

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6434131号
(P6434131)

(45) 発行日 平成30年12月5日(2018.12.5)

(24) 登録日 平成30年11月16日(2018.11.16)

(51) Int.Cl.		F I			
G06F	9/50	(2006.01)	G06F	9/50	150C
G06F	11/16	(2006.01)	G06F	11/16	666
G06F	11/20	(2006.01)	G06F	11/20	664

請求項の数 13 (全 25 頁)

(21) 出願番号	特願2017-512128 (P2017-512128)	(73) 特許権者	000005108
(86) (22) 出願日	平成27年4月15日 (2015.4.15)		株式会社日立製作所
(86) 国際出願番号	PCT/JP2015/061594		東京都千代田区丸の内一丁目6番6号
(87) 国際公開番号	W02016/166844	(74) 代理人	110001689
(87) 国際公開日	平成28年10月20日 (2016.10.20)		青稜特許業務法人
審査請求日	平成29年9月14日 (2017.9.14)	(72) 発明者	愛甲 和秀
			東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		(72) 発明者	木下 雅文
			東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		(72) 発明者	小島 剛
			東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内

最終頁に続く

(54) 【発明の名称】 分散処理システム、タスク処理方法、記憶媒体

(57) 【特許請求の範囲】

【請求項1】

それぞれにワーカーノード識別子が設定され、それぞれがデータをキャッシュするメモリを有する複数のワーカーノードを含む分散処理システムは、

前記ワーカーノード識別子の中の第1のワーカーノード識別子が設定され、第1のタスクの実行結果の第1の出力データを、キャッシュ要求に応じて自ワーカーノードに有される第1のメモリにキャッシュし、他ワーカーノードへ前記第1の出力データをコピーさせる、前記複数のワーカーノードの中の第1のワーカーノードと、

前記ワーカーノード識別子の中の第2のワーカーノード識別子が設定され、前記第1の出力データのコピーを、自ワーカーノードに有される第2のメモリにキャッシュする、前記複数のワーカーノードの中の第2のワーカーノードと、

前記第1のワーカーノード識別子の情報を有し、前記第1の出力データを入力する第2のタスクを割り当てるワーカーノードを前記第1のワーカーノードであると、前記第1のワーカーノード識別子の情報に基づき選択するマスタノードとを備えたことを特徴とする分散処理システム。

【請求項2】

前記マスタノードは、前記第1の出力データと前記第1のワーカーノードの識別子の情報を管理するキャッシュ管理部と、アプリケーションからのタスクを含む処理要求を登録し、

前記管理された情報に基づき前記第1の出力データが前記第1のワーカーノードにキャッシュされているかを判定し、キャッシュされていると判定した場合は、前記第1のワーカーノードを選択し、キャッシュされていないと判定した場合は、任意のワーカーノードを選択し、

前記タスクの処理状態を登録し、

前記処理要求に含まれる全てのタスクの前記処理状態が処理済みとなるまで、前記選択と前記登録を繰り返し、全てのタスクが処理済みとなると前記アプリケーションへ応答するスケジュール管理部と、

前記選択されたワーカーノードへ前記第1の出力データを入力とする第2のタスクの実行を要求し、

10

前記選択されたワーカーノードと前記第2のタスクとを関連付けて登録し、

前記第1のワーカーノードの障害を検知して、タスクの再実行が必要かを判定し、

前記第1のワーカーノードの復旧通知を受けて、前記第1のワーカーノードに関する情報を登録する

タスク配置管理部と、

を有し、

前記第1のワーカーノードは、

出力データを格納するストレージ装置へのインタフェースと、

前記タスク配置管理部の前記第2のタスクの実行の要求を受け付けて登録し、

前記第1の出力データを取得し、

20

前記第2のタスクを実行し、

前記第2のタスクの実行結果の第2の出力データをキャッシュするか判定し、

前記第2の出力データを前記ストレージ装置へ格納させ、

前記タスク配置管理部へ前記第2のタスクの実行の完了を通知し、

前記第1のワーカーノードの復旧時に前記復旧通知を送る

タスク実行部と、

前記第1の出力データの取得として、前記第1の出力データが前記第1のワーカーノードにキャッシュされているか判定し、キャッシュされていないと判定した場合は、前記ストレージ装置から前記第1の出力データを取得し、

前記第2の出力データをキャッシュすると判定した場合は、前記第2の出力データのサイズが前記第1のメモリの残容量未満であると、前記第2の出力データをキャッシュ登録処理するデータ配置管理部と、

30

前記キャッシュ登録処理に対応して、前記第1の出力データまたは前記第2の出力データを前記第1のメモリへキャッシュし、

前記第1の出力データまたは前記第2の出力データのコピー要求を発行する

データ保管部と、

前記キャッシュ登録処理に対応して前記第1のメモリへキャッシュされた前記第1の出力データまたは前記第2の出力データの情報を登録する

ローカルキャッシュ管理部と

を有し、

40

前記第2のワーカーノードは、

前記第1のワーカーノードのデータ保管部からのコピー要求に基づいて前記第1の出力データまたは前記第2の出力データのコピーを前記第2のメモリへキャッシュし、

前記第1のワーカーノードの障害を検知し、前記第1の出力データのコピーの情報を前記第1の出力データの情報として登録し、前記第1の出力データの情報の登録をデータ配置管理部へ通知し、前記第2のメモリにキャッシュしたコピーデータのマスタへの変更又は削除する

データ保管部と、

前記第1の出力データの情報の登録通知を受信し、前記スケジュール管理部から前記第2のタスクの前記処理状態を取得して判定し、前記処理状態が処理中か処理済みの場合は、

50

前記第1の出力データのコピーを削除処理、前記処理状態が処理中と処理済みのいずれでもない場合は、前記第1の出力データのコピーの情報を前記第1の出力データの情報として登録するよう前記データ保管部に要求する

データ配置管理部と、

前記第1の出力データのコピーの情報を前記第1の出力データの情報として登録し、前記第1の出力データの情報と前記第2のワーカースレッドの識別子の情報を前記キャッシュ管理部へ更新要求する

ローカルキャッシュ管理部と、

を有すること

を特徴とする請求項1に記載の分散処理システム。

10

【請求項3】

前記マスタノードは、

前記選択された前記第1のワーカースレッドへ前記第2のタスクの実行を要求し、

前記第1のワーカースレッドは、

前記第2のタスクの実行の要求に応じて、前記第2のタスクを実行し、前記第2のタスクの実行結果の第2の出力データをキャッシュすると判定した場合、前記第2の出力データを前記第1のメモリにキャッシュし、他ワーカースレッドへ前記第2の出力データをコピーさせ、前記第1のワーカースレッドで前記第2の出力データをキャッシュしたことを前記マスタノードへ通知すること

を特徴とする請求項1に記載の分散処理システム。

20

【請求項4】

前記第1のワーカースレッドは、

コンシステントハッシングによる選択、ラウンドロビンによる選択、ランダムによる選択、あるいはメモリ容量使用率が少ないワーカースレッドを選択のいずれかにより、前記複数のワーカースレッドの中から前記第2の出力データのコピー先である他ワーカースレッドを選択し、前記選択した他ワーカースレッドへ前記第2の出力データをコピーさせることを特徴とする請求項3に記載の分散処理システム。

【請求項5】

前記第1のワーカースレッドは、

前記第2のタスクの実行の要求に応じて、キャッシュされた前記第1の出力データを前記第1のメモリから取得し、前記取得した第1の出力データを入力として前記第2のタスクを実行すること

を特徴とする請求項4に記載の分散処理システム。

30

【請求項6】

前記第1のワーカースレッドは、

前記第2のタスクの実行結果の第2の出力データをキャッシュすると判定した場合、前記第2の出力データのサイズが前記第1のメモリの残容量未満であるとさらに判定して、前記第2の出力データを前記第1のメモリにキャッシュすること

