

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6885193号
(P6885193)

(45) 発行日 令和3年6月9日(2021.6.9)

(24) 登録日 令和3年5月17日(2021.5.17)

(51) Int.Cl.		F I			
G06F	9/50	(2006.01)	G06F	9/50	150Z
G06F	9/48	(2006.01)	G06F	9/48	300Z
G06F	15/80	(2006.01)	G06F	15/80	
G06F	15/177	(2006.01)	G06F	15/177	A

請求項の数 5 (全 20 頁)

(21) 出願番号	特願2017-95200 (P2017-95200)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成29年5月12日 (2017.5.12)		神奈川県川崎市中原区上小田中4丁目1番1号
(65) 公開番号	特開2018-194875 (P2018-194875A)	(74) 代理人	100104190 弁理士 酒井 昭徳
(43) 公開日	平成30年12月6日 (2018.12.6)	(72) 発明者	小久保 良輔 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
審査請求日	令和2年2月13日 (2020.2.13)	(72) 発明者	橋本 剛 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	井上 宏一

最終頁に続く

(54) 【発明の名称】 並列処理装置、ジョブ管理方法、およびジョブ管理プログラム

(57) 【特許請求の範囲】

【請求項1】

実行待ちの各ジョブの実行に使用されるノード数と、前記各ジョブの実行にかかる実行予定時間とに基づいて、前記各ジョブの実行規模を算出し、

算出した前記実行規模が大きいジョブから順に、複数のノードが配置された領域を区分けして分割された複数のエリアのうち、故障可能性が高い問題ノードの数が少ないエリアを選択し、選択した前記エリア内のノード群にジョブを割り当てる、

制御部を有することを特徴とする並列処理装置。

【請求項2】

前記制御部は、

前記ジョブの割り当てを行う際に、前記複数のノードそれぞれの使用状態を示す情報を参照して、前記ジョブの実行に使用されるノード数に基づいて、選択した前記エリアにおいて、前記問題ノードを含まない、前記ジョブを割り当て可能なノード群を探索し、

前記ノード群が探索された場合に、当該ノード群を選択して前記ジョブの割り当てを行う、

ことを特徴とする請求項1に記載の並列処理装置。

【請求項3】

前記制御部は、

前記複数のエリアの全てについて前記問題ノードを含まないノード群を選択した前記ジョブの割り当てができないときは、選択した前記エリアにおいて、前記問題ノードの数が

最小となるように、前記ジョブを割り当て可能なノード群を探索し、前記ノード群が探索された場合に、当該ノード群を選択して前記ジョブの割り当てを行う、

ことを特徴とする請求項 2 に記載の並列処理装置。

【請求項 4】

実行待ちの各ジョブの実行に使用されるノード数と、前記各ジョブの実行にかかる実行予定時間とに基づいて、前記各ジョブの実行規模を算出し、

算出した前記実行規模が大きいジョブから順に、複数のノードが配置された領域を区分けして分割された複数のエリアのうち、故障可能性が高い問題ノードの数が少ないエリアを選択し、選択した前記エリア内のノード群にジョブを割り当てる、

処理をコンピュータが実行することを特徴とするジョブ管理方法。

【請求項 5】

実行待ちの各ジョブの実行に使用されるノード数と、前記各ジョブの実行にかかる実行予定時間とに基づいて、前記各ジョブの実行規模を算出し、

算出した前記実行規模が大きいジョブから順に、複数のノードが配置された領域を区分けして分割された複数のエリアのうち、故障可能性が高い問題ノードの数が少ないエリアを選択し、選択した前記エリア内のノード群にジョブを割り当てる、

処理をコンピュータに実行させることを特徴とするジョブ管理プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、並列処理装置、ジョブ管理方法、およびジョブ管理プログラムに関する。

【背景技術】

【0002】

従来、コンピュータシステムを用いて科学技術計算などの大規模な計算を行う場合、複数の計算機を用いた並列計算が行われる。並列計算が可能なコンピュータシステムは、並列計算機システムと呼ばれる。大規模な並列計算機システムは、並列計算を行う多数の計算機と、管理用計算機とが含まれる。管理用計算機は、計算機に実行させるジョブを管理する。以下の説明では、並列計算を行う計算機を「計算ノード」と表記し、管理用計算機を「管理ノード」と表記する場合がある。

【0003】

大規模な並列計算機システムにおいては、複数の計算ノードを管理し並列に動作させることで、システム全体の演算性能を高めている。システム全体の演算性能を向上するためには、大量の計算ノードが必要となる。また、大規模な並列計算機システムの管理ノードでは、ジョブスケジューラ機能が、計算ノード群にユーザのジョブを割り当てる制御を実施する。

【0004】

先行技術としては、例えば、障害影響度の高いジョブをリスク度の低い実行サーバで実行し、障害影響度の高いジョブを実行中または実行が予定されている実行サーバの多重度を下げて、高負荷状態にせず障害リスク度が低い状態に保つものがある。また、ジョブに関する情報を基にジョブの形状毎の影響度を求め、影響度が高い順に所定数のジョブの形状を計算対象形状として決定し、計算対象形状及び影響度を基に、計算ノードそれぞれへのジョブの割り当て方であるジョブの事前配置を決定する技術がある。投入されたジョブが計算対象形状のいずれかに一致する場合、事前配置にしたがい投入されたジョブが計算ノードへ割り当てられる。

【先行技術文献】

【特許文献】

【0005】

【特許文献 1】特開 2011 - 215661 号公報

【特許文献 2】特開 2015 - 69577 号公報

10

20

30

40

50

【発明の概要】

【発明が解決しようとする課題】

【0006】

しかしながら、従来技術では、大規模な並列計算機システムの実行稼働率を低下させないように、ジョブを計算ノードに割り当てるのが難しい場合がある。

【0007】

一つの側面では、本発明は、複数のノードを含むシステムの実行稼働率を向上させることを目的とする。

【課題を解決するための手段】

【0008】

1つの実施態様では、実行待ちの各ジョブの実行に使用されるノード数と、前記各ジョブの実行にかかる実行予定時間とに基づいて、前記各ジョブの実行規模を算出し、算出した前記実行規模が大きいジョブから順に、複数のノードが配置された領域を区分けして分割された複数のエリアのうち、故障可能性が高い問題ノードの数が少ないエリアからジョブを割り当てる、並列処理装置が提供される。

