

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5176558号  
(P5176558)

(45) 発行日 平成25年4月3日(2013.4.3)

(24) 登録日 平成25年1月18日(2013.1.18)

(51) Int.Cl. F I  
G 0 6 F 9/50 (2006.01) G 0 6 F 9/46 4 6 5 A

請求項の数 6 (全 25 頁)

(21) 出願番号	特願2008-8355 (P2008-8355)	(73) 特許権者	000005223
(22) 出願日	平成20年1月17日 (2008.1.17)		富士通株式会社
(65) 公開番号	特開2009-169756 (P2009-169756A)		神奈川県川崎市中原区上小田中4丁目1番1号
(43) 公開日	平成21年7月30日 (2009.7.30)	(74) 代理人	100104190
審査請求日	平成22年8月20日 (2010.8.20)		弁理士 酒井 昭徳
		(72) 発明者	繁田 聡一
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	清水 智弘
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		審査官	清木 泰

最終頁に続く

(54) 【発明の名称】 分散処理プログラム、分散処理装置、および分散処理方法

(57) 【特許請求の範囲】

【請求項1】

複数の計算機に複数のジョブを分散処理させる分散処理装置を、  
 前記複数の計算機から前記ジョブの割当先を決定する決定手段、  
 前記割当先により計測された、前記割当先に先に送信されたジョブ群の処理要求が受信されてから前記先に送信されたジョブ群の実行が開始されるまでの待ち時間、および前記先に送信されたジョブ群の実行時間を受信する受信手段、  
 前記割当先に先に送信されたジョブ群の処理要求を送信してから、前記先に送信されたジョブ群の処理結果を前記割当先から受信するまでの経過時間を計測する計測手段、  
 前記受信手段によって受信された前記待ち時間と前記実行時間と、前記計測手段によって計測された前記経過時間と、前記割当先に先に送信されたジョブ群を形成するジョブの数とに基づいて、前記割当先との通信にかかる通信時間のパラメータを算出し、前記実行時間と前記先に送信されたジョブ群を形成するジョブの数とに基づいて、ジョブ単位の実行時間を算出し、前記割当先に割り当てるジョブ群の実行時間と前記通信時間との対応情報と、算出した前記通信時間のパラメータと前記算出したジョブ単位の実行時間とに基づいて、前記割当先に割り当てる前記ジョブの数を算出する算出手段、  
 前記算出手段によって算出された前記ジョブの数に基づいて、前記割当先に割り当てるジョブ群を生成する生成手段、  
 前記生成手段によって生成された前記ジョブ群の処理要求を、前記割当先に送信する送信手段、

10

20

として機能させることを特徴とする分散処理プログラム。

【請求項 2】

前記決定手段は、

前記各計算機の使用状態に基づいて、前記複数の計算機から前記割当先を決定することを特徴とする請求項 1 に記載の分散処理プログラム。

【請求項 3】

前記決定手段は、

前記各計算機の処理性能に基づいて、前記複数の計算機から前記割当先を決定することを特徴とする請求項 1 または 2 に記載の分散処理プログラム。

【請求項 4】

前記決定手段は、

前記各計算機が実行可能なジョブのジョブタイプに基づいて、前記複数の計算機から前記割当先を決定することを特徴とする請求項 1 ~ 3 のいずれか一つに記載の分散処理プログラム。

【請求項 5】

複数の計算機に複数のジョブを分散処理させる分散処理装置であって、

前記複数の計算機から前記ジョブの割当先を決定する決定手段と、

前記割当先により計測された、前記割当先に先に送信されたジョブ群の処理要求が受信されてから前記先に送信されたジョブ群の実行が開始されるまでの待ち時間、および前記先に送信されたジョブ群の実行時間を受信する受信手段と、

前記割当先に先に送信されたジョブ群の処理要求を送信してから、前記先に送信されたジョブ群の処理結果を前記割当先から受信するまでの経過時間を計測する計測手段と、

前記待ち時間と前記実行時間と、計測した前記経過時間と、前記割当先に先に送信されたジョブ群を形成するジョブの数とに基づいて、前記割当先との通信にかかる通信時間のパラメータを算出し、前記実行時間と前記先に送信されたジョブ群を形成するジョブの数とに基づいて、ジョブ単位の実行時間を算出し、前記割当先に割り当てるジョブ群の実行時間と前記通信時間との対応情報と算出した前記通信時間のパラメータと前記算出したジョブ単位の実行時間とに基づいて、前記割当先に割り当てる前記ジョブの数を算出する算出手段と、

算出した前記ジョブの数に基づいて、前記割当先に割り当てるジョブ群を生成する生成手段と、

生成した前記ジョブ群の処理要求を前記割当先に送信する送信手段と、

を備えることを特徴とする分散処理装置。

【請求項 6】

複数の計算機に複数のジョブを分散処理させる分散処理装置が、

前記複数の計算機から前記ジョブの割当先を決定し、

前記割当先により計測された、前記割当先に先に送信されたジョブ群の処理要求が受信されてから前記先に送信されたジョブ群の実行が開始されるまでの待ち時間、および前記先に送信されたジョブ群の実行時間を受信し、

前記割当先に先に送信されたジョブ群の処理要求を送信してから、前記先に送信されたジョブ群の処理結果を前記割当先から受信するまでの経過時間を計測し、

受信した前記待ち時間と前記実行時間と、計測した前記経過時間と、前記割当先に先に送信されたジョブ群を形成するジョブの数とに基づいて、前記割当先との通信にかかる通信時間のパラメータを算出し、前記実行時間と前記先に送信されたジョブ群を形成するジョブの数とに基づいて、ジョブ単位の実行時間を算出し、前記割当先に割り当てるジョブ群の実行時間と前記通信時間との対応情報と、算出した前記通信時間のパラメータと前記算出したジョブ単位の実行時間とに基づいて、前記割当先に割り当てる前記ジョブの数を算出し、

算出した前記ジョブの数に基づいて、前記割当先に割り当てるジョブ群を生成し、

生成した前記ジョブ群の処理要求を、前記割当先に送信する、

10

20

30

40

50

処理を実行することを特徴とする分散処理方法。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、マスタ計算機（以下、単に「マスタ」という）が複数のワーカ計算機（以下、単に「ワーカ」という）に処理群を分散処理させるグリッドコンピューティングにおける分散処理プログラム、分散処理装置、および分散処理方法に関する。

【背景技術】

【0002】

従来、ネットワークを介して通信可能なマスタ/ワーカ間でやり取りされるジョブの流れでは、まず、マスタがジョブおよびその処理に必要なデータをワーカに投入する。つぎに、投入されたワーカは、ジョブの処理を実行する。そして、そのワーカは、ジョブの処理結果をマスタに返す。マスタは、これらを複数のワーカに対しておこなうことにより、ジョブ全体をワーカに分散処理させている。

10

【0003】

一般に、マスタに入力されるジョブは、ジョブ単位の実行時間がマスタ/ワーカ間の通信遅延時間に比べて、十分に大きいジョブ（粗粒度なジョブ）であると想定されている。粗粒度なジョブを分散処理する技術は既実用化されており、例えば、大量のバッチジョブを分散処理するためのグリッド・ミドルウェアである System Walker CyberGRIP（登録商標）などがある。

20

【0004】

また、下記特許文献1には、複数のタスクから構成されるタスクグループ単位で処理の割り当てをおこなう技術が記載されている。具体的には、まず、上位管理装置から下位管理装置にタスクグループの処理を割り当てる。つぎに、下位管理装置からタスク実行装置にタスクグループに含まれるタスクを割り当て、タスク実行装置によりタスクを実行し、その実行結果を下位管理装置に送信する。そして、下位管理装置により、タスクの実行結果を収集し、その収集結果を上位管理装置に送信する。最後に、上位管理装置により、収集結果を集約して出力する。

【0005】

また、下記特許文献2には、マスタからワーカに割り当てるタスク量を、ワーカからマスタへのタスク要求に応じて動的に調整する技術が記載されている。具体的には、まず、マスタが各ワーカにて処理すべきタスクを割り当てる。このとき、最初に同一量のタスクを各ワーカに割り当て、所定時間内の各ワーカからのタスク要求回数に応じて、一度のタスク要求に対して割り当てるタスク量を変化させる。

30

【0006】

【特許文献1】特開2004-110318号公報

【特許文献2】特開平11-195007号公報

【発明の開示】

【発明が解決しようとする課題】

【0007】

40

上述した従来技術では、マスタに入力されたジョブをジョブ単位でワーカに投入することで、ジョブ全体をワーカに分散処理させている。これは、ジョブ単位の実行時間がマスタ/ワーカ間の通信時間に比べて十分に大きい場合に有効な手法である。ところが、ジョブ単位の実行時間がマスタ/ワーカ間の通信時間よりも短いジョブを大量に処理することが要求されるアプリケーションもある。

【0008】

しかしながら、実行時間の短いジョブをジョブ単位でワーカに投入すると、ジョブ処理中におけるマスタ/ワーカ間の通信時間やワーカのアイドル時間などのオーバーヘッドが顕在化してしまう。これでは、ネットワーク上のトラフィックが増大化し、ひいてはジョブ処理にかかる所要時間の長期化を招くという問題がある。

50

## 【 0 0 0 9 】

この発明は、上述した従来技術による問題点を解消するため、実行時間の短いジョブを束ねたジョブ群を適切に割り当てることにより、ジョブ処理中におけるマスタ/ワーカ間の通信トラフィックの低減を図り、効率的な分散処理を実現することができる分散処理プログラム、分散処理装置、および分散処理方法を提供することを目的とする。

## 【課題を解決するための手段】

## 【 0 0 1 0 】

上述した課題を解決し、目的を達成するため、この分散処理プログラム、分散処理装置、および分散処理方法は、通信可能なワーカ計算機群の中からジョブの割当先を決定し、決定された割当先のワーカ計算機の処理性能と、前記割当先との通信にかかる通信時間とに基づいて、前記割当先に割り当てる前記ジョブのジョブ数を算出し、算出されたジョブ数に基づいて、前記割当先に割り当てるジョブ群を生成し、生成されたジョブ群の処理要求を、前記割当先に送信することを要件とする。