を特徴とする請求項5に記載の分散処理システム。

【請求項7】

前記第2のワーカースレッドは、

他ワーカースレッドの障害を監視し、前記第1のワーカースレッドの障害を検知した場合、前記第2のメモリにキャッシュした前記第1の出力データのコピーの属性をスレーブからマスタへ更新し、前記第2のタスクが処理中であると登録されるより前の状態であることを判定し、前記マスタノードが有する前記第1のワーカースレッド識別子の情報を前記第2のワーカースレッドの識別子の情報へ更新する要求を前記マスタノードへ発行し、

前記マスタノードは、

前記第2のワーカースレッドからの前記更新する要求に応じて、前記マスタノードが有する前記第1のワーカースレッド識別子の情報を前記第2のワーカースレッドの識別子の情報へ更新すること

40

50

を特徴とする請求項 6 に記載の分散処理システム。

【請求項 8】

前記マスタノードは、

前記第 1 のワーカーノード識別子の情報に基づき前記第 1 のワーカーノードを選択し、前記第 1 のワーカーノードの障害を検知すると、第 3 のワーカーノードへ前記第 2 のタスクの実行要求を発行し、前記第 2 のタスクが処理中であると登録し、

前記第 2 のワーカーノードは、

前記第 1 のワーカーノードの障害を検知し、前記第 2 のタスクが前記第 3 のワーカーノードで処理中であると登録されていると、前記第 2 のメモリに格納している前記第 1 の出力データを削除すること

10

を特徴とする請求項 7 に記載の分散処理システム。

【請求項 9】

前記マスタノードは、

前記第 1 のワーカーノードの障害を検知し、前記第 1 のワーカーノード以外のワーカーノードへ前記第 2 のタスクの実行の要求を発行すること

を特徴とする請求項 8 に記載の分散処理システム。

【請求項 10】

ワーカーノードは、

実行の要求される第 1 のタスクの識別子と、前記第 1 のタスクに入力される第 1 のデータの識別子を受け付け、

20

前記受け付けた第 1 のデータの識別子に基づき、前記ワーカーノードのメモリに前記第 1 のデータがキャッシュされているかを判定し、

前記キャッシュされているかを判定した結果がキャッシュされていると判定された場合、前記第 1 のデータを前記メモリから取得し、

前記取得した第 1 のデータを入力として前記第 1 のタスクを実行し、

前記第 1 のタスクの実行結果の第 2 のデータをキャッシュするかを判定し、

前記キャッシュするかを判定した結果がキャッシュすると判定された場合、前記第 2 のデータを前記メモリにキャッシュして前記ワーカーノードで前記第 2 のデータをキャッシュしたことを通知し、他ワーカーノードへ前記第 2 のデータをコピーさせること

を特徴とするワーカーノードのタスク処理方法。

30

【請求項 11】

前記ワーカーノードは、

他ワーカーノードの障害を検知し、

前記障害の検知された他ワーカーノードでキャッシュされたデータのコピーである第 3 のデータを有する場合、前記第 3 のデータを入力とする第 3 のタスクの実行状態の情報を取得し、

前記取得した実行状態の情報を判定し、前記実行状態の情報が実行中である情報と判定した場合、前記第 3 のデータを前記メモリから削除し、

前記実行状態の情報が実行前である情報と判定した場合、前記ワーカーノードで前記第 3 のデータをキャッシュしたことを通知すること

40

を特徴とする請求項 10 に記載のワーカーノードのタスク処理方法。

【請求項 12】

CPU を有するワーカーノードのプログラムを格納する記憶媒体は、

前記 CPU が

実行の要求される第 1 のタスクの識別子と、前記第 1 のタスクに入力される第 1 のデータの識別子を受け付け、

前記受け付けた第 1 のデータの識別子に基づき、前記ワーカーノードのメモリに前記第 1 のデータがキャッシュされているかを判定し、

前記キャッシュされているかを判定した結果がキャッシュされていると判定された場合、前記第 1 のデータを前記メモリから取得し、

50

前記取得した第1のデータを入力として前記第1のタスクを実行し、
 前記第1のタスクの実行結果の第2のデータをキャッシュするかを判定し、
 前記キャッシュするかを判定した結果がキャッシュすると判定された場合、前記第2のデータを前記メモリにキャッシュして前記ワーカーノードで前記第2のデータをキャッシュしたことを通知し、他ワーカーノードへ前記第2のデータをコピーさせるプログラムを格納したことを特徴とする記憶媒体。

【請求項13】

前記CPUが、
 他のワーカーノードの障害を検知し、
 前記障害の検知された他のワーカーノードでキャッシュされたデータのコピーである第3のデータを有する場合、前記第3のデータを入力とする第3のタスクの実行状態の情報を取得し、
 前記取得した実行状態の情報を判定し、前記実行状態の情報が実行中である情報と判定した場合、前記第3のデータを前記メモリから削除し、
 前記実行状態の情報が実行前である情報と判定した場合、前記ワーカーノードで前記第3のデータをキャッシュしたことを通知するプログラムをさらに格納したことを特徴とする請求項12に記載の記憶媒体。

10

【発明の詳細な説明】

【技術分野】

20

【0001】

本発明は、分散処理システム、タスク処理方法、記憶媒体に関するものである。

【背景技術】

【0002】

大量データの分析処理を行うシステムにおける性能向上を目的とし、分散処理化による処理性能向上、サーバDRAMメモリ上にデータを配置し処理対象のデータが配置されたサーバにジョブスケジューリングすることによるI/O性能向上を実現するような、データ配置を考慮した分散処理環境が知られている。

【0003】

このような分散処理環境のデータ複製処理において、障害発生時にジョブスケジューラに複製データの配置先を通知することにより、障害復旧時間を短縮する技術が知られている。(特許文献1)

30

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2012-073975号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかし、特許文献1に記載の技術では、データ更新に伴うデータ複製処理の度にジョブスケジューラへの通知処理が発生し、ジョブスケジューラへの負荷集中及び、データ配置情報の更新による定常時の性能劣化及び同期処理にタイムラグが発生する。

40

【0006】

本発明の目的は、定常時の性能を維持しつつ、障害発生時に短時間で復旧することにある。

【課題を解決するための手段】

【0007】

本発明に係る代表的な分散処理システムは、それぞれにワーカーノード識別子が設定され、それぞれがデータをキャッシュするメモリを有する複数のワーカーノードを含む分散処理システムであって、前記ワーカーノード識別子の中の第1のワーカーノード識別子が

50

設定され、第1のタスクの実行結果の第1の出力データを、キャッシュ要求に応じて自ワーカーノードに有される第1のメモリにキャッシュし、他ワーカーノードへ前記第1の出力データをコピーさせる、前記複数のワーカーノードの中の第1のワーカーノードと、前記ワーカーノード識別子の中の第2のワーカーノード識別子が設定され、前記第1の出力データのコピーを、自ワーカーノードに有される第2のメモリにキャッシュする、前記複数のワーカーノードの中の第2のワーカーノードと、前記第1のワーカーノード識別子の情報を有し、前記第1の出力データを入力する第2のタスクを割り当てるワーカーノードを前記第1のワーカーノードであると、前記第1のワーカーノード識別子の情報に基づき選択するマスタノードとを備えたことを特徴とする。

【発明の効果】

10

【0008】

本発明によれば、定常時の性能を維持することが可能となり、障害発生時に短時間で復旧することができる。また、その復旧時間は、データサイズやタスク実行時間に依存せず、所定の時間内とすることが可能である。