【発明の効果】

【0009】

本発明の一側面によれば、複数のノードを含むシステムの実行稼働率を向上させることができる。

【図面の簡単な説明】

【0010】

【図1A】図1Aは、実施の形態にかかるジョブ管理方法の一実施例を示す説明図（その1）である。

【図1B】図1Bは、実施の形態にかかるジョブ管理方法の一実施例を示す説明図（その2）である。

【図2】図2は、並列計算機システム200のシステム構成例を示す説明図である。

【図3】図3は、並列処理装置101のハードウェア構成例を示すブロック図である。

【図4】図4は、ノード管理テーブル220の記憶内容の一例を示す説明図である。

【図5】図5は、ジョブ管理テーブル230の記憶内容の一例を示す説明図である。

【図6】図6は、問題ノード一覧情報600の具体例を示す説明図である。

【図7】図7は、並列処理装置101の機能的構成例を示すブロック図である。

【図8】図8は、ノード管理テーブル220の記憶内容の更新例を示す説明図である。

【図9】図9は、ジョブ管理テーブル230の記憶内容の更新例を示す説明図である。

【図10】図10は、並列処理装置101のジョブ管理処理手順の一例を示すフローチャートである。

【図11】図11は、ジョブ割当処理の具体的処理手順の一例を示すフローチャート（その1）である。

【図12】図12は、ジョブ割当処理の具体的処理手順の一例を示すフローチャート（その2）である。

【発明を実施するための形態】

【0011】

以下に図面を参照して、本発明にかかる並列処理装置、ジョブ管理方法、およびジョブ管理プログラムの実施の形態を詳細に説明する。

【0012】

（実施の形態）

図1Aおよび図1Bは、実施の形態にかかるジョブ管理方法の一実施例を示す説明図である。図1において、並列処理装置101は、複数のノードNに実行させるジョブを管理するコンピュータ（いわゆる、管理ノード）である。ノードNは、並列計算機システムの構成要素であり、並列計算を行うコンピュータ（いわゆる、計算ノード）である。ジョブは、ユーザがコンピュータに依頼する仕事の単位である。ジョブとしては、例えば、科学

10

20

30

40

50

技術計算などの大規模な計算を行うジョブが挙げられる。

【 0 0 1 3 】

ここで、大規模な並列計算機システムにおけるジョブは、特定1ノードに割り当てられるのではなく、同時に複数のノードを占有して実行する場合が多い。また、メッシュなしトラスネットワークを持つシステムでは、1ジョブへの割り当て範囲の部分ネットワークが、サブメッシュなしサブトラス（ n 次元直方体状）であることが必要な場合が多い。例えば、トラスネットワークを有する並列計算機システムのジョブスケジューラでは、計算ノードにジョブを「 n 次元直方体の形に割り当てる」のように割り当てる。

【 0 0 1 4 】

一方で、大規模な並列計算機システムでは、計算ノードの台数増加に比例して計算ノードの故障率が高くなる傾向がある。例えば、ユーザのジョブを実行している計算ノードがハードウェア故障により停止してしまうと、該当ノード上で実行されているジョブは異常終了してしまう。

【 0 0 1 5 】

このため、大規模な並列計算機システムにおいては、各計算ノードのシステムログとしてハードウェア故障を予見するログが出力されたことを事前検知し、該当計算ノードをジョブ実行で利用しないよう運用から動的に切り離すシステム監視機能が知られている。システム監視機能により運用から切り離された計算ノードは、管理ノードのジョブスケジューラ機能により、当該計算ノードに新規にジョブを割り当てないよう制御される。

【 0 0 1 6 】

ところが、ハードウェアが故障する確率が高い計算ノードを特定できたとしても、必ず故障する計算ノードを特定するのは困難である。故障確率が高いとはいえ、まだ健全な計算ノードを運用から切り離してしまうと、並列計算機システムの実行稼働率（スループット）が低下してしまう。なお、並列計算機システムの実行稼働率は、例えば、下記式（1）によって表すことができる。

【 0 0 1 7 】

並列計算機システムの実行稼働率 = （各計算ノードで正常終了したジョブが実行されていた時間） / （各計算ノードの電源が投入されていた時間）・・・（1）

【 0 0 1 8 】

そこで、本実施の形態では、故障確率が高いノード N を避けながら、できるだけ並列計算機システムの実行稼働率を向上させるジョブ管理方法について説明する。以下、並列処理装置101の処理例について説明する。ここでは、複数のノード N として、「ノード N 1～ N 60」を例に挙げて説明する。また、実行待ちのジョブとして、「ジョブ J 1～ J 3」を例に挙げて説明する。また、本実施の形態では、複数のノード N が配置された領域として、2次元の領域を例に挙げて説明するが、3次元以上の n 次元の領域にも適用可能である。

【 0 0 1 9 】

（1）並列処理装置101は、実行待ちの各ジョブ J の実行ノード数 C と実行予定時間 T とに基づいて、各ジョブ J の実行規模 S を算出する。ここで、実行ノード数 C は、実行待ちの各ジョブ J の実行に使用されるノード数である。実行予定時間 T は、各ジョブの実行にかかる予定時間である。ジョブ J の実行ノード数 C および実行予定時間 T は、例えば、ジョブ J を投入するユーザにより指定される。

【 0 0 2 0 】

また、実行規模 S は、ジョブ J が異常終了した際に並列計算機システムの実行稼働率に与える影響度合いが高いほど大きくなる指標である。例えば、実行ノード数 C が多いジョブ J ほど、実行中に多くのノード N を占有することになり、異常終了した際に実行稼働率に与える影響度合いは高いといえる。また、実行予定時間 T が長いジョブ J ほど、実行中に長い時間ノード N を占有することになり、異常終了した際に実行稼働率に与える影響度合いは高いといえる。

【 0 0 2 1 】

10

20

30

40

50

このため、並列処理装置 101 は、例えば、実行待ちの各ジョブ J の実行ノード数 C と実行予定時間 T とを乗算することにより、各ジョブ J の実行規模 S を算出してもよい。図 1 A の例では、各ジョブ J 1 ~ J 3 の実行規模 S 1 ~ S 3 が算出された結果、各ジョブ J 1 ~ J 3 が実行規模 S 1 ~ S 3 の大きい順にソートされている (J 1 J 2 J 3)。

【 0 0 2 2 】

(2) 並列処理装置 101 は、複数のノード N が配置された領域を区分けして複数のエリア A に分割する。ここで、領域は、複数のノード N が配置された平面あるいは空間のことである。以下の説明では、複数のノード N が配置された領域を「ノードエリア A R」と表記する場合がある。

【 0 0 2 3 】

具体的には、例えば、並列処理装置 101 は、ノードエリア A R を四角形 (あるいは、n 次元直方体形状) に均等に区分けして複数のエリア A に分割する。分割数は、例えば、並列計算機システムのシステムサイズに応じて設定される。図 1 A の例では、ノードエリア A R がエリア A 1 ~ A 4 に分割されている。また、エリア A 1 ~ A 4 内に存在する問題ノードの数が少ない順にエリア A 1 ~ A 4 がソートされている。

【 0 0 2 4 】

問題ノードとは、故障可能性が高いノード N である。問題ノードは、例えば、ハードウェア故障を予見するログが出力されたノード N であってもよく、また、複数のノード N のうち使用年数等をもとに相対的に故障可能性が高いと判断されたノード N であってもよい。なお、図 1 A および図 1 B 中、問題ノードは、白抜きの四角で表す。

【 0 0 2 5 】

(3) 並列処理装置 101 は、算出した実行規模 S が大きいジョブ J から順に、ノードエリア A R を区分けして分割された複数のエリア A のうち、問題ノードの数が少ないエリア A からジョブ J を割り当てる。具体的には、例えば、並列処理装置 101 は、ジョブ J の割り当てを行う際に、問題ノードを含まないノード N 群を選択してジョブ J の割り当てを行う。

【 0 0 2 6 】

図 1 B の例では、まず、ジョブ J 1 ~ J 3 のうち実行規模 S が最大のジョブ J 1 が、問題ノードの数が最小のエリア A 2 内の空き領域に割り当てられる。つぎに、2 番目に実行規模 S が大きいジョブ J 2 が、問題ノードの数が 2 番目に少ないエリア A 1 内の空き領域に割り当てられる。最後に、実行規模 S が最小のジョブ J 3 が、問題ノードの数が 3 番目に少ないエリア A 3 内の空き領域に割り当てられる。

