせる割り込みが演算プロセッサ上に発生し、これを契機にスケジューラが起動され、エラー処理を行い得るものとなる。

#### 【0022】

#### 【発明の実施の形態】

30

以下、図1～図31を用いて、本発明の一実施形態による並列ジョブ多重スケジューリング方法及び装置について説明する。

最初に、図1を用いて、本発明の一実施形態による計算機システムの構成について説明する。

図1は、本発明の一実施形態による計算機システムの全体構成図である。

#### 【0023】

演算プロセッサ群1000は、例えば、8台の演算プロセッサ1001, ..., 1008によって構成されている。ここでは、演算プロセッサ1001～1008の台数は、8台としているがこれに限られるものでない。演算プロセッサ1001～1008は、それぞれ、共有記憶装置2000に接続され、データの授受を行う。共有記憶装置2000は、演算プロセッサ間データ授受領域2101, ..., 2108を備えている。演算プロセッサ1001～1008間のデータをやり取りは、演算プロセッサ間データ授受領域2101, ..., 2108を介して行われる。例えば、送り元である演算プロセッサ1001から受け手である演算プロセッサ1002にデータを送るには、送り元の演算プロセッサ1001が受け手の演算プロセッサのデータ授受領域2102にデータを書き込み、このデータを受け手の演算プロセッサ1002が参照することで行われる。ユーザのプログラムコードは、外部記憶装置3000に保存され、実行時には共有記憶装置上に配置されて実行される。

40

#### 【0024】

ここで、本実施形態においては、演算プロセッサ群1000を構成する演算プロセッサ1

50

001～1008は、全て同等の演算機能を有するものである。並列ジョブの逐次処理部分及び並列処理部分を実行する演算プロセッサを、逐次演算プロセッサと称する。また、並列処理部分のみを実行する演算プロセッサを、並列演算プロセッサと称する。そして、演算プロセッサ群1000は、複数の逐次演算プロセッサと、複数の並列演算プロセッサから構成される並列演算プロセッサ群に論理的に分割されている。

#### 【0025】

並列演算プロセッサ群を構成する複数の並列演算プロセッサは、ジョブの並列処理部分を1群として並列的に実行するため、これらの並列演算プロセッサの集合体を、並列演算プロセッサ群と称している。複数の並列演算プロセッサは、それぞれ、別々の並列ジョブに割り当てられる。一方、複数の逐次演算プロセッサは、それぞれ、独立してジョブの逐次

10

#### 【0026】

ここで、本実施形態においては、例えば、演算プロセッサ1001，…，1004が、複数の逐次演算プロセッサを構成し、演算プロセッサ1005，…，1008が、並列演算プロセッサ群を構成するように、論理分割されているものとする。なお、複数の逐次演算プロセッサと並列演算プロセッサ群の論理分割は自由であり、例えば、演算プロセッサ1001，1002が、複数の逐次演算プロセッサを構成し、演算プロセッサ1003，…，1008が、並列演算プロセッサ群を構成するように、論理分割してもよいものである

20

#### 【0027】

演算プロセッサ間同期機構4000は、演算プロセッサ1001～1008に接続され、演算プロセッサ1001～1008間の同期処理を行うものである。演算プロセッサ間同期機構4000は、同期範囲指定機構4100を備えている。同期範囲指定機構4100は、並列演算プロセッサ群を、複数並列ジョブで共用使用するために、各逐次演算プロセッサに対応して上記並列演算プロセッサ群が使用可能かをプログラム制御可能な並列演算機組合せ機構とするものである。同期範囲指定機構4100によって指定された計算機システム内の演算プロセッサ1000間のみで同期処理を行うことが可能である。同期範囲指定機構4100の詳細については、図2を用いて後述する。

30

#### 【0028】

ここで、並列処理部分を実行する機構が、従来の並列専用機構とは異なり、逐次演算プロセッサと同等の複数個の演算プロセッサによって論理的に構成することが重要である。これによって、計算機システム内の演算プロセッサの逐次及び並列演算プロセッサへの割り当て比率を変更することが可能になる。また、必要に応じて演算プロセッサを物理的に増加させることによって、演算性能を向上させることも可能である。

40

#### 【0029】

各演算プロセッサ1001～1008は、演算プロセッサ間同期機構4000に対して、同期到達通知信号線SRLと同期完了通知信号線SCLによって接続される。同期処理に際しては、同期処理に達した演算プロセッサ1000は、同期到達通知命令を実行することによって、同期到達通知信号線SRLに信号を出力し、同期完了待ち命令を実行して、同期完了通知信号線SCLから信号が来るのを待つ。演算プロセッサ間同期機構4000は、同期範囲指定機構4100によって同期の範囲内と指定された全ての演算プロセッサが同期到達信号線SRLに対して信号を出した時点で同期処理が完了したとして、同期範囲内の演算プロセッサの同期完了信号線SCLに対して信号を出して同期処理の完了を通知する。同期完了信号を受けた演算プロセッサは、同期完了待ち状態が解除され、同期処

50

理に引き続く処理の実行に移る。

【0030】

演算プロセッサ間割り込み制御機構5000は、演算プロセッサ間同期機構4000に接続されており、同期範囲指定機構4100によって同期範囲外と設定された演算プロセッサからの同期要求に対しては、この演算プロセッサに割り込みを発生させて、共有記憶装置2000内のスケジューラを起動する。また、同期範囲指定機構4100によって同期範囲内と設定された演算プロセッサが同期要求を行い、その同期の成立前に同期範囲外に設定され直した場合にも、この演算プロセッサに割り込みを発生させてスケジューラを起動する。

【0031】

次に、図2を用いて、本発明の一実施形態による計算機システムに用いる同期範囲指定機構4100の構成について説明する。

図2は、本発明の一実施形態による同期範囲指定機構の構成図である。

【0032】

図1において説明したように、計算機システムが8台の演算プロセッサによって構成される場合においては、同期範囲指定機構4100は、8ビットの同期範囲指定レジスタ4110を備えている。

【0033】

一般的に、n台の演算プロセッサから成る計算機システムでは、同期処理に参加する演算プロセッサは、同期範囲指定機構が備えるnビットレジスタ（同期範囲指定レジスタ）によって指定される。このレジスタの第kビットが演算プロセッサkに対応するものとし、このビットに1が設定されていれば、対応する演算プロセッサは同期処理に参加する。同期範囲内と設定された演算プロセッサからの同期要求に対してはそのまま同期処理を続行する。

【0034】

一方、同期範囲外と設定された演算プロセッサからの同期要求に対してはこの演算プロセッサに割り込みを発生させてスケジューラを起動する。また、同期範囲内と設定された演算プロセッサが同期要求を行い、その同期の成立前に同期範囲外に設定され直した場合にも、同演算プロセッサに割り込みを発生させてスケジューラを起動する。

【0035】

ここで、同期範囲外の演算プロセッサからの同期要求は無視されるのではなく、同期要求に対して、要求元の演算プロセッサ上に割り込みを発生させること、及び、同期成立前に同期範囲外に設定し直された演算プロセッサ上に割り込みを発生させることが重要である。

【0036】

図2に示す例では、8ビットからなる同期範囲指定レジスタ4110のビット1, 5, 6, 7及び8に"1"を設定し、8台の演算プロセッサの内、図1に示した演算プロセッサ1001, 1005, 1006, 1007, 1008が同期処理に参加するように指示している。ここで、演算プロセッサ1001は、逐次演算プロセッサであり、演算プロセッサ1005, ..., 1008は、並列演算プロセッサ群を構成する並列演算プロセッサである。演算プロセッサ1005, ..., 1008が、ジョブの並列処理部分を実行する際には、演算プロセッサ1001も並列処理部分を実行する。それ以外のときには、演算プロセッサ1001は、逐次処理部分を実行する。

【0037】

次に、図3を用いて、共有記憶装置2000中に用意されるデータについて説明する。

図3は、本発明の一実施形態による共有記憶装置中に用意されるデータの説明図である。

【0038】

共有記憶装置2000は、ジョブ実行コード2100が各ジョブに対応して配置され、各演算プロセッサ1001~1008は自らに割り当てられたジョブ実行コード2100を実行する。また、OS実行コード2200は、ユーザジョブのコードとは別に配置される

10

20

30

40

50

。これらのジョブ及びハードウェア資源管理のための情報を保持する管理情報 2 3 0 0 が、共有記憶装置 2 0 0 0 上に確保されている。管理情報 2 3 0 0 の詳細については、図 4 を用いて後述する。

【 0 0 3 9 】

次に、図 4 を用いて、共有記憶装置 2 0 0 0 上に確保されている管理情報 2 3 0 0 の内容について説明する。

図 4 は、本発明の一実施形態における共有記憶装置上の管理情報の構成図である。

【 0 0 4 0 】

ジョブ管理情報先頭アドレス設定フィールド 2 3 0 1 は、ジョブ管理情報が存在するメモリアドレスを保持する領域である。ユーザ実行ジョブは起動されると、ジョブ毎に、図 5 に示すジョブ管理情報 2 4 0 0 が共有記憶装置 2 0 0 0 上に確保され、ユーザ実行ジョブは、ジョブ管理情報 2 4 0 0 に基づいて管理される。

10

【 0 0 4 1 】

階級設定フィールド 2 3 0 2 は、現在実行中の階級を設定するフィールドである。逐次並列間データ授受領域対応付けフィールド 2 3 0 3 については、図 1 3 を用いて後述する。構成変更情報設定フィールド 2 3 0 4 は、ここに設定された情報に基づいて階級切り替え及びシステム内の演算プロセッサ構成を変更するのに用いられ、詳細については図 2 2 を用いて後述する。同期範囲指定機構設定フィールド 1 3 0 0 5 は、このフィールドの設定によって同期範囲指定レジスタ 4 1 1 0 ( 図 2 ) の内容を指定する。

【 0 0 4 2 】

同期タイムアウトレジスタ設定フィールド 2 3 0 6 は、同期処理が失敗したと判断するまでの時間を設定する特別なレジスタ ( 同期タイムアウトレジスタ ) の設定に用いられるものである。同期範囲指定フィールド 2 3 0 7 及び同期到達フラグ 2 3 0 9 は、後述するように、演算プロセッサ間の同期処理をソフトウェアで実現するときのみ用いるものである。

20

【 0 0 4 3 】

単一実行フラグ 2 3 0 8 については、図 3 0 を用いて後述する。待ち状態フラグ 2 3 1 0 は、演算プロセッサの数だけ用意されているものであり、詳細については後述する。

【 0 0 4 4 】

次に、図 5 を用いて、共有記憶装置 2 0 0 0 上に確保されているユーザ実行ジョブを管理するジョブ管理情報 2 4 0 0 の構成について説明する。

30

図 5 は、本発明の一実施形態において用いるジョブ管理情報の構成図である。

【 0 0 4 5 】

並列ジョブ識別番号フィールド 2 4 0 1 は、並列ジョブ識別番号を保持するフィールドであり、並列ジョブ識別番号によって、ジョブは一意に指定することができる。所属階級識別フィールド 2 4 0 2 は、ジョブが所属するジョブ階級を保持するフィールドであり、詳細については後述する。逐次実行時間フィールド 2 4 0 3 は、ジョブの起動からの逐次実行時間を記録するフィールドである。

【 0 0 4 6 】

並列実行時間フィールド 2 4 0 4 は、ジョブの起動からの並列実行時間を記録するフィールドである。先行逐次情報フィールド 2 4 0 5 は、先行逐次情報を記録するフィールドであり、詳細については図 2 9 を用いて後述する。ジョブ状態保存フィールド 2 4 0 6 は、ジョブ状態を保存するフィールドであり、ジョブ自らの状態を記録する。

40

【 0 0 4 7 】

レジスタ未退避フラグ 2 4 0 7 は、レジスタが未退避であることを示すフラグである。レジスタ退避範囲 2 4 0 8 は、レジスタの退避範囲を示すフラグであり、その詳細については図 1 2 を用いて後述する。レジスタ内容退避エリア 2 4 0 9 は、ジョブ切り替え時にレジスタの内容を退避するエリアである。

【 0 0 4 8 】

次に、図 6 を用いて、共有記憶装置 2 0 0 0 上に確保されている階級管理情報 2 5 0 0 に

50



について説明する。

図 6 は、本発明の一実施形態において用いる階級管理情報の構成図である。

【 0 0 4 9 】

後述するように、ジョブは階級分けされ、これらの階級を切り替えて実行される。各階級は、図 6 に示すようなフィールドから成る階級管理情報 2 5 0 0 を共有記憶装置 2 0 0 0 上にもち、階級管理情報 2 5 0 0 に従って管理される。

【 0 0 5 0 】

階級毎に、階級管理情報 2 5 0 0 は、階級識別番号フィールド 2 5 0 1 と、キュー先頭アドレス設定フィールド 2 5 0 2 と、同期範囲指定レジスタ退避フィールド 2 5 0 3 と、所属ジョブ識別番号フィールド 2 5 0 4 とから構成され、共有記憶装置 2 0 0 0 上に保持される。

10

【 0 0 5 1 】

図 1 に示す例では、演算プロセッサ 1 0 0 0 を共有記憶装置 2 0 0 0 によって接続しているが、演算プロセッサ間のデータ授受方法は演算プロセッサ間に備えられた通信装置によるメッセージ送達によるものでもよいものである。

【 0 0 5 2 】

本実施形態において重要なことは、逐次演算プロセッサと同等の複数の演算プロセッサによって並列演算機構を論理的に構成することと、同期範囲指定機構を備え、同期範囲外の演算プロセッサからの同期要求に対しては割り込みを発生することである。

【 0 0 5 3 】

20

次に、図 7 を用いて、本実施形態による並列ジョブ多重スケジューリング方法が対象とする並列ジョブの形式について説明する。

図 7 は、本発明の一実施形態によって並列ジョブを複数の演算プロセッサを用いて実行する場合の 1 方式の概略図である。

【 0 0 5 4 】

図 7 において、矩形部 U 1 ~ U 1 3 , P 1 ~ P 1 3 の 1 つ 1 つが情報処理プログラムが実行する計算の 1 単位を表すものとする。一般に、情報処理プログラムは計算機システム上で実行する処理の列を記述したものであり、単一の演算プロセッサから成る計算機システムで情報処理プログラムを実行する場合は、その記述された処理の列を逐次的に実行する。例えば、図 7 ( A ) に示すように、1 情報処理プログラムが処理 U 1 ~ U 1 3 から成る列を記述していたとすると、単一演算プロセッサから成る計算機システムでは、処理 U 1 から順に処理 U 1 3 までが順番に行われる。ここでは、時間は図中の左から右に流れているとして表している。

30

【 0 0 5 5 】

一方、処理 U 1 ~ U 1 3 の内、例えば、処理 U 2 ~ U 6 の 5 単位の処理の順序が任意であるとすれば、これらは複数の演算プロセッサ上で並列に実行することができる並列処理可能部分である。同様に処理 U 8 ~ U 1 2 の 5 単位の処理の順序が任意であるとすれば、これらは複数の演算プロセッサ上で並列に実行することができる並列処理可能部分である。残りの処理 U 1 , U 7 , U 1 3 は、逐次処理部分である。

【 0 0 5 6 】

40

そこで、本実施形態においては、処理 U 1 ~ U 1 3 の全体処理は、処理 P 1 ~ P 1 3 のように 5 台の演算プロセッサ上で実行することができる。ここで、並列処理可能部分である処理 U 2 ~ U 6 は、並列処理部分として処理 P 2 ~ P 6 のように時間的に同時に処理することが可能である。並列処理可能部分である処理 U 8 ~ U 1 2 も同様に、並列処理部分として処理 P 8 ~ P 1 2 のように並列に実行される。逐次処理部分 U 1 , U 7 , U 1 3 は、逐次処理部分として処理 P 1 , P 7 , P 1 3 のように実行される。ここでは、時間は図中で、左から右に流れているとして表している。また、上下方向に重ねられた矩形は時間的に同時に別々の演算プロセッサ上で実行されることを表している。

【 0 0 5 7 】

図 7 ( B ) において、処理 P 3 ~ P 6 及び処理 P 9 ~ P 1 2 は、それぞれ、例えば、図 1

50

に示した並列演算プロセッサ群を構成する並列演算プロセッサ1005～1008によって実行される。また、逐次処理部分である処理P1、P7、P13及び並列処理部分の一部である処理P2、P8は、例えば、図1に示した複数の逐次演算プロセッサの1つである逐次演算プロセッサ1001によって実行される。