10

## 【 0 0 1 1 】

この分散処理プログラム、分散処理装置、および分散処理方法によれば、ジョブ単位の実行時間がマスタ/ワーカ間の通信時間よりも短いジョブを束ねて、1つのジョブ群単位でワーカ群に分散処理させることができる。このとき、マスタ/ワーカ間の通信時間およびワーカの処理性能に応じたジョブ数で束ねることで、ワーカ間における終了時刻を平準化し、全体の所要時間の短縮を図ることができる。

20

## 【 0 0 1 2 】

また、この分散処理プログラム、分散処理装置、および分散処理方法において、前記割当先により計測された、前記処理要求よりも先に送信された一のジョブ群の処理要求が受信されてから当該一のジョブ群の実行が開始されるまでの待ち時間、および前記一のジョブ群の実行時間に関する情報を受信し、前記割当先に前記一のジョブ群の処理要求を送信してから、前記一のジョブ群の処理結果を前記割当先から受信するまでの経過時間を計測し、受信された待ち時間および実行時間に関する情報と、計測された経過時間とを用いて、前記通信時間を算出することとしてもよい。

## 【 0 0 1 3 】

この分散処理プログラム、分散処理装置、および分散処理方法によれば、割当先にジョブ群を割り当てることで得られた待ち時間、実行時間および経過時間に関する過去の計測値を用いて、マスタと割当先との間の通信時間を算出することができる。

30

## 【 0 0 1 4 】

また、この分散処理プログラム、分散処理装置、および分散処理方法において、さらに、前記一のジョブ群の実行時間に関する情報を用いて、前記処理性能を算出することとしてもよい。

## 【 0 0 1 5 】

この分散処理プログラム、分散処理装置、および分散処理方法によれば、割当先にジョブ群を割り当てることで得られた実行時間に関する過去の計測値を用いて、割当先の処理性能を算出することができる。

## 【 0 0 1 6 】

また、この分散処理プログラム、分散処理装置、および分散処理方法において、前記各ワーカ計算機の使用状態に基づいて、前記ワーカ計算機群の中から前記割当先を決定することとしてもよい。

40

## 【 0 0 1 7 】

この分散処理プログラム、分散処理装置、および分散処理方法によれば、管理下にあるワーカ群のうち、マスタから割り当てられたジョブ群を処理中、または、機能を停止中のワーカを排除することで、割当先候補を絞り込むことができる。

## 【 0 0 1 8 】

また、この分散処理プログラム、分散処理装置、および分散処理方法において、前記各ワーカ計算機の処理性能に基づいて、前記ワーカ計算機群の中から前記割当先を決定する

50

こととしてもよい。

【 0 0 1 9 】

この分散処理プログラム、分散処理装置、および分散処理方法によれば、ジョブ群を高速に処理可能なワーカを割当先に決定することができる。

【 0 0 2 0 】

また、この分散処理プログラム、分散処理装置、および分散処理方法において、前記各ワーカ計算機が実行可能なジョブのジョブタイプに基づいて、前記ワーカ計算機群の中から前記割当先を決定することとしてもよい。

【 0 0 2 1 】

この分散処理プログラム、分散処理装置、および分散処理方法によれば、ワーカWが受け入れ可能なジョブのジョブタイプを考慮して、ジョブの割当先を決定することで、ジョブ群を適切なワーカWに割り当てることができる。

10

【発明の効果】

【 0 0 2 2 】

この分散処理プログラム、分散処理装置、および分散処理方法によれば、実行時間の短いジョブを束ねたジョブ群を適切に割り当てることにより、ジョブ処理中におけるマスタ/ワーカ間の通信トラフィックの低減を図り、効率的な分散処理を実現することができるという効果を奏する。

【発明を実施するための最良の形態】

【 0 0 2 3 】

20

以下に添付図面を参照して、この分散処理プログラム、分散処理装置、および分散処理方法の好適な実施の形態を詳細に説明する。なお、本明細書において、分散処理装置とは、グリッドコンピューティングシステムを構成する計算機（マスタまたはワーカ）であり、分散処理プログラムとは、分散処理装置にインストールされたプログラムである。

【 0 0 2 4 】

（実施の形態1）

（グリッドコンピューティングシステムのシステム構成）

まず、本実施の形態にかかるグリッドコンピューティングシステム100のシステム構成について説明する。図1は、グリッドコンピューティングシステムおよび分散処理装置のシステム構成図である。図1において、グリッドコンピューティングシステム100は、インターネット、LAN、WANなどのネットワーク110を介して通信可能なマスタMとワーカW1～Wm群とから構成される。

30

【 0 0 2 5 】

各ワーカW1～Wmは、異なる処理能力を持つこととしてもよく、また、OSやハードウェア・アーキテクチャなど異なる構造を持つこととしてもよい。さらに、ネットワーク110の通信品質は、一定または画一的である必要もない。

【 0 0 2 6 】

このグリッドコンピューティングシステム100では、マスタMが、実行時間の短いジョブ（例えば、解析用プログラム）を束ねたジョブ群を生成し、そのジョブ群を適切なワーカW1～Wmに投入する。そして、ワーカW1～Wmが、投入されたジョブ群を実行し、その処理結果をマスタMに返す。このとき、ジョブの処理結果は、ジョブ単位で返すのではなく、ジョブ群単位で返す。

40

【 0 0 2 7 】

なお、マスタ/ワーカ間でやり取りされる各ジョブの実行時間は一様である。すなわち、各ジョブの粒度（処理量）は揃っている、あるいは、無視できる程度のばらつきである。さらに、各ジョブの実行時間は、ジョブ処理中におけるマスタ/ワーカ間の通信時間（転送時間）よりも短い時間である。

【 0 0 2 8 】

ここで処理対象となるジョブとしては、例えば、金融機関などのオプション・リスク計算において、モンテカルロ法を用いて確率的にシナリオをシミュレーションする際に数万

50

～数千万個処理することが要求される、実行時間が短い細粒度なジョブが挙げられる。

【 0 0 2 9 】

このような実行時間の短いジョブを大量に処理するケースでは、ジョブ単位でマスタMからワーカW1～Wmにジョブ投入すると、ジョブ処理中におけるマスタノワーカ間の通信時間やワーカW1～Wmのアイドル時間といったオーバーヘッドが顕在化してしまう。そこで、複数のジョブを束ねたジョブ群を1つの処理単位とする。

【 0 0 3 0 】

図2は、ワーカWにおけるジョブ処理過程の概要を示す説明図である。図2において、グラフ210は、実行時間の短いジョブを1つの処理単位として各ワーカW1～Wmに分散処理させる従来の形態のジョブ処理過程を示している。グラフ220は、実行時間の短いジョブを束ねて(図2では、3個)1つの処理単位として各ワーカW1～Wmに分散処理させる本実施の形態のジョブ処理過程を示している。

10

【 0 0 3 1 】

グラフ210に示すように、実行時間の短いジョブを1つの処理単位として分散処理させた場合には、ジョブ処理中における、マスタノワーカ間の通信時間やワーカW1～Wmがつぎのジョブを受け取るまでのアイドル時間などのオーバーヘッドが顕在化してしまう。

【 0 0 3 2 】

一方、グラフ220に示すように、複数のジョブを束ねて1つの処理単位として分散処理させた場合には、ジョブ処理中における、マスタノワーカ間の通信頻度および通信時間を削減し、上記オーバーヘッドを隠蔽する(全所要時間に対する通信時間およびワーカW1～Wmのアイドル時間の割合を小さくする)ことができる。

20

【 0 0 3 3 】

ここで、図1に示したワーカW1を例に挙げて、グリッドコンピューティング100におけるジョブの投入例を説明する。まず、マスタMにおいて、入力されたジョブJ1, J2, J3の3個を束ねたジョブ群JG1を生成する。つぎに、その割当先に決定されたワーカW1に生成されたジョブ群JG1を投入(ネットワーク110経由で送信)する。

【 0 0 3 4 】

このあと、ワーカW1において、マスタMから投入されたジョブ群JG1を実行し、ジョブ群JG1を構成するすべてのジョブJ1, J2, J3の実行が完了したあと、ジョブJ1, J2, J3の処理結果R1, R2, R3を束ねた処理結果RG1をマスタMに返す。

30

【 0 0 3 5 】

このように、複数のジョブJ1, J2, J3を束ねたジョブ群JG1を1つの処理単位として分散処理させることにより、ジョブ処理中におけるマスタMとワーカW1との間の通信時間を低減し、通信時間やアイドル時間などのオーバーヘッドを隠蔽する。

【 0 0 3 6 】

さらに、割当先のワーカW1で処理された処理結果は、すべてのジョブJ1, J2, J3の処理が完了したあとに、割当先のワーカW1からマスタMに返すことで、ジョブJ1, J2, J3の処理中におけるネットワーク110上のトラフィックの低減を図る。

40

【 0 0 3 7 】

(マスタMおよびワーカWのハードウェア構成)

つぎに、実施の形態1にかかるマスタMおよびワーカW1～Wmのハードウェア構成について説明する。なお、以降において、特に指定する場合を除いて「ワーカW1～Wm」を「ワーカW」と表記する。図3は、マスタMおよびワーカWのハードウェア構成を示すブロック図である。

【 0 0 3 8 】

図3において、マスタMおよびワーカWは、CPU301と、ROM302と、RAM303と、HDD(ハードディスクドライブ)304と、HD(ハードディスク)305と、FDD(フレキシブルディスクドライブ)306と、着脱可能な記録媒体の一例とし

50

てのFD（フレキシブルディスク）307と、ディスプレイ308と、I/F（インターフェース）309と、キーボード310と、マウス311と、スキャナ312と、プリンタ313とを備えている。また、各構成部は、バス300によってそれぞれ接続されている。

#### 【0039】

ここで、CPU301は、マスタMおよびワーカWの全体の制御を司る。ROM302は、ブートプログラムなどのプログラムを記録している。RAM303は、CPU301のワークウェアとして使用される。HDD304は、CPU301の制御にしたがってHD305に対するデータのリード/ライトを制御する。HD305は、HDD304の制御で書き込まれたデータを記憶する。