【 0 0 2 7 】

なお、各ジョブ J の割当先となる空き領域は、例えば、四角形 (あるいは、n 次元直方体形状) のサブトラスを形成するノード N 群であって、他のジョブ J の実行に使用されていない未使用のノード N を、少なくとも各ジョブ J の実行ノード数 C 分含むノード N 群を含む領域である。

【 0 0 2 8 】

このように、並列処理装置 101 によれば、実行に使用されるノード数が多く、実行に時間がかかるジョブ J に、故障可能性が高い問題ノードができるだけ割り当てられないように、ジョブ J を実行するノード N を効率よく選定することができる。このため、異常終了時の影響度合いが大きいジョブ J が問題ノードに割り当たる確率を下げ、並列計算機システムの実行稼働率 (スループット) を向上させることができる。また、ノードエリア A R を分割したエリア単位でジョブ J の割当先となるノード N 群を探索することができるため、ジョブ J の割当先を決める際の処理時間を短縮して、ジョブ J の開始時間が遅延するのを防ぐことができる。

【 0 0 2 9 】

(並列計算機システム 200 のシステム構成例)

つぎに、図 1 に示した並列処理装置 101 を含む並列計算機システム 200 のシステム構成例について説明する。

10

20

30

40

50

【0030】

図2は、並列計算機システム200のシステム構成例を示す説明図である。図2において、並列計算機システム200は、並列処理装置101と、ノードN1～Nn（n：2以上の自然数）と、クライアント装置201と、を含む。並列計算機システム200において、並列処理装置101、ノードN1～Nnおよびクライアント装置201は、有線または無線のネットワーク210を介して接続される。ネットワーク210は、例えば、LAN（Local Area Network）、WAN（Wide Area Network）、インターネットなどである。

【0031】

並列処理装置101は、ノード管理テーブル220およびジョブ管理テーブル230を有し、ノードN1～Nnに実行させるジョブを管理する。ノード管理テーブル220およびジョブ管理テーブル230の記憶内容については、図4および図5を用いて後述する。並列処理装置101は、例えば、サーバである。

10

【0032】

ノードN1～Nnは、並列計算を行うコンピュータである。各ノードN1～Nnは、例えば、サーバである。ノードN1～Nnは、例えば、ノード間的高速通信を可能にするトラスネットワークを形成する。図1Aに示したノードN1～N60は、例えば、ノードN1～Nnに相当する（n=60）。

【0033】

以下の説明では、ノードN1～Nnのうちの任意のノードを「ノードN」と表記する場合がある。また、ノードN1～Nnが配置された領域を「ノードエリアAR」と表記する場合がある。

20

【0034】

クライアント装置201は、並列計算機システム200のユーザ（管理者を含む）が使用するコンピュータである。クライアント装置201は、例えば、PC（Personal Computer）である。なお、図2の例では、クライアント装置201を1台のみ表記したが、これに限らない。例えば、クライアント装置201は、並列計算機システム200のユーザごとに設けられる。

【0035】

（並列処理装置101のハードウェア構成例）

30

図3は、並列処理装置101のハードウェア構成例を示すブロック図である。図3において、並列処理装置101は、CPU（Central Processing Unit）301と、メモリ302と、I/F（Interface）303と、ディスクドライブ304と、ディスク305と、を有する。また、各構成部は、バス300によってそれぞれ接続される。

【0036】

ここで、CPU301は、並列処理装置101の全体の制御を司る。メモリ302は、例えば、ROM（Read Only Memory）、RAM（Random Access Memory）およびフラッシュROMなどを有する。具体的には、例えば、フラッシュROMやROMが各種プログラムを記憶し、RAMがCPU301のワークエリアとして使用される。メモリ302に記憶されるプログラムは、CPU301にロードされることで、コーディングされている処理をCPU301に実行させる。

40

【0037】

I/F303は、通信回線を通じてネットワーク210に接続され、ネットワーク210を介して外部のコンピュータ（例えば、図2に示したノードN1～Nn、クライアント装置201）に接続される。そして、I/F303は、ネットワーク210と装置内部とのインターフェースを司り、外部のコンピュータからのデータの入出力を制御する。I/F303には、例えば、モデムやLANアダプタなどを採用することができる。

【0038】

ディスクドライブ304は、CPU301の制御に従ってディスク305に対するデー

50

タのリード/ライトを制御する。ディスク305は、ディスクドライブ304の制御で書き込まれたデータを記憶する。ディスク305としては、例えば、磁気ディスク、光ディスクなどが挙げられる。

【0039】

なお、並列処理装置101は、上述した構成部のほかに、例えば、SSD(Solid State Drive)、入力装置、ディスプレイ等を有することにもよい。また、図2に示したノードN1~Nnおよびクライアント装置201についても、並列処理装置101と同様のハードウェア構成により実現することができる。ただし、クライアント装置201は、上述した構成部のほかに、入力装置、ディスプレイを有する。

【0040】

(ノード管理テーブル220の記憶内容)

つぎに、並列処理装置101が有するノード管理テーブル220の記憶内容について説明する。ノード管理テーブル220は、例えば、図3に示したメモリ302、ディスク305などの記憶装置により実現される。

【0041】

図4は、ノード管理テーブル220の記憶内容の一例を示す説明図である。図4において、ノード管理テーブル220は、ノードID、位置(x、y)、エリアID、故障可能性フラグおよび使用中フラグのフィールドを有する。各フィールドに情報を設定することで、ノード管理情報(例えば、ノード管理情報400-1~400-n)がレコードとして記憶される。

【0042】

ここで、ノードIDは、並列計算機システム200に含まれるノードNを一意に識別する識別子である。位置(x、y)は、ノードNの位置を示す座標である。なお、ここではノードエリアARとして2次元の領域を例に挙げて説明するが、ノードエリアARが3次元以上のn次元の空間の場合には、位置フィールドには、n次元座標系におけるノードNの位置を示す座標が設定される。

【0043】

エリアIDは、ノードNが属するエリアAを一意に識別する識別子である。エリアAは、ノードN1~Nnが配置されたノードエリアARを区分けして分割されたエリアである。故障可能性フラグは、ノードNが、故障可能性が高い問題ノードであるか否かを示すフラグである。故障可能性フラグ「0」は、ノードNが問題ノードではないことを示す。故障可能性フラグ「1」は、ノードNが問題ノードであることを示す。

【0044】

使用中フラグは、ノードNが、ジョブJの実行に使用されているか否かを示すフラグである。使用中フラグ「0」は、ノードNがジョブJの実行に使用されていない空きノードであることを示す。使用中フラグ「1」は、ノードNがジョブJの実行に使用されている使用中ノードであることを示す。

【0045】

(ジョブ管理テーブル230の記憶内容)

つぎに、並列処理装置101が有するジョブ管理テーブル230の記憶内容について説明する。ジョブ管理テーブル230は、例えば、図3に示したメモリ302、ディスク305などの記憶装置により実現される。