なお、以降の説明でも、同様の記法を用いて説明するものとする。

#### 【0058】

次に、図8を用いて、本発明の一実施形態による並列ジョブの多重実行の概略について説明する。

図8は、本発明の一実施形態による並列ジョブの多重実行の概略説明図である。

#### 【0059】

図8においては、複数の演算プロセッサのジョブに対する割り当ての方法及び演算プロセッサ構成の動的な変更の概略を説明している。

#### 【0060】

上述したように、演算プロセッサ1001～1008は、逐次演算プロセッサと並列演算プロセッサ群とに論理的に分割されている。ここでは、演算プロセッサ1001～1004が、逐次演算プロセッサを構成し、演算プロセッサ1005～1008が、並列演算プロセッサ群を構成しているものとする。

#### 【0061】

第1のジョブJ10と第2のジョブJ20は、並列的に実行されている並列ジョブである。第1のジョブJ10は、逐次処理部分である処理P101と、並列処理部分である処理P102～P106と、逐次処理部分である処理P107から構成されている。また、第2のジョブJ20は、逐次処理部分である処理P201と、並列処理部分である処理P202～P206と、逐次処理部分である処理P207から構成されている。

#### 【0062】

並列ジョブJ10、J20は、時刻T1において、逐次演算プロセッサ1003、1004上でそれぞれ起動される。第1のジョブJ10の並列処理部分の処理P102～P106の方が、第2のジョブJ20の並列処理部分の処理P202～P206よりも早い時刻T2において起動されるとすると、並列処理部分の処理の起動時には、第1のジョブJ10の逐次処理部分の処理P101を実行した逐次演算プロセッサ1004が、並列演算プロセッサ群1005～1008を獲得して実行を進める。並列処理部分の処理の起動時には、逐次演算プロセッサ1004と並列演算プロセッサ群1005～1008との間で、同期処理を行う。この同期処理は、ある所定の時期においては、唯一の逐次演算プロセッサからのみ可能なようにすることで並列演算プロセッサ群の割り当て制御を行う。即ち、時刻T2に逐次演算プロセッサ1004が同期処理を開始すると、この時刻T2以降においては、逐次演算プロセッサ1003は同期処理を行うことができない。逐次演算プロセッサ1003の同期処理は、逐次演算プロセッサ1004と並列演算プロセッサ1005～1008の並列処理が終了するまで行えないようになっている。なお、この詳細な処理については、図12を用いて後述する。

#### 【0063】

逐次演算プロセッサ1004による同期処理が完了すると、並列処理部分の処理P102～P106が、逐次演算プロセッサ1004と並列演算プロセッサ群1005～1008によって並列的に実行される。

#### 【0064】

並列処理部分の処理P102～P106の終結時には、並列演算プロセッサ群1005～1008は、終了通知を発行し、これを契機にして他の並列実行待ちジョブの並列処理部分の実行に切り替わる。この際、並列処理部分の実行が終了しているため、並列演算プロセッサのレジスタの退避が大幅に削減されることによってオーバーヘッドの少ない切り替えを実現することができる。一般に、並列演算機構は、多くのレジスタを使って高速実行を実現しているため、従来の時分割による実行並列ジョブの切り替え方式では、並列機構使用中にジョブを切り替えると多量のレジスタの退避・回復が必要であり、性能劣化の大

10

20

30

40

50

きな原因となっているのに対して、本実施形態では、オーバーヘッドを低減できるものである。

【0065】

時刻T3以降において、並列プロセッサ1005～1008によるジョブJ10の並列処理部分の実行がそれぞれ終了し、終了通知がなされ、そして、時刻T4において、全てのジョブJ10の並列処理部分の実行の終了が確認されたものとする。時刻T4において、ジョブJ20の逐次処理部分の処理J201が実行終了しており、並列処理部分P202～P206の実行待ちになっていたすると、逐次演算プロセッサ1003は、並列演算プロセッサ1004～1008に対して同期処理を行った後、並列処理部分の処理P202～P206が、逐次演算プロセッサ1003と並列演算プロセッサ群1005～1008によって並列的に実行される。なお、ジョブJ10の並列処理部分の実行が終了すると、逐次演算プロセッサ1004は、ジョブJ10の残りの逐次処理部分の処理P107を実行する。

10

【0066】

時刻T5における並列処理部分の処理P202～P206の終結時には、並列演算プロセッサ群1005～1008は、終了通知を発行し、逐次演算プロセッサ1003は、ジョブJ20の残りの逐次処理部分の処理P207を実行する。なお、このとき、他の並列実行待ちジョブがあれば、そのジョブの並列処理部分の実行に切り替わる。

【0067】

以上のようにして、本実施形態においては、複数のジョブを並列的に実行できるため、計算機システム内の全プロセッサの稼働率を高めることができる。即ち、図8に示す例では、2つのジョブJ10、J20が並列的に実行されている。

20

【0068】

ここで、従来の方式においては、例えば、時刻T1からT6の間においては、ジョブJ10だけしか実行できないものである。その結果、並列プロセッサ1005～1008が稼働できる時間は、時刻T2からT3の間に限られ、残りの時間(T1～T2、T3～T6)は、並列プロセッサ1005～1008は稼働していない遊休状態となっているものである。

【0069】

それに対して、本実施形態の方式においては、並列プロセッサ1005～1008が稼働できる時間は、時刻T2からT3と、時刻T4からT5の間となり、並列プロセッサ1005～1008は稼働していない遊休状態となっている時間(T1～T2、T3～T4、T5～T6)を短くすることができるため、並列演算プロセッサ1005～1008の稼働率を向上することができるものである。

30

【0070】

さらに、本実施形態においては、並列演算プロセッサ群は、逐次演算プロセッサと同等の演算プロセッサから構成されているため、逐次演算プロセッサと並列演算プロセッサ群の構成比率を動的に変更することが可能である。

【0071】

図8において、時刻T6以降は、演算プロセッサの構成比率を変更した状態を示している。即ち、時刻T6以降においては、例えば、演算プロセッサ1001が、逐次演算プロセッサを構成し、演算プロセッサ1002～1008が、並列演算プロセッサ群を構成している。従って、時刻T7において、新しいジョブJ30が起動されると、逐次処理部分の処理P301の実行が、逐次プロセッサ1001上で起動される。時刻T8において、逐次処理部分の処理P301が終了すると、次の並列処理部分の処理の起動時に、逐次演算プロセッサ1001と並列演算プロセッサ群1002～1008との間で、同期処理を行う。逐次演算プロセッサ1001による同期処理が完了すると、並列処理部分の処理P302～P309が、逐次演算プロセッサ1001と並列演算プロセッサ群1002～1008によって並列的に実行される。さらに、時刻T9において、並列処理が終了すると、逐次演算プロセッサ1001は、逐次処理部分の処理P310を実行する。

40

50

## 【 0 0 7 2 】

以上説明したように、並列演算量が多い場合には、全体の演算プロセッサの内の1台を逐次演算プロセッサとして構成し、残りの7台の演算プロセッサを並列演算プロセッサ群として構成することにより、並列処理部分の処理時には、8台の演算プロセッサをすべて使用して実行することが可能となる。例えば、日中のように、並行ジョブを多数実行したい場合には、図8の左側に示すように、演算プロセッサ1001～1004によって逐次演算プロセッサを構成し、演算プロセッサ1005～1008によって並列演算プロセッサ群を構成するようにして、複数の逐次演算プロセッサとして構成する演算プロセッサの台数と、並列演算プロセッサ群を構成する演算プロセッサの台数とを等しくする。また、夜間のように、処理するジョブが少なく、しかも、並列演算量が多い場合には、全体の演算プロセッサの内の1台を逐次演算プロセッサとして構成し、残りの7台の演算プロセッサを並列演算プロセッサ群として構成することができる。

10

## 【 0 0 7 3 】

従って、図8に示す例では、例えば、時刻T6までは、日中のシステム運用の形態を示し、時刻T7以降が夜間のシステム運用とする。時刻T6～T7における演算プロセッサの構成比率の変更は、システムオペレータによって行われる。

## 【 0 0 7 4 】

このようにして、実行対象ジョブの逐次処理と並列処理の比率に応じて計算機システム内の演算プロセッサの構成比率を変更することができる。

## 【 0 0 7 5 】

従来の並列ジョブ実行を目的とした計算機システムにおいては、逐次処理を行う逐次専用演算プロセッサと、並列処理を行う並列専用演算プロセッサを固定的な形態で備えている。並列ジョブ毎に逐次処理部分と並列処理部分の比率が異なっているにも拘わらず、固定的な構成の計算機システムのみが提供され続けている。しかも、従来方式の計算機システムでは、並列演算の性能向上のために特別な並列演算プロセッサを採用するため、並列演算プロセッサは特別な演算しか行えない上、並列演算プロセッサ全体は1体でしか扱えず、分割使用は不能であった。すなわち、逐次演算量、並列演算量に応じて、逐次演算プロセッサと並列演算プロセッサの構成を替えることができないものである。

20

## 【 0 0 7 6 】

それに対して、本実施形態においては、計算機システムを構成する演算プロセッサを全て同等の機能の演算プロセッサにより構成するようにしているので、実行並列ジョブの性質、即ち、逐次処理部分と並列処理部分の比率に応じて、柔軟に逐次演算処理能力と並列演算処理能力の比率を動的に変化させて、計算機システムの稼働率を高めることが可能となる。

30

## 【 0 0 7 7 】

次に、図9を用いて、図7(B)に示したように並列化されたジョブの実行するための計算機システムの構成について説明する。

図9は、本発明の一実施形態による並列ジョブ多重スケジューリング装置のシステム構成図である。

## 【 0 0 7 8 】

本実施形態においては、計算機システムが8台の演算プロセッサによって構成されているものとして説明する。また、演算プロセッサ1001～1004を逐次演算プロセッサとして用い、演算プロセッサ1005～1008を並列演算プロセッサ群として用いるものとする。

40

## 【 0 0 7 9 】

1並列ジョブは、1つの逐次演算プロセッサ上で起動される。ここでは、ある並列ジョブが、逐次演算プロセッサ1001上で起動されたとする。ジョブが起動されると、逐次処理部分の処理(図7;処理P1)の実行を逐次演算プロセッサ1001上で開始する。この後、並列処理部分の処理(図7;処理P2～P6)に達すると、逐次演算プロセッサ1001と並列演算プロセッサ群1005～1008を用いて、これら5単位の処理(図7

50

；処理 P 2 ～ P 6 ）を 5 台の演算プロセッサ 1 0 0 1 ， 1 0 0 5 ～ 1 0 0 8 上で並列に行う。並列演算プロセッサ群を構成する演算プロセッサ 1 0 0 5 ～ 1 0 0 8 は一括して 1 並列ジョブに割り当てられる。

#### 【 0 0 8 0 】

1 並列ジョブの並列処理部分の実行の開始時と終了時にはそのジョブの逐次処理部分を実行する 1 台の逐次演算プロセッサ 1 0 0 1 と並列演算プロセッサ群を構成する演算プロセッサ 1 0 0 5 ～ 1 0 0 8 の間で同期処理を行う。この同期処理に参加する演算プロセッサ 1 0 0 1 ， 1 0 0 5 ～ 1 0 0 8 は、同期範囲指定機構 4 1 0 0 （図 1 ）によって指示される。並列処理部分開始時に行われる同期処理を契機として、各演算プロセッサ上では、並列処理部分の内、自分に割り当てられた処理を一斉に開始し、また、並列処理部分の終了時にも同期処理を行い、全ての演算プロセッサの処理が終了したことを確認して次の処理に進む。

10

#### 【 0 0 8 1 】

各演算プロセッサ 1 0 0 1 ～ 1 0 0 8 は、演算プロセッサ間データ授受領域 2 1 0 1 ～ 2 1 0 8 を共有記憶装置 2 0 0 0 中に備え、並列処理部分の開始時には逐次演算プロセッサ 1 0 0 1 上で実行されているプログラムが、各並列演算プロセッサに割り当てられた処理の先頭アドレスを演算プロセッサ間データ授受領域 2 1 0 1 ～ 2 1 0 8 に書き込み、そのデータを各並列演算プロセッサ 1 0 0 5 ～ 1 0 0 8 が読み出すことによって、1 並列処理部分の複数の演算プロセッサへの割り当てを実現する。

#### 【 0 0 8 2 】

20

各演算プロセッサ 1 0 0 1 ～ 1 0 0 8 が実行するジョブ実行コードは、ユーザコード領域であるジョブ実行コード 2 1 0 0 に保持される。図 7 に示す例では、処理 P 1 ， P 2 ， P 7 ， P 8 ， P 1 3 が逐次演算プロセッサ 1 0 0 1 上で実行される。処理 P 3 ， P 9 が演算プロセッサ 1 0 0 5 上で、処理 P 4 ， P 1 0 が演算プロセッサ 1 0 0 6 上で、処理 P 5 ， P 1 1 が演算プロセッサ 1 0 0 7 上で、処理 P 6 ， P 1 2 が演算プロセッサ 1 0 0 8 上でそれぞれ実行される。

#### 【 0 0 8 3 】

次に、図 1 0 を用いて、本実施形態における 1 並列処理部分の起動・終結時の実行の詳細について説明する。

図 1 0 は、本発明の一実施形態による並列ジョブ多重スケジューリング方法における並列処理部分の起動・終結時の実行について説明するフローチャートである。

30

#### 【 0 0 8 4 】

並列処理部分の起動と終結は、演算プロセッサ間に備えられたバリア線による同期処理を用いて行われる。図 1 において説明したように、演算プロセッサ 1 0 0 0 は、自らが同期に達したことを他の演算プロセッサに知らせるため同期到達通知信号線 S R L を介して伝達される同期到達通知命令及び、他の演算プロセッサが同期に達するまで待つための同期完了待ち命令を備えている。これらの命令を用いて演算プロセッサ間の同期処理が行なわれる。

#### 【 0 0 8 5 】

また、並列演算プロセッサ群 1 0 0 5 ～ 1 0 0 8 上で並列処理部分の実行を開始するに当たって必要となる初期データは、逐次演算機構と並列演算機構の間で互いに参照可能な演算プロセッサ間データ授受領域 2 1 0 5 ～ 2 1 0 8 を起動情報格納領域として用いて授受することにより実現する。

40

#### 【 0 0 8 6 】

図 1 0 は、1 並列処理部分の逐次演算プロセッサからの起動、並列演算プロセッサ群上での処理、並列処理部分の終結の概略を示している。図 1 0 の左側は、逐次演算プロセッサ 1 0 0 1 の処理を示しており、図 1 0 の右側は、並列演算プロセッサ群を構成する個々の演算プロセッサ 1 0 0 5 ～ 1 0 0 8 の処理を示している。

ステップ S 1 0 0 1 において、並列処理部分の起動を待っている並列演算プロセッサ 1 0 0 5 ～ 1 0 0 8 は、同期到達通知命令を実行し、さらに、ステップ S 1 0 0 2 において、

50

同期完了待ち命令を実行して待ち状態になっている。この時点で、この同期処理に対して、同期に達していないのは逐次演算プロセッサ１００１のみである。

【００８７】

一方、ステップＳ１０１１において、逐次演算プロセッサ１００１は、並列処理部分の処理を並列演算プロセッサ群１００５～１００８上で起動するに当たって、並列演算プロセッサ群上で実行されるコードの先頭番地を、演算プロセッサ間データ授受領域２１０５～２１０８に、起動情報として設定する。

【００８８】

ステップＳ１０１２において、起動情報の設定が完了すると、逐次演算プロセッサ１００１は、同期到達通知命令を実行する。さらに、ステップＳ１０１３において、自らが並列処理部分の起動のための準備が整ったことを並列演算プロセッサ群側に知らせ、同期完了待ち命令を実行する。その後、同期処理の完了を待つ。

10

【００８９】

ステップ１０２１においては、逐次演算プロセッサ１００１上での同期到達通知命令の実行によって、逐次演算プロセッサ１００１と並列演算プロセッサ群１００５～１００８との間の同期処理が完了する。

【００９０】

ステップＳ１００３において、並列演算プロセッサ群１００５～１００８は、起動のための同期処理の終了後、それぞれ演算プロセッサ間データ授受領域２１０５～２１０８から逐次演算プロセッサ１００１が設定した並列処理部分の起動のための情報であるステップＳ１０１１において設定された実行コード先頭番地）を取り出し、ステップＳ１００４において、それぞれの並列処理部分の処理の実行が開始される。