10

#### 【0040】

FDD306は、CPU301の制御にしたがってFD307に対するデータのリード/ライトを制御する。FD307は、FDD306の制御で書き込まれたデータを記憶したり、FD307に記憶されたデータをマスタMおよびワーカWに読み取らせたりする。

#### 【0041】

また、着脱可能な記録媒体として、FD307のほか、CD-ROM（CD-R、CD-RW）、MO、DVD（Digital Versatile Disk）、メモ리카ードなどであってもよい。ディスプレイ308は、カーソル、アイコンあるいはツールボックスをはじめ、文書、画像、機能情報などのデータを表示する。このディスプレイ308には、たとえば、CRT、TFT液晶ディスプレイ、プラズマディスプレイなどを採用

20

#### 【0042】

I/F309は、通信回線を通じてインターネットなどのネットワーク110に接続され、このネットワーク110を介して他の装置に接続される。そして、I/F309は、ネットワーク110と内部のインターフェースを司り、ネットワーク110からのデータの入出力を制御する。I/F309には、たとえばモデムやLANアダプタなどを採用することができる。

#### 【0043】

キーボード310は、文字、数字、各種指示などの入力のためのキーを備え、データの入力をおこなう。また、タッチパネル式の入力パッドやテンキーなどであってもよい。マウス311は、カーソルの移動や範囲選択、あるいはウィンドウの移動やサイズの変更などをおこなう。ポインティングデバイスとして同様の機能を備えるものであれば、トラックボールやジョイスティックなどであってもよい。

30

#### 【0044】

スキャナ312は、画像を光学的に読み取り、装置内に画像データを読み込む。なお、スキャナ312は、OCR機能を持たせてもよい。また、プリンタ313は、画像データや文書データを印刷する。プリンタ313には、たとえば、レーザプリンタやインクジェットプリンタなどを採用することができる。

#### 【0045】

（ワーカ管理テーブルの記憶内容）

40

ここで、ジョブの割当先となるワーカW1～WmのIPアドレスおよび使用状態を特定する場合に用いられるワーカ管理テーブルについて説明する。図4は、ワーカ管理テーブルの記憶内容を示す説明図である。図4において、ワーカ管理テーブル400には、ワーカ識別子、IPアドレスおよび状態がワーカW1～Wmごとに記憶されている。

#### 【0046】

状態とは、各ワーカW1～Wmの使用状態である。この状態は、マスタMからのジョブ投入前は「空き」となる。そして、マスタMからジョブが投入されると「空き」から「使用中」に変更される。また、ワーカW1～Wmから処理結果が返ってくると「使用中」から「空き」に変更される。さらに、ワーカW1～Wmの機能が停止している場合には「停止中」となる。このワーカ管理テーブル400は、図3に示したRAM303やHD30

50

5などの記憶部によりその機能を実現する。

【0047】

(スループットテーブルの記憶内容)

ここで、ワーカW1～Wmの処理性能を特定する場合に用いられるスループットテーブルについて説明する。図5は、スループットテーブルの記憶内容を示す説明図である。図5において、スループットテーブル500には、ワーカW1～Wmごとのスループット値が記憶されている。

【0048】

スループット値は、各ワーカW1～Wmの処理性能を表わす指標であり、単位時間当たり処理されたジョブ数を表現する。このスループット値は、各ワーカW1～WmのCPU使用率など(例えば、マスタMとは異なる他のコンピュータ装置から投入されたジョブを処理中)によって動的に変化する。

【0049】

また、未だジョブの割当先となっていないワーカW1～Wm(例えば、W2)のスループット値は記憶されていない(図5中「-」)。このスループットテーブル500は、図3に示したRAM303やHD305などの記憶部によりその機能を実現する。

【0050】

(分散処理装置の機能的構成)

つぎに、分散処理装置の機能的構成について説明する。まず、マスタMの機能的構成について説明する。図6は、マスタMの機能的構成を示すブロック図である。図6において、マスタMは、決定部601と、算出部602と、生成部603と、送信部604と、受信部605と、計測部606と、から構成されている。

【0051】

これら各機能601～606は、マスタMの記憶部に記憶された当該機能601～606に関するプログラムをCPUに実行させることにより、または、入出力I/Fにより、当該機能を実現することができる。また、各機能601～606からの出力データは上記記憶部に保持される。また、図6中矢印で示した接続先の機能は、接続元の機能からの出力データを記憶部から読み込んで、当該機能に関するプログラムをCPUに実行させるものとする。

【0052】

まず、決定部601は、ワーカW群の中からジョブの割当先を決定する機能を有する。具体的には、例えば、各ワーカWの使用状態に基づいて、ワーカW群の中から割当先を決定することとしてもよい。例えば、ワーカ管理テーブル400からワーカWごとの状態を読み出して、使用状態が「空き」のワーカW群の中から割当先を決定する。これは、管理下にあるワーカW群のうち、マスタMから割り当てられたジョブ群を処理中、または、機能を停止中のワーカWを排除することで、割当先候補を絞り込むことを意味する。

【0053】

また、各ワーカWの処理性能に基づいて、ワーカW群の中から割当先を決定することとしてもよい。例えば、スループットテーブル500からワーカWごとのスループット値Tを読み出して、スループット値Tが最大のワーカWを割当先に決定する。これは、ジョブ群を高速に処理可能なワーカWを割当先に決定することを意味する。さらに、使用状態が「空き」のワーカW群の中から、スループット値Tが最大のワーカWを割当先に決定してもよく、また、ランダムに割当先を決定してもよい。

【0054】

算出部602は、決定部601によって決定された割当先のワーカWの処理性能と、当該割当先との通信にかかる通信時間とに基づいて、割当先に割り当てるジョブのジョブ数を算出する機能を有する。具体的には、例えば、ジョブ数nのジョブを束ねたジョブ群の実行時間が、ジョブ処理中におけるマスタM/ワーカW間の通信時間の1～2倍程度になるようにジョブ数nを算出する。

【0055】

10

20

30

40

50

すなわち、マスタM/ワーカW間でのジョブ群のジョブ処理中における、通信時間やアイドル時間などのオーバーヘッドが顕在化することを回避する。なお、算出部602による算出処理の具体例については後述する。

**【0056】**

生成部603は、算出部602によって算出された算出結果に基づいて、割当先に割り当てるジョブ群を生成する機能を有する。具体的には、例えば、ジョブキューにある一連のジョブを算出部602によって算出されたジョブ数nで束ねることにより1つのジョブ群を生成する。このとき、ジョブ数n分のジョブがジョブキューにない場合は、ジョブキューにあるすべてのジョブを束ねることとなる。あるいは、待ち時間の上限を設定して、ジョブキューにジョブ数n分のジョブが溜まるのを待つこととしてもよい。

10

**【0057】**

各ジョブには固有のジョブID（例えば、図1に示したJ1, J2, J3）が割り付けられており、さらに、1つのジョブ群を構成する各ジョブには共通のジョブ群ID（例えば、JG1）を割り付ける。このようにジョブIDおよびジョブ群IDを割り付けることにより、複数のジョブからなるジョブ群を認識することができる。

**【0058】**

なお、ジョブキューとは、マスタM内の記憶部に構築されるデータ構造であり、マスタMに入力されたジョブをキューイングする機能を有する。また、分散処理対象となるジョブは、マスタMに直接入力することとしてもよく、また、ネットワーク110を介して外部のコンピュータ装置から取得することとしてもよい。

20

**【0059】**

送信部604は、生成部603によって生成されたジョブ群の処理要求を、割当先に送信する機能を有する。具体的には、例えば、決定部601によって決定されたワーカWのIPアドレスをワーカ管理テーブル400から読み出して、そのIPアドレスを宛先に設定することでジョブ群の処理要求を割当先に送信することができる。また、送信部604により処理要求が送信されると、図4に示したワーカ管理テーブル400の記憶内容のうち、割当先の状態が「空き」から「使用中」に書き換えられる。

**【0060】**

ここで、上記算出部602による算出処理の具体例について説明する。受信部605は、割当先により計測された、送信部604によって上記処理要求よりも先に送信された一のジョブ群の処理要求が受信されてから当該一のジョブ群の実行が開始されるまでの待ち時間、および一のジョブ群の実行時間に関する情報を受信する機能を有する。

30

**【0061】**

待ち時間とは、例えば、割当先において、一のジョブ群の処理要求の受信が検出されてから、一のジョブ群の準備（実行に必要な関数の呼び出しなど）が整うまでの時間である。実行時間とは、一のジョブ群の実行を開始してから、該一のジョブ群を構成するすべてのジョブの実行が完了するまでに要した時間である。受信部605によって受信された情報は、例えば、後述するパラメータテーブル700に記憶される。

**【0062】**

また、計測部606は、割当先に一のジョブ群の処理要求を送信してから、一のジョブ群の処理結果を割当先から受信するまでの経過時間を計測する機能を有する。具体的には、例えば、一のジョブ群の処理要求が送信された送信日時と、一のジョブ群の処理結果が受信された受信日時とから経過時間を計測することができる。計測部606によって計測された計測結果は、例えば、後述するパラメータテーブル700に記憶される。

40

**【0063】**

そして、算出部602は、受信部605によって受信された待ち時間および実行時間に関する情報と、計測部606によって計測された経過時間とを用いて、通信時間を算出することとしてもよい。すなわち、割当先にジョブ群を割り当てることで得られた待ち時間、実行時間および経過時間に関する過去の計測値を用いて、マスタと割当先との間の通信時間を算出する。

50

【 0 0 6 4 】

また、算出部 6 0 2 は、一のジョブ群の実行時間に関する情報を用いて、割当先の処理性能を算出することとしてもよい。すなわち、割当先にジョブ群を割り当てることで得られた実行時間に関する過去の計測値を用いて、割当先の処理性能を算出する。