【図面の簡単な説明】

【0009】

【図1】分散処理システムの構成の例を示す図である。

【図2】マスタノードの構成の例を示す図である。

【図3】ワーカーノードの構成の例を示す図である。

【図4】タスク管理テーブルの例を示す図である。

20

【図5】タスク配置管理テーブルの例を示す図である。

【図6】キャッシュ管理テーブルの例を示す図である。

【図7】ローカルタスク管理テーブルの例を示す図である。

【図8】ローカルキャッシュ管理テーブルの例を示す図である。

【図9】コピーデータ管理テーブルの例を示す図である。

【図10】タスク登録の処理フローの例を示す図である。

【図11】キャッシュ登録の処理フローの例を示す図である。

【図12】障害タスク復旧の処理フローの例を示す図である。

【図13】データ配置更新処理の情報の更新の例を示す図である。

【図14】データ削除処理の情報の更新の例を示す図である。

30

【図15】ワーカーノード復旧の処理フローの例を示す図である。

【図16】実施例2の分散処理システムの構成の例を示す図である。

【図17】実施例2のワーカーノードの構成の例を示す図である。

【図18】実施例2のデータノードの構成の例を示す図である。

【図19】実施例3のキャッシュ登録の処理フローの例を示す図である。

【発明を実施するための形態】

【0010】

以下、図面を参照しながら、実施の形態を説明する。なお、以後の説明では、「xxxテーブル」等の表現にて情報を説明することがあるが、これら情報はテーブルのデータ構造以外で表現されていてもよい。そのため、データ構造に依存しないことを示すために「xxxテーブル」等について「xxx情報」と呼ぶことがある。各情報の内容を説明する際に、「番号」、「名称」という表現の識別情報が採用されるが、他の種類の識別情報が使用されて良い。以後の説明における「xxx処理」や「xxx部」は、「xxxプログラム」であってもよい。以後の説明における「処理」や「部」を主語とした説明は、プロセッサを主語とした説明としてもよい。プロセッサの処理の一部または全ては、専用ハードウェアによって実現されてもよい。各種プログラムは、プログラム配布サーバや、計算機が読み取り可能な記憶媒体によって各計算機にインストールされてもよい。

40

【実施例1】

【0011】

図1は、分散処理システムの構成の例を示す図である。分散処理システムは、全体とし

50

てデータセンタに設置された計算機システム1であってもよい。この計算機システム1は、マスタノード2及びクライアント端末3がネットワークスイッチ4を介して複数のワーカノード5A、5B、5C(以下、ワーカノード5A、5B、5Cを特に区別する必要のない場合は代表的にワーカノード5と記載し、他の符号も同じ表記とする)と接続されると共に、マスタノード2及びワーカノード5がそれぞれストレージスイッチ6を介してストレージ装置7と接続されることにより構成されている。

【0012】

マスタノード2は、パーソナルコンピュータ又はワークステーション等から構成され、制御プログラム群211と管理テーブル群212とアプリケーション213を有する。マスタノード2は具体的には図2の例に示すように、メモリ21、CPU(Central Processing Unit)22、ネットワークインタフェース23及びディスクインタフェース24を備え、例えばOS(Operating System)の管理下で動作する。

10

【0013】

CPU22は、マスタノード2全体の動作制御を司るCPUであり、メモリ21に格納された後述の制御プログラム群211及び管理テーブル群212に基づいて必要な処理を実行する。メモリ21は、各内容を後述する制御プログラム群211及び管理テーブル群212を記憶するために用いられるほか、CPU22のワークメモリとしても用いられる。

【0014】

ネットワークインタフェース23は、図1に示したネットワークスイッチ4に対応した通信インタフェースであり、マスタノード2が各ワーカノード5と通信する際のプロトコル制御を行う。またディスクインタフェース24は、ストレージスイッチ6に対応した通信インタフェースであり、マスタノード2がストレージ装置7と通信する際のプロトコル制御を行う。

20

【0015】

図1に戻り、ワーカノード5は、マスタノード2と同様にパーソナルコンピュータ又はワークステーション等から構成される。ワーカノード5Aは「W1」で識別され、制御プログラム群511Aと管理テーブル群512Aを有し、ワーカノード5Bは「W2」で識別され、制御プログラム群511Bと管理テーブル群512Bを有し、ワーカノード5Cは「W3」で識別され、制御プログラム群511Cと管理テーブル群512Cを有する。

30

【0016】

制御プログラム群511A、511B、511Cのそれぞれは、記憶されているワーカノード5が異なるだけで、同じ内容のプログラム群である。管理テーブル群512A、512B、512Cのそれぞれは、記憶されているワーカノード5が異なり、同じ項目のテーブル群であって、各項目の内容はそれぞれで異なってもよいし、同じであってもよい。

【0017】

ワーカノード5は、図3に示すように、CPU52、メモリ51、ネットワークインタフェース53及びディスクインタフェース54を備え、例えばOSの管理下で動作する。これらは、それぞれマスタノード2のCPU22、メモリ21、ネットワークインタフェース23及びディスクインタフェース24と同様のものであるため、それらの詳細については説明を省略する。なお、後述するキャッシュにメモリ51が利用されても良く、ワーカノード5の制御プログラム群511と管理テーブル群512の内容については後述する。

40

【0018】

図1に戻り、ストレージ装置7は、複数のディスク装置71を備えて構成される。ストレージ装置7は、複数のディスク装置71を用いてRAID(Redundant Arrays of Inexpensive Disks)グループを構成してもよい。ストレージ装置7は、RAIDグループ上に複数のボリュームを構成してもよい。複数のディスク装置71は、記憶媒体の異なるデ

50

ディスク装置、例えば、HDD (Hard Disk Drive) と、SSD (Solid State Drive) でもよい。

【0019】

ディスク装置は、例えばFC (Fibre Channel) ディスク、SCSI (Small Computer System Interface) ディスク、SATA (Serial ATA) ディスク、ATA (AT Attachment) ディスク又はSAS (Serial Attached SCSI) ディスク等であり、大容量のデータを記憶することのできる記憶媒体である。

【0020】

マスタノード2は、図2に示すように、メモリ21内に制御プログラム群211として、スケジュール管理部2111、タスク配置管理部2112、及びキャッシュ管理部2113を有する。これらの各部はプログラムであり、説明の分かり易さのために分けてあるが、一つに纏めて実現されても良いし、実装上の都合により任意に分けてもよい。

10

【0021】

スケジュール管理部2111は、アプリケーション213からの処理要求に対して、タスク管理テーブル2121を用いて、各タスクの進捗を管理する。タスク配置管理部2112は、スケジュール管理部2112またはタスク実行部5111からの処理要求に対して、タスク配置管理テーブル2122を用いて、各ワーカーノード5へのタスク割り当てを管理する。

【0022】

キャッシュ管理部2113は、スケジュール管理部2111またはデータ配置管理部5113からの処理要求に対して、キャッシュ管理テーブル2123を用いて、各ワーカーノード5に対するキャッシュデータの配置を管理する。

20

【0023】

ワーカーノード5は、図3に示すように、メモリ51内に制御プログラム群511として、タスク実行部5111、ローカルキャッシュ管理部5112、データ配置管理部5113、及びデータ保管部5114を有する。これらの各部はプログラムであり、説明の分かり易さのために分けてあるが、一つに纏めて実現されても良いし、実装上の都合により任意に分けてもよい。

【0024】

タスク実行部5111は、タスク配置管理部2112からの処理要求に対して、ローカルタスク管理テーブル5121を用いて、割り当てられた各タスクの実行ならびに進捗を管理する。ローカルキャッシュ管理部5112は、データ配置管理部5113または他のワーカーノード5上で動作しているローカルキャッシュ管理部5112からの処理要求に対して、ローカルキャッシュ管理テーブル5122を用いて、割り当てられたキャッシュデータを管理する。

30

【0025】

データ配置管理部5113は、タスク実行部5111またはローカルキャッシュ管理部5113からの処理要求に対して、マスタノード2と各ワーカーノード5の間のタスク割り当て情報及びキャッシュデータ情報の整合性を管理する。データ保管部5114は、そのデータ保管部5114を有するワーカーノード5のメモリ51を管理し、コピーデータ管理テーブル5123を用いて、キャッシュデータのメモリ上の配置及び他のワーカーノード5上のデータ保管部5114と連携してコピー処理を実行する。

40

【0026】

なお、各部を制御プログラム群211、511の一部のプログラムとして説明したが、各部というプログラムをCPU22、52が実行することにより物としての各部を構成してもよい。例えば、タスク実行部5111というプログラムをCPU52が実行することにより、タスク実行部という物を構成してもよい。