20

【００９１】

一方、ステップＳ１０１４において、逐次演算プロセッサ１００１も、自らに割り当てられた並列処理部分の処理の実行を開始する。

【００９２】

ステップＳ１０１５において、逐次演算プロセッサ１００１は、自らに割り当てられた並列処理部分の処理の実行が終了すると、並列演算プロセッサ群上の並列処理部分の終結を待つために、同期到達通知命令を実行し、さらに、ステップＳ１０１６において、同期完了待ち命令を実行する。

30

【００９３】

他方では、ステップＳ１００５において、並列演算プロセッサ群を構成する演算プロセッサ１００５～１００８は、それぞれ、自らに割り当てられた並列処理部分の処理の実行が終了すると、並列演算プロセッサ群上の他の並列処理部分の終結を待つために、同期到達通知命令を実行し、さらに、ステップＳ１００６において、同期完了待ち命令を実行する。

【００９４】

ステップＳ１０２２において、演算プロセッサ１００１は、全ての演算プロセッサの処理の実行が終了し、同期に達した時点で、１並列処理部分の処理が終了する。並列処理部分の終結に際しては、並列処理部分の起動時に用いた演算プロセッサ間データ授受領域２１０５～２１０８に計算結果等のデータを書き込み、逐次演算プロセッサ１００１がこれらのデータを読み込むことにより、並列演算プロセッサ群１００５～１００８から逐次演算プロセッサ１００１へデータを渡すことができる。

40

【００９５】

次に、図１１を用いて、本発明の一実施形態による図８の左側に示したような複数の並列ジョブの多重実行方式について説明する。

図１１は、本発明の一実施形態による並列ジョブの多重実行の概略説明図である。

【００９６】

図１１においては、複数の演算プロセッサのジョブに対する割り当ての方法の概略を説明している。

50

本実施形態における複数の並列ジョブの多重実行方式においては、複数の並列ジョブは、複数の逐次演算プロセッサ上で互いに独立に起動される。それぞれの並列ジョブは、並列処理部分に達すると、並列演算プロセッサ群を用いてその並列処理部分の実行を行う。本実施形態においては、1プログラムの逐次処理部分に専用の逐次演算プロセッサを割り付けるため、1プログラムに一括して割り付けられる並列演算プロセッサ群の使用要求が衝突しない限り、各プログラムは停止することなく、互いに並行に実行することができる。

【0097】

図11に示す例は、図9に示した計算機システム上で複数の情報処理プログラムを多重に実行した場合について示している。図9に示した計算機システムにおいて、演算プロセッサ1001～1004が逐次演算プロセッサとして、演算プロセッサ1005～1008が並列演算プロセッサ群として用いられている。

10

【0098】

第1のジョブJ10と第2のジョブJ20は、並列的に実行されている並列ジョブである。第1のジョブJ10は、逐次処理部分である処理P101と、並列処理部分である処理P102～P106と、逐次処理部分である処理P107から構成されている。また、第2のジョブJ20は、逐次処理部分である処理P201と、並列処理部分である処理P202～P206と、逐次処理部分である処理P207から構成されている。

【0099】

並列ジョブJ10、J20は、時刻T1において、逐次演算プロセッサ1001、1002上でそれぞれ起動される。このように、2つのジョブはそれぞれの逐次処理部分を別々の演算プロセッサで並行に処理する。第1のジョブJ10の並列処理部分の処理P102～P106の方が、第2のジョブJ20の並列処理部分の処理P202～P206よりも早い時刻T2において起動されるとすると、並列処理部分の処理の起動時には、第1のジョブJ10の逐次処理部分の処理P101を実行した逐次演算プロセッサ1002が、並列演算プロセッサ群1005～1008を獲得して実行を進める。並列処理部分の処理の起動時には、逐次演算プロセッサ1002と並列演算プロセッサ群1005～1008との間で、同期処理SYN1を行う。

20

【0100】

逐次演算プロセッサ1002による同期処理が完了すると、並列処理部分の処理P102～P106が、逐次演算プロセッサ1002と並列演算プロセッサ群1005～1008によって並列的に実行される。

30

【0101】

並列処理部分の処理P102～P106の終結時には、並列演算プロセッサ群1005～1008は、終了通知END1を発行し、これを契機にして他の並列実行待ちジョブの並列処理部分の実行に切り替わる。この際、並列処理部分の実行が終了しているため、並列演算プロセッサのレジスタの退避が大幅に削減されることによってオーバーヘッドの少ない切り替えを実現することができる。

【0102】

時刻T3以降において、並列プロセッサ1005～1008によるジョブJ10の並列処理部分の実行がそれぞれ終了し、終了通知END1がなされる。

40

【0103】

一方、ジョブJ20は、並列演算プロセッサ群1005～1008を使用しなければ、時刻T1～T4に示すように、ジョブJ10の処理と並行して処理P210を実行することができる。

【0104】

時刻T4において、ジョブJ20の逐次処理部分の処理J201が実行終了しており、並列処理部分P202～P206の実行待ちになっていたすると、逐次演算プロセッサ1001は、並列演算プロセッサ1005～1008に対して同期処理SYN2を行った後、並列処理部分の処理P202～P206が、逐次演算プロセッサ1001と並列演算プロセッサ群1005～1008によって並列的に実行される。なお、ジョブJ10の並列処

50

理部分の実行が終了すると、逐次演算プロセッサ 1002 は、ジョブ J10 の残りの逐次処理部分の処理 P107 を実行する。

【0105】

時刻 T5 における並列処理部分の処理 P202 ~ P206 の終結時には、並列演算プロセッサ群 1005 ~ 1008 は、終了通知 END2 を発行し、逐次演算プロセッサ 1001 は、ジョブ J20 の残りの逐次処理部分の処理 P207 を実行する。

【0106】

このように複数のジョブが同時に並列処理部分を実行しなければ、1ジョブの逐次処理部分の実行中に遊休する並列演算プロセッサ群を他のジョブの並列処理部分の実行に割り当てることができ、システム全体の稼働率を高めることが可能になる。

10

【0107】

また、性能測定を行うような特別なジョブを除くと、並列ジョブの並列化率は必ずしも高くないことが知られており、1並列ジョブの処理中には並列演算機構の遊休時間帯が多数存在する。このため、この遊休時間帯を有効活用する本発明の方法によって、計算機システムのスループットの改善が可能になる。

【0108】

次に、図12を用いて、図11に示したように複数のプログラムが同時に並列実行部分を処理しようとした場合の並列演算プロセッサ群のプログラムへの割り当て制御方式に関して説明する。

図12は、本発明の一実施形態による並列ジョブ多重スケジューリング方法における並列演算プロセッサ群のプログラムへの割り当て制御方式の説明図である。

20

【0109】

最初に、並列処理部分の起動時の処理について説明する。

計算機システムは、図9に示したように、8台の演算プロセッサから構成されている場合の処理として説明する。

【0110】

第1のジョブ J10 と第2のジョブ J20 は、図11において説明したように、並列的に実行されている並列ジョブである。並列ジョブ J10, J20 は、時刻 T1 において、逐次演算プロセッサ 1001, 1002 上でそれぞれ起動される。

【0111】

30

図中の符号 4111, 4112 は、それぞれ、図2に示した8ビットからなる同期範囲指定レジスタ 4110 のビット1, 2の状態を示している。演算プロセッサ 1001 と演算プロセッサ 1002 に対応する同期範囲指定レジスタ 4110 のビットの値は、図示する例では、実行開始当初は、逐次演算プロセッサの内、逐次演算プロセッサ 1001 に対するビットに "1" が設定されているものとする。逐次演算プロセッサ 1002 に対するビットには "0" が設定されている。

【0112】

並列処理部分の起動のためには、逐次演算プロセッサ 1001 ~ 1004 及び並列演算プロセッサ群 1005 ~ 1008 は、同期到達通知命令を用いて同期処理を行う。同期到達通知命令は、同期範囲指定機構 4100 中のレジスタ 4110 の内容によって処理の場合分けを行う。同期到達通知命令を実行した演算プロセッサに対応するビットに "1" が設定されていれば、そのまま同期処理を行う。一方、ビットに "0" に設定されていれば、その逐次演算プロセッサには並列演算プロセッサ群の使用が許可されていないとして、その旨を伝える割り込みを発生させる。

40

【0113】

図12に示した初期状態から、2つのジョブ J10, J20 が実行を開始し、ジョブ J10 が先に並列処理部分の処理 P102 ~ P106 に達したとする。この例では、同期範囲指定機構 4100 のレジスタ 4110 のビット1のビット状態を見ると、"1" が設定されており、ジョブ J10 は並列処理が許可されているので、そのまま起動 SYN1 のための同期到達通知命令 SRL1 は実行され、並列処理部分の処理 P102 ~ P106 が開

50



始される。

【 0 1 1 4 】

一方、ジョブ J 1 0 の並列処理部分の処理 P 1 0 2 ~ P 1 0 6 の実行中にジョブ J 2 0 が、並列処理部分の起動 S Y N 2 のために、同期到達通知命令 S R L 2 を実行したとすると、同期範囲指定機構 4 1 0 0 のレジスタ 4 1 1 2 のビット 2 に " 0 " が設定されているため、演算プロセッサ上には割り込み I N T 1 が発生する。この割り込み I N T 1 を契機として、共有記憶装置 2 0 0 0 内のスケジューラ S C H 1 が起動され、並列処理要求 R E Q 1 は並列実行待ちキュー C U E に入れられ、並列演算プロセッサ群が割り当てられるのを待つことになる。

【 0 1 1 5 】

並列実行待ちキュー C U E は、図 6 において説明したように、共有記憶装置 2 0 0 0 上に保持された階級管理情報 2 5 0 0 の中のキュー先頭アドレス設定フィールド 2 5 0 2 の中に、並列処理要求 R E Q が発生した順に順次入れられる。

【 0 1 1 6 】

このように、高々 1 つの逐次演算プロセッサのみのビットに " 1 " を設定し、並列処理部分の起動時に実行される同期到達通知命令が自らを実行した演算プロセッサに対応する同期範囲指定機構 4 1 0 0 の中のレジスタ 4 1 1 0 に対応するビットを調べ、自分が並列演算プロセッサ群の使用が可能か否かの場合分けに従った処理を行うことによって、排他的に並列演算プロセッサ群を割り付けることが可能である。

【 0 1 1 7 】

1 時期には、同期範囲指定機構 4 1 0 0 内のレジスタ 4 1 1 0 は、並列演算プロセッサ群に対応するビットと、並列演算プロセッサ群が割り当てられた逐次演算プロセッサに対応するビットに " 1 " が設定されている。例えば、図 2 に示した例では、並列演算プロセッサ群 1 0 0 5 ~ 1 0 0 8 に対応するビット 5 ~ 8 と、並列演算プロセッサ群が割り当てられた逐次演算プロセッサ 1 0 0 1 に対応するビット 1 に " 1 " が設定されている。

【 0 1 1 8 】

従って、並列演算プロセッサ群が割り当てられている逐次演算プロセッサの同期到達通知命令を用いた並列処理要求に対しては、そのまま並列処理の開始を許可し、一方、並列演算プロセッサが割り当てられていない逐次演算プロセッサの並列処理要求に対しては割り込みを発生させてスケジューラを起動して待ち状態にし、起動要求を並列実行待ちキューに入れることが可能になる。

【 0 1 1 9 】

なお、並列演算プロセッサ群上でのジョブの切替え時には、旧ジョブに割り当てられた逐次演算プロセッサに対応するビットに " 0 " を設定してその演算プロセッサの並列処理の起動を禁止し、新ジョブに割り当てられた逐次演算プロセッサに対応するビットには " 1 " を設定して起動を許可にする。

【 0 1 2 0 】

このように同期範囲指定機構 4 1 0 0 を、逐次演算プロセッサと並列演算プロセッサ群の組み合わせを指定する並列演算機組み合わせ機構として用いることが、本実施形態の方法の大きな特徴である。

【 0 1 2 1 】

次に、図 1 2 に戻って、並列処理部分の実行完了時の計算機システムの動作について説明する。図 1 2 の右側の部分は、並列演算プロセッサ群上での実行プログラムの変更の概略を示している。

【 0 1 2 2 】

本実施形態においては、並列演算プロセッサ群を構成する演算プロセッサは、1 つのプログラムの並列処理部分に一括して割り付けられ、1 つの並列処理部分の起動時と終結時に並列演算プロセッサ群と逐次演算プロセッサ間で同期処理を行う。

【 0 1 2 3 】

上述したように、ジョブ J 1 0 が並列処理部分の処理 P 1 0 2 ~ P 1 0 6 を並列演算プロ

10

20

30

40

50

セッサ群 1 0 0 5 ~ 1 0 0 8 上で処理している間に、ジョブ J 2 0 が並列演算の起動に失敗し、並列実行待ちキューに入れられているものとする。

【 0 1 2 4 】

並列処理部分の終結のための同期処理の完了後に、 1 引き続き逐次処理部分の実行時間が一定時間より長いとコンパイラが判断できたとき、或いは、 2 プログラムの記述中にユーザが並列処理部分の最終部分で指示を行ったときに終結時の同期完了後に、並列演算プロセッサ群を即座に他のジョブの並列処理部分の実行に割り当て替える特別な終了通知（並列部分終了通知）を行う。この並列部分終了通知が行われると、並列演算プロセッサ群を構成する演算プロセッサ上に割り込みが発生する。

【 0 1 2 5 】

即ち、ジョブ J 1 0 の並列処理部分の処理 P 1 0 2 ~ P 1 0 6 が終結するに当たって、並列演算プロセッサ群 1 0 0 5 ~ 1 0 0 8 は、並列演算プロセッサ群を並列実行待ちキューに入れられているジョブに譲渡することを指示する並列部分終了通知 E N D 1 が発行される。この終了通知 E N D 1 によって、並列演算プロセッサ群 1 0 0 5 ~ 1 0 0 8 上では割り込み I N T 2 が発生する。また、この割り込み I N T 2 を契機として共有記憶装置 2 0 0 0 内のスケジューラ S C H 2 が起動され、終結した先行情報処理プログラムの状態の内の規定の最小限のレジスタのみを退避する。

【 0 1 2 6 】

この時、ジョブ J 1 0 の並列演算プロセッサ群 1 0 0 5 ~ 1 0 0 8 上での処理は完了しているため、演算プロセッサの全ての状態を退避する必要はなく、規定の最小限のレジスタの内容のみを退避する。即ち、並列処理部分の処理の終結時には、並列演算プロセッサ群を構成する演算プロセッサ 1 0 0 5 ~ 1 0 0 8 のレジスタ上の非常に少ない量のデータの退避のみを行えばよいため、並列演算プロセッサ群の割り当てられたジョブの変更を高速に行うことができる。

【 0 1 2 7 】

ここで、必要最小限のレジスタのみが退避されていることをレジスタ退避範囲フラグ 2 4 0 8（図 5）を設定することによって記録し、並列実行待ちの他のジョブの処理を開始する。

【 0 1 2 8 】

なお、終了通知が行われなかったときは、切り替えを行わない。ここで、後続する逐次処理部分の実行時間が一定時間より長い場合にのみ終了通知を発行することが重要である。これによって、 1 並列ジョブの逐次処理時間が、他の並列ジョブの並列処理部分の実行に十分な長さをもっている時のみ、並列演算プロセッサの切り替えが行われ、並列演算プロセッサ群の遊休時間に、他の並列ジョブの並列処理部分を割り当てることが可能になる。もちろん、後続の逐次処理部分の実行時間が短い場合に終了通知を発行して、並列演算プロセッサ群の割り当て先を変更しても計算結果には影響を与えない。

【 0 1 2 9 】

ところが、例えば、 1 並列ジョブが実行時間が短い並列処理部分と逐次処理部分を頻繁に繰り返す場合に、並列処理部分の終結毎に終了通知を発行して並列演算プロセッサ群の割り当て先を変更すると、スケジューリングのオーバーヘッドが実質的な計算に費やされる計算時間に対して大きくなり、計算機システム全体の効率を低下させる。

【 0 1 3 0 】

以上のように、本実施形態においては、単純に並列処理部分の終了を契機として並列演算プロセッサの割り当て先を切り替えるのではなく、実行ジョブの逐次処理及び並列処理の実行時間の長さと分布を解析し、その結果に基づいて切り替えを行うという機能をもつことによって、 1 ジョブの処理中に遊休する並列演算プロセッサ群を低オーバーヘッドで他のジョブに割り当てることが可能となる。