【 0 0 6 5 】

ここで、算出部 6 0 2 による算出処理に用いられるパラメータテーブルについて説明する。パラメータテーブルは、例えば、マスタ M の管理下にあるワーカ W ごとに、ワーカ識別子と関連付けて保持されている。ここでは、あるワーカ W のパラメータテーブルを例に挙げて説明する。図 7 は、パラメータテーブルの記憶内容を示す説明図である。図 7 において、パラメータテーブル 7 0 0 には、マスタ / ワーカ間の通信時間に関するパラメータ  $k, c$  と、ワーカ W に割り当てられたジョブ群  $J G_1 \sim J G_p$  ごとにパラメータ情報 7 0 0 - 1 ~ 7 0 0 - p が記憶されている。

10

【 0 0 6 6 】

具体的には、パラメータ情報 7 0 0 - 1 ~ 7 0 0 - p は、パラメータ  $n, t, w, k, c$  を有している。パラメータ  $n$  は、ジョブ群に含まれているジョブの数を表わす値である。パラメータ  $t$  は、マスタ M において、ジョブ群の処理要求が送信されてからそのジョブ群の処理結果が受信されるまでの経過時間を表わす値である。パラメータ  $w$  は、ワーカ W において、ジョブ群の処理要求が受信されてからそのジョブ群の実行が開始されるまでの待ち時間を表わす値である。

20

【 0 0 6 7 】

パラメータ  $t$  は、ワーカ W において一のジョブ群の処理に要した実行時間を表わす値である。パラメータ  $k, c$  は、マスタ M / ワーカ W 間の通信時間に関する値である。ここで、マスタ M / ワーカ W 間の通信時間とは、例えば、マスタ M からワーカ W への処理要求の送信、およびワーカ W からマスタ M への処理結果の送信にかかる時間である。

【 0 0 6 8 】

具体的には、例えば、マスタ M / ワーカ W 間の通信時間は、下記式 ( 1 ) を用いて求めることができる。ただし、通信時間を  $K$ 、ジョブ群を構成するジョブ数を  $n$  とする。

【 0 0 6 9 】

$$K = k \times n + c \quad \dots ( 1 )$$

30

【 0 0 7 0 】

また、パラメータ  $k, c$  は、パラメータ  $t, w, n_x, n_y$  を用いて算出することができる。具体的には、例えば、パラメータ  $k, c$  は、下記式 ( 2 ) を用いて求めることができる。

【 0 0 7 1 】

$$t = w + \frac{K}{n} + ( k \times n + c ) \quad \dots ( 2 )$$

【 0 0 7 2 】

より詳細に説明すると、上記式 ( 2 ) を用いて、 $n_x, n_y$  である  $x, y$  (  $x, y$  は自然数 ) について、パラメータ  $k, c$  の連立方程式 ( 下記式 ( 3 ) および ( 4 ) ) をたてることで、パラメータ  $k, c$  を求めることができる。

【 0 0 7 3 】

$$t_x = w_x + \frac{K_x}{n_x} + ( k \times n_x + c ) \quad \dots ( 3 )$$

40

【 0 0 7 4 】

$$t_y = w_y + \frac{K_y}{n_y} + ( k \times n_y + c ) \quad \dots ( 4 )$$

【 0 0 7 5 】

このあと、パラメータテーブル 7 0 0 から、 $n_x, n_y$  である 2 つのパラメータ情報 7 0 0 - 1 ~ 7 0 0 - p を読み出す。そして、読み出したパラメータ情報 7 0 0 - 1 ~ 7 0 0 - p に含まれるパラメータ  $n_x, n_y, t_x, t_y, w_x, w_y, K_x, K_y$  を上記式 ( 3 ) および ( 4 ) に代入して、パラメータ  $k, c$  を求める。

【 0 0 7 6 】

なお、 $n_x, n_y$  である 2 つのパラメータ情報 7 0 0 - 1 ~ 7 0 0 - p を読み出す場合には、現在時刻から遡って最も直近の時刻に記憶された 2 つのパラメータ情報 7 0 0 - 1 ~

50

700 - pを読み出すこととしてもよい。

【0077】

また、パラメータ  $\tau$  , nを用いて、1ジョブ当たりの実行時間を表わすパラメータ  $\tau$  を算出することができる。具体的には、例えば、パラメータ  $\tau$  は、下記式(5)を用いて求めることができる。ただし、jは1からpの自然数である。

【0078】

【数1】

$$\tau = \frac{\sum_{j=1}^p \delta_j}{\sum_{j=1}^p n_j} \quad \dots(5) \quad 10$$

【0079】

そして、求めたパラメータ  $\tau$  , k , cを下記式(6)に代入することにより、割当先に割り当てる一のジョブ群として束ねるジョブ数nを算出することができる。ただし、sは、ジョブ数nのジョブ群の実行時間が通信時間Kのs倍であることを示すパラメータである。パラメータsの値は、任意に設定可能(例えば、s = 1)である。

【0080】

$$n \times \tau = s (k \times n + c) \quad \dots (6) \quad 20$$

【0081】

なお、上記式(6)を用いて算出されたジョブ数nがn > 0であった場合には、例えば、nに数%のゆらぎを持たせたものをジョブ数とすることとしてもよい。また、上記式(6)を用いて算出されたジョブ数nがn = 0であった場合には、例えば、予め設定されている数(例えば、n = 1)をジョブ数とすることとしてもよい。

【0082】

上記パラメータsの値は、例えば、図3に示したキーボード310やマウス311などをユーザが操作することで、任意に設定可能である。さらに、マスタMに入力されたジョブの数が多い場合には、パラメータsの値を大きくし(例えば、s = 3)、ジョブキューにあるジョブの数が少なくなるにつれて、パラメータsの値を小さくする(例えば、s = 1)こととしてもよい。 30

【0083】

また、図5に示したスループットテーブル500に記憶されているスループット値Tは、上記式(5)を用いて求めたパラメータ  $\tau$  の逆数(1 /  $\tau$ )、すなわち、単位時間当たりに処理可能なジョブ数によって表現することができる。

【0084】

さらに、パラメータテーブル700の記憶内容は、管理下にあるワーカWから処理結果を受信する都度、更新される。さらに、上記スループット値Tは、例えば、パラメータテーブル700の記憶内容が更新される都度、合わせて更新されることとなる。パラメータテーブル700を更新する更新処理については後述する。 40

【0085】

なお、パラメータk , cを求めるために必要となるパラメータ情報700 - 1 ~ 700 - pがパラメータテーブル700に記憶されていない場合には(例えば、初回の分散処理時)、予め設定されている初期値を用いることとしてもよい。つまり、上記パラメータ  $n_x$  ,  $n_y$  ,  $t_x$  ,  $t_y$  ,  $w_x$  ,  $w_y$  ,  $\tau_x$  ,  $\tau_y$  に相当する初期値を上記式(3)および(4)に代入して、パラメータk , cを求める。

【0086】

つぎに、ワーカWの機能的構成について説明する。図8は、ワーカWの機能的構成を示すブロック図である。図8において、ワーカWは、受信部801と、実行部802と、送信部803と、計測部804と、から構成されている。 50

## 【 0 0 8 7 】

これら各機能 8 0 1 ~ 8 0 4 は、ワーカ W の記憶部に記憶された当該機能 8 0 1 ~ 8 0 4 に関するプログラムを CPU 3 0 1 に実行させることにより、または、入出力 I / F により、当該機能を実現することができる。また、各機能 8 0 1 ~ 8 0 4 からの出力データは上記記憶部に保持される。また、図 8 中矢印で示した接続先の機能は、接続元の機能からの出力データを記憶部から読み込んで、当該機能に関するプログラムを CPU に実行させるものとする。

## 【 0 0 8 8 】

まず、受信部 8 0 1 は、ジョブ群の処理要求をマスタ M から受信する機能を有する。また、実行部 8 0 2 は、マスタ M から割り当てられたジョブ群の処理を実行する機能を有する。ジョブ群の処理は、ワーカ W の CPU によって実行され、その処理結果はワーカ W 内の記憶部に保持される。

10

## 【 0 0 8 9 】

送信部 8 0 3 は、実行部 8 0 2 によってジョブ群を構成するすべてのジョブの処理の実行が完了した結果、ジョブ群の処理結果をマスタ M に送信する機能を有する。なお、ジョブ群の実行完了は、例えば、以下の手順で検出することができる。具体的には、実行部 8 0 2 は、ジョブ群の処理開始時に、予めジョブ群に含まれるジョブ数 n を送信部 8 0 3 に通知する。

## 【 0 0 9 0 】

さらに、実行部 8 0 2 は、ジョブ群に含まれる個々のジョブの実行完了ごとに実行結果を送信部 8 0 3 へ通知する。送信部 8 0 3 は、実行部 8 0 2 から通知されるジョブの実行結果の個数を数えることによって該ジョブ群に含まれるすべてのジョブの実行が完了したことを検出することができる。すなわち、実行部 8 0 2 から予め通知されたジョブ数 n に等しい個数の実行結果を受け取ったときに、ジョブ群の実行が完了したことを検出する。

20

## 【 0 0 9 1 】

そして、送信部 8 0 3 は、すべてのジョブの実行が完了したことが検出された場合、ジョブ群の処理結果をマスタ M に送信する。このとき、ジョブ群の処理要求から特定される IP アドレスを宛先に設定することで処理結果をマスタ M に送信することができる。

## 【 0 0 9 2 】

計測部 8 0 4 は、受信部 8 0 1 によってジョブ群の処理要求が受信されてから、実行部 8 0 2 によってジョブ群の実行が開始されるまでの待ち時間を計測する機能を有する。具体的には、例えば、処理要求が受信された受信日時と、ジョブ群の実行が開始された開始日時とから待ち時間を計測することができる。

30

## 【 0 0 9 3 】

また、計測部 8 0 4 は、実行部 8 0 2 によるジョブ群の実行時間を計測する機能を有する。具体的には、例えば、ジョブ群の実行が開始された開始日時と、ジョブ群の実行が完了した完了日時とから実行時間を計測することができる。

## 【 0 0 9 4 】

また、送信部 8 0 3 は、計測部 8 0 4 によって計測された上記待ち時間または / および実行時間に関する情報をマスタ M に送信する機能を有する。この情報は、マスタ M の受信部 6 0 5 で受信され、算出部 6 0 2 による算出処理に用いられる。