【0046】

図5は、ジョブ管理テーブル230の記憶内容の一例を示す説明図である。図5において、ジョブ管理テーブル230は、ジョブID、実行ノード数、実行予定時間および実行規模のフィールドを有し、各フィールドに情報を設定することで、ジョブ管理情報(例えば、ジョブ管理情報500-1~500-3)をレコードとして記憶する。

【0047】

ここで、ジョブIDは、実行待ちのジョブJを一意に識別する識別子である。実行ノード数は、ジョブJの実行に使用されるノード数である。実行予定時間は、ジョブJの実行

10

20

30

40

50

にかかる予定時間である。実行規模は、ジョブJが異常終了した際に並列計算機システム200の実行稼働率に与える影響度合いを表す指標である。

【0048】

(問題ノード一覧情報600の具体例)

つぎに、並列処理装置101が用いる問題ノード一覧情報600の具体例について説明する。

【0049】

図6は、問題ノード一覧情報600の具体例を示す説明図である。図6において、問題ノード一覧情報600は、ノードN1~Nnのうちの故障可能性が高い問題ノードを識別するノードIDを示す情報である。問題ノード一覧情報600は、例えば、並列処理装置101において作成されてもよく、また、並列処理装置101とは異なる他のコンピュータにおいて作成されることにしてもよい。

10

【0050】

(並列処理装置101の機能的構成例)

図7は、並列処理装置101の機能的構成例を示すブロック図である。図7において、並列処理装置101は、取得部701と、受付部702と、算出部703と、分割部704と、割当制御部705と、を含む構成である。取得部701~割当制御部705は制御部となる機能であり、具体的には、例えば、図3に示したメモリ302、ディスク305などの記憶装置に記憶されたプログラムをCPU301に実行させることにより、または、I/F303により、その機能を実現する。各機能部の処理結果は、例えば、メモリ302、ディスク305などの記憶装置に記憶される。より具体的には、各機能部は、例えば、並列処理装置101のジョブスケジューラにより実現することができる。

20

【0051】

取得部701は、ノードNの位置情報を取得する。ここで、ノードの位置情報は、ノードNの位置を示す情報であり、例えば、ノードエリアARにおけるノードNの位置を示す座標である。ノードNの位置情報には、例えば、ノードNを識別するノードIDが含まれる。ノードIDとしては、例えば、ノードNのMAC(Media Access Control)アドレスを用いることができる。

【0052】

具体的には、例えば、取得部701は、ネットワーク210(図2参照)を介して、他のコンピュータ(例えば、クライアント装置201)からノードNの位置情報を受信することにより、ノードNの位置情報を取得することにもよい。また、取得部701は、例えば、不図示の入力装置を用いたユーザの操作入力により、ノードNの位置情報を取得することにもよい。

30

【0053】

取得されたノードNの位置情報は、例えば、図4に示したノード管理テーブル220に記憶される。ここで、ノード管理テーブル220の記憶内容の更新例について説明する。

【0054】

図8は、ノード管理テーブル220の記憶内容の更新例を示す説明図である。図8の(8-1)において、ノード管理テーブル220のノードIDおよび位置(x、y)の各フィールドに情報が設定された結果、ノード管理情報(例えば、ノード管理情報400-1~400-3)がレコードとして記憶される。ただし、この時点では、各ノード管理情報のエリアIDフィールドは「-(Null)」である。また、各ノード管理情報の故障可能性フラグおよび使用中フラグの各フィールドは初期状態「0」である。

40

【0055】

図7の説明に戻り、取得部701は、問題ノードを示す情報を取得する。ここで、問題ノードは、故障可能性が高いノードNである。具体的には、例えば、取得部701は、ネットワーク210を介して、他のコンピュータ(例えば、クライアント装置201)から問題ノード一覧情報600を受信することにより、問題ノードを示す情報を取得することにもよい。また、取得部701は、例えば、不図示の入力装置を用いたユーザの操作

50

入力により、問題ノード一覧情報 600 を取得することにしてもよい。

【0056】

また、取得部 701 は、各ノード N のシステムログを監視して、問題ノード一覧情報 600 を作成することにしてもよい。より詳細に説明すると、取得部 701 は、例えば、ノード N のシステムログとしてハードウェア故障を予見するログを検出すると、当該ノード N を問題ノードとして問題ノード一覧情報 600 に登録する。

【0057】

問題ノードを示す情報が取得されると、例えば、ノード管理テーブル 220 内の対応するノード管理情報の故障可能性フラグが「1」に更新される。例えば、問題ノード一覧情報 600 が示すノード ID「N15」を例に挙げると、図 8 の(8-2)に示すように、ノード管理情報 400-15 の故障可能性フラグが「1」に更新される。

10

【0058】

受付部 702 は、ジョブ J の実行ノード数 C と実行予定時間 T とを受け付ける。ここで、実行ノード数 C は、ジョブ J の実行に使用されるノード数である。実行予定時間 T は、ジョブの実行にかかる予定時間である。実行予定時間 T の単位は、任意に設定可能であり、例えば、「分」や「時間」に設定される。

【0059】

具体的には、例えば、並列計算機システム 200 のユーザが、クライアント装置 201 において、ジョブ J を投入する際に、ジョブ J の実行ノード数 C と実行予定時間 T とを指定する。この場合、受付部 702 は、クライアント装置 201 において指定されたジョブ J の実行ノード数 C と実行予定時間 T とを受け付ける。また、受付部 702 は、例えば、不図示の入力装置を用いたユーザの操作入力により、ジョブ J の実行ノード数 C と実行予定時間 T とを受け付けることにしてもよい。

20

【0060】

受け付けられたジョブ J の実行ノード数 C と実行予定時間 T は、例えば、図 5 に示したジョブ管理テーブル 230 に記憶される。ここで、ジョブ管理テーブル 230 の記憶内容の更新例について説明する。

【0061】

図 9 は、ジョブ管理テーブル 230 の記憶内容の更新例を示す説明図である。図 9 の(9-1)において、ジョブ管理テーブル 230 のジョブ ID、実行ノード数および実行予定時間の各フィールドに情報が設定された結果、ジョブ管理情報(例えば、ジョブ管理情報 500-1~500-3)がレコードとして記憶される。ただし、この時点では、各ジョブ管理情報の実行規模フィールドは「-」である。

30

【0062】

図 7 の説明に戻り、算出部 703 は、実行待ちの各ジョブ J の実行ノード数 C と実行予定時間 T とに基づいて、各ジョブ J の実行規模 S を算出する。ここで、実行規模 S は、ジョブ J が異常終了した際に並列計算機システム 200 の実行稼働率に与える影響度合いを表す指標である。

【0063】

具体的には、例えば、算出部 703 は、ジョブ管理テーブル 230 を参照して、実行待ちの各ジョブ J の実行ノード数 C と実行予定時間 T とを乗算することにより、各ジョブ J の実行規模 S を算出する。算出された各ジョブ J の実行規模 S は、例えば、図 9 の(9-2)に示すように、各ジョブ J のジョブ ID と対応付けて、ジョブ管理テーブル 230 の実行規模フィールドに記憶される。