【 0 1 3 1 】

割り込み I N T 2 の発生後、スケジューラ S C H 2 は、並列実行待ちキュー C U E に入っている先頭の並列処理要求 R E Q 1 を読みだして、並列実行待ち状態であったジョブ J 2

10

20

30

40

50

0に並列演算プロセッサ群1005~1008が割り付けられ、ジョブJ20の並列処理部分の処理P202~P206の実行が開始される。

【0132】

並列演算プロセッサ群上でのジョブの切替え時には、スケジューラSCH2は、同期範囲指定機構4100内のレジスタ4110の設定を変更し、旧ジョブであるジョブJ10に割り当てられた逐次演算プロセッサ1001に対応するビット4111に"0"を設定して、その演算プロセッサの並列処理の起動を禁止し、新ジョブであるジョブJ20に割り当てられた逐次演算プロセッサ1002に対応するビット4112には"1"を設定して起動を許可にする。なお、スケジューラSCH2の動作の詳細については、図14~図21を用いて後述する。

10

【0133】

ここで、スケジューラが並列プロセッサ群上のみで起動されることが重要である。ジョブJ10の並列処理部分の処理P102~P106が終了すると、演算プロセッサ1001は、ジョブJ10の引き続く逐次処理部分の処理P107の実行を即座に開始する。

【0134】

一方、これと並行して、並列演算プロセッサ群1005~1008上では、ジョブ切り替えのための処理が行われるが、この処理はジョブJ10の逐次処理部分の処理P107の実行に影響を及ぼすことなく行われるため、ジョブ切り替えのオーバーヘッドを隠蔽することが可能になる。

【0135】

20

以上説明したように、本実施形態による並列ジョブ多重スケジューリング方法によれば、並列処理部分は並列演算プロセッサ群と逐次演算プロセッサが一体となって実行するため、逐次演算プロセッサが遊休することはない。これは、並列演算プロセッサ群を逐次演算プロセッサと同等の演算プロセッサで構成し、並列演算プロセッサ群を複数の逐次演算プロセッサで共有して使用するという本実施形態の方法を用いることで初めて可能になるものである。

【0136】

また、従来の方法では、切り替えのタイミングはタイマ割り込みによるものであるため、切り替えに際しては、一般に多量のレジスタの内容の退避・回復が必要になる。それに対して、本実施形態の方法では、並列処理部分の終了を並列演算プロセッサ群の割り当て替えの契機として用いることによって、切りのよいタイミングでジョブを切り替え、退避・回復するレジスタの数を大幅に削減することが可能になり、ジョブ切り替えのオーバーヘッドを取り除いている。

30

【0137】

ここで、本実施形態による同期処理の失敗検出方法について説明する。

1つの同期処理に参加する複数の並列演算プロセッサの内、予め指定した時間(所定時間)内に同期ポイントに達しない並列演算プロセッサが存在したとき、その同期処理は失敗したと定義する。本実施形態では、同期処理が失敗したと判断するまでの時間を設定する同期タイムアウトレジスタが用意されている。

【0138】

40

同期処理が開始され、この指定された時間が経過する前にすべての演算プロセッサが同期ポイントに達しない場合、同期失敗を知らせる割り込みが演算プロセッサ上に発生し、これを契機にスケジューラが起動され、エラー処理を行う。同期タイムアウトレジスタの設定は、共有記憶装置2000上に備えられた同期タイムアウトレジスタ設定フィールド2307(図4)によって行う。

【0139】

次に、図13を用いて、本実施形態によるジョブ切り替え時のデータ授受領域の管理方法について説明する。

図13は、本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるデータ授受領域の管理方法の説明図である。

50

## 【 0 1 4 0 】

逐次演算プロセッサと並列演算プロセッサとの間のデータの授受方法については既に説明したとおりである。本実施形態においては、複数のジョブを並行に動作させるため、逐次演算プロセッサと並列演算プロセッサとの間のデータの授受に用いる演算プロセッサ間データ授受領域 2 1 0 1 ~ 2 1 0 8 の内容は、ジョブの切替え毎に、退避・回復する必要がある。そこで、並列に実行させようとするジョブの数だけ演算プロセッサ間データ授受領域を実記憶空間上に用意し、これらを同一の仮想記憶空間上の番地に対応づけるようにしている。演算プロセッサ間データ授受領域は、演算プロセッサ間データ授受領域対応付けフィールド 2 3 0 3 ( 図 4 ) に設定される。ジョブの切り替えに際して、両者の対応づけを変更することによって、演算プロセッサ間データ授受領域の状態を退避・回復する必要がなくなり、並列演算機構の割り当て変更を高速に行うことができる。

10

## 【 0 1 4 1 】

ここで、図 1 3 を用いて、例えば、2 つのジョブ A , B を並行に実行する場合について説明する。例えば、ジョブ A は、図 1 1 に示したジョブ J 1 0 であり、ジョブ B は、図 1 1 に示したジョブ J 2 0 とする。演算プロセッサ 1 0 0 5 に対する演算プロセッサ間データ授受領域 2 1 0 5 を例にとって説明する。

## 【 0 1 4 2 】

2 つのジョブ A , B に対応して、それぞれ専用の 2 つの演算プロセッサ間データ授受領域が、実記憶空間上の異なった記憶ページ P A 1 , P A 2 上に用意される。それぞれの記憶ページ P A 1 , P A 2 の先頭番地を、実記憶番地 a と実記憶番地 b とする。

20

## 【 0 1 4 3 】

一方、仮想記憶空間上では、仮想記憶番地 x に、演算プロセッサ間データ授受領域 2 1 0 5 として用意される。逐次演算プロセッサ 1 0 0 1 及び並列演算プロセッサ群の中の並列演算プロセッサ 1 0 0 5 は、見かけ上、この仮想記憶空間上の領域 2 1 0 5 を介してデータのやり取りを行う。一方、仮想記憶番地 x には、現在実行されているジョブ専用の領域の番地が対応づけられている。この対応付けは、逐次並列間データ授受領域対応付けフィールド C F によって指定される。

## 【 0 1 4 4 】

図示の例は、ジョブ A の専用領域の先頭番地 a が仮想記憶番地 x に対応付けられている状況を示している。これによって、ジョブ A の実行中はこの領域に対する逐次演算プロセッサ 1 0 0 1 からのデータの書き込み及び並列演算プロセッサ 1 0 0 5 からのデータの読み出しによって、逐次演算プロセッサと並列演算プロセッサ群との間のデータの授受を実現する。

30

## 【 0 1 4 5 】

実行ジョブの切替え時、例えば、ジョブ A からジョブ B に実行ジョブが切り替わる時は、仮想記憶番地 x に対応づけられた実記憶番地を " a " から " b " に変更することによって、ジョブ B の専用領域への書き込みと読み出しができるようになる。

## 【 0 1 4 6 】

次に、図 1 4 ~ 図 2 1 を用いて、本実施形態による並列演算プロセッサ割り当て切り替えのためのスケジューラの処理手順について説明する。

40

図 1 4 は、本発明の一実施形態による図 9 で示した並列演算プロセッサ群上で起動されるスケジューラ S C H 2 ( 図 1 2 ) の全体の動作の概略を示すフローチャートであり、図 1 5 ~ 図 2 1 は、図 1 4 に示した処理の細部を説明するフローチャートである。

## 【 0 1 4 7 】

図 1 5 は、図 1 4 のステップ S 1 4 5 1 の処理の詳細を説明したものである。図 1 6 は、図 1 4 のステップ 1 4 1 2 の処理の詳細を説明したものである。図 1 7 は、図 1 4 のステップ S 1 4 1 3 の処理の詳細を説明したものである。図 1 8 は、図 1 4 のステップ S 1 4 0 2 , S 1 4 0 3 の処理の詳細を説明したものである。図 1 9 は、図 1 4 のステップ S 1 4 0 4 の詳細を説明したものである。図 2 0 は、図 1 4 のステップ S 1 4 1 4 の処理の詳細を説明したものである。図 2 1 は、図 1 4 のステップ S 1 4 3 1 , S 1 4 3 2 の処理の

50

詳細を説明したものである。

【0148】

図14において、ステップS1401において、並列演算プロセッサ1005～1007における並列処理部分が終了し、ステップS1411において、並列演算プロセッサ1008における並列処理部分が終了し、ステップS1421において、逐次演算プロセッサ1001における並列処理部分が終了すると、ステップS1441において、並列終結時の同期処理が行われる。

【0149】

並列終結時の同期処理に引き続いて、ステップS1451において、終了通知とスケジューラの起動処理が実行される。即ち、図15に詳細に示すように、ステップS14511 10  
において、並列演算プロセッサ1005～1007は、並列演算プロセッサ群の譲渡を指示する終了通知を発行し、また、ステップS14516において、並列演算プロセッサ1008は、並列演算プロセッサ群の譲渡を指示する終了通知を発行する。

【0150】

次に、ステップS14512において、並列演算プロセッサ1005～1007では、割り込みが発生し、また、ステップS14517において、並列演算プロセッサ1008では、割り込みが発生する。これらの割り込みを契機として、ステップS14513において、並列演算プロセッサ1005～1007では、スケジューラが起動し、また、ステップS14518において、並列演算プロセッサ1008では、スケジューラが起動する。

【0151】

スケジューラが起動すると、ステップS14514において、並列演算プロセッサ1005～1007では、以降で必要となるワークレジスタを退避し、また、ステップS14519 20  
において、並列演算プロセッサ1008では、以降で必要となるワークレジスタを退避する。

【0152】

次に、図14に戻って、ステップS1442において、演算プロセッサ1001，1005～1008は、全てスケジューリング可能な状態に達したことを保証するために同期処理を行う。

【0153】

次に、演算プロセッサ1008が親プロセッサとして、スケジューラの中で振る舞う。ステップS1412において、演算プロセッサ1008は、図6のキュー先頭アドレス設定フィールド2502を参照して、並列実行待ちキューCUE（図12）が空か否かを調べる。即ち、図16に詳細に示すように、ステップS14121において、並列演算プロセッサ1008は、並列実行待ちキューCUEが空か否かを調べる。

【0154】

その結果、キューが空の場合には、ステップS14122において、演算プロセッサ1008は、現ジョブの並列起動を再許可した後、ステップS14123において、変数no\_\_switchに"#f"（fail）を書き込む。

【0155】

キューが空でない場合には、ステップS14124において、演算プロセッサ1008は 40  
、変数no\_\_switchに"#t"（true）を書き込む。これらの変数no\_\_switchへの書き込みは、各演算プロセッサに伝達される。

【0156】

さらに、ステップS14125において、演算プロセッサ1008は、キューから切り替え先新ジョブの並列起動要求REQ1（図12）を取り出し、演算プロセッサ1005～1007が参照可能な位置に格納する。さらに、ステップS14126において、新ジョブの逐次演算プロセッサに対応する同期範囲指定レジスタのビットに"1"を設定しての並列演算プロセッサ群の使用を許可する。即ち、例えば、図12に示したように、ジョブJ20の逐次演算プロセッサ1002に対する同期範囲指定レジスタのビット4112に"1"を設定する。

10

20

30

40

50

## 【0157】

図14に戻って、ステップS1412の処理が終了すると、ステップS1443において、同期命令を実行する。この同期処理は、キューが空であった場合には、変数no\_switchへの変更と現ジョブの並列起動再許可の終了の保証を行うため、空でなかった場合は、変数no\_switchへの変更及び、新ジョブの設定、新ジョブの逐次演算プロセッサに対する並列演算プロセッサ群の使用の許可の設定の完了を保証するために行うものである。

## 【0158】

続けて、演算プロセッサ1008の処理について説明すると、ステップS1412における並列待ちキューが空の場合には、ステップS1413において、現ジョブ続行する。即ち、図17に詳細に示すように、ステップS14131において、演算プロセッサ1008は、ワークレジスタを回復し、その後、ステップS14132において、現ジョブを続行する。元のジョブでは、並列処理部分は終了しているため、並列演算プロセッサ群1008上では、次の並列演算実行要求が到着するまで、待ち状態に入ることになる。

## 【0159】

一方、演算プロセッサ1005～1007は、ステップS1443における同期処理の完了後、ステップS1402において、並列待ちキューが空か否かを判断し、ステップS1403における現ジョブの続行か、若しくは、ステップS1404における新ジョブへの切替を実行する。ここで、ステップS1402における判断は、図18のステップS1402に示すように、演算プロセッサ1008)が、ステップS14123, S14124 (図16)において設定した変数no\_switchの値に従って場合分けして、処理を進める。

## 【0160】

図18において、変数no\_switchが#fであった場合は、ステップS14031において、演算プロセッサ1005～1007は、実行ジョブの変更は行わず、先に退避したワークレジスタを回復し、ステップS14032において、元のジョブの実行を再開する。元のジョブでは、並列処理部分は終了しているため、並列演算プロセッサ群1005～1007上では、次の並列演算実行要求が到着するまで、待ち状態に入ることになる。

## 【0161】

一方、演算プロセッサ8がno\_switchに#tを設定している場合、ステップS1404において、新ジョブへの切替を実行する。即ち、図19に詳細に示したように、演算プロセッサ1005～1007は、ステップS14041において、旧ジョブの実行の状態を規定された最小限のレジスタのみを退避する。そして、レジスタ退避範囲フラグ2408 (図5)を設定して、最小限のレジスタのみが退避されたことを記録した後、新しいジョブの実行を開始する。

## 【0162】

即ち、ステップS14042において、図13において説明したように、新ジョブの仮想空間の切り替え設定を行い、ステップS14043において、データ授受領域の対応付けの変更を行う。さらに、ステップ14044において、これから回復するジョブのレジスタ退避範囲フラグを参照する。レジスタ退避範囲フラグが設定されていれば、ステップS14045において、必要最低限のレジスタのみを回復し、設定されていなければ、ステップS14046において、全レジスタを回復して、ステップS14047において、新ジョブを再開する。

## 【0163】

一方、図14のステップS1414における演算プロセッサ1008の新ジョブへの切替は、図20に詳細を示す処理フローに従って行われる。即ち、図20のステップS14141において、最低限のレジスタを退避し、ステップS14142において、最低限のレジスタのみが退避されたことをレジスタ退避範囲フラグを設定し、ステップS14143において、図13において説明したように、新ジョブの仮想空間の切り替え設定を行い

10

20

30

40

50

、ステップS 1 4 1 4 4において、データ授受領域の対応付けの変更を行う。さらに、ステップ1 4 1 4 5において、これから回復するジョブのレジスタ退避範囲フラグを参照する。レジスタ退避範囲フラグが設定されていれば、ステップS 1 4 1 4 6において、必要最低限のレジスタのみを回復し、ステップS 1 4 1 4 7において、レジスタ退避範囲フラグをリセットする。一方、レジスタ退避範囲フラグが設定されていなければ、ステップS 1 4 1 4 8において、全レジスタを回復して、ステップS 1 4 1 4 9において、新ジョブを再開する。

#### 【0 1 6 4】

次に、図1 4における演算プロセッサ1 0 0 2の処理について説明する。

ステップS 1 4 3 1において、並列起動要求がなされた場合、図1 2において説明したように、既に他の並列ジョブが実行中であるため、並列起動に失敗し、ステップS 1 4 3 2において、並列演算プロセッサの割当待ちとなる。即ち、図2 1に詳細に説明するように、ステップS 1 4 3 1において、並列起動に失敗すると、ステップS 1 4 3 2 1において、起動禁止割り込みが発生する。これを契機に、ステップS 1 4 3 2 2において、スケジューラが起動され、ステップS 1 4 3 2 3において、スケジューラは、起動要求を並列実行待ちキューへ登録し、ステップS 1 4 3 2 4において、並列演算プロセッサ群の割り当て待ちとなっている。以上が、並列演算プロセッサ上で起動されたスケジューラの動作の概略である。

#### 【0 1 6 5】

さらに、図1 4において、ステップS 1 4 0 4 , S 1 4 1 4における新ジョブへの切替え処理が行われると、ステップS 1 4 4 4において、演算プロセッサ1 0 0 2 , 1 0 0 5 ~ 1 0 0 8は、並列部分起動時の同期処理を行う。同期処理が終了すると、ステップS 1 4 0 5において、演算プロセッサ1 0 0 5 ~ 1 0 0 7は、並列処理部分の処理を開始し、ステップS 1 4 1 5において、演算プロセッサ1 0 0 8は、並列処理部分の処理を開始し、ステップS 1 4 3 3において、演算プロセッサ1 0 0 2は、並列処理部分の処理を開始する。以上のようにして、並列演算プロセッサ群上で起動されるスケジューラを用いた並列演算プロセッサ割り当て切り替えが実行される。