40

## 【 0 0 9 5 】

(分散処理装置の分散処理手順)

つぎに、分散処理装置の分散処理手順について説明する。まず、マスタ M における分散処理手順について説明する。図 9 は、マスタ M における分散処理手順の一例を示すフローチャートである。図 9 のフローチャートにおいて、まず、ジョブキューにジョブがあるか否かを判断する (ステップ S 9 0 1 )。

## 【 0 0 9 6 】

ここで、ジョブがマスタ M に入力されるのを待って (ステップ S 9 0 1 : N o )、ジョブキューにジョブがある場合 (ステップ S 9 0 1 : Y e s )、決定部 6 0 1 により、管理

50

下にあるワーカW群の中からジョブの割当先を決定する割当先決定処理を実行する（ステップS902）。

【0097】

このあと、算出部602により、ステップS902において決定された割当先のワーカWの処理性能と、割当先との通信にかかる通信時間とに基づいて、割当先に割り当てる一のジョブ群として束ねるジョブ数nを算出する（ステップS903）。

【0098】

つぎに、ジョブキューへのジョブ入力を待つ待ち時間の上限値を設定する（ステップS904）。このあと、ステップS903において算出されたジョブ数nがジョブキューにある全ジョブ数N以下であるか否かを判断する（ステップS905）。

10

【0099】

ここで、ジョブ数n 全ジョブ数Nの場合（ステップS905：Yes）、生成部603により、割当先に割り当てるジョブ数nのジョブ群を生成する（ステップS906）。そして、送信部604により、ステップS906において生成されたジョブ群の処理要求を割当先に送信して（ステップS907）、本フローチャートによる一連の処理を終了する。

【0100】

また、ステップS905において、ジョブ数n > 全ジョブ数Nの場合（ステップS905：No）、ステップS904において設定された待ち時間の上限値に達したか否かを判断する（ステップS908）。ここで、上限値に達していない場合（ステップS908：No）は、ステップS905に戻る。

20

【0101】

一方、上限値に達した場合（ステップS908：Yes）、生成部603により、ジョブキューにある全ジョブ数Nのジョブを用いて、割当先に割り当てるジョブ群を生成する（ステップS906）。

【0102】

なお、ステップS904において設定される待ち時間の上限値は、例えば、予め規定された値（例えば、1分）を設定することとしてもよく、また、複数の上限値候補の中から、ジョブ数nに応じた上限値を自動選択することとしてもよい。

【0103】

具体的には、例えば、ジョブ数nと上限値とを関連付けて表わす上限値テーブルを参照することで、ステップS903において算出されたジョブ数に応じた上限値を選択し、待ち時間の上限値に設定することとしてもよい。この上限値テーブルは、例えば、マスタM内の記憶部に予め保持されている。

30

【0104】

また、ステップS904において待ち時間の上限値が設定されると、待ち時間の計測が開始され、このあと、ステップS908において上限値が設定されてからの経過時間が待ち時間の上限値に達したか否かを判断することとしてもよい。

【0105】

つぎに、図9に示したステップS902における割当先決定処理の詳細な処理手順について説明する。図10は、割当先決定処理手順の一例を示すフローチャートである。図10において、まず、ワーカ管理テーブル400に基づいて、未使用（空き）のワーカWがあるか否かを判断する（ステップS1001）。ここで、未使用のワーカWがなかった場合には（ステップS1001：No）、いずれかのワーカWが使用可能となるのを待つ。

40

【0106】

一方、未使用のワーカWがあった場合（ステップS1001：Yes）、スループットテーブル500に基づいて、未使用のワーカWの中から、スループット値Tが最大のワーカWを選択する（ステップS1002）。そして、ステップS1002において選択されたワーカWを割当先に決定し（ステップS1003）、図9に示したステップS903に移行する。

50

## 【0107】

つぎに、ワーカWにおける分散処理手順について説明する。図11は、ワーカWにおける分散処理手順の一例を示すフローチャートである。図11のフローチャートにおいて、まず、受信部801により、ジョブ群の処理要求をマスタMから受信したか否かを判断する(ステップS1101)。

## 【0108】

ここで、処理要求を受信するのを待って(ステップS1101:No)、受信した場合(ステップS1101:Yes)、計測部804により、ジョブ群の処理要求が受信されてから、ジョブ群の実行が開始されるまでの待ち時間の計測を開始する(ステップS1102)。

10

## 【0109】

このあと、実行部802により、マスタMから割り当てられたジョブ群の処理を実行する(ステップS1103)。このとき、計測部804により、待ち時間の計測を終了するとともに、ジョブ群の実行時間の計測を開始する(ステップS1104)。

## 【0110】

このあと、ジョブ群を構成するすべてのジョブの実行が完了したことが検出されると(ステップS1105)、計測部804により、ジョブ群の実行時間の計測を終了する(ステップS1106)。そして、送信部803により、計測部804による計測結果を含むジョブ群の処理結果をマスタMに送信して(ステップS1107)、本フローチャートによる一連の処理を終了する。

20

## 【0111】

つぎに、マスタMにおける、図7に示したパラメータテーブル700の記憶内容を更新する更新処理手順について説明する。図12は、パラメータテーブルの更新処理手順の一例を示すフローチャートである。図12のフローチャートにおいて、まず、受信部605により、ワーカWからジョブ群の処理結果を受信したか否かを判断する(ステップS1201)。

## 【0112】

ここで、ジョブ群の処理結果を受信するのを待って(ステップS1201:No)、受信された場合(ステップS1201:Yes)、ジョブ群の処理結果からパラメータ $n$ ,  $w$ を取得するとともに、計測部606による計測結果からパラメータ $t$ を取得する(ステップS1202)。このとき、上記ジョブ群のジョブ群IDから特定される計測結果からパラメータ $t$ を取得する。

30

## 【0113】

このあと、ワーカ管理テーブル400に基づいて、更新対象のエントリがあるか否かを判断する(ステップS1203)。これは、ステップS1201において受信された処理結果が、マスタMの管理下にあるワーカWからのものか否かを判断するステップである。

## 【0114】

ここで、更新対象のエントリがある場合(ステップS1203:Yes)、そのエントリのパラメータテーブル700に記憶済みのパラメータがあるか否かを判断する(ステップS1204)。これは、ステップS1201における処理結果の受信が、初回の受信であるか否かを判断するステップである。

40

## 【0115】

ここで、記憶済みのパラメータがある場合(ステップS1204:Yes)、算出部602により、パラメータテーブル700の記憶内容、およびステップS1202において取得されたパラメータ $n$ ,  $w$ ,  $t$ を用いて、パラメータ $k$ ,  $c$ を算出する(ステップS1205)。

## 【0116】

このあと、ステップS1202において取得されたパラメータ $n$ ,  $w$ ,  $t$ 、およびステップS1205において算出されたパラメータ $k$ ,  $c$ を用いてパラメータテーブル700の記憶内容を更新して(ステップS1206)、本フローチャートによる一連の処理

50

を終了する。

【0117】

一方、ステップS1204において、記憶済みのパラメータがない場合（ステップS1204：No）、ステップS1202において取得されたパラメータ $n$ 、 $w$ 、 $t$ を用いてパラメータテーブル700の記憶内容を更新して（ステップS1206）、本フローチャートによる一連の処理を終了する。

【0118】

また、ステップS1203において、更新対象のエントリがない場合（ステップS1203：No）、更新処理を強制終了するエラー処理を実行して（ステップS1208）、本フローチャートによる一連の処理を終了する。さらに、ステップS1206において、パラメータテーブル700の記憶内容の更新に合わせて、ワーカ管理テーブル400の記憶内容（スループット値 $T$ ）を更新することとしてもよい。

10

【0119】

以上説明したように、実施の形態1によれば、ジョブ単位の実行時間がマスタ $M$ /ワーカ $W$ 間の通信時間よりも短いジョブを束ねて、ジョブ群単位でワーカ群に分散処理させることができる。このとき、割当先のワーカ $W$ の処理性能に応じたジョブ数で束ねることで、ワーカ $W$ 間における終了時刻を平準化することができる。

【0120】

このように、実行時間の短いジョブを束ねたジョブ群を適切にワーカ $W$ に割り当てることで、ジョブ処理中におけるマスタ $M$ /ワーカ $W$ 間の通信トラフィックの低減を図り、効率的な分散処理を実現することができる。

20

【0121】

ここで、複数のジョブを、処理性能が異なるワーカ $W_a$ 、 $W_b$ 、 $W_c$ に分散処理させた場合の効果を具体的に説明する。図13は、ワーカ $W$ ごとの所要時間を示すガントチャートである。図13において、(1)は、ワーカ $W_a$ 、 $W_b$ 、 $W_c$ の処理性能（スループット値 $T$ ）を考慮することなく、一定の数（ここでは、3個）のジョブを束ねた場合の所要時間を表わしている。

【0122】

(2)は、本実施の形態で説明したように、ワーカ $W_a$ 、 $W_b$ 、 $W_c$ の処理性能（スループット値 $T$ ）に応じて、束ねるジョブ数を変化させた場合の所要時間を表わしている。(1)では、処理性能が低いワーカ $W_b$ に律速され、全ジョブの終了時刻が遅くなっている。

30

【0123】

一方、(2)では、ワーカ $W_a$ 、 $W_b$ 、 $W_c$ の処理性能に応じて求めたジョブ数からなるジョブ群を割り当てているため、各ワーカ $W_a$ 、 $W_b$ 、 $W_c$ のジョブ処理にかかる所要時間が平準化され、全ジョブの終了時刻が(1)の場合に比べて短縮されている。

【0124】

(実施の形態2)

つぎに、実施の形態2にかかる分散処理装置について説明する。実施の形態2では、図1に示したグリッドコンピューティングシステム100を構成するワーカ $W_1 \sim W_m$ ごとに受け入れ可能（実行可能）なジョブタイプが異なる場合（同一の場合を含む）の分散処理について説明する。なお、実施の形態1において説明した箇所と同一箇所については、同一符号を付して図示および説明を省略する。