40

【0064】

分割部 704 は、ノード N1~Nn が配置されたノードエリア AR を区分けして複数のエリア A に分割する。例えば、ノードエリア AR が 2 次元の平面の場合、各エリア A は、四角形の領域となる。例えば、ノードエリア AR が n 次元の空間の場合、各エリア A は、n 次元直方体の領域となる。具体的には、例えば、分割部 704 は、ノードエリア AR を四角形(あるいは、n 次元直方体形状)に均等に区分けして複数のエリア A に分割する。

50

分割数は、例えば、並列計算機システム 200 のシステムサイズに応じて適宜設定される。

【0065】

また、分割部 704 は、各エリア A の探索開始位置を設定する。ここで、探索開始位置とは、各エリア A において、ジョブ J を割り当てる空き領域を探索する際の開始位置となる位置である。空き領域とは、ジョブ J の実行に使用されていない未使用のノード N 群を含む領域である。各エリア内のどの位置を探索開始位置とするかは任意に設定可能である。具体的には、例えば、分割部 704 は、ノードエリア AR を四角形に区分けした各エリア A の左下の位置を探索開始位置に設定することにしてもよい。

【0066】

一例として、図 1A に示したようにノードエリア AR を 4 分割する場合、ノードエリア AR の左下を原点とすると、左下のエリア A1 の探索開始位置は、「 $(x, y) = (0, 0)$ 」となる。また、右下のエリア A2 の探索開始位置は、「 $(x, y) = (x \text{ 軸最大値} \div 2, 0)$ 」となる。また、左上のエリア A3 の探索開始位置は、「 $(x, y) = (0, y \text{ 軸最大値} \div 2)$ 」となる。また、右上のエリア A4 の探索開始位置は、「 $(x, y) = (x \text{ 軸最大値} \div 2, y \text{ 軸最大値} \div 2)$ 」となる。

【0067】

また、分割部 704 は、ノード N が属するエリア A を特定する。具体的には、例えば、分割部 704 は、ノード管理テーブル 220 を参照して、各ノード N が属するエリア A を特定する。特定された結果（エリア A のエリア ID）は、例えば、図 8 の（8-3）に示すように、各ノード N のノード ID と対応付けて、ノード管理テーブル 220 のエリア ID フィールドに記憶される。

【0068】

割当制御部 705 は、実行待ちのジョブ J を割り当てる制御を行う。具体的には、例えば、割当制御部 705 は、ノード管理テーブル 220 を参照して、各エリア A の問題ノード数 p を算出する。問題ノード数 p は、各エリア A に属する問題ノードの数である。一例として、エリア A1 の問題ノード数 p1 を算出するとする。この場合、割当制御部 705 は、エリア ID フィールドに「A1」が設定されたノード管理情報のうち、故障可能性フラグに「1」が設定されたノード管理情報の数を、エリア A1 の問題ノード数 p1 として算出する。

【0069】

そして、割当制御部 705 は、ジョブ管理テーブル 230 を参照して、算出された実行規模 S が大きいジョブ J から順に、複数のエリア A のうち、算出した問題ノード数 p が少ないエリア A からジョブ J を割り当てる。この際、割当制御部 705 は、例えば、問題ノードを含まないノード N 群を選択してジョブ J の割り当てを行う。

【0070】

より詳細に説明すると、例えば、まず、割当制御部 705 は、ジョブ管理テーブル 230 を参照して、実行待ちのジョブ J1 ~ J3 を実行規模 S が大きい順にソートする。ここで、実行規模 S1 ~ S3 の大小関係を「 $S1 > S2 > S3$ 」とする。この場合、ジョブ J1 ~ J3 を実行規模 S が大きい順にソートすると、{J1, J2, J3} となる。なお、実行規模 S が同一のジョブ J が存在する場合には、割当制御部 705 は、例えば、それらジョブ J について、キューに入れられた順にソートすることにしてもよい。

【0071】

また、割当制御部 705 は、複数のエリア A を問題ノード数 p が少ない順にソートする。ここで、複数のエリア A を「エリア A1 ~ A4」とし、エリア A1 ~ A4 の問題ノード数 p1 ~ p4 の大小関係を「 $p4 > p3 > p1 > p2$ 」とする。この場合、エリア A1 ~ A4 を問題ノード数 p が少ない順にソートすると、{A2, A1, A3, A4} となる。なお、問題ノード数 p が同一のエリア A が存在する場合には、割当制御部 705 は、例えば、それらエリア A について、探索開始位置に近い問題ノードの数が少ないエリアを上位にソートすることにしてもよい。

10

20

30

40

50

【 0 0 7 2 】

つぎに、割当制御部 7 0 5 は、{ J 1 , J 2 , J 3 } から、実行規模 S が最大のジョブ J 1 を選択する。また、割当制御部 7 0 5 は、{ A 2 , A 1 , A 3 , A 4 } から、問題ノード数 p が最小のエリア A 2 を選択する。そして、割当制御部 7 0 5 は、ノード管理テーブル 2 2 0 を参照して、選択したエリア A 2 から、問題ノードを含まない、ジョブ J 1 を割り当て可能なノード N 群を探索する。

【 0 0 7 3 】

ここで、ジョブ J 1 を割り当て可能なノード N 群は、例えば、サブトラスを形成するノード N の集合であって、他のジョブ J の実行に使用されていない未使用のノード N を、少なくともジョブ J 1 の実行ノード数 C 1 分含むノード N の集合である。

10

【 0 0 7 4 】

具体的には、例えば、割当制御部 7 0 5 は、エリア A 2 の探索開始位置から徐々に範囲を広げながら、問題ノードを含まない、ジョブ J 1 を割り当て可能なノード N 群を探索する。この際、割当制御部 7 0 5 は、例えば、ノード単位、あるいは、シャーシ単位で範囲を広げることにもよい。シャーシとは、サブトラスを形成するノード N の集合である。そして、ノード N 群の探索に成功すると、割当制御部 7 0 5 は、探索したノード N 群を選択してジョブ J 1 を割り当てる。

【 0 0 7 5 】

一方、ノード N 群の探索に失敗すると、割当制御部 7 0 5 は、{ A 2 , A 1 , A 3 , A 4 } から、問題ノード数 p が次に少ないエリア A 1 を選択する。そして、割当制御部 7 0 5 は、選択したエリア A 2 から、問題ノードを含まない、ジョブ J 1 を割り当て可能なノード N 群を探索する。割当制御部 7 0 5 は、ノード N 群の探索に成功する、あるいは、未選択のエリア A がなくなるまで、上述した一連の処理を繰り返す。

20

【 0 0 7 6 】

また、ジョブ J 1 の割り当てが完了すると、割当制御部 7 0 5 は、{ J 1 , J 2 , J 3 } から、実行規模 S が次に大きいジョブ J 2 を選択して、ジョブ J 1 と同様の処理を行う。そして、ジョブ J 2 の割り当てが完了すると、割当制御部 7 0 5 は、{ J 1 , J 2 , J 3 } から、実行規模 S が次に大きいジョブ J 3 を選択して、ジョブ J 1 , J 2 と同様の処理を行う。

【 0 0 7 7 】

ただし、複数のエリア A の全てについて、問題ノードを含まないノード N 群を選択したジョブ J の割り当てができないときがある。この場合、割当制御部 7 0 5 は、問題ノード数 p が少ないエリア A から、問題ノードの数が最小となるようにノード N 群を選択してジョブ J の割り当てを行うことにしてもよい。