#### 【0 1 6 6】

次に、図2 2を用いて、本実施形態による計算機システム内の演算プロセッサの逐次演算プロセッサ及び並列演算プロセッサ群への割り付け個数の動的な変更について説明する。図2 2は、本発明の一実施形態による並列ジョブ多重スケジューリング方法により演算プロセッサの割り付けを動的変更する例の説明図である。

#### 【0 1 6 7】

上述したように、本実施形態においては、複数の演算プロセッサを論理的に逐次演算プロセッサと並列演算プロセッサに分割し、1時期には1つの逐次演算プロセッサが複数の演算プロセッサから構成される並列演算プロセッサ群を一括して獲得して、この並列演算プロセッサ群を用いて並列処理部分の実行を行っている。ここで、逐次演算プロセッサに用いられる演算プロセッサと並列演算プロセッサ群に用いられる演算プロセッサは同等のものであり、且つ同等の装置によって接続されるため、計算機システム内の演算プロセッサを、逐次演算プロセッサに割り当てる個数と並列演算プロセッサ群に割り当てる個数の比率を動的に変更することが可能であり、実行ジョブが要求する計算機資源の大小に応じた運用を行えるものである。この割り付け個数の動的変更は、共有記憶装置2 0 0 0の管理情報2 3 0 0の中の構成変更情報フィールド2 3 0 4 (図4)に保持されたデータに基づいて、構成比率、頻度、契機等が指定される。

#### 【0 1 6 8】

図2 2に示す例は、時間Tの経過毎に演算プロセッサにタイマ割り込みを発生させ、このタイマ割り込みを契機として、演算プロセッサ上でスケジューラを起動する。起動されたスケジューラは、直前まで実行されていたプログラムの状態を退避し、計算機システム内の逐次演算プロセッサに割り当てる演算プロセッサ数と並列演算プロセッサ群に割り当てる演算プロセッサ数の構成比率を変更するように運用される。なお、スケジューラの処理

10

20

30

40

50

動作の詳細については、図 2 3 を用いて後述する。

【 0 1 6 9 】

例えば、図 9 に示した 8 台の演算プロセッサから構成される計算機システムにおいて、8 台の演算プロセッサの内の 4 台の演算プロセッサを逐次演算プロセッサとし、残りの 4 台の演算プロセッサを並列演算プロセッサ群として用いることによって複数のジョブを並行実行し、システム全体のスループットを向上させるシステム構成や、8 台の演算プロセッサの内の 1 台を逐次演算プロセッサとして用い、残り 7 台を並列演算プロセッサ群として用いる 1 ジョブの最高性能を実現するシステム構成など、異なった演算プロセッサの使用方法を同一の枠組内で実現することが可能であり、且つこれらのシステム構成を時間 T の経過毎に切り替えることでさまざまな要求を満たすことが可能になる。

10

【 0 1 7 0 】

ここで、図 2 2 のシステム構成の動的変更について具体的に説明する。

図 9 に示した 8 台の演算プロセッサから構成される計算機システムにおいて、演算プロセッサ 1 0 0 1 ~ 1 0 0 8 は、逐次演算プロセッサと並列演算プロセッサ群とに論理的に分割されている。階級 B<sup>T</sup> の実行時には、演算プロセッサ 1 0 0 1 ~ 1 0 0 4 が、逐次演算プロセッサを構成し、演算プロセッサ 1 0 0 5 ~ 1 0 0 8 が、並列演算プロセッサ群を構成しているものとする。また、階級 A<sup>T</sup> の実行時には、演算プロセッサ 1 0 0 1 が、逐次演算プロセッサを構成し、演算プロセッサ 1 0 0 2 ~ 1 0 0 8 が、並列演算プロセッサ群を構成しているものとする。

【 0 1 7 1 】

20

第 1 のジョブ J 1 0 , 第 2 のジョブ J 2 0 , 第 4 のジョブ J 4 0 , 第 5 のジョブ J 5 0 は、並列的に実行されている並列ジョブである。第 1 のジョブ J 1 0 は、逐次処理部分である処理 P 1 0 1 と、並列処理部分である処理 P 1 0 2 ~ P 1 0 6 と、逐次処理部分である処理 P 1 0 7 から構成されている。また、第 2 のジョブ J 2 0 は、逐次処理部分である処理 P 2 0 1 と、並列処理部分である処理 P 2 0 2 ~ P 2 0 6 と、逐次処理部分である処理 P 2 0 7 から構成されている。第 4 のジョブ J 4 0 は、逐次処理部分である処理 P 4 0 1 のみから構成されており、第 5 のジョブ J 5 0 は、逐次処理部分である処理 P 5 0 1 のみから構成されている。

【 0 1 7 2 】

最初に、並列ジョブ J 1 0 , J 2 0 , J 4 0 , J 5 0 が、逐次演算プロセッサ 1 0 0 4 , 1 0 0 3 , 1 0 0 2 , 1 0 0 1 上でそれぞれ起動される。第 1 のジョブ J 1 0 の逐次処理部分の処理 P 1 0 1 が終了すると、第 1 のジョブ J 1 0 の逐次処理部分の処理 P 1 0 1 を実行した逐次演算プロセッサ 1 0 0 4 が、並列演算プロセッサ群 1 0 0 5 ~ 1 0 0 8 を獲得して、ジョブ J 1 0 の並列処理部分の処理 P 1 0 2 ~ P 1 0 6 が実行を進める。

30

【 0 1 7 3 】

並列処理部分の処理 P 1 0 2 ~ P 1 0 6 が終結すると、逐次演算プロセッサ 1 0 0 4 は、続く逐次処理部分 P 1 0 7 の処理を開始する。また、並列プロセッサ 1 0 0 5 ~ 1 0 0 8 によるジョブ J 1 0 の並列処理部分の実行の終了時に、ジョブ J 2 0 の並列処理部分 P 2 0 2 ~ P 2 0 6 の実行待ちになっていたすると、逐次演算プロセッサ 1 0 0 3 は、並列演算プロセッサ 1 0 0 4 ~ 1 0 0 8 に対して同期処理を行った後、並列処理部分の処理 P 2 0 2 ~ P 2 0 6 が、逐次演算プロセッサ 1 0 0 3 と並列演算プロセッサ群 1 0 0 5 ~ 1 0 0 8 によって並列的に実行される。並列処理部分の処理 P 2 0 2 ~ P 2 0 6 の終結時には、並列演算プロセッサ群 1 0 0 5 ~ 1 0 0 8 は、終了通知を発行し、逐次演算プロセッサ 1 0 0 3 は、ジョブ J 2 0 の残りの逐次処理部分の処理 P 2 0 7 を実行する。

40

【 0 1 7 4 】

共有記憶装置 2 0 0 0 の管理情報 2 3 0 0 の中の構成変更情報フィールド 2 3 0 4 ( 図 4 ) に保持されたデータに基づいて、例えば、時間 T の経過毎に演算プロセッサ 1 0 0 1 ~ 1 0 0 8 上に割り込みを発生する。この割り込みを契機として、演算プロセッサ 1 0 0 1 ~ 1 0 0 8 上ではスケジューラ S C H 1 0 , S C H 1 1 が起動され、システムの構成と実行対象となるプログラムを変更して、逐次演算プロセッサと並列演算プロセッサ群の構成

50



比率，即ち、割り付け個数の動的変更を行う。

【0175】

即ち、階級 $B^T$ の実行時が終了すると、割り込みが発生し、スケジューラSCH10が起動して、演算プロセッサ1001が、逐次演算プロセッサを構成し、演算プロセッサ1002～1008が、並列演算プロセッサ群を構成するシステム構成に変更して、階級 $A^T$ を実行可能とする。

【0176】

階級 $A^T$ を実行時には、例えば、第3のジョブJ30が実行される。新しいジョブJ30が起動されると、逐次処理部分の処理P301の実行が、逐次プロセッサ1001上で起動される。逐次処理部分の処理P301が終了すると、次の並列処理部分の処理の起動時に、逐次演算プロセッサ1001と並列演算プロセッサ群1002～1008との間の同期処理が完了すると、並列処理部分の処理P302～P309が、逐次演算プロセッサ1001と並列演算プロセッサ群1002～1008によって並列的に実行される。さらに、並列処理が終了すると、逐次演算プロセッサ1001は、逐次処理部分の処理P310を実行する。

10

【0177】

逐次処理部分の処理P310が終了すると、次の並列処理部分の処理の起動時に、逐次演算プロセッサ1001と並列演算プロセッサ群1002～1008との間の同期処理が完了すると、並列処理部分の処理P311～P318が、逐次演算プロセッサ1001と並列演算プロセッサ群1002～1008によって並列的に実行される。さらに、並列処理が終了すると、逐次演算プロセッサ1001は、逐次処理部分の処理P319を実行する。

20

【0178】

以上のようにして、実行ジョブが要求する計算機資源の大小に応じて、逐次演算プロセッサと並列演算プロセッサ群の構成比率を変更した運用が可能となる。従って、実行並列ジョブの性質，即ち、逐次処理部分と並列処理部分の比率に応じて、柔軟に逐次演算処理能力と並列演算処理能力の比率を動的に変化させて、計算機システムの稼働率を高めることが可能となる。

【0179】

次に、図23を用いて、本実施形態によるジョブ階級切り替え時のスケジューラの動作の概略について説明する。

30

図23は、本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるジョブ階級切り替え時のスケジューラの動作を示すフローチャートである。

【0180】

ステップS2301において、共有記憶装置2000の中の管理情報2300の構成変更情報設定フィールド2304（図4）に設定されている構成変更情報に従って、逐次演算プロセッサ1001に対してタイマ割り込みが発生する。この割り込みを契機として、ステップS2302において、演算プロセッサ1001は、スケジューラSCHを起動する。

【0181】

スケジューラSCHが起動されると、ステップS2303において、演算プロセッサ1001は、以降の処理で用いるワークレジスタの内容を退避する。次に、ステップS2304において、構成変更情報を参照して切り替え先の階級を取得する。切り替え先の階級とは、図22において説明したような階級 $B^T$ や階級 $A^T$ のことである。ステップS2305において、切り替え先階級があったか否かを判断し、ない場合には、ステップS2306において、退避しておいたワークレジスタを回復し、ステップS2307において、現行階級を続行する。

40

【0182】

一方、切り替え先階級が存在した場合、ステップS2308において、演算プロセッサ1001は、演算プロセッサ1002～1008に対して割り込みを発生させる。

50

## 【0183】

他方、演算プロセッサ1002～1008では、ステップS2321において、演算プロセッサ1001からの割り込みが発生する。割り込みが発生した演算プロセッサ1002～1008は、ステップS2322において、スケジューラを起動し、ステップS2323において、以降の処理で用いるワークレジスタの内容を退避する。

## 【0184】

その後、ステップS2331において、演算プロセッサ1001, 1002～1008は、同期処理を行う。

## 【0185】

この同期処理が完了すると、ステップS2309において、演算プロセッサ1001は、レジスタを退避し、ステップS2310において、現階級での同期範囲指定機構の状態を退避する。その後、ステップS2311において、切り替え先の階級に所属するジョブのレジスタを回復する。

10

## 【0186】

次に、ステップS2312において、新階級の同期範囲指定機構の状態を回復し、ステップS2313において、新階級でのジョブの仮想空間を設定し(14017)、さらに、ステップS2314において、データ授受領域の内容を階級切り替えの反映したものに変更する。

## 【0187】

一方、演算プロセッサ1002～1008は、ステップS2324において、レジスタの退避を行った後、ステップS2325において、切り替え先階級のジョブのレジスタ内容を回復する。さらに、ステップS2326において、新ジョブの仮想空間の設定し、ステップS2327において、データ授受領域の設定の変更を行う。

20

## 【0188】

その後、ステップS2332において、演算プロセッサ1001, 1002～1008は、同期処理を行う。

## 【0189】

この同期処理が完了すると、ステップS2315において、演算プロセッサ1001は、全ての演算プロセッサが新階級の実行の準備ができたことを確認した後、新階級の実行に移り、また、ステップS2328において、演算プロセッサ1002～1008は、新階級の実行に移る。

30

## 【0190】

本実施形態におけるジョブ階級切り替えの運用に当たっては次に行っている。即ち、前述したように、計算機システムの最高速度を実現しようとすれば、全ての演算プロセッサをそのジョブに割り付ける必要が生じる。一方、最終的には最高速度を目指すジョブでも、デバッグの段階では計算機システム全てを占有する必要はなく、複数のデバッグ作業を並行して進行させることが望ましい場合や、1つのジョブの性能をある程度犠牲にしても、複数のジョブを並行に動作させ、計算機システム全体の稼働率を高めることが要求される場合もある。

## 【0191】

そこで、本実施形態においては、これらの要求に従って、並列ジョブは階級分けされる。計算機システム上では、構成変更情報設定フィールド2304(図4)に設定された構成変更情報に基づいて階級を切り替え、それぞれの階級に割り付けられた時間内で階級に属するジョブをその階級の実行方式に従って処理する。

40

## 【0192】

さらに、計算機システムは、各階級毎にスケジューリングのためのジョブスケジューリングキュー(図6)を備える。計算機上で起動されたジョブは自分の属する階級のジョブスケジューリングキューに入れられ、その階級の処理の順番が到来し次第、処理される。

## 【0193】

同一階級に属するジョブは、並列実行待ちキュー(図6)を用いて同一階級のプログラム

50

間で共有される並列演算プロセッサ群を切り替えて処理を進める。階級内でこれらのジョブがどのように処理されるかは、その階級の定義に従うものである。

【0194】

ここで、図22を参照して、以上述べたジョブ階級切り替えの運用例について説明する。図22に示した例では、階級Aは全システムを占有する並列ジョブの階級とし、階級Bは4台の逐次演算プロセッサと4台の演算プロセッサから成る並列演算プロセッサ群を備えるシステム構成で実行するジョブの階級と定めるとする。図22は、階級Bの実行から階級Aの実行への切替えの様子を示している。

【0195】

図22において、階級Bに属する4つの並列ジョブJ10、J20、J40、J50を階級Bに割り付けられた計算時間Tの中で処理している。階級Bでは、4台の演算プロセッサを逐次演算プロセッサとして用い、同時に4つの並列ジョブを実行する。5つ以上のジョブが同時に起動されている場合は、5つ目以降のジョブは、階級Bのジョブスケジューリングキューの中で先行する4つのジョブの実行の終了を待つ。

10

【0196】

階級Bに割り当てられた時間Tの経過後に起動されたスケジューラSCH10は、階級Bの実行の状態を退避し、他方の階級Aの実行のためにシステム内のプロセッサ構成の変更と階級Aに属するジョブの処理を開始する。

【0197】

階級Aでは、同時には1つのプログラムしか処理しないので、同時に起動されている2つ目以降のジョブは、階級Aのジョブスケジューリングキューの中で実行中のジョブの終了を待つ。

20

【0198】

同様にして、階級Aに割り当てられた時間Tが経過すると、再びスケジューラSCH11が起動され、階級Aの実行状態を退避し、システム内の演算プロセッサの構成の変更を行い、前回退避した階級Bの実行の状態を回復し、再開する。

以降、同様にして交互に異なった階級を時分割で実行することになる。

【0199】

なお、以上の説明では、階級Aと階級Bに割り付ける計算時間Tを等しくしたが、階級毎に異なった時間を割り付けることも可能である。また、階級A、階級Bの2つの階級のみではなく、他の演算プロセッサの構成をもつ階級を加えることも可能である。

30

【0200】

また、同一の演算プロセッサの構成で複数の階級を作成することも可能である。さらに、同一の演算プロセッサの構成で、異なる計算時間の割り当てられた階級や、時分割された時間の割り当ての優先度を異ならせたものを作成することも可能である。

【0201】

また、1階級内では、同時に処理され得るジョブ数の上限をその階級の逐次演算機構の数としたが、各階級が備えるジョブスケジューリングキューを用いて、それ以上の数のジョブを階級に割り振られた計算時間内で切り替えて実行することも可能である。

【0202】

40

なお、構成変更情報内に保持された時間Tと切り替えの契機の設定によって、例えば、昼間は前述の階級Bで運用し、夜間は階級Aで運用するといった方式を実現することも可能である。

【0203】

次に、図24を用いて、本実施形態による待ち状態の発生を契機とした計算機システム内の演算プロセッサの逐次演算プロセッサ及び並列演算プロセッサ群への割り付け個数の動的な変更、即ちスケジューリングについて説明する。