40

【0125】

まず、実施の形態2にかかるマスタ $M$ の機能的構成について説明する。決定部601は、各ワーカが実行可能なジョブのジョブタイプに基づいて、ワーカ群の中から割当先を決定する機能を有する。ジョブタイプとは、例えば、ジョブの実行時に呼び出される関数名で分類されるジョブの種別である。

【0126】

より具体的には、例えば、ジョブの実行時に起動されるアプリケーションによってジョ

50

ブタイプを分類することができる。すなわち、各ワーカWが持つ機能によっては、マスタMから投入されるジョブを実行できる場合と、実行できない場合とがある。そこで、マスタMが、各ワーカWが実行可能なジョブのジョブタイプを判断し、ジョブの適切な割当先を決定する。

**【0127】**

ここで、実施の形態2にかかるマスタMの分散処理手順の概要について説明する。図14は、分散処理手順の概要を示す説明図である。図14において、ジョブタイプの異なる3つのジョブJ<sub>a</sub>、J<sub>b</sub>、J<sub>c</sub>がマスタMに入力されたとする。ここでは、ジョブJ<sub>a</sub>のジョブタイプをA、ジョブJ<sub>b</sub>のジョブタイプをB、ジョブJ<sub>c</sub>のジョブタイプをCと表記する。

10

**【0128】**

この場合、各ジョブJ<sub>a</sub>、J<sub>b</sub>、J<sub>c</sub>のジョブタイプを判断して、ジョブタイプごとに分類されたジョブキューA、B、Cに対応するジョブJ<sub>a</sub>、J<sub>b</sub>、J<sub>c</sub>を配置する。ここでは、ジョブJ<sub>a</sub>がジョブキューAに配置され、ジョブJ<sub>b</sub>がジョブキューBに配置され、ジョブJ<sub>c</sub>がジョブキューCに配置される。

**【0129】**

このあと、決定部601は、各ワーカWが実行可能なジョブのジョブタイプに基づいて、ワーカW群の中から各ジョブキューA、B、Cに配置されたジョブJ<sub>a</sub>、J<sub>b</sub>、J<sub>c</sub>の割当先を決定することとなる。

**【0130】**

20

さらに、決定部601は、あるジョブタイプのジョブを実行可能なワーカWが複数存在する場合には、複数のワーカWの中からスループット値が最大のワーカWを割当先に決定することとしてもよい。このとき、例えば、ジョブタイプAのジョブの割当先を決定する場合には、スループット値が最大のワーカWを決定し、ジョブタイプBのジョブの割当先を決定する場合には、スループット値が2以上のワーカWをランダムに決定するなどのポリシーを設定することとしてもよい。

**【0131】**

ここで、各ワーカWが実行可能なジョブのジョブタイプを特定する場合に用いられるワーカ管理テーブルについて説明する。図15は、ワーカ管理テーブルの記憶内容を示す説明図である。図15において、ワーカ管理テーブル1500には、ワーカW<sub>1</sub>～W<sub>m</sub>ごとに、IPアドレス、状態および実行可能タイプが記憶されている。実行可能タイプとは、各ワーカWが実行可能なジョブのジョブタイプである。

30

**【0132】**

つぎに、ワーカW<sub>1</sub>～W<sub>m</sub>のジョブタイプごとの処理性能を特定する場合に用いられるスループットテーブルについて説明する。図16は、スループットテーブルの記憶内容を示す説明図である。図16において、スループットテーブル1600には、ワーカW<sub>1</sub>～W<sub>m</sub>ごとに、実行可能なジョブのジョブタイプごとのスループット値が記憶されている。

**【0133】**

ワーカW<sub>i</sub>を例に挙げると、ジョブタイプAのジョブに関するスループット値T<sub>ia</sub>、ジョブタイプBのジョブに関するスループット値T<sub>ib</sub>、およびジョブタイプCのジョブに関するスループット値T<sub>ic</sub>を有している。

40

**【0134】**

ここで、決定部601による割当先決定処理の具体例について説明する。まず、決定部601は、ワーカ管理テーブル1500からワーカWごとの状態および実行可能タイプを読み出して、ジョブキューにあるジョブを実行可能かつ使用状態が「空き」であるワーカWをワーカW群の中から特定する。そして、スループットテーブル1600から特定されたワーカWのスループット値を読み出し、スループット値が最大のワーカWを割当先に決定する。

**【0135】**

(分散処理装置の分散処理手順)

50

つぎに、実施の形態 2 にかかる分散処理装置の分散処理手順について説明する。なお、分散処理装置の分散処理手順のうち、決定部 601 による割当先決定処理（図 9 で示したステップ S 902 に相当）以外は、実施の形態 1 と同様のため説明を省略する。図 17 は、割当先決定処理手順の他の一例を示すフローチャートである。

【0136】

図 17 において、まず、ジョブキューにあるジョブのジョブタイプを取得する（ステップ S 1701）。このあと、ワーカ管理テーブル 1500 に基づいて、ステップ S 1701 において取得されたジョブタイプのジョブを実行可能なワーカ W があるか否かを判断する（ステップ S 1702）。

【0137】

ここで、実行可能なワーカ W がある場合（ステップ S 1702：Yes）、ワーカ管理テーブル 1500 に基づいて、上記ジョブタイプのジョブを実行可能なワーカ W のうち、未使用（空き）のワーカ W があるか否かを判断する（ステップ S 1703）。ここで、未使用のワーカ W がなかった場合には（ステップ S 1703：No）、いずれかのワーカ W が使用可能となるのを待つ。

【0138】

一方、未使用のワーカ W があった場合（ステップ S 1703：Yes）、スループットテーブル 1600 に基づいて、上記未使用のワーカ W の中から、スループット値 T が最大のワーカ W を選択する（ステップ S 1704）。そして、選択されたワーカ W を割当先に決定し（ステップ S 1705）、図 9 に示したステップ S 903 に移行する。

【0139】

また、ステップ S 1702 において、実行可能なワーカ W がなかった場合には（ステップ S 1702：No）、入力されたジョブが実行不可能である旨を示すメッセージを出力するエラー処理を実行して（ステップ S 1706）、処理を終了する。

【0140】

なお、実施の形態 2 にかかるマスタ M における分散処理は、ワーカ M 内に構築されたジョブキューごとに実行することとしてもよい。また、実施の形態 1 で説明したパラメータテーブル 700（図 7 参照）は、各ワーカ W が実行可能なジョブのジョブタイプごとに作成することとしてもよい。これにより、各ワーカ W の処理性能（スループット値）を、ジョブタイプごとに算出することができる。

【0141】

以上説明したように、実施の形態 2 によれば、ワーカ W が受け入れ可能なジョブのジョブタイプを考慮して、ジョブの割当先を決定することで、ジョブ群を適切なワーカ W に割り当てることができる。これにより、割り当てられたジョブ群を実行できない場合のエラー処理や、他のワーカ W へのジョブ群の再割り当てにかかる処理などを削減し、マスタ M / ワーカ W 間の通信トラフィックの低減を図ることができる。

【0142】

また、ワーカ W に入力されたジョブをジョブタイプごとに設けたジョブキューに配置することで、同じ関数を読み出すジョブを束ねることとなり、ワーカ W のキャッシュヒット率を向上させることができる。

【0143】

なお、本実施の形態で説明した分散処理方法は、予め用意されたプログラムをパーソナル・コンピュータやワークステーションなどのコンピュータで実行することにより実現することができる。このプログラムは、ハードディスク、フレキシブルディスク、CD-ROM、MO、DVD などのコンピュータで読み取り可能な記録媒体に記録され、コンピュータによって記録媒体から読み出されることによって実行される。またこのプログラムは、インターネットなどのネットワーク（伝送媒体）を介してコンピュータに配布することが可能な形態であってもよい。

【0144】

上述した実施の形態に関し、さらに以下の付記を開示する。

10

20

30

40

50

## 【 0 1 4 5 】

( 付記 1 ) 通信可能なワーカ計算機群に複数のジョブを分散処理させるマスタ計算機を、  
前記ワーカ計算機群の中から前記ジョブの割当先を決定する決定手段、  
前記決定手段によって決定された割当先のワーカ計算機の処理性能と、前記割当先との  
通信にかかる通信時間とに基づいて、前記割当先に割り当てる前記ジョブのジョブ数を算  
出する算出手段、  
前記算出手段によって算出された算出結果に基づいて、前記割当先に割り当てるジョブ  
群を生成する生成手段、  
前記生成手段によって生成されたジョブ群の処理要求を、前記割当先に送信する送信手  
段、  
として機能させることを特徴とする分散処理プログラム。

10

## 【 0 1 4 6 】

( 付記 2 ) 前記マスタ計算機を、  
前記割当先により計測された、前記送信手段によって前記処理要求よりも先に送信され  
た一のジョブ群の処理要求が受信されてから当該一のジョブ群の実行が開始されるまでの  
待ち時間、および前記一のジョブ群の実行時間に関する情報を受信する受信手段、  
前記割当先に前記一のジョブ群の処理要求を送信してから、前記一のジョブ群の処理結  
果を前記割当先から受信するまでの経過時間を計測する計測手段として機能させ、  
前記算出手段は、  
前記受信手段によって受信された待ち時間および実行時間に関する情報と、前記計測手  
段によって計測された経過時間とを用いて、前記通信時間を算出することを特徴とする付  
記 1 に記載の分散処理プログラム。

20

## 【 0 1 4 7 】

( 付記 3 ) 前記算出手段は、  
さらに、前記一のジョブ群の実行時間に関する情報を用いて、前記処理性能を算出する  
ことを特徴とする付記 2 に記載の分散処理プログラム。

## 【 0 1 4 8 】

( 付記 4 ) 前記決定手段は、  
前記各ワーカ計算機の使用状態に基づいて、前記ワーカ計算機群の中から前記割当先を  
決定することを特徴とする付記 1 ~ 3 のいずれか一つに記載の分散処理プログラム。

30

## 【 0 1 4 9 】

( 付記 5 ) 前記決定手段は、  
前記各ワーカ計算機の処理性能に基づいて、前記ワーカ計算機群の中から前記割当先を  
決定することを特徴とする付記 1 ~ 4 のいずれか一つに記載の分散処理プログラム。