【0027】

タスク管理テーブル2121は、図4に示すように、タスク識別子欄21211、処理欄21212、入力データ欄21213、出力データ欄21214、キャッシュ要求欄2

50

1 2 1 5、状態欄 2 1 2 1 6、及び優先度欄 2 1 2 1 7を有する。これらの欄の情報は、例えば、1つ目のタスクT1が、ストレージ装置 7 上に配置された file 1 を load 処理によりアプリケーション 2 1 3 が扱うデータオブジェクトD1に変換するタスクであり、既に処理が完了しており、さらに他のタスクと比較して優先度が高いことを表す。

【 0 0 2 8 】

図 4 の例では、タスク管理テーブル 2 1 2 1 の上からタスクが実行されるため、出力データ欄 2 1 2 1 4 のデータオブジェクトが入力データ欄 2 1 2 1 3 の次の行のデータオブジェクトとなっている。また、キャッシュ要求欄 2 1 2 1 5 は、出力データ欄 2 1 2 1 4 のデータオブジェクトをキャッシュへ格納するか否かを示す情報であり、予め設定される。状態欄 2 1 2 1 6 はタスクの状態を示す情報であり、その内容については後述する。

10

【 0 0 2 9 】

タスク配置管理テーブル 2 1 2 2 は、図 5 に示すように、タスク識別子欄 2 1 2 2 1、及びワーカーノード識別子欄 2 1 2 1 3 を有する。タスク識別子欄 2 1 2 2 1 は、タスク管理テーブル 2 1 2 1 のタスク識別子欄 2 1 2 1 1 の情報に対応する情報である。ワーカーノード識別子欄 2 1 2 2 2 は、タスクが割り当てられたワーカーノード 5 を識別する情報である。このため、例えば、1つ目のタスクT1がワーカーノードW2で処理中であることを表す。

【 0 0 3 0 】

キャッシュ管理テーブル 2 1 2 3 は、図 6 に示すように、データ識別子欄 2 1 2 3 1、及びワーカーノード識別子欄 2 1 2 3 2 を有する。データ識別子欄 2 1 2 3 1 は、タスク管理テーブル 2 1 2 1 の入力データ欄 2 1 2 1 3 と出力データ欄 2 1 2 1 4 のデータオブジェクトに対応する情報である。ワーカーノード識別子欄 2 1 2 3 2 は、データオブジェクトがキャッシュに配置されたワーカーノード 5 を識別する情報である。このため、例えば、1つ目のデータオブジェクトD1がワーカーノードW2上に配置されていることを表す。

20

【 0 0 3 1 】

ローカルタスク管理テーブル 5 1 2 1 は、図 7 に示すように、タスク識別子欄 5 1 2 1 1、処理欄 5 1 2 1 2、入力データ欄 5 1 2 1 3、出力データ欄 5 1 2 1 4、キャッシュ要求欄 5 1 2 1 5、状態欄 5 1 2 1 6、及び優先度欄 5 1 2 1 7を有する。これらの情報それぞれは、タスク管理情報テーブル 2 1 2 1 の各欄と同様であるため、その詳細については説明を省略する。

30

【 0 0 3 2 】

ローカルキャッシュ管理テーブル 5 1 2 2 は、図 8 に示すように、データ識別子欄 5 1 2 2 1 を有する。この情報は、例えば、データオブジェクトD2が、このローカルキャッシュ管理テーブル 5 1 2 2 を有するワーカーノード 5 のローカルメモリにキャッシュされていることを表す。

【 0 0 3 3 】

コピーデータ管理テーブル 5 1 2 3 は、図 9 に示すように、データ識別子欄 5 1 2 3 1、コピーデータ識別子欄 5 1 2 3 2、ワーカーノード識別子欄 5 1 2 3 3、及び状態欄 5 1 2 3 4 を有する。データ識別子欄 5 1 2 3 1 は、タスク管理テーブル 2 1 2 1 の入力データ欄 2 1 2 1 3 と出力データ欄 2 1 2 1 4 等のデータオブジェクトに対応する情報である。コピーデータ識別子欄 5 1 2 3 2 は、データ識別子欄 5 1 2 3 1 のデータオブジェクトの複数のコピーそれぞれを識別する情報である。

40

【 0 0 3 4 】

ワーカーノード識別子欄 5 1 2 3 3 は、コピーデータ識別子欄 5 1 2 3 2 のコピーであるデータオブジェクトそれぞれが配置されたワーカーノード 5 を識別する情報である。状態欄 5 1 2 3 4 は、コピーデータ識別子欄 5 1 2 3 2 のコピーであるデータオブジェクトそれぞれがマスタ(M)であるかスレーブ(S)であるかの属性情報である。

【 0 0 3 5 】

これにより、例えば、データ保管部 5 1 1 4 はデータオブジェクトD2に対して独自の識別子を付けたデータオブジェクトM2、C2として管理しており、それぞれワーカーノードW1

50

、W2上に配置し、データオブジェクトM2をマスタ(M)、データオブジェクトC2をスレーブ(S)として、データオブジェクトD2に対するアクセス要求を受信した際には、マスタ(M)であるデータオブジェクトM2をワーカーノードW1から取得し応答することを表す。

【0036】

データ保管部5114による独自の識別子を付加したデータ管理により、タスク実行処理とデータコピー処理を独立して動作させ、タスク実行処理からはデータのコピー管理処理を隠蔽することが可能となり、コピー処理をタスク実行処理に影響を与えることなく実現することができる。なお、マスタ(M)とスレーブ(S)の設定については後述する。

【0037】

図10は、タスク登録の処理フローの例を示す図である。まず、スケジュール管理部2111は、アプリケーション213からの処理要求のタスクあるいはOSを経由しての処理要求のタスクを受け取ると、タスク管理テーブル2121に登録する(SP301)。その際、状態欄21216の値は「未割り当て」とする。

【0038】

スケジュール管理部2111は、キャッシュ管理部2113に対してキャッシュ判定要求を発行し、要求を受けたキャッシュ管理部2113はタスク管理テーブル2121の入力データ欄21213を参照することにより特定した入力データのデータ識別子が、キャッシュ管理テーブル2123のデータ識別子欄21231上に存在するかどうかの情報と、存在する場合は対応するワーカーノード識別欄21232の情報を応答する(SP303)。

【0039】

スケジュール管理部2111は、その応答の情報により入力データが既にキャッシュされているかどうかを判定する(SP304)。スケジュール管理部2111は、ステップSP304の判定で既にキャッシュされている場合、キャッシュ管理テーブル2123のワーカーノード識別子欄21232の情報よりキャッシュデータが配置されているワーカーノード5を特定する。

【0040】

そして、特定したワーカーノード5に対するタスク割り当ての処理要求をタスク配置管理部2112へ発行し、タスク管理テーブル2121の状態欄21216を「割り当て済」に更新する(SP305)。

【0041】

一方、スケジュール管理部2111は、ステップSP304の判定でキャッシュされていない場合、任意のワーカーノード5を選択してタスク配置管理部2112へタスク割り当ての処理要求を発行し、タスク管理テーブル2121の状態欄21216を「割り当て済」に更新する(SP306)。なお、任意のワーカーノード5を選択する際、下記のような選択でもよい。

【0042】

- ・ランダムに選択
- ・ラウンドロビンで選択
- ・負荷の低いワーカーノード5を優先して選択
- ・処理中のタスク数が少ないワーカーノード5を優先して選択
- ・同じ入力データを必要とする複数のタスクが、1つのワーカーノード5に配置されるように選択

また、1つの選択では1つのワーカーノード5に決定できない場合、他の選択との組合せに基づいて選択してもよい。

【0043】

次に、処理要求を受けたタスク配置管理部2112は、処理要求で指定されたワーカーノード5上のタスク実行部5111に対してタスク処理要求を発行した後(SP307)、スケジュール管理部2111へ処理要求に対する応答をする。応答を受けたスケジュール管理部2111は、タスク配置管理テーブル2122に新しい行を追加してタスク処理