図24は、本発明の一実施形態による並列ジョブ多重スケジューリング方法による待ち状態の発生を契機としたスケジューリングの説明図である。

【0204】

50

計算機システム内の演算プロセッサ1001～1008は、それぞれ、共有記憶装置2000の中の管理情報2300（図4）に、自らが待ち状態になっているか否かの情報を保持する待ち状態フラグ2310（図4）を有している。並列演算プロセッサ群を構成する複数の演算プロセッサの内の1台が待ち状態になると、待ち状態フラグ2310によって、並列演算プロセッサ群に参加する他の演算プロセッサ及び、逐次演算プロセッサが待ち状態になったか否かを調べる。ここで、全ての並列演算プロセッサ群中の全ての演算プロセッサ及び逐次演算プロセッサが待ち状態になったことが判明した場合、スケジューラは並列演算機構の割り当てを待ち状態となった現在のジョブから、並列実行待ちキューに入れられた他のジョブに切り替える。

【0205】

10

図24に示す例において、計算機システムは、図9に示したように、8台の演算プロセッサから構成されているものとする。第1のジョブJ60と第2のジョブJ70は、並列的に実行されている並列ジョブである。並列ジョブJ60、J70は、逐次演算プロセッサ1001、1002上でそれぞれ起動される。

【0206】

2つのジョブJ60、J70が実行を開始し、ジョブJ60の処理P601及びジョブJ70の処理P701が実行される。ジョブJ60の処理P601が終了し、並列処理部分の処理P602～P606に達したとすると同期処理SYN6を行って、演算プロセッサ1001及び演算プロセッサ1005～1008を用いて、並列処理部分の処理P602～P606が実行される。

20

【0207】

一方、ジョブJ60の並列処理部分の処理P602～P606の実行中にジョブJ70が、並列処理部分の起動SYN6を行うと、他の並行ジョブが実行中であるため、割り込みINT6が発生する。この割り込みINT6を契機として、共有記憶装置2000内のスケジューラSCH6が起動され、並列処理要求REQ6は並列実行待ちキューCUEに入れられ、並列演算プロセッサ群が割り当てられるのを待つことになる。

【0208】

この状態から、並列演算プロセッサ群を構成する全ての演算プロセッサ1005～1008及び演算プロセッサ1001が、ディスクの読み書き等を行うことにより、待ち状態WAIT1になったとすると、スケジューラSCH6が並列実行待ちキューCUEの中からジョブJ70の並列実行要求REQ6を取り出す。スケジューラSCH6の詳細な処理については、図25を用いて後述する。

30

【0209】

並列実行待ち状態であったジョブJ70に並列演算プロセッサ群1005～1008が割り付けられ、ジョブJ70の並列処理部分の処理P702～P706の実行が開始される。

【0210】

待ち状態を生じた事象WAT-RELが終了し、ジョブJ60の待ち状態が解除されると、スケジューラSCH6は、実行中の並列ジョブJ70の終了を待って、再び、ジョブJ60に並列演算プロセッサ群を割り付け、待ち状態に引き続く処理P602～P606を再開する。ジョブJ60の並列処理部分の処理P602～P606が終了すると、演算プロセッサ1001は、ジョブJ60の引き続く逐次処理部分の処理P607の実行を即座に開始する。

40

【0211】

また、ジョブJ70の並列処理部分の処理P702～P706が終了すると、演算プロセッサ1002は、ジョブJ70の引き続く逐次処理部分の処理P707の実行を即座に開始する。

【0212】

以上が並列演算プロセッサ群上での待ち状態の発生を契機としたスケジューリングの方法である。

50

## 【0213】

次に、同じく、図24を用いて、逐次演算プロセッサ上での待ち状態の発生を契機としたスケジューリングの方法について説明する。

逐次演算プロセッサ上で待ち状態になった場合は、その待ち状態になったジョブの並列処理部分が並列演算プロセッサ群上で処理されていなければ、他のジョブの逐次処理部分に逐次演算プロセッサを割り当て替える。

## 【0214】

即ち、図24において、演算プロセッサ1002上のジョブJ70がその逐次処理部分の処理P707の実行中に待ち状態WAIT2となったとし、且つこの時、並列演算プロセッサ群1005～1008上では、ジョブJ70の並列処理部分は実行されていないものとする。すると、演算プロセッサ1002上では、スケジューラSCH7が、他のジョブの逐次処理部分への切り替えCHAN1を行う。

10

## 【0215】

以上のようにして、逐次演算プロセッサ上での待ち状態発生を契機としてスケジューリングを行うことができる。

## 【0216】

ここで、図25を用いて、本実施形態による待ち状態の発生を契機としたスケジューリング時の演算プロセッサの処理動作について説明する。

図25は、本発明の一実施形態による並列ジョブ多重スケジューリング方法における待ち状態の発生によるスケジューリングの処理を説明するフローチャートである。

20

## 【0217】

図25のステップS2501において、例えば、システムコールが発行され、ステップS2502において、待ち状態が発生する。次に、ステップS2503において、スケジューラSCHが起動され、待ち状態フラグが設定される。

## 【0218】

ここで、ステップS2504において、待ち状態フラグが設定された演算プロセッサが同期範囲内であるか否かが判断され、同期範囲内であれば、ステップS2505において、同期範囲内の全ての演算プロセッサが待ち状態か否かが判断され、他の演算プロセッサが待ち状態にならなければ、ステップS2506において、そのまま続行する。並列実行待ちの他の同期範囲内の全ての演算プロセッサが待ち状態になったら、ステップS2507において、並列実行待ちの他のジョブに切り替わる。即ち、図24に示す例では、ジョブJ60が待ち状態となり、ジョブJ70に切り替わる。

30

## 【0219】

また、ステップS2504における判断で、待ち状態になった演算プロセッサが同期範囲内でないとき、すなわち、並列処理部分を実行中でないジョブに割り当てられた逐次演算プロセッサならば、ステップS2508において、スケジューリングキュー内の他のジョブの逐次処理部分に切り替える。即ち、図24に示す例では、ジョブJ70が待ち状態となり、他のジョブに切り替わる。

## 【0220】

次に、図26を用いて、本実施形態による複数の並列演算プロセッサ群を備えた計算機システムの動作について説明する。

40

図26は、本発明の一実施形態による並列ジョブ多重スケジューリング装置における複数の並列演算プロセッサ群を備えた計算機システムの動作の説明図である。

## 【0221】

本実施形態においては、計算機システムが、合計12台の演算プロセッサ1001～1012によって構成されているものとして説明する。演算プロセッサ1001～1004の4台は逐次演算プロセッサとして構成され、演算プロセッサ1005～1008の4台が第1の並列演算プロセッサ群として構成され、演算プロセッサ1009～1012の4台が第2の並列演算プロセッサ群として構成されているものとする。

## 【0222】

50

4 台の逐次演算プロセッサ 1 0 0 1 ~ 1 0 0 4 では、4 つのジョブ J 1 0 , J 2 0 , J 8 0 , J 9 0 が起動され、それぞれの処理 P 1 0 1 , P 2 0 1 , P 8 0 1 , P 9 0 1 が実行される。最初に並列実行部分に達したジョブ J 8 0 が、最初に第 1 の並列演算プロセッサ群 1 0 0 5 ~ 1 0 0 8 を獲得し、並列処理部分の処理 P 2 0 2 ~ P 2 0 6 の実行を開始する。次に並列処理部分に達したジョブ J 1 0 は、第 2 の並列演算プロセッサ群 1 0 0 9 ~ 1 0 1 2 を獲得して、ジョブ J 8 0 の並列処理部分の処理 P 8 0 2 ~ P 8 0 6 と並行して、自らの並列処理部分の処理 P 1 0 2 ~ P 1 0 6 の実行を開始する。

#### 【 0 2 2 3 】

ここで、並列演算プロセッサ群が 1 つのみの場合には、同期範囲指定機構 4 1 0 0 中の同期範囲指定レジスタ 4 1 1 0 のレジスタのビットを、図 2 において示したように、演算プロセッサに対応付けて同期範囲の指定及び逐次演算プロセッサの並列演算プロセッサ群の使用の許可・禁止を制御した。しかしながら、並列演算プロセッサ群が複数個用意された環境では、逐次演算プロセッサと並列演算プロセッサ群との組み合わせを指示・制御する同期範囲指定レジスタを全演算プロセッサ個数分のビット数だけ用意し、並列演算プロセッサが割り当てられている逐次演算プロセッサは、自らの同期範囲指定レジスタに自らの演算プロセッサ番号を設定し、並列演算プロセッサは自分が割り当てられている先の逐次演算プロセッサの演算プロセッサ番号を、同期範囲指定レジスタに設定する。並列演算プロセッサ群が割り付けられていない逐次演算プロセッサに対応するビットは " 0 " を設定する。自らの演算プロセッサ番号が設定されている演算プロセッサに並列演算プロセッサ群の使用が許可される。並列処理部分の起動時の同期処理では、このビットの値によって、同期処理を続行するか、割り込みを発生してスケジューラを起動し、並列実行待ちキューに登録するかが判断される。

#### 【 0 2 2 4 】

ここで、図 2 8 を用いて、本実施形態による同期範囲指定レジスタの構成について説明する。

図 2 8 は、本発明の一実施形態による複数の並列演算プロセッサ群を備えた計算機システムにおける同期範囲指定レジスタの構成図である。

#### 【 0 2 2 5 】

1 2 台の演算プロセッサから構成される計算機システムでは、1 2 ビットの同期範囲指定レジスタ 3 1 1 0 A から構成されている。逐次演算プロセッサ 1 0 0 2 上のジョブ J 2 0 には、第 2 の並列演算プロセッサ群を構成する演算プロセッサ 1 0 0 9 ~ 1 0 1 2 が割り当てられることを示すため、同期範囲指定レジスタ 3 1 1 0 A のビット 2 , 9 ~ 1 2 には、" 2 " が設定されている。また、逐次演算プロセッサ 1 0 0 4 上のジョブ J 9 0 には、第 1 の並列演算プロセッサ群を構成する演算プロセッサ 1 0 0 5 ~ 1 0 0 8 が割り当てられることを示すため、同期範囲指定レジスタ 3 1 1 0 A のビット 4 , 5 ~ 8 には、" 4 " が設定されている。並列演算プロセッサ群の使用が禁止された演算プロセッサ 1 0 0 1 と演算プロセッサ 1 0 0 3 に対応する同期範囲指定レジスタ 3 1 1 0 A のビット 1 , 3 には、" 0 " が設定されている。

#### 【 0 2 2 6 】

並列演算プロセッサ群 1 0 0 1 上のジョブ J 1 0 の並列処理部分の処理 P 1 0 2 ~ P 1 0 6 の実行が完了すると、引き続いて起動されたジョブ J 2 0 の並列処理部分の処理 P 2 0 2 ~ P 2 0 6 が、第 2 の並列演算プロセッサ群 1 0 0 9 ~ 1 0 1 2 で処理される。一方、第 1 の並列演算プロセッサ群では、ジョブ J 8 0 の処理の終了後に起動されたジョブ J 9 0 の並列処理部分の処理 P 9 0 2 ~ P 9 0 6 の実行が開始される。

#### 【 0 2 2 7 】

この処理の終了後、第 2 の並列演算プロセッサ群では、ジョブ J 8 0 の並列処理部分の処理 P 8 0 8 ~ P 8 1 2 の実行が開始される。ジョブ J 8 0 の 1 つ前の並列処理部分の処理 P 8 0 2 ~ P 8 0 6 は、第 1 の並列演算プロセッサ群で実行されたが、同一のジョブ中の並列処理部分に前回の並列処理部分が実行された並列演算プロセッサ群とは異なる並列演算プロセッサ群が割り当てられることも有り得る。

10

20

30

40

50

## 【0228】

次に、図27を用いて、演算プロセッサ間の同期処理をソフトウェア的に実現する方法について説明する。

図27は、本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるソフトウェアの同期処理を説明するフローチャートである。

## 【0229】

上述した例では、演算プロセッサ間の同期処理を専用のハードウェア機構を用いていたが、本実施形態においては、その機能をソフトウェア的に実現するようにしている。

## 【0230】

演算プロセッサは、バリア同期機構を保持する代わりに、共有記憶装置2000上に、各演算プロセッサに対応した同期到達フラグ2309（図4）を備える。同期ポイントに達した演算プロセッサは、同期到達フラグ2309を設定することで同期要求を行う。全ての演算プロセッサの同期到達フラグ2309が設定された時点で同期処理が完了したことをソフトウェアによって検知する。

10

## 【0231】

同期範囲指定機構4100の代わりに、共有記憶装置2000上に、同期範囲指定フィールド2307（図4）を備え、同期範囲内の演算プロセッサを同期範囲指定フィールド2307に登録することで、同期範囲の指定をソフトウェアで行う。各演算プロセッサは同期要求に際して、同期範囲指定フィールド2307の内容によって自らが同期の範囲内か否かを判定し、同期範囲内であれば、そのまま同期処理を続行し、同期範囲外であれば自らを実行する演算プロセッサ上に割り込みを発生させる。さらに、同期範囲内と設定されていた演算プロセッサが同期要求を行い、要求後で同期成立前に、プログラムにて同期範囲指定機構に対して該演算プロセッサを同期範囲外と設定した場合は、同設定を行った演算プロセッサが、新たに同期範囲外に指定された演算プロセッサに対して割り込みを発生させることによって、ハードウェアによって実現していた同期範囲指定機構4100をソフトウェアによって実現することができる。

20

## 【0232】

また、同期タイムアウト機構もソフトウェア的に実現可能である。同期処理に到達した演算プロセッサは、同期範囲内の他の演算プロセッサが同期に達するまで無限に待つのではなく、一定の時間が経過しても同期に達しない演算プロセッサが存在した場合、同期処理の失敗として、同期範囲内の演算プロセッサに割り込みを発生させスケジューラを起動してエラー処理を行う。

30

## 【0233】

図27のステップS2801において、同期範囲内の演算プロセッサを同期範囲外に設定し直すときは、演算プロセッサが同期フラグを設定しているか否かを調べる。設定していなければ、ステップS2802において、そのまま同期範囲外に設定する。一方、同期フラグが設定されていれば、ステップS2803において、演算プロセッサ間割り込みを同演算プロセッサに発生させる。割り込みを受けた演算プロセッサは、ステップS2804において、スケジューラを起動し、同演算プロセッサ上で実行されているジョブを並列実行待ちキューに登録すると同時に同演算プロセッサを同期範囲外に設定する。

40

## 【0234】

次に、図29を用いて、終了通知を動的な判断によって発行する方法について説明する。図29は、本発明の一実施形態による並列ジョブ多重スケジューリング方法における終了通知の動的判断処理を示すフローチャートである。

## 【0235】

静的な終了通知の発行は、コンパイラによって自動的に引き続く逐次処理部分の実行時間の解析を行い、終了通知を発行することにより行われる。

## 【0236】

本方式の適用の可否は、1並列処理の終了時にどの程度の頻度で終了通知を発行し、退避・回復のオーバーヘッドが低い並列演算プロセッサ群上の実行ジョブの切り替えが行えるか

50

に依存するものである。少しでも多くの場合に終了通知が発行されるように、ユーザは並列処理の終了時に引き続く逐次処理部分が長いと予めプログラムの構造として明らかな時は、並列演算プロセッサ群の割り当て変更を促進するために、動的に終了通知を発行する命令を実行することができる。

#### 【0237】

本実施形態においては、実行時に並列処理部分と逐次処理部分の実行時間を、共有記憶装置上のジョブ管理情報2400中の先行逐次情報フィールド2405（図5）に記録し、そのデータに基づいて統計的な解析を加えて、確率的に1並列処理部分に引き続く逐次処理部分の実行時間が一定時間よりも長いと判断されたときは、終了通知をスケジューラが自動的に発行するといった、実行時の動的な解析に基づいた処理を行う。

10

#### 【0238】

図29のステップS2901において、先行逐次処理部分の実行時間を共有記憶装置2000中のジョブ管理情報2400中に新たに備える先行逐次情報フィールド2405（図5）に保存する。次に、ステップS2902において、演算プロセッサは、並列処理部分の処理を実行する。

#### 【0239】

ステップS2903において、スケジューラは、並列処理部分の終了時には、先行逐次情報フィールド2405の内容に基づいて後続逐次処理部分の実行時間が一定時間より長いかな否かを判定する。長いと判定された場合には、ステップS2904において、終了通知を行い、短いと判定された場合には、ステップS2905において、終了通知は行わないようにする。