30

【 0 0 7 8 】

すなわち、割当制御部 7 0 5 は、問題ノードを含むことを許容して、ジョブ J を割り当て可能なノード N 群を探索する。この際、割当制御部 7 0 5 は、例えば、問題ノードの数が最小となるように、ジョブ J を割り当て可能なノード N 群をエリア A から探索する。そして、ノード N 群の探索に成功すると、割当制御部 7 0 5 は、探索したノード N 群を選択してジョブ J を割り当てる。ただし、複数のエリア A の全てについて、問題ノードを含むことを許容してもノード N 群の探索に失敗した場合、割当制御部 7 0 5 は、ジョブ J をキューに戻すことにしてもよい。

40

【 0 0 7 9 】

なお、ジョブ J の割り当てが完了すると、割当制御部 7 0 5 は、ジョブ J を割り当てたノード N に対応する、ノード管理テーブル 2 2 0 内の使用中フラグを「 1 」に変更する。また、ジョブ J の実行が終了すると、割当制御部 7 0 5 は、ジョブ J が割り当てられていたノード N に対応する、ノード管理テーブル 2 2 0 内の使用中フラグを「 0 」に変更する。

【 0 0 8 0 】

(並列処理装置 1 0 1 のジョブ管理処理手順)

50

つぎに、並列処理装置 101 のジョブ管理処理手順について説明する。ジョブ管理処理は、例えば、定期的に行われることにしてもよく、新たなジョブ J が投入される、あるいは、投入済みのいずれかのジョブ J の実行が完了したことに応じて実行されることにしてもよい。また、ノード N の位置情報は、ノード管理テーブル 220 に記憶されているとする。

【0081】

図 10 は、並列処理装置 101 のジョブ管理処理手順の一例を示すフローチャートである。図 10 のフローチャートにおいて、まず、並列処理装置 101 は、問題ノード一覧情報 600 を取得する (ステップ S1001)。そして、並列処理装置 101 は、取得した問題ノード一覧情報 600 に基づいて、ノード管理テーブル 220 内の故障可能性フラグを更新する (ステップ S1002)。

10

【0082】

つぎに、並列処理装置 101 は、ジョブ J の実行ノード数 C と実行予定時間 T とを受け付ける (ステップ S1003)。受け付けられたジョブ J の実行ノード数 C と実行予定時間 T は、ジョブ管理テーブル 230 に記憶される。

【0083】

そして、並列処理装置 101 は、ジョブ管理テーブル 230 を参照して、実行待ちの各ジョブ J の実行ノード数 C と実行予定時間 T とを乗算することにより、各ジョブ J の実行規模 S を算出する (ステップ S1004)。算出された各ジョブ J の実行規模 S は、ジョブ管理テーブル 230 に記憶される。

20

【0084】

つぎに、並列処理装置 101 は、ジョブ管理テーブル 230 を参照して、実行待ちのジョブ J を実行規模 S が大きい順にソートする (ステップ S1005)。そして、並列処理装置 101 は、ノード N1 ~ Nn が配置されたノードエリア AR を区分けして複数のエリア A に分割する (ステップ S1006)。この際、並列処理装置 101 は、各エリア A の探索開始位置を設定する。

【0085】

つぎに、並列処理装置 101 は、ノード管理テーブル 220 を参照して、各ノード N が属するエリア A を特定する (ステップ S1007)。特定された結果 (エリア A のエリア ID) は、ノード管理テーブル 220 に記憶される。そして、並列処理装置 101 は、ノード管理テーブル 220 を参照して、各エリア A の問題ノード数 p を算出する (ステップ S1008)。

30

【0086】

つぎに、並列処理装置 101 は、複数のエリア A を問題ノード数 p が少ない順にソートする (ステップ S1009)。つぎに、並列処理装置 101 は、実行規模 S が大きい順にソートした実行待ちのジョブ J の先頭から未選択のジョブ J を選択する (ステップ S1010)。

【0087】

そして、並列処理装置 101 は、選択したジョブ J を割り当てるジョブ割当処理を実行する (ステップ S1011)。ジョブ割当処理の具体的な処理手順については、図 11 および図 12 を用いて後述する。つぎに、並列処理装置 101 は、実行規模 S が大きい順にソートした実行待ちのジョブ J のうち選択されていない未選択のジョブ J があるか否かを判断する (ステップ S1012)。

40

【0088】

ここで、未選択のジョブ J がある場合 (ステップ S1012: Yes)、並列処理装置 101 は、ステップ S1010 に戻る。一方、未選択のジョブ J がいない場合 (ステップ S1012: No)、並列処理装置 101 は、本フローチャートによる一連の処理を終了する。これにより、実行待ちのジョブ J の割り当てを行うことができる。

【0089】

つぎに、ステップ S1011 のジョブ割当処理の具体的な処理手順について説明する。

50

【 0 0 9 0 】

図 1 1 および図 1 2 は、ジョブ割当処理の具体的処理手順の一例を示すフローチャートである。図 1 1 のフローチャートにおいて、まず、並列処理装置 1 0 1 は、問題ノード数 p が少ない順にソートした複数のエリア A の先頭から未選択のエリア A を選択する（ステップ $S 1 1 0 1$ ）。

【 0 0 9 1 】

つぎに、並列処理装置 1 0 1 は、選択したエリア A から、問題ノードを含まない、選択したジョブ J を割り当て可能なノード N 群を探索する（ステップ $S 1 1 0 2$ ）。なお、ジョブ J を割り当て可能なノード N 群は、例えば、サブトラスを形成するノード N の集合であって、未使用のノード N を実行ノード数 C 分含むノード N の集合である。

10

【 0 0 9 2 】

そして、並列処理装置 1 0 1 は、ノード N 群が探索されたか否かを判断する（ステップ $S 1 1 0 3$ ）。ここで、ノード N 群が探索された場合（ステップ $S 1 1 0 3$: Yes）、並列処理装置 1 0 1 は、探索したノード N 群を選択してジョブ J を割り当てて（ステップ $S 1 1 0 4$ ）、ジョブ割当処理を呼び出したステップに戻る。

【 0 0 9 3 】

一方、ノード N 群が探索されなかった場合（ステップ $S 1 1 0 3$: No）、並列処理装置 1 0 1 は、問題ノード数 p が少ない順にソートした複数のエリア A のうち、ステップ $S 1 1 0 1$ において選択されていない未選択のエリア A があるか否かを判断する（ステップ $S 1 1 0 5$ ）。

20

【 0 0 9 4 】

ここで、未選択のエリア A がある場合（ステップ $S 1 1 0 5$: Yes）、並列処理装置 1 0 1 は、ステップ $S 1 1 0 1$ に戻る。一方、未選択のエリア A がない場合には（ステップ $S 1 1 0 5$: No）、並列処理装置 1 0 1 は、図 1 2 に示すステップ $S 1 2 0 1$ に移行する。

【 0 0 9 5 】

図 1 2 のフローチャートにおいて、まず、並列処理装置 1 0 1 は、問題ノード数 p が少ない順にソートした複数のエリア A の先頭から未選択のエリア A を選択する（ステップ $S 1 2 0 1$ ）。そして、並列処理装置 1 0 1 は、選択したエリア A から、問題ノードを含むことを許容して、問題ノードの数が最小となるように、ジョブ J を割り当て可能なノード N 群を探索する（ステップ $S 1 2 0 2$ ）。