10

20

30

40

50

要求を発行したタスク及び要求先のワーカーノード5を登録し、タスク管理テーブル2121の状態欄21216を「処理中」に更新する(S P 3 0 8)。

【0044】

一方、タスク実行部5111は、タスク配置管理部2112からのタスク処理要求を受信すると、その内容をローカルタスク管理テーブル5121に登録する(S P 3 0 9)。次に、ローカルタスク管理テーブル5121からタスクを1つ選択し、選択したタスクの入力データ欄51213の入力データの取得要求をデータ配置管理部5113に発行する(S P 3 1 0)。

【0045】

入力データの取得要求を受信したデータ配置管理部5113は、ローカルキャッシュ管理部5112を介してローカルキャッシュ管理テーブル5122を参照し、指定された入力データがキャッシュされているかを判定する(S P 3 1 1)。ステップS P 3 1 1の判定でキャッシュされていない場合は、ストレージ装置7から該当データを取得した後(S P 3 1 2)、タスク実行部5111に応答する(S P 3 1 3)。

【0046】

この際、タスク管理テーブル2121を参照し、ストレージ装置7から取得したデータがキャッシュ要求のあったデータだった場合、後述するキャッシュ登録処理を行うとしてもよい。

【0047】

キャッシュされていた場合は、ローカルキャッシュ管理部5112及びデータ保管部5114を介してキャッシュデータを取得し、タスク実行部5111に応答する(S P 3 1 4)。応答を受信したタスク実行部5111は、応答に含まれる入力データを用いてステップS P 3 0 9で登録したタスクを実行した後(S P 3 1 5)、ローカルタスク管理テーブル5121のキャッシュ要求欄51215を参照し、処理結果のキャッシュ要求の有無を判定する(S P 3 1 6)。

【0048】

タスク実行部5111は、ステップS P 3 1 6の判定でキャッシュ要求がなかった場合はキャッシュへ登録せずにステップS P 3 1 8に進み、キャッシュ要求があった場合は同一ワーカーノード5内のデータ配置管理部5113にキャッシュ要求を発行する。

【0049】

その要求を受信したデータ配置管理部5113は、キャッシュ登録処理(後述)を行ってタスク実行部5111に応答する(S P 3 1 7)。応答を受信したタスク実行部5111は、ストレージ装置7にタスクの処理結果の出力データを格納した後、実行結果をマスタノード2のタスク配置管理部2112に応答する(S P 3 1 8)。

【0050】

次に、マスタノード2において、ワーカーノード5から受信したタスク実行処理の実行結果に基づいて、タスク配置管理テーブル2122から該当タスクの行を削除した後、実行結果をスケジュール管理部2111に応答する(S P 3 1 9)。応答を受信したスケジュール管理部2111は、タスク管理テーブル2121の状態欄21216のうち、結果を受信したタスクの項目を「処理済」に更新する(S P 3 2 0)。

【0051】

その結果、すべてのタスクの処理が完了したか、すなわちタスク管理テーブル2121の状態欄21216のすべての値が「処理済」になっているかどうかを判定する(S P 3 2 1)。ステップS P 3 2 1の判定で、すべてのタスクが「処理済」だった場合は、タスク登録処理を完了し、処理要求したアプリケーション213あるいはOSへ応答する(S P 3 2 2)。「処理済」でないタスクがあった場合はステップS P 3 0 1の直後に戻る。

【0052】

図11は、キャッシュ登録の処理フローの例を示す図である。この処理は、図10におけるキャッシュ登録処理S P 3 1 7の処理手順である。まず、データ配置管理部5113は、キャッシュ要求として指定されたデータのサイズが、例えばメモリ51等のキャッシ

10

20

30

40

50

メモリの残容量を超過するかどうかを判定する (S P 4 0 1)。

【 0 0 5 3 】

データ配置管理部 5 1 1 3 は、ステップ S P 4 0 1 の判定で、残容量超過だった場合にキャッシュ登録の処理を終了し、残容量以下だった場合にローカリキャッシュ管理部 5 1 1 2 を介して同一ワーカーノード内のデータ保管部 5 1 1 4 へデータ格納要求を発行する (S P 4 0 2)。

【 0 0 5 4 】

データ格納要求を受信したデータ保管部 5 1 1 4 は、自身が管理するメモリ空間にキャッシュ対象のデータを格納して、別のワーカーノード 5 を選択し、コピーデータをキャッシュさせる (S P 4 0 3)。キャッシュしたデータ及びコピーしたデータの情報をコピーデータ管理テーブル 5 1 2 3 に登録した後、ローカルキャッシュ管理部 5 1 1 2 にデータの格納結果を含む応答を発行する (S P 4 0 4)。

10

【 0 0 5 5 】

応答を受信したローカルキャッシュ管理部 5 1 1 2 は、自身が持つローカルキャッシュ管理テーブル 5 1 2 2 にキャッシュデータを登録した後、マスタノード 2 のキャッシュ管理部 2 1 1 3 にキャッシュしたデータ及び配置先のワーカーノード 5 の識別子を含むキャッシュ情報更新要求を発行する (S P 4 0 5)。

【 0 0 5 6 】

情報更新要求を受信したキャッシュ管理部 2 1 1 3 は、情報更新要求の中で指定されたデータ及びワーカーノード 5 の情報をキャッシュ管理テーブル 2 1 2 3 に登録し、ローカルキャッシュ管理部 5 1 1 3 を介してデータ配置管理部 5 1 1 3 に応答した後、キャッシュ登録の処理を終了する (S P 4 0 6)。

20

【 0 0 5 7 】

この際、キャッシュ登録するデータは、マスタデータ、すなわちコピーデータ管理テーブル 5 1 2 3 の状態欄 5 1 2 3 4 が「M」のデータの情報のみであり、コピーデータ、すなわちコピーデータ管理テーブル 5 1 2 3 の状態欄 5 1 2 3 4 が「S」のデータの情報は通知しない。

【 0 0 5 8 】

このように、コピーデータの管理をワーカーノード 5 で実行し、マスタノード 2 から隠ぺいすることにより、マスタノード 2 の負荷を低減することができる。

30

【 0 0 5 9 】

なお、コピーデータを配置するワーカーノード 5 の選定する際、下記のような選択でもよく、1つの選択では所定数のワーカーノード 5 に決定できない場合、他の選択との組合せに基づいて選択してもよい。

【 0 0 6 0 】

- ・コンシステントハッシング法で選択
- ・ラウンドロビンで選択
- ・ランダムに選択
- ・メモリ容量使用率が少ないワーカーノード 5 を優先して選択。

【 0 0 6 1 】

また、ステップ S P 4 0 1 の判定で容量超過だった場合はキャッシュ登録処理を直ちに終了するとしたが、ステップ S P 4 0 1 実行時点でキャッシュしているデータのうち、L R U (Least Recently Used) 等の基準に基づいて選択した参照頻度の低いデータと、新たにキャッシュする対象データとを入れ替えるとしてもよい。

40

【 0 0 6 2 】

図 1 2 は、障害タスク復旧の処理フロー示す図である。まず、タスク配置管理部 2 1 1 2 は、タスク実行部 5 1 1 1 に対するハートビート等の手段を用いて、障害が発生したワーカーノード 5 を検知する (S P 5 0 1)。そして、タスク配置管理部 2 1 1 2 はタスク管理テーブル 2 1 2 1 とキャッシュ管理テーブル 2 1 2 3 を参照し、障害が発生したワーカーノード 5 上で実行中のタスク及びキャッシュされているデータを特定し、再実行が必要

50

なタスクがあるかどうかを判定する (S P 5 0 2)。

【 0 0 6 3 】

タスク配置管理部 2 1 1 2 は、ステップ S P 5 0 2 の判定において、再実行の必要なタスクがなかった場合、障害タスク復旧の処理を終了し、再実行の必要なタスクがあった場合、スケジュール管理部 2 1 1 1 に対して、図 1 0 を用いて説明したタスク登録処理の再実行要求を発行する (S P 5 0 3)。