20

#### 【0240】

なお、1並列ジョブの実行開始から終了までの全ての逐次処理部分の実行時間を記録しなくてもよい。ある定数Aを定め、過去A回分の逐次処理部分の実行時間の履歴を保持し、これに基づいて判定を行うことも可能である。例えば、並列処理部分が繰り返し構造の中に存在する場合、実行時には同並列処理部分が繰り返し回数分周期的に出現する。ここで、仮に、同並列処理部分間の逐次処理部分の実行時間が一定時間よりも長いにも拘わらず、終結時に終了通知を発するコードが生成されていなかったとする。この時、繰り返しで実行される並列処理部分は過去A回の履歴を参照することにより、終結時に終了通知を発行することが可能になる。

30

#### 【0241】

本実施形態によれば、コンパイラによる静的な解析による終了通知を発行が難しい場合にも、この動的な終了通知の発行によって、並列演算プロセッサ群の割り当て変更が可能となる。

#### 【0242】

次に、図30及び図31を用いて、本実施形態における実行ジョブが、唯一つのときの実行方法について説明する。

図30は、本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるジョブが単一で実行されているときのレジスタ退避時の動作を示すフローチャートであり、図31は、本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるジョブが単一で実行されているときのレジスタ回復時の動作を示すフローチャートである。

40

#### 【0243】

1個のジョブのみが実行されている状態では、すべての並列ジョブの多重実行のための処理を省くことが可能になる。これまでの説明では、並列処理部分の起動と終結においては、レジスタ退避範囲フラグ2408（図5）を設定することによって、最小限のレジスタの退避・回復を行うものとして説明している。しかし、計算機システム内で実行されるジョブが1つである場合は、この最小限のレジスタの退避・回復も不要である。

#### 【0244】

実行されるジョブが1つである状態を検出するために、計算機システム内にジョブ単一実行フラグ2308（図4）を備え、ジョブの起動及び終結に伴って、実行中のジョブの数

50



が1となった場合は、スケジューラがジョブ単一実行フラグを設定する。ジョブ単一実行フラグを参照し、起動中のジョブ数が1であるときは、並列演算プロセッサ群のレジスタの内容の退避・回復を行わず、並列ジョブの多重実行のための一切のオーバーヘッドを生じさせずに、並列演算プロセッサ群のスケジューリングを可能にする。

【0245】

最初に、図30に基づいて、レジスタの退避処理について説明する。

終了通知に引き続いて、ステップS3001において、スケジューラは、並列演算プロセッサ群上のジョブを切り替えるときに、ジョブ単一実行フラグ2308（図4）が設定されているかを調べる。設定されていれば、ステップS3002において、レジスタの退避は行わず、ステップS3003において、レジスタが未退避であることをジョブ管理情報中のレジスタ未退避フラグ2407（図5）を設定することで記録する。

10

【0246】

単一実行フラグが設定されていないときは、ステップS3004において、規定の最小限のレジスタのみを退避し、さらに、ステップS3005において、最小限のレジスタのみが退避されたことをレジスタ退避範囲フラグ2408（図5）を設定して記録する。

【0247】

次に、図31をに基づいて、レジスタの回復処理について説明する。

ステップS3101において、スケジューラは、ジョブ単一実行フラグ2308（図4）が設定されているかを調べる。設定されていれば、ステップS3102において、レジスタの回復は行わず、ステップS3103において、ジョブ管理情報中のレジスタ未退避フラグ2407（図5）を解除する。

20

【0248】

単一実行フラグが設定されていないときは、ステップS3104において、切り替え先の新ジョブのレジスタ退避範囲フラグ2408（図5）が設定されているかを調べる。設定されていれば、ステップS3105において、必要なレジスタのみ回復し、さらに、ステップS3106において、新ジョブのレジスタ退避範囲フラグ2408を解除する。設定されていなければ、ステップS3107において、全てのレジスタを回復する。

【0249】

以上説明したように、本実施形態によれば、複数のジョブが同時に並列処理部分を実行しなければ、1ジョブの逐次処理部分の実行中に遊休する並列演算プロセッサ群を他のジョブの並列処理部分の実行に割り当てることができ、複数のジョブを並列的に実行できるため、計算機システム内の全プロセッサの稼働率を高めることができる。

30

【0250】

また、並列処理部分を実行する並列演算プロセッサ群は、従来の並列専用機構とは異なり、逐次演算プロセッサと同等の複数個の演算プロセッサによって論理的に構成しているため、計算機システム内の演算プロセッサの逐次及び並列演算プロセッサへの割り当て比率を動的に変更することができる。従って、実行対象ジョブの逐次処理と並列処理の比率に応じて計算機システム内の演算プロセッサの構成比率を変更することができるので、計算機システムの稼働率を高めることが可能となる。

【0251】

40

また、同期範囲指定機構を用いて、同期範囲内の演算プロセッサによる並列処理を実行し、同期範囲外の演算プロセッサからの同期要求に対して、要求元の演算プロセッサ上に割り込みを発生させ、また、同期成立前に同期範囲外に設定し直された演算プロセッサ上に割り込みを発生させてスケジューラを起動して待ち状態にし、起動要求を並列実行待ちキューに入れるようにしているため、1時期には、1つの逐次演算プロセッサに対して並列演算プロセッサ群を容易に割り当てることができる。

【0252】

また、並列処理部分の起動時に実行される同期到達通知命令が自らを実行した演算プロセッサに対応する同期範囲指定機構の中のレジスタに対応するビットを調べることで、容易に、自分が並列演算プロセッサ群の使用が可能か否かの場合分けに従った処理を行うこと

50

によって、排他的に並列演算プロセッサ群を割り付けることが可能である。

【0253】

また、並列処理部分の処理の終結時には、並列演算プロセッサ群は終了通知を発行し、これを契機にして他の並列実行待ちジョブの並列処理部分の実行に切り替える。この際、並列処理部分の実行が終了しているため、並列演算プロセッサのレジスタの退避が大幅に削減されることによってオーバーヘッドの少ない切り替えを実現することができる。従って、ジョブの変更を高速に行うことができる。

【0254】

また、後続する逐次処理部分の実行時間が一定時間より長い場合にのみ終了通知を発行するようにしているため、1並列ジョブの逐次処理時間が、他の並列ジョブの並列処理部分の実行に十分な長さをもっている時のみ、並列演算プロセッサの切り替えが行われ、並列演算プロセッサ群の遊休時間に、他の並列ジョブの並列処理部分を割り当てることが可能になるため、スケジューリングのオーバーヘッドを少なくして、計算機システム全体の効率を向上できる。

10

【0255】

また、スケジューラは並列プロセッサ群上のみで起動されて、ジョブ切り替えの処理を行うようにしているため、この処理は、他のジョブの逐次処理部分の処理の実行に影響を及ぼすことなく行われるため、ジョブ切り替えのオーバーヘッドを隠蔽できる。

【0256】

また、並列処理部分の終了を並列演算プロセッサ群の割り当て替えの契機として用いることによって、切りのよいタイミングでジョブを切り替え、退避・回復するレジスタの数を大幅に削減することが可能になり、ジョブ切り替えのオーバーヘッドを低減できる。

20

【0257】

また、同期タイムアウトレジスタを用いて、予め指定した時間内に同期ポイントに達しない並列演算プロセッサが存在したとき、その同期処理は失敗したものとして、同期失敗を知らせる割り込みが演算プロセッサ上に発生し、これを契機にスケジューラが起動され、エラー処理を行える。

【0258】

また、ジョブの切り替えに際しては、演算プロセッサ間データ授受領域を実記憶空間上に用意し、これらを演算プロセッサ間データ授受領域対応つけフィールドを用いて同一の仮想記憶空間上の番地に対応づけ、ジョブの切り替えに際して、両者の対応づけを変更することによって、演算プロセッサ間データ授受領域の状態を退避・回復する必要がなくなり、並列演算機構の割り当て変更を高速に行うことができる。

30

【0259】

また、並列ジョブは階級分けされ、構成変更情報設定フィールドに設定された構成変更情報に基づいて階級を切り替え、それぞれの階級に割り付けられた時間内で階級に属するジョブをその階級の実行方式に従って処理することができる。

【0260】

また、各階級毎にスケジューリングのためのジョブスケジューリングキューを備えており、ジョブは自分の属する階級のジョブスケジューリングキューに入れるようにしているため、各階級毎に、処理の順番に応じた処理が行える。

40

【0261】

また、同一階級に属するジョブは、並列実行待ちキューを用いて同一階級のプログラム間で共有される並列演算プロセッサ群を切り替えて処理を進めることができる。

【0262】

また、並列演算プロセッサ群上での待ち状態の発生を契機として、ジョブ切り替えのためのスケジューリングを実行できる。

【0263】

また、複数の並列演算プロセッサ群を備えることにより、複数の並列処理を並行して実行できる。

50

## 【 0 2 6 4 】

また、演算プロセッサ間の同期処理は、ソフトウェア的に実現することも可能である。

## 【 0 2 6 5 】

また、並列処理の終了時に引き続く逐次処理部分が長いと予めプログラムの構造として明らかな時には、動的な終了通知の発行によって、並列演算プロセッサ群の割り当て変更ができる。

## 【 0 2 6 6 】

## 【 発明の効果 】

本発明によれば、並列ジョブ多重スケジューリング方法及び装置における複数の並列化されたジョブの同時実行時に、システム内の全プロセッサの稼働率を向上することができる。

10

## 【 図面の簡単な説明 】

【 図 1 】 本発明の一実施形態による計算機システムの全体構成図である。

【 図 2 】 本発明の一実施形態による同期範囲指定機構の構成図である。

【 図 3 】 本発明の一実施形態による共有記憶装置中に用意されるデータの説明図である。

【 図 4 】 本発明の一実施形態における共有記憶装置上の管理情報の構成図である。

【 図 5 】 本発明の一実施形態において用いるジョブ管理情報の構成図である。

【 図 6 】 本発明の一実施形態において用いる階級管理情報の構成図である。

【 図 7 】 本発明の一実施形態によって並列ジョブを複数の演算プロセッサを用いて実行する場合の 1 方式の概略図である。

20

【 図 8 】 本発明の一実施形態による並列ジョブの多重実行の概略説明図である。

【 図 9 】 本発明の一実施形態による並列ジョブ多重スケジューリング装置のシステム構成図である。

【 図 1 0 】 本発明の一実施形態による並列ジョブ多重スケジューリング方法における並列処理部分の起動・終結時の実行について説明するフローチャートである。

【 図 1 1 】 本発明の一実施形態による並列ジョブの多重実行の概略説明図である。

【 図 1 2 】 本発明の一実施形態による並列ジョブ多重スケジューリング方法における並列演算プロセッサ群のプログラムへの割り当て制御方式の説明図である。

【 図 1 3 】 本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるデータ授受領域の管理方法の説明図である。

30

【 図 1 4 】 本発明の一実施形態による図 9 で示した並列演算プロセッサ群上で起動されるスケジューラ S C H 2 ( 図 1 2 ) の全体の動作の概略を示すフローチャートである。

【 図 1 5 】 図 1 4 に示した処理の細部を説明するフローチャートである。

【 図 1 6 】 図 1 4 に示した処理の細部を説明するフローチャートである。

【 図 1 7 】 図 1 4 に示した処理の細部を説明するフローチャートである。

【 図 1 8 】 図 1 4 に示した処理の細部を説明するフローチャートである。

【 図 1 9 】 図 1 4 に示した処理の細部を説明するフローチャートである。

【 図 2 0 】 図 1 4 に示した処理の細部を説明するフローチャートである。

【 図 2 1 】 図 1 4 に示した処理の細部を説明するフローチャートである。

40

【 図 2 2 】 本発明の一実施形態による並列ジョブ多重スケジューリング方法により演算プロセッサの割り付けを動的変更する例の説明図である。

【 図 2 3 】 本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるジョブ階級切り替え時のスケジューラの動作を示すフローチャートである。

【 図 2 4 】 本発明の一実施形態による並列ジョブ多重スケジューリング方法による待ち状態の発生を契機としたスケジューリングの説明図である。

【 図 2 5 】 本発明の一実施形態による並列ジョブ多重スケジューリング方法における待ち状態の発生によるスケジューリングの処理を説明するフローチャートである。

【 図 2 6 】 本発明の一実施形態による並列ジョブ多重スケジューリング装置における複数の並列演算プロセッサ群を備えた計算機システムの動作の説明図である。

50

【図 27】 本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるソフトウェア的同期処理を説明するフローチャートである。

【図 28】 本発明の一実施形態による複数の並列演算プロセッサ群を備えた計算機システムにおける同期範囲指定レジスタの構成図である。

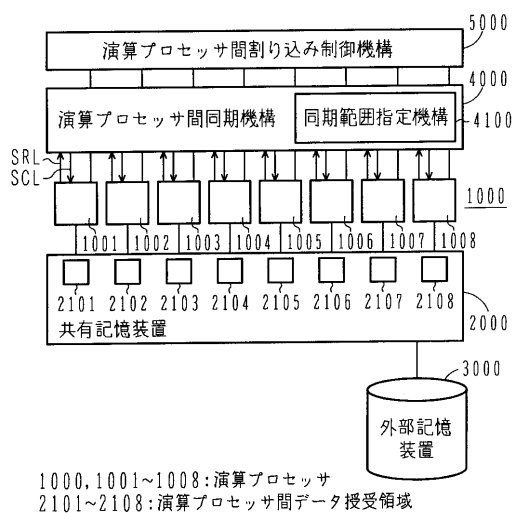
【図 29】 本発明の一実施形態による並列ジョブ多重スケジューリング方法における終了通知の動的判断処理を示すフローチャートである。

【図 30】 本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるジョブが単一で実行されているときのレジスタ退避時の動作を示すフローチャートである。

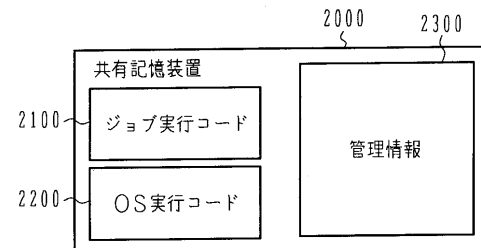
【図 31】 本発明の一実施形態による並列ジョブ多重スケジューリング方法におけるジョブが単一で実行されているときのレジスタ回復時の動作を示すフローチャートである。