30

【 0 0 9 6 】

つぎに、並列処理装置 1 0 1 は、ノード N 群が探索されたか否かを判断する（ステップ $S 1 2 0 3$ ）。ここで、ノード N 群が探索された場合（ステップ $S 1 2 0 3$: Yes）、並列処理装置 1 0 1 は、探索したノード N 群を選択してジョブ J を割り当てて（ステップ $S 1 2 0 4$ ）、ジョブ割当処理を呼び出したステップに戻る。

【 0 0 9 7 】

一方、ノード N 群が探索されなかった場合（ステップ $S 1 2 0 3$: No）、並列処理装置 1 0 1 は、問題ノード数 p が少ない順にソートした複数のエリア A のうち、ステップ $S 1 2 0 1$ において選択されていない未選択のエリア A があるか否かを判断する（ステップ $S 1 2 0 5$ ）。

40

【 0 0 9 8 】

ここで、未選択のエリア A がある場合（ステップ $S 1 2 0 5$: Yes）、並列処理装置 1 0 1 は、ステップ $S 1 2 0 1$ に戻る。一方、未選択のエリア A がない場合には（ステップ $S 1 2 0 5$: No）、並列処理装置 1 0 1 は、選択したジョブ J をキューに入れて（ステップ $S 1 2 0 6$ ）、ジョブ割当処理を呼び出したステップに戻る。

【 0 0 9 9 】

これにより、実行規模 S が大きいジョブ J に、故障可能性が高い問題ノードができるだけ割り当てられないように制御することができる。

【 0 1 0 0 】

50

以上説明したように、実施の形態にかかる並列処理装置 101 によれば、実行待ちの各ジョブ J の実行ノード数 C と実行予定時間 T とに基づいて、各ジョブ J の実行規模 S を算出することができる。そして、並列処理装置 101 によれば、実行規模 S が大きいジョブ J から順に、ノード N1 ~ Nn が配置されたノードエリア AR を区分けして分割された複数のエリア A のうち、問題ノード数 p が少ないエリア A からジョブ J を割り当てることができる。

【0101】

これにより、実行に使用されるノード数が多く、実行に時間がかかるジョブ J に、故障可能性が高い問題ノードができるだけ割り当てられないように、ジョブ J を実行するノード N を効率よく選定することができる。このため、異常終了時の影響度合いが大きいジョブ J が問題ノードに割り当たる確率を下げ、並列計算機システム 200 の実行稼働率（スループット）を向上させることができる。また、ジョブ J の割当先を決める際の処理時間を短縮することができ、ジョブ J の開始時間が遅延するのを防ぐことができる。

10

【0102】

また、並列処理装置 101 によれば、ジョブ J の割り当てを行う際に、問題ノードを含まないノード N 群を選択してジョブ J の割り当てを行うことができる。これにより、実行中のジョブ J が異常終了してしまうのを防ぐことができ、無駄な処理が生じて並列計算機システム 200 の実行稼働率が低下するのを抑制することができる。

【0103】

また、並列処理装置 101 によれば、複数のエリア A の全てについて問題ノードを含まないノード N 群を選択したジョブ J の割り当てができないときは、問題ノードの数が最小となるようにノード N 群を選択してジョブ J の割り当てを行うことができる。これにより、問題ノードを避けたジョブ J の割り当てができない場合は、問題ノード数を最小化することで、実行中のジョブ J が異常終了する確率を下げるることができる。

20

【0104】

また、並列処理装置 101 によれば、問題ノード数 p として、各ノード N で記録されるシステムログからハードウェア故障が予見されたノード N の数を計数することができる。これにより、ハードウェア故障の可能性が高い問題ノードにジョブ J ができるだけ割り当てられないように、ジョブ J を実行するノード N を効率よく選定することができる。

【0105】

これらのことから、並列処理装置 101 によれば、トラスネットワークを有するような大規模な並列計算機システム 200 において、実行稼働率を下げないように、部分ネットワーク（例えば、2次元平面状や n次元直方体状のサブトラス）にジョブ J を割り当てることが可能となる。

30

【0106】

なお、本実施の形態で説明したジョブ管理方法は、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等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。本ジョブ管理プログラムは、ハードディスク、フレキシブルディスク、CD (Compact Disc) - ROM、MO (Magneto - Optical disk)、DVD (Digital Versatile Disk)、USB (Universal Serial Bus) メモリ等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、本ジョブ管理プログラムは、インターネット等のネットワークを介して配布してもよい。

40

【0107】

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。

【0108】

(付記 1) 実行待ちの各ジョブの実行に使用されるノード数と、前記各ジョブの実行にかかる実行予定時間とに基づいて、前記各ジョブの実行規模を算出し、

算出した前記実行規模が大きいジョブから順に、複数のノードが配置された領域を区分

50

けて分割された複数のエリアのうち、故障可能性が高い問題ノードの数が少ないエリアからジョブを割り当てる、

制御部を有することを特徴とする並列処理装置。

【0109】

(付記2) 前記制御部は、

前記ジョブの割り当てを行う際に、前記問題ノードを含まないノード群を選択して前記ジョブの割り当てを行う、

ことを特徴とする付記1に記載の並列処理装置。

【0110】

(付記3) 前記制御部は、

前記複数のエリアの全てについて前記問題ノードを含まないノード群を選択した前記ジョブの割り当てができないときは、前記問題ノードの数が最小となるようにノード群を選択して前記ジョブの割り当てを行う、

ことを特徴とする付記2に記載の並列処理装置。

【0111】

(付記4) 前記複数のノードは、トラスネットワークを形成するノードである、ことを特徴とする付記1～3のいずれか一つに記載の並列処理装置。

【0112】

(付記5) 前記問題ノードは、前記複数のノードそれぞれで記録されるシステムログからハードウェア故障が予見されたノードである、ことを特徴とする付記1～4のいずれか一つに記載の並列処理装置。

【0113】

(付記6) 実行待ちの各ジョブの実行に使用されるノード数と、前記各ジョブの実行にかかる実行予定時間とに基づいて、前記各ジョブの実行規模を算出し、

算出した前記実行規模が大きいジョブから順に、複数のノードが配置された領域を区分けして分割された複数のエリアのうち、故障可能性が高い問題ノードの数が少ないエリアからジョブを割り当てる、

処理をコンピュータが実行することを特徴とするジョブ管理方法。

【0114】

(付記7) 実行待ちの各ジョブの実行に使用されるノード数と、前記各ジョブの実行にかかる実行予定時間とに基づいて、前記各ジョブの実行規模を算出し、

算出した前記実行規模が大きいジョブから順に、複数のノードが配置された領域を区分けして分割された複数のエリアのうち、故障可能性が高い問題ノードの数が少ないエリアからジョブを割り当てる、

処理をコンピュータに実行させることを特徴とするジョブ管理プログラム。

【符号の説明】

【0115】

- 101 並列処理装置
- 200 並列計算機システム
- 201 クライアント装置
- 210 ネットワーク
- 220 ノード管理テーブル
- 230 ジョブ管理テーブル
- 300 バス
- 301 CPU
- 302 メモリ
- 303 I/F
- 304 ディスクドライブ
- 305 ディスク
- 600 問題ノード一覧情報