【 0 0 6 4 】

一方、データ保管部 5 1 1 4 は、ワーカーノード 5 の相互にハートビート等の手段を用いて、障害が発生したワーカーノード 5 を検知し (S P 5 0 4)、コピーデータ管理テーブル 5 1 2 3 を参照し、障害が発生したワーカーノード 5 上に配置されたデータのコピーデータのの中から 1 つをマスタとして昇格させ (S P 5 0 5)、コピーデータ管理テーブル 5 1 2 3 を更新する (S P 5 0 6)。

10

【 0 0 6 5 】

そして、データ配置変更を通知するイベントを、マスタ昇格したデータが配置されているワーカーノード 5 上で実行中のデータ配置管理部 5 1 1 3 に対して発行する (S P 5 0 6)。なお、タスク配置管理部 2 1 1 2 の用いるハートビート信号とワーカーノード 5 の相互に用いるハートビート信号とは同じであってもよいし、異なってもよい。

【 0 0 6 6 】

マスタとして昇格させられるワーカーノード 5 は予め設定されていてもよいし、昇格させるルールが予め設定され、そのルールにしたがって昇格させてもよい。特に、コピーデータが 1 つである場合はそのコピーデータをスレーブとし、そのコピーデータを有するワーカーノード 5 をマスタへ昇格させてもよい。

20

【 0 0 6 7 】

また、マスタ昇格処理の S P 5 0 5 において、障害で消失したデータを一度に復旧するとしたが、マスタノード 2 のスケジュール管理部 2 1 1 1 に問い合わせを行うことで、タスク管理テーブル 2 1 2 1 の入力データ欄 2 1 2 1 3 及び優先度欄 2 1 2 1 7 を参照することにより消失したデータを入力するタスクの優先度を確認し、優先度の高いタスクのデータから逐次的にマスタ昇格させるとしてもよい。

【 0 0 6 8 】

次に、データ配置変更イベントを受信したデータ配置管理部 5 1 1 3 は、スケジュール管理部 2 1 1 1 からタスク配置情報を収集し (S P 5 0 7)、マスタ昇格したデータを使用するタスクが既に別のワーカーノード 5 に割り当てられているかどうか判定する (S P 5 0 8)。ステップ S P 5 0 8 の判定において、「その他」すなわち「処理済」または「処理中」だった場合、ステップ S P 5 1 3 に進む。

30

【 0 0 6 9 】

データ配置管理部 5 1 1 3 は、ステップ S P 5 0 8 の判定において、「別ノードに割り当て済」だった場合、キャッシュ配置管理部 5 1 1 3 はタスク配置管理部 2 1 1 2 を経由してスケジュール管理部 2 1 1 1 に対して該当するタスクの再割り当て要求を発行する。要求を受け取ったスケジュール管理部 2 1 1 1 は、データ配置管理部 5 1 1 3 へタスク再割り当て要求に回答した後、図 1 0 を用いて説明したタスク登録処理を実行する (S P 5 0 3)。

40

【 0 0 7 0 】

一方、スケジュール管理部 2 1 1 1 からの応答を受け取ったデータ配置管理部 5 1 1 3 はステップ S P 5 0 9 に進む。また、ステップ S P 5 0 8 の判定において、「未割り当て」または「自ノードに割り当て済」だった場合、データ配置更新処理として、ローカルキャッシュ管理部 5 1 1 2 に対して、マスタ昇格したデータを含むデータ配置更新要求を発行する (S P 5 0 9)。

【 0 0 7 1 】

更新要求を受信したローカルキャッシュ管理部 5 1 1 2 は、自身が管理するローカルキャッシュ管理テーブル 5 1 2 2 にマスタ昇格したデータを登録した後マスタノード 2 のキ

50

キャッシュ管理部 2 1 1 3 に対して、マスタ昇格したデータ及びそのデータが配置されているワーカーノード 5 の識別情報を含むデータ配置更新要求を発行する (S P 5 1 0)。

【 0 0 7 2 】

要求を受信したキャッシュ管理部 2 1 1 3 は、キャッシュ管理テーブル 2 1 2 3 におけるデータ識別子欄 2 1 2 3 1 を参照し、受信した要求において指定されたデータのワーカーノード識別子欄 2 1 2 3 2 の情報を、受信したワーカーノード 5 の情報に更新した後、ローカルキャッシュ管理部 5 1 1 2 を介してデータ配置管理部 5 1 1 3 に処理結果を応答する (S P 5 1 1)。

【 0 0 7 3 】

応答を受信したデータ配置管理部 5 1 1 3 は、応答内容に基づいて、データ更新処理の成功/失敗を判定する (S P 5 1 2)。ステップ S P 5 1 2 の判定で、処理成功だった場合は、障害タスク復旧処理を終了し、処理失敗だった場合は、データ保管部 5 1 1 4 に対してデータ削除処理要求を発行する (S P 5 1 3)。

【 0 0 7 4 】

要求を受信したデータ保管部 5 1 1 4 は、削除対象データ及び別のワーカーノード 5 上に配置されているコピーデータをメモリ上から削除し (S P 5 1 4)、削除した内容に基づいてコピーデータ管理テーブル 5 1 2 3 を更新した後、データ配置管理部 5 1 1 3 に応答し (S P 5 1 5)、障害タスク復旧処理を終了する。

【 0 0 7 5 】

このように、タスク確認処理 (S P 5 0 7) の結果に基づいて、コピーデータの使用可否を判定することにより、マスタ昇格処理 (S P 5 0 5) よりも先にタスク再実行要求 (S P 5 0 2) が実行されて「処理中」の状態に遷移してしまい、キャッシュデータが別のワーカーノード 5 上で生成されてしまうことにより、ローカルキャッシュ管理テーブル 5 1 2 2 上に同一のキャッシュデータが複数登録されてしまうという状態の発生を抑制する。

【 0 0 7 6 】

すなわち、コピー目的ではないデータが複数のワーカーノード 5 上に配置されてしまうことによるメモリ利用効率の悪化に対し、これを抑制してメモリ利用効率の向上を実現することができる。

【 0 0 7 7 】

さらに、「割り当て済」の状態に遷移してしまっていた場合においては、タスク再割り当て要求を発行することにより、無駄なキャッシュデータ再構築処理を回避することにより復旧時間を短縮することができる。

【 0 0 7 8 】

図 1 3 を用いて、障害処理におけるキャッシュ管理テーブル 1 2 1 3、ローカルキャッシュ管理テーブル 5 1 2 2、及びコピーデータ管理テーブル 5 1 2 3 のデータ配置更新処理 (S P 5 0 4 ~ S P 5 0 6、S P 5 0 9 ~ S P 5 1 1) の具体例について述べる。図 1 3 は、ワーカーノード 5 (W1) に障害が発生し、それを検知したワーカーノード 5 (W2、W3) 及びマスタノード 2 のデータ配置情報を更新する例である。

【 0 0 7 9 】

(1) まず、ワーカーノード 5 (W1) の障害を検知したワーカーノード 5 (W2) は、コピーデータ管理テーブル 5 1 2 3 を参照することにより、ワーカーノード 5 (W1) が保持していたデータ D2 (M2) が消失したために、自身が保持するデータ C2 をマスタ昇格させ、データ D2 として扱う必要があることを判定する。

【 0 0 8 0 】

(2) マスタ昇格処理が必要と判断したワーカーノード 5 (W2) は、コピーデータ管理テーブル 5 1 2 3 の D2 の行の状態欄 5 1 2 3 4 のうち、M2 の値をマスター (M) からスレーブ (S) に、C2 の値をスレーブ (S) からマスタ (M) に更新することにより、データ C2 をデータ D2 のマスタデータとして扱うように変更し、データ配置の変更が発生したことをワーカーノード 5 (W2) のデータ配置管理部 5 1 1 3 にイベント通知する。

【 0 0 8 1 】

10

20

30

40

50

なお、ワーカーノード5 (W1) の障害発生を検知したワーカーノード5 (W3) も同様に、コピーデータ管理テーブル5 1 2 3を更新することにより、ワーカーノード5 (W2) にデータ配置変更が発生したことを検知する。この更新処理は、ワーカーノード5 (W1) の障害を検知したタイミングで実施、または、マスタ昇格処理を行ったワーカーノード5 (W2) からのマスタ昇格イベント通知を受けて更新を実施するとしてもよい。