10

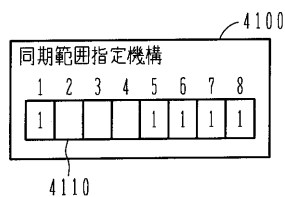
【図 1】



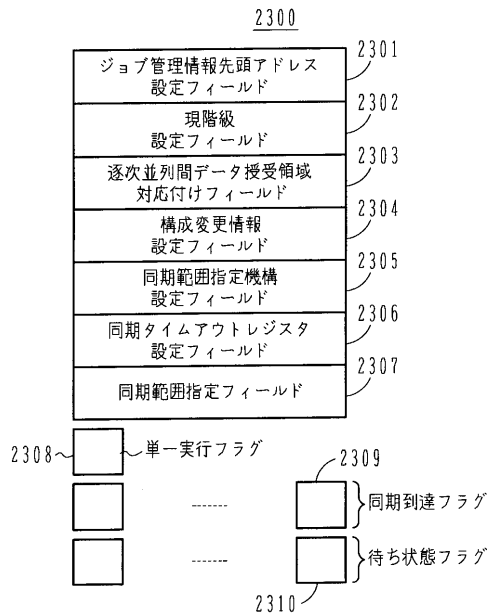
【図 3】



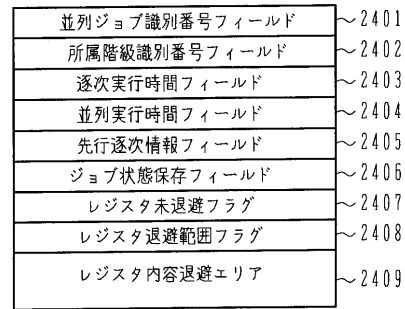
【図 2】



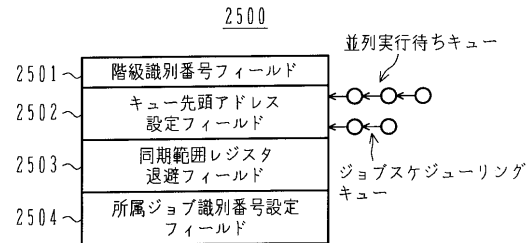
【図 4】



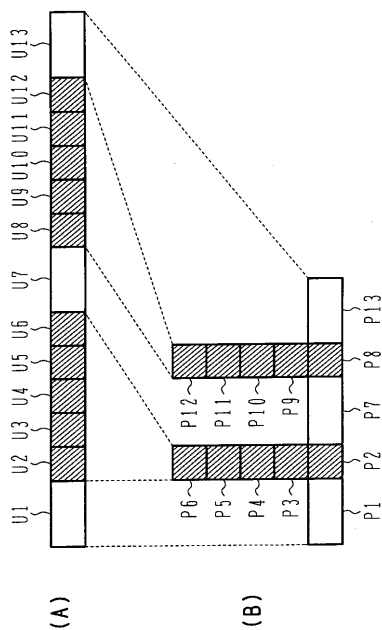
【図 5】



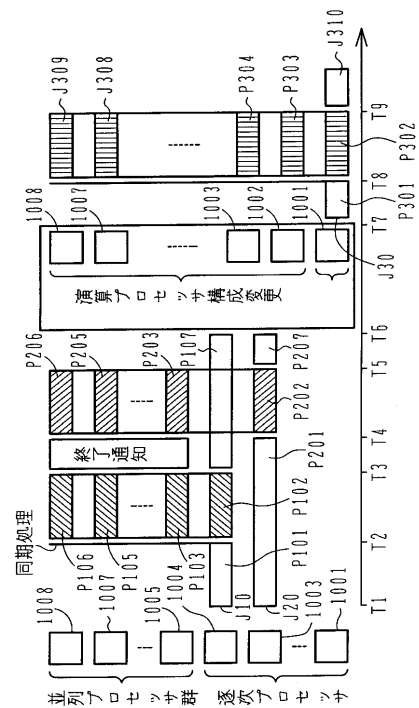
【図 6】



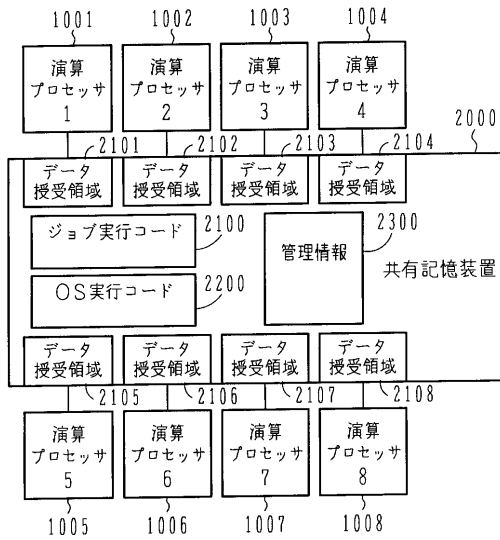
【図 7】



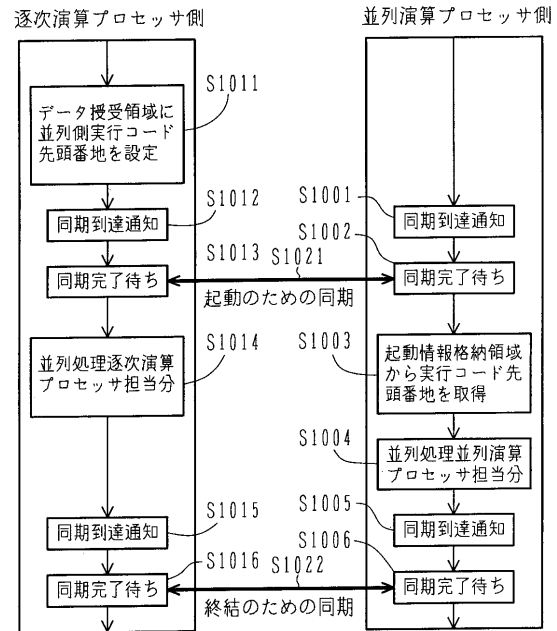
【図 8】



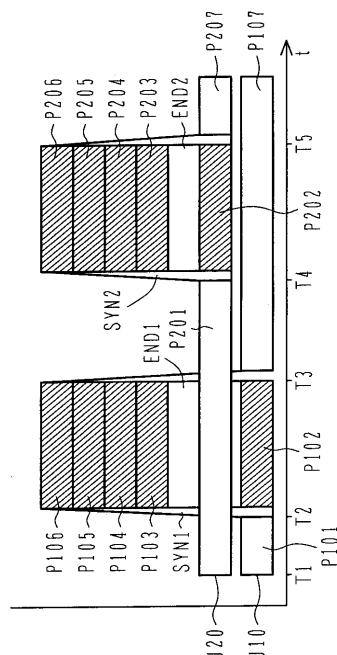
【図 9】



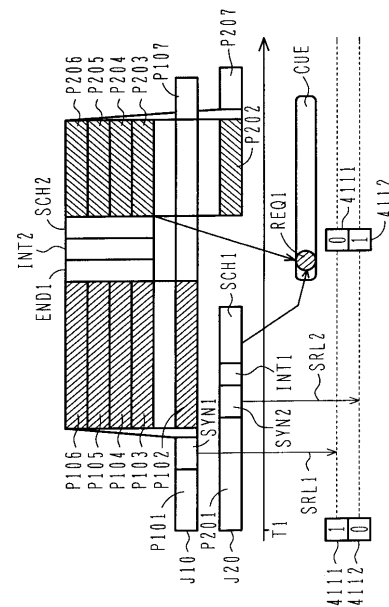
【図 10】



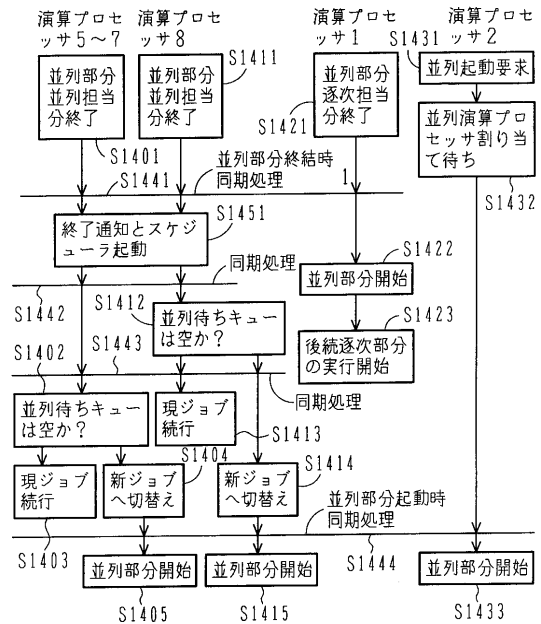
【図 11】



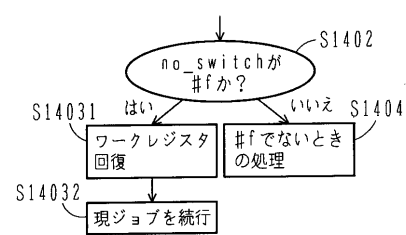
【図 12】



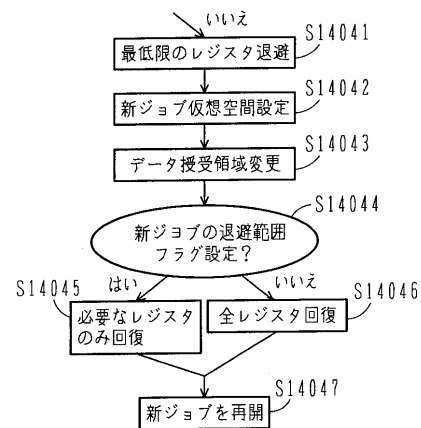
【 図 1 4 】



【 図 1 8 】



【 図 1 9 】

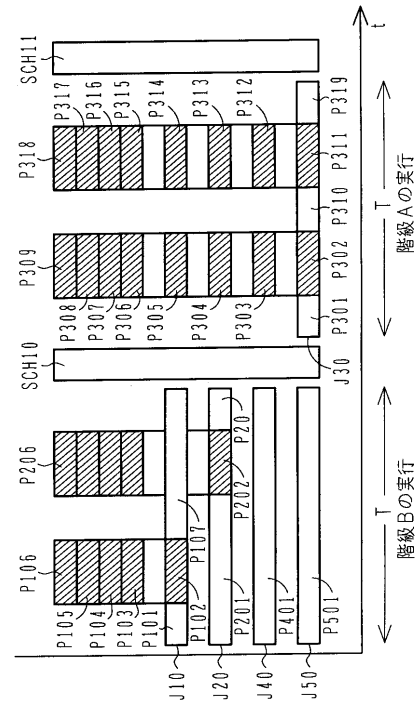


```

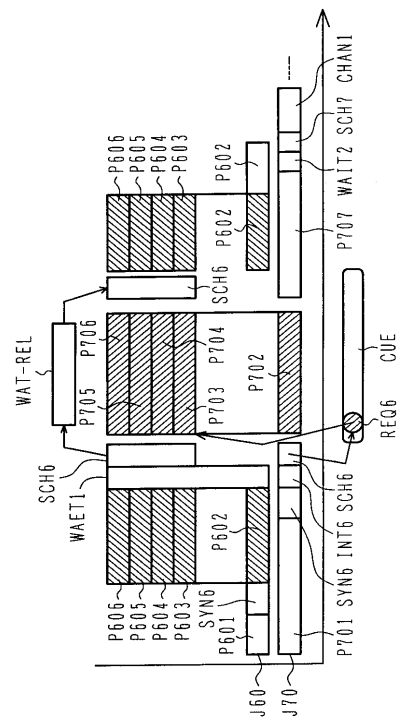
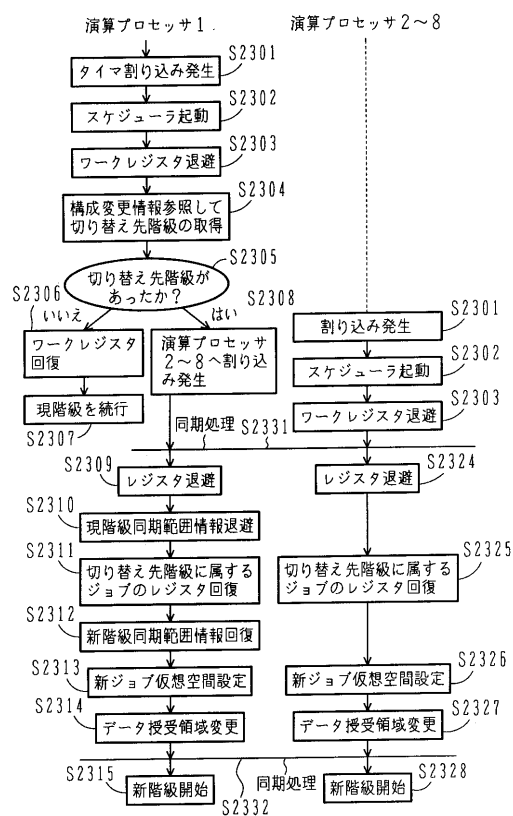
graph TD
    S14131[ワークレジスタ回復] --> S14132[現ジョブを続行]

```

【 ㄨ 2 2 】

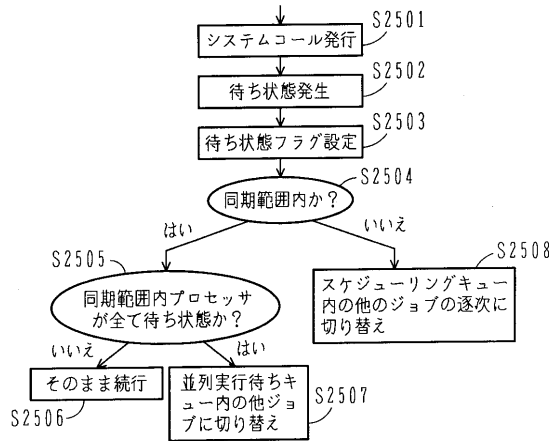


【 図 2 4 】

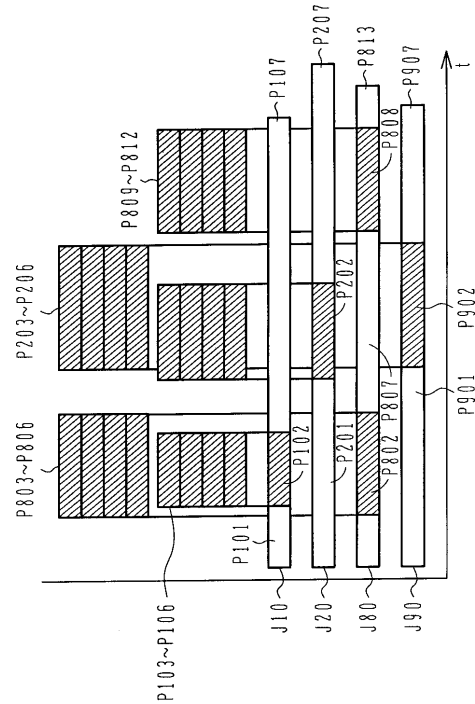




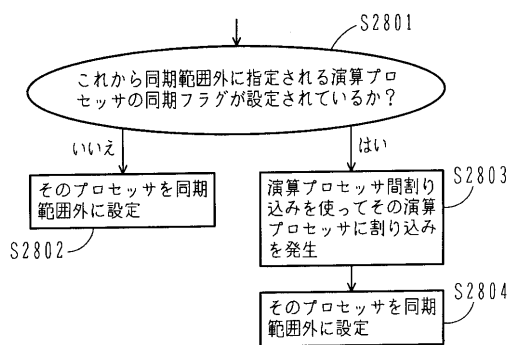
【図 25】



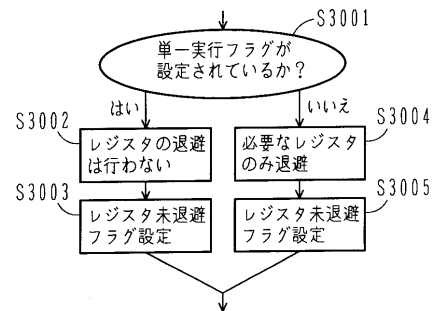
【図 26】



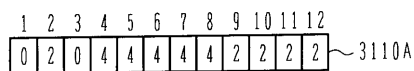
【図 27】



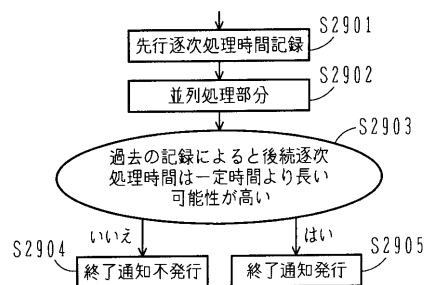
【図 30】



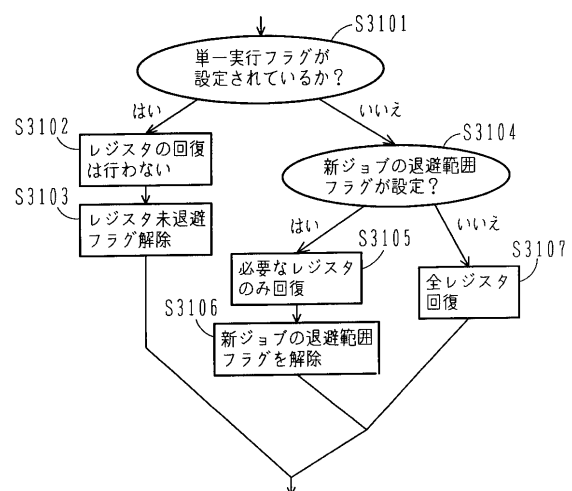
【図 28】



【図 29】



【図 31】



---

フロントページの続き

- (72)発明者 熊崎 裕之  
神奈川県横浜市戸塚区戸塚町5030番地  
本社内 株式会社 日立製作所 ソフトウェア開発
- (72)発明者 助川 直伸  
東京都国分寺市東恋ヶ窪一丁目280番地 株式会社 日立製作所 中央研究所内
- (72)発明者 中島 恵  
神奈川県横浜市戸塚区戸塚町5030番地  
本社内 株式会社 日立製作所 ソフトウェア開発
- (72)発明者 深川 正一  
神奈川県秦野市堀山下1番地  
事業部内 株式会社 日立製作所 汎用コンピュータ

審査官 殿川 雅也

(56)参考文献 特開平02-144657(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 9/46- 9/54

G06F 15/16- 15/177