10

20

30

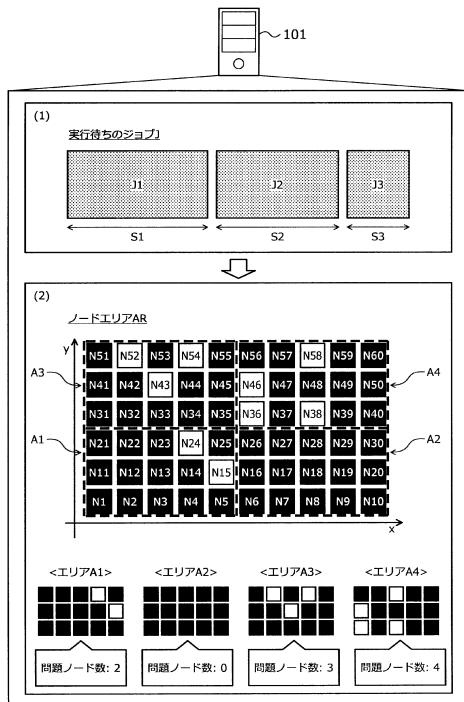
40

50

- 701 取得部
- 702 受付部
- 703 算出部
- 704 分割部
- 705 割当制御部
- A エリア
- AR ノードエリア
- N ノード

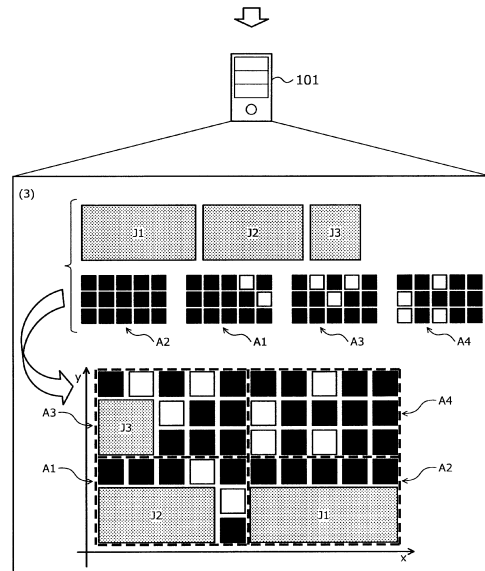
【図1A】

実施の形態にかかるジョブ管理方法の一実施例を示す説明図(その1)



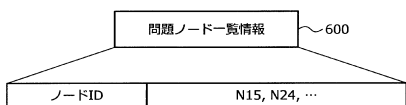
【図1B】

実施の形態にかかるジョブ管理方法の一実施例を示す説明図(その2)



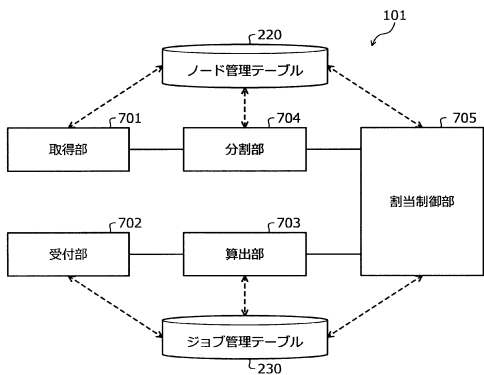
【図6】

問題ノード一覧情報600の具体例を示す説明図



【図7】

並列処理装置101の機能的構成例を示すブロック図



【図8】

ノード管理テーブル2.0の記憶内容の更新例を示す説明図

ノードID	位置 (x, y)	エリアID	故障可能性フラグ	使用中フラグ
400-1	(x1, y1)	-	0	0
400-2	(x2, y2)	-	0	0
400-3	(x3, y3)	-	0	0
...

(8-1)

ノードID	位置 (x, y)	エリアID	故障可能性フラグ	使用中フラグ
...
N15	(x15, y15)	-	1	0
...

(8-2)

ノードID	位置 (x, y)	エリアID	故障可能性フラグ	使用中フラグ
N1	(x1, y1)	A1	0	0
N2	(x2, y2)	A1	0	0
N3	(x3, y3)	A1	0	0
...

(8-3)

【図9】

ジョブ管理テーブル2.3.0の記憶内容の更新例を示す説明図

ジョブID	実行ノード数	実行予定時間	実行規模
500-1	C1	T1	-
500-2	C2	T2	-
500-3	C3	T3	-
...

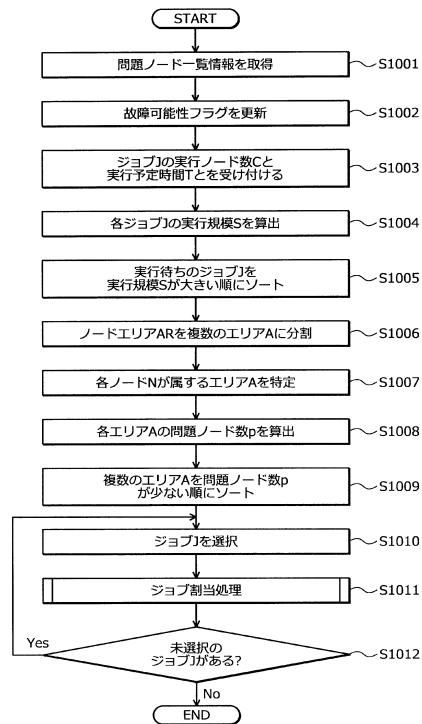
(9-1)

ジョブID	実行ノード数	実行予定時間	実行規模
500-1	C1	T1	S1
500-2	C2	T2	S2
500-3	C3	T3	S3
...

(9-2)

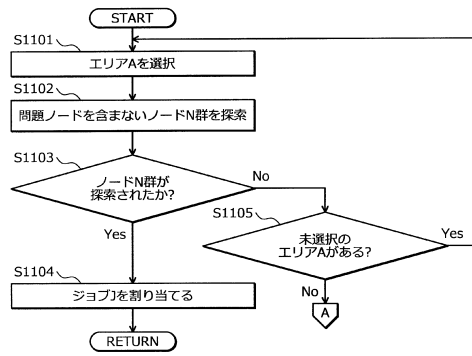
【図10】

並列処理装置101のジョブ管理処理手順の一例を示すフローチャート



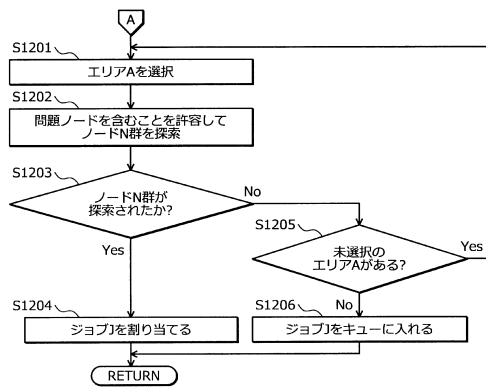
【図 11】

ジョブ割当処理の具体的な処理手順の一例を示すフローチャート（その1）



【図 12】

ジョブ割当処理の具体的な処理手順の一例を示すフローチャート（その2）



フロントページの続き

(56)参考文献 特開2003 - 58518 (JP, A)
特開2015 - 69577 (JP, A)
特表2007 - 533031 (JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 9/455 - 9/54
G06F 15/177
G06F 15/80