【0082】

また、障害発生によりワーカーノード5 (W1) のデータを引き継いだワーカーノード5 (W2) は、障害によって消失したデータの管理だけでなく、メモリ容量そのものの管理も引き継ぐ。すなわち、新たなデータD5のキャッシュ要求を受け取った際に、ワーカーノード5 (W2) 自身が管理するメモリ領域は容量を超えるが、ワーカーノード5 (W1) のコピーデータの格納領域として管理している領域には格納可能だった場合、コピーデータの格納領域に対するデータ格納処理として、データD5のコピーデータ識別子であるM5の配置先をワーカーノード5 (W1) とし、データD5のコピーデータ識別子であるC5の配置先をワーカーノード5 (W2) とし、C5の状態をマスタ(M)とするように実行する。

10

【0083】

(3) データ配置変更イベントを受け取ったワーカーノード5 (W2) のデータ配置管理部5 1 1 3は、データ配置更新処理として、同じくワーカーノード5 (W2) で動作しているローカルキャッシュ管理部5 1 1 2にデータD2が新たにワーカーノード5 (W2) に配置されたことを通知する。

【0084】

通知を受け取ったローカルキャッシュ管理部5 1 1 2は、ローカルキャッシュ管理テーブル5 1 2 2のデータ識別子5 1 2 2 1欄にデータD2を追加し、その結果をマスタノード2のキャッシュ管理部2 1 1 3に通知する。通知を受け取ったキャッシュ管理部2 1 1 3は、キャッシュ管理テーブル2 1 2 3のデータD2のワーカーノード識別子欄2 1 2 2 2をW2に更新する。

20

【0085】

以上の一連の処理結果により、マスタノード5 (W2) 及び他のすべてのワーカーノード5が、ワーカーノード5 (W1) の障害によって消失したデータD2がワーカーノード5 (W2) 上に配置されたことを検知することができる。

【0086】

次に、図14を用いて、障害処理におけるデータ配置更新処理よりも前に、タスクが再配置されてしまった場合 (SP307 ~ SP311、SP401 ~ SP406) のデータ削除処理 (SP513 ~ SP515) の具体例について述べる。

30

【0087】

図14は、図13で例示したワーカーノード5 (W1) の障害発生によるデータD2のデータ配置更新処理よりも前に、データD2を入力とするタスクT3が別のワーカーノード5上で処理を開始してしまった場合のワーカーノード5 (W2、W3) 及びマスタノード2のデータ配置情報を更新する例である。

【0088】

(1) まず、マスタノード2のスケジュール管理部2 1 1 1は、タスク配置管理部2 1 1 2を介してワーカーノード5 (W3) に対するタスクT3の処理要求を発行し、タスクT3の状態を「処理中」に更新する。

40

【0089】

タスクT3の処理要求を受け取ったワーカーノード5 (W3) は、タスクT3をローカルタスク管理テーブル5 1 2 1に登録した後、タスクT3の入力データ欄5 1 2 1 3を参照しデータD2が必要であることを検知し、さらに、ローカルキャッシュ管理テーブル5 1 2 2を参照し、自身のメモリ上にデータD2がないことを検知する。

【0090】

(2) 入力データD2がメモリ上にないことを検知したデータ配置管理部5 1 1 3は、ストレージ装置7からデータD2を取得し、データ保管部5 1 1 4にデータD2の登録要求を発

50

行する。要求を受け取ったデータ保管部 5 1 1 4 は、データD2を新たに登録し、マスターデータをワーカーノード 5 (W3) に、コピーデータをワーカーノード 5 (W1) に配置する。なお、コピーデータの配置先は、障害で停止しているワーカーノード 5 (W1) を除外して正常動作しているワーカーノード 5 から選択するとしてもよい。

【 0 0 9 1 】

データ登録の完了を確認したデータ配置管理部 5 1 1 3 は、ローカルキャッシュ管理部 5 1 1 2 にデータD2のキャッシュ登録要求を発行し、要求を受け取ったローカルキャッシュ管理部 5 1 1 2 は、ローカルキャッシュ管理テーブル 5 1 2 2 のデータ識別子 5 1 2 2 1 欄にデータD2を追加し、その結果をマスタノード 2 のキャッシュ管理部 2 1 1 3 に通知する。通知を受け取ったキャッシュ管理部 2 1 1 3 は、キャッシュ管理テーブル 2 1 2 3 のデータD2のワーカーノード識別子欄 2 1 2 3 2 をW3に更新する。

10

【 0 0 9 2 】

この状態で図 1 3 に示したデータ配置更新処理におけるマスタ昇格が完了すると、データD2のマスタデータがワーカーノード 5 (W2、W3) の 2 つ存在することになるが、マスタノード 2 のキャッシュ管理テーブル 2 1 2 3 のワーカーノード識別子欄 2 1 2 3 2 欄には 1 つのワーカーノード 5 しか登録できないため、どちらか一方は、使用されないデータとなってしまう。そこで、処理中のタスクT3を中断させないために、ワーカーノード 5 (W2) 上に配置されたデータD2を削除する。

【 0 0 9 3 】

(3) ワーカーノード 5 (W2) で処理中のデータ配置更新処理におけるタスク確認処理において、マスタノードのタスク管理テーブル 2 1 2 1 を参照すると、マスタ昇格したデータD2を入力とするタスクT3が処理中、すなわちデータD2はすでに他のワーカーノード 5 上に配置されていると判定することができる。

20

【 0 0 9 4 】

(4) したがって、マスタ昇格したデータD2のデータ削除処理として、データ保管部 5 1 1 4 に、データD2の削除要求を発行する。要求を受け取ったデータ保管部 5 1 1 4 は、コピーデータ管理テーブル 5 1 2 3 のワーカーノード識別子欄 5 1 2 3 3 及び状態欄 5 1 2 3 4 を参照し、自身すなわちワーカーノード 5 (W2) がマスタとなっているD2としてデータC2を管理していることを特定し、メモリ上からデータC2を削除した後、他のワーカーノード 5 に対して、D2に関連するコピーデータの削除要求を発行する。

30

【 0 0 9 5 】

以上の一連の処理結果により、障害処理におけるデータ配置更新処理よりも前にタスクが再配置されてしまった場合においても、マスタノード 2 と各ワーカーノード 5 のデータ配置情報の整合性を保ちつつ、無駄なデータを削除し、メモリ領域を効率よく使用することができる。

【 0 0 9 6 】

図 1 5 は、ワーカーノード復旧の処理フローの例を示す図である。まず、復旧したワーカーノード 5 上で動作するタスク実行部 5 1 1 1 は、起動時にタスク配置管理部 2 1 1 2 に対し復旧通知を発行し (S P 6 0 1)、通知を受信したタスク配置管理部 2 1 1 2 は、自身の管理対象として復旧したワーカーノード 5 を登録する (S P 6 0 2)。

40

【 0 0 9 7 】

また、復旧したワーカーノード 5 上で動作するデータ保管部 5 1 1 4 は、他のワーカーノード 5 上で動作するデータ保管部 5 1 1 4 に対し復旧通知を発行し、他のワーカーノード 5 のデータ保管部 5 1 1 4 は復旧を検知する (S P 6 0 4)。復旧したワーカーノード 5 のデータ保管部 5 1 1 4 は、障害発生前に自身が担当していたキャッシュデータを他のワーカーノード 5 から収集するデータ転送処理を行う (S P 6 0 5)。

【 0 0 9 8 】

このデータ転送処理で障害発生前のデータ配置に戻すことにより、復旧したワーカーノード 5 も含めたワーカーノード 5 のメモリ使用量を素早く均等にすることができる。なお、担当していたキャッシュデータがストレージ装置 7 上に配置されていた場合は、ストレ