## 【 0 1 5 0 】

( 付記 6 ) 前記決定手段は、  
前記各ワーカ計算機が実行可能なジョブのジョブタイプに基づいて、前記ワーカ計算機  
群の中から前記割当先を決定することを特徴とする付記 1 ~ 5 のいずれか一つに記載の分  
散処理プログラム。

## 【 0 1 5 1 】

( 付記 7 ) マスタ計算機と通信可能なワーカ計算機を、  
前記マスタ計算機から割り当てられたジョブ群の処理を実行する実行手段、  
前記マスタ計算機から前記ジョブ群の処理要求を受信する受信手段、  
前記受信手段によって前記処理要求が受信されてから、前記実行手段によって前記ジョ  
ブ群の実行が開始されるまでの待ち時間を計測する計測手段、  
前記計測手段によって計測された待ち時間に関する情報を前記マスタ計算機に送信する  
送信手段、  
として機能させることを特徴とする分散処理プログラム。

40

## 【 0 1 5 2 】

( 付記 8 ) 前記計測手段は、

50

前記実行手段による前記ジョブ群の実行時間を計測し、  
前記送信手段は、  
前記計測手段によって計測された実行時間に関する情報を前記マスタ計算機に送信することを特徴とする付記 7 に記載の分散処理プログラム。

【 0 1 5 3 】

(付記 9) 通信可能なワーカ計算機群に複数のジョブを分散処理させる分散処理装置であって、

前記ワーカ計算機群の中から前記ジョブの割当先を決定する決定手段と、

前記決定手段によって決定された割当先のワーカ計算機の処理性能と、当該割当先との通信にかかる通信時間とに基づいて、前記割当先に割り当てる前記ジョブのジョブ数を算出する算出手段と、

前記算出手段によって算出された算出結果に基づいて、前記割当先に割り当てるジョブ群を生成する生成手段と、

前記生成手段によって生成されたジョブ群の処理要求を、前記割当先に送信する送信手段と、

を備えることを特徴とする分散処理装置。

【 0 1 5 4 】

(付記 10) マスタ計算機から割り当てられたジョブ群の処理を実行する実行手段と、

前記マスタ計算機から前記ジョブ群の処理要求を受信する受信手段と、

前記受信手段によって前記処理要求が受信されてから、前記実行手段によって前記ジョブ群の実行が開始されるまでの待ち時間を計測する計測手段と、

前記計測手段によって計測された待ち時間に関する情報を前記マスタ計算機に送信する送信手段と、

を備えることを特徴とする分散処理装置。

【 0 1 5 5 】

(付記 11) 通信可能なワーカ計算機群に複数のジョブを分散処理させる分散処理方法であって、

前記ワーカ計算機群の中から前記ジョブの割当先を決定する決定工程と、

前記決定工程によって決定された割当先のワーカ計算機の処理性能と、当該割当先との通信にかかる通信時間とに基づいて、前記割当先に割り当てる前記ジョブのジョブ数を算出する算出工程と、

前記算出工程によって算出された算出結果に基づいて、前記割当先に割り当てるジョブ群を生成する生成工程と、

前記生成工程によって生成されたジョブ群の処理要求を、前記割当先に送信する送信工程と、

を含んだことを特徴とする分散処理方法。

【 0 1 5 6 】

(付記 12) マスタ計算機から割り当てられたジョブ群の処理を実行する実行工程と、

前記マスタ計算機から前記ジョブ群の処理要求を受信する受信工程と、

前記受信工程によって前記処理要求が受信されてから、前記実行工程によって前記ジョブ群の実行が開始されるまでの待ち時間を計測する計測工程と、

前記計測工程によって計測された待ち時間に関する情報を前記マスタ計算機に送信する送信工程と、

を含んだことを特徴とする分散処理方法。

【 図面の簡単な説明 】

【 0 1 5 7 】

【 図 1 】 グリッドコンピューティングシステムおよび分散処理装置のシステム構成図である。

【 図 2 】 ワーカ W におけるジョブ処理過程の概要を示す説明図である。

【 図 3 】 マスタ M およびワーカ W のハードウェア構成を示すブロック図である。

10

20

30

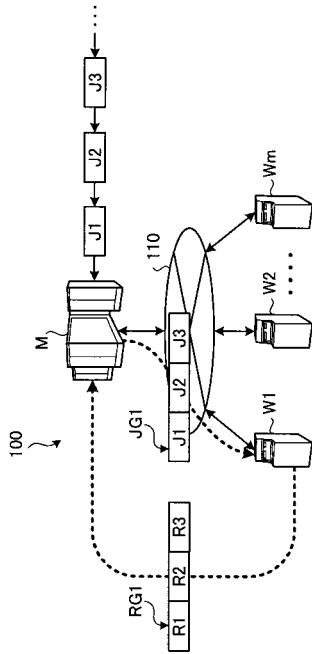
40

50

- 【図4】ワーカ管理テーブルの記憶内容を示す説明図(その1)である。
- 【図5】スループットテーブルの記憶内容を示す説明図(その1)である。
- 【図6】マスタMの機能的構成を示すブロック図である。
- 【図7】パラメータテーブルの記憶内容を示す説明図である。
- 【図8】ワーカWの機能的構成を示すブロック図である。
- 【図9】マスタMにおける分散処理手順の一例を示すフローチャートである。
- 【図10】割当先決定処理手順の一例を示すフローチャートである。
- 【図11】ワーカWにおける分散処理手順の一例を示すフローチャートである。
- 【図12】パラメータテーブルの更新処理手順の一例を示すフローチャートである。
- 【図13】ワーカWごとの所要時間を示すガントチャートである。 10
- 【図14】分散処理手順の概要を示す説明図である。
- 【図15】ワーカ管理テーブルの記憶内容を示す説明図(その2)である。
- 【図16】スループットテーブルの記憶内容を示す説明図(その2)である。
- 【図17】割当先決定処理手順の他の一例を示すフローチャートである。
- 【符号の説明】
- 【0158】
- 100 グリッドコンピューティングシステム
- 210, 220 グラフ
- 400, 1500 ワーカ管理テーブル
- 500, 1600 スループットテーブル 20
- 601 決定部
- 602 算出部
- 603 生成部
- 604 送信部
- 605 受信部
- 606 計測部
- 700 パラメータテーブル
- 700 - 1 ~ 700 - p パラメータ情報
- 801 受信部
- 802 実行部 30
- 803 送信部
- 804 計測部
- M マスタ
- W, W1 ~ Wm ワーカ

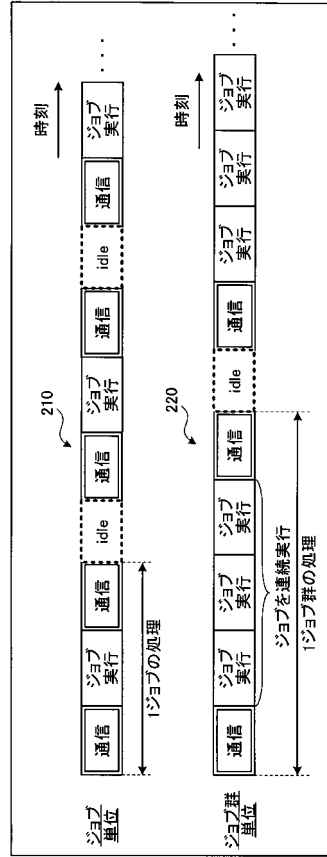
【図1】

グリッドコンピューティングシステムおよび分散処理装置のシステム構成図



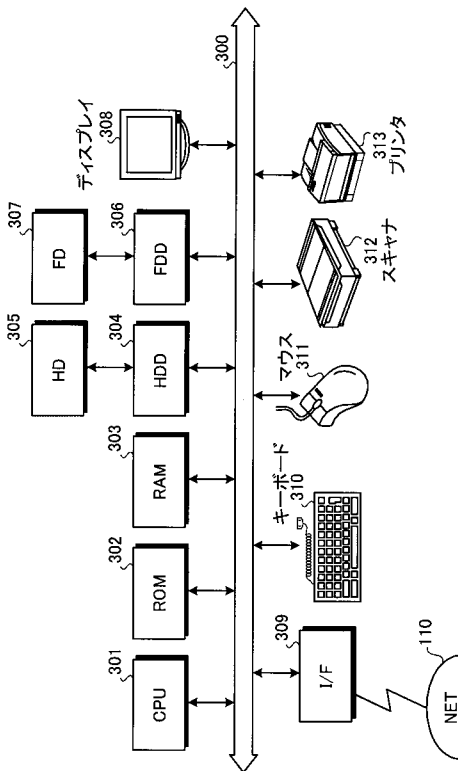
【図2】

ワーカWにおけるジョブ処理過程の概要を示す説明図



【図3】

マスタMおよびワーカWのハードウェア構成を示すブロック図



【図4】

ワーカ管理テーブルの記憶内容を示す説明図(その1)

ワーカ識別子	IPアドレス	状態
W1	xxx.yyy.z.t	使用中
W2	xxx.yyy.c.d	停止中
⋮	⋮	⋮
Wi	xxx.yyy.g.j	空き
⋮	⋮	⋮
Wm	xxx.yyy.k.l	使用中

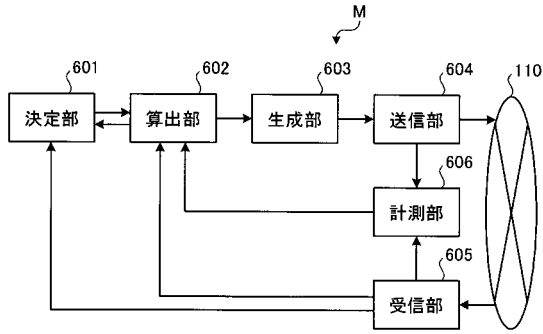
【図5】

スループットテーブルの記憶内容を示す説明図(その1)

ワーカ識別子	スループット値(T)
W1	T1
W2	—
⋮	⋮
Wi	Ti
⋮	⋮
Wm	Tm

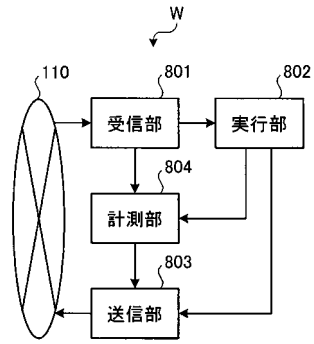
【図6】

マスタMの機能的構成を示すブロック図



【図8】

ワーカWの機能的構成を示すブロック図



【図7】

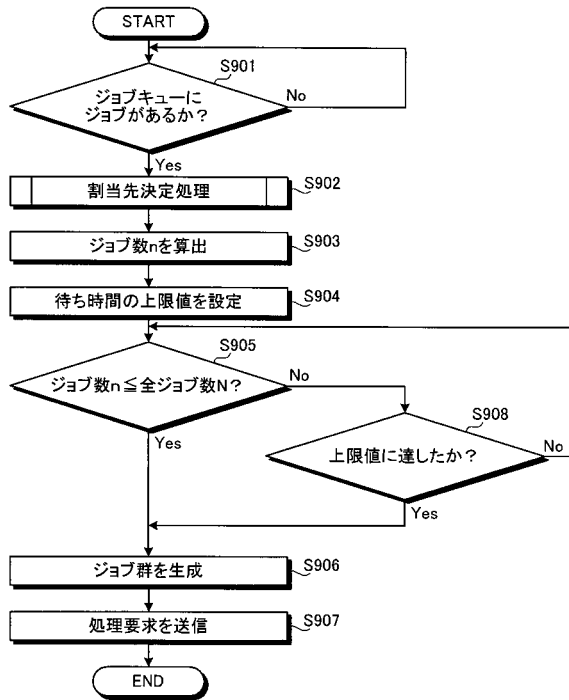
パラメータテーブルの記憶内容を示す説明図

700

通信時間にかかるパラメータ					
k	c				
ジョブ群ID	ジョブ群の実行時間にかかるパラメータ				
	ジョブ数n	経過時間t	待ち時間w	実行時間 $\delta$	
700-1	JG <sub>1</sub>	n <sub>1</sub>	t <sub>1</sub>	w <sub>1</sub>	$\delta_1$
700-2	JG <sub>2</sub>	n <sub>2</sub>	t <sub>2</sub>	w <sub>2</sub>	$\delta_2$
⋮	⋮	⋮	⋮	⋮	⋮
700-j	JG <sub>j</sub>	n <sub>j</sub>	t <sub>j</sub>	w <sub>j</sub>	$\delta_j$
⋮	⋮	⋮	⋮	⋮	⋮
700-p	JG <sub>p</sub>	n <sub>p</sub>	t <sub>p</sub>	w <sub>p</sub>	$\delta_p$

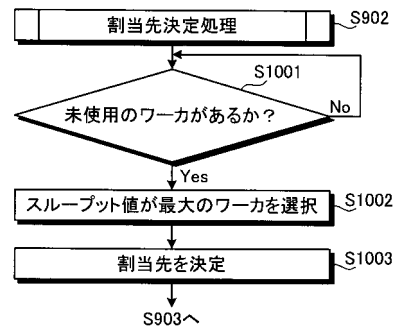
【図9】

マスタMにおける分散処理手順の一例を示すフローチャート



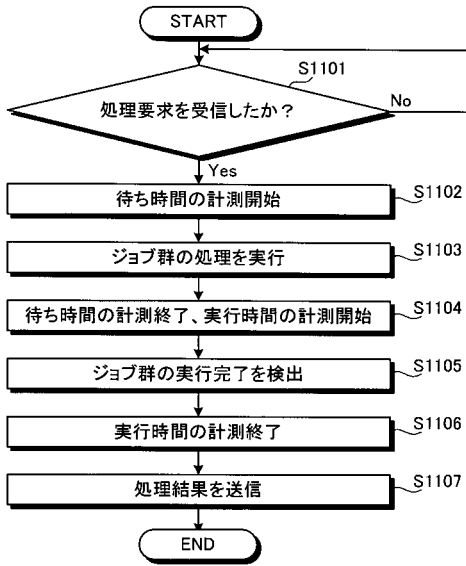
【図10】

割当先決定処理手順の一例を示すフローチャート



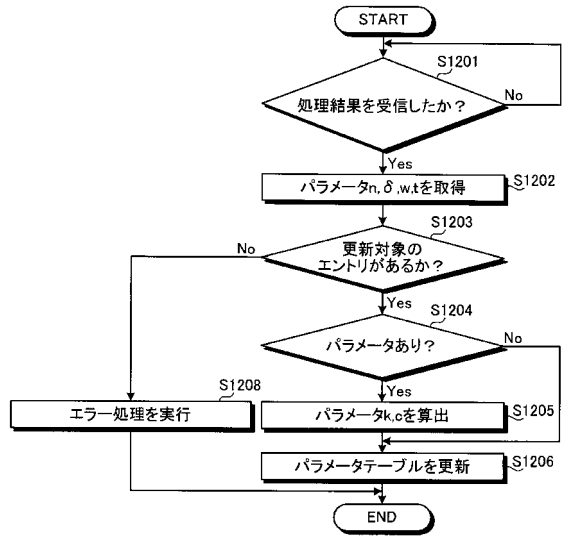
【図11】

ワーカWにおける分散処理手順の一例を示すフローチャート



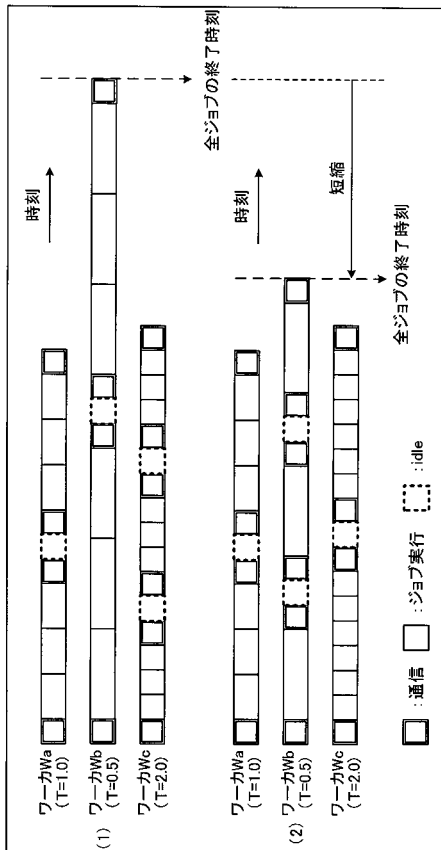
【図12】

パラメータテーブルの更新処理手順の一例を示すフローチャート



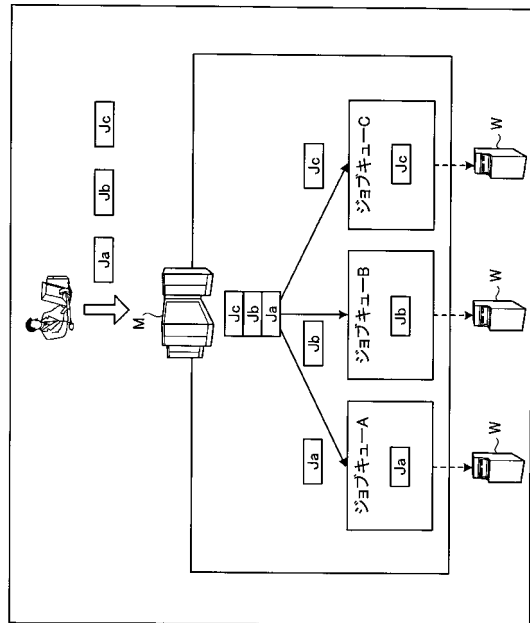
【図13】

ワーカWごとの所要時間を示すガントチャート



【図14】

分散処理手順の概要を示す説明図



【図15】

ワーカ管理テーブルの記憶内容を示す説明図(その2)

1500

ワーカ識別子	IPアドレス	状態	実行可能タイプ
W1	xxx.yyy.z.t	使用中	A,C
W2	xxx.yyy.o.d	停止中	A,B
⋮	⋮	⋮	⋮
Wi	xxx.yyy.g.j	空き	A,B,C
⋮	⋮	⋮	⋮
Wm	xxx.yyy.k.l	使用中	A

【図16】

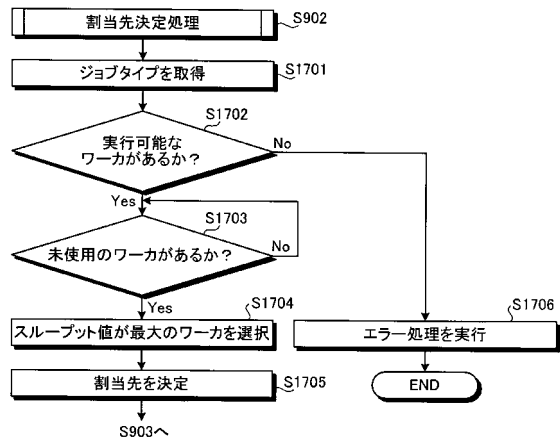
スループットテーブルの記憶内容を示す説明図(その2)

1600

ワーカ識別子	スループット値(T)		
	A	B	C
W1	T1a	—	T1c
W2	—	—	—
⋮	⋮	⋮	⋮
Wi	Tia	Tib	Tic
⋮	⋮	⋮	⋮
Wm	Tma	—	—

【図17】

割当先決定処理手順の他の一例を示すフローチャート



---

フロントページの続き

- (56)参考文献 特開2000-242614(JP,A)  
特開2003-208414(JP,A)  
特開平04-223548(JP,A)  
特開平05-012228(JP,A)  
特開2009-169757(JP,A)  
特開2005-148911(JP,A)  
特開2004-038226(JP,A)  
特開2001-312485(JP,A)  
米国特許第06757730(US,B1)  
米国特許第07093250(US,B1)

(58)調査した分野(Int.Cl., DB名)

G06F 9/46 - 9/54  
G06F 15/16 - 15/177