50

ージ装置 7 から取得してもよい。

【0099】

次に、ステップ SP 606 ~ SP 616 の処理を実行するが、それぞれ、図 12 におけるステップ SP 505 ~ SP 515 と同様の処理であるため、その詳細については説明を省略する。

【0100】

以上で説明したように、コピーデータの管理をワーカーノード 5 のローカルキャッシュ管理部 5112 にオフロードしマスタノード 2 への負荷を低減することにより定常時のアプリケーション 213 の応答性能を維持しつつ、ワーカーノード 5 の障害発生時に、コピーデータの配置をマスタノード 2 のキャッシュ管理部 2113 に通知しキャッシュデータを引き継ぐことにより、障害発生時のシステム停止及び性能劣化を抑止することができる。

10

【0101】

また、障害からの復旧時間は、障害検知時間 (SP 504)、系切り替え処理時間 (SP 505 ~ 506)、イベント通知時間 (SP 509 ~ SP 511) 等を含む。ここで、障害検知時間は、データ保管部 5114 が、SP 504 で障害が発生したワーカーノード 5 を検知するのに要する時間である。障害検知時間は、監視周期などの設定値に依存して変動するが、システム稼働後は一定となる値である。一方、SP 505 や SP 506 の実行に要する系切り替え処理時間や、SP 509 ~ SP 511 の実行に要するイベント通知時間、またデータ復旧のための他の処理の実行に要する時間は、障害検知時間 (SP 504) に対して、十分小さい値となる。したがって、復旧時間は、主にシステム稼働後に一定となる障害検知時間となる。

20

【0102】

一般に、障害時にデータ再構築する場合、障害により消失したデータを再計算又はディスク装置から再取得するため、データ再構築時間は、タスクを再実行する時間 (タスク実行時間) やデータサイズに依存して変動する時間になり、復旧時間が増加する要因となっていた。これに対し、本実施形態では、図 11 のキャッシュ登録処理の SP 403 において、コピーデータを別ワーカーノード 5 に配置しておくので、図 12 の障害タスク復旧処理において、データ配置更新のイベント通知 (SP 509 ~ SP 511) を実行すれば良く、データ再構築時間を削減可能となる。本実施形態の復旧時間は、データサイズやタスク実行時間に依存せず、主にシステム稼働後に一定となる障害検知時間となるため、所定の時間内とすることが可能である。

30

【実施例 2】

【0103】

実施例 2 は、ワーカーノード 5 のデータ保管部 5114 とコピーデータ管理テーブル 5123 をワーカーノード 5 以外の装置に含め、その装置をワーカーノード 5 とストレージ装置 7 (ストレージスイッチ 6) との間に配置する構成とした点で実施例 1 とは異なる。以下の実施例 2 の説明において、実施例 1 と同じ構成、同じ処理に関しては同一の符号を付してその説明を省略し、実施例 1 と異なる点について説明する。

【0104】

図 16 は、実施例 2 における分散処理システムの構成の例を示す図である。また、図 17、図 18 のそれぞれは、実施例 2 におけるワーカーノード 5、データノード 8 の構成の例を示す図である。図 16 ~ 図 18 から明らかなように、実施例 1 において説明したデータ配置管理部 5113 とデータ保管部 5114 を分離し、データ保管部 5114 をデータノード 8 のデータ保管部 8111 とすることにより、データコピー処理の負荷がワーカーノード 5 におけるタスク実行処理に影響を与えず、通常時のタスク処理性能を向上するものである。

40

【0105】

まず、図 16 を参照して、分散処理システムの構成について説明する。マスタノード 2 及びクライアント端末 3 がネットワークスイッチ 4 を介して複数のワーカーノード 5 と接

50

続され、複数のワーカーノード5がネットワークスイッチ9を介し複数のデータノード8と接続されると共に、マスタノード2及び各データノード8がそれぞれストレージスイッチ6を介してストレージ装置7と接続されることにより構成されている。なお、ネットワークスイッチ4とネットワークスイッチ9は共用であって、1つのネットワークスイッチであってよい。

【0106】

ワーカーノード5は、図17に示すように、メモリ51内に制御プログラム群511として、タスク実行部5111、ローカルキャッシュ管理部5112、及びデータ配置管理部5113を有する。これらの各部はプログラムであり、説明の分かり易さのために分けてあるが、一つに纏めて実現されても良いし、実装上の都合により任意に分けてもよい。

また、ワーカーノード5は、ネットワークスイッチ9と接続するためのネットワークインタフェース55を有する。ネットワークインタフェース55はネットワークインタフェース53と共用されてもよい。

10

【0107】

ワーカーノード5は、メモリ51内に管理テーブル群512として、ローカルタスク管理テーブル5121、ローカルキャッシュ管理テーブル5122を有する。ワーカーノード5の各構成については、既に説明したとおりであるが、データ保管部5114とコピーデータ管理テーブル5123がなく、ローカルキャッシュ管理部5112とデータ配置管理部5113はインタフェース55を介してデータノード8と通信可能に構成されている。

20

【0108】

データノード8は、図18に示すように、メモリ81内に制御プログラム群811としてデータ保管部8111を有する。データ保管部8111はデータノード8という物理的な格納場所が異なるため、異なる符号を付しているが、データ保管部5114と同じ処理をするためのプログラムである。実施例1におけるローカルキャッシュ管理部5112及びデータ配置管理部5113とデータ保管部5114との間のやり取りを、実施例2のデータ保管部8111はネットワークインタフェース83経由でワーカーノード5と通信する。

【0109】

また、メモリ81内に管理テーブル群812としてコピーデータ管理テーブル8121を有する。コピーデータ管理テーブル8121もデータノード8という物理的な格納場所が異なるため、異なる符号を付しているが、コピーデータ管理テーブル5123と同じ構成の情報を格納する。

30

【0110】

以上で説明したように、ワーカーノード5上にデータ配置管理部5113を配置し、ワーカーノード5とは物理的に異なるデータノード8上にデータ保管部8111を分離して配置し、ネットワークを介してデータ配置制御させることにより、データコピー処理の負荷をデータノード8が担うことができるようになり、ワーカーノード5におけるタスク実行処理に影響を与えないため、通常時のタスク処理性能を向上させることができる。

【実施例3】

40

【0111】

実施例3は、データをキャッシュとしてメモリ上に登録する際に、キャッシュデータを生成したワーカーノード5とは別のワーカーノード5に登録する点において、実施例1及び実施例2とは異なる。以下では、このキャッシュ登録処理に絞って詳細に説明する。なお、実施例3の説明において、実施例1あるいは実施例2と同じ構成、同じ処理に関しては同一の符号を付してその説明を省略する。

【0112】

図19は、実施例3におけるキャッシュ登録の処理フローの例を示す図である。既に説明したワーカーノード5のキャッシュ登録処理において、メモリ容量不足によりキャッシュデータをメモリ上に配置できない場合に、別のワーカーノード5上にキャッシュデータ

50

を配置することにより、ストレージ装置 7 へのアクセス数の低減し、通常時のタスク処理性能を向上するものである。

【 0 1 1 3 】

図 1 9 に示した処理フローは、図 1 1 を用いて説明したキャッシュ登録処理フローにおける S P 4 0 1 に該当する。データ配置管理部 5 1 1 3 は、自身が動作するワーカーノード 5 に搭載されているメモリ容量、現在の使用量、登録要求で指定されたデータのサイズを比較することにより、データの登録可否を判定する (S P 7 0 1)。

【 0 1 1 4 】

ステップ S P 7 0 1 の判定において登録可能だった場合、データ配置管理部 5 1 1 3 は、同一ワーカーノード 5 上で動作するローカルキャッシュ管理部 5 1 1 2 を介してデータ保管部 5 1 1 4 にキャッシュ登録要求を送付し (S P 7 0 2)、キャッシュ登録の処理を終了する。これにより、自身の動作するワーカーノード 5 のメモリ上にキャッシュデータを配置する。なお、この処理により得られる結果は、図 1 9 を用いて説明した結果と同一の内容になる。

【 0 1 1 5 】

ステップ S P 7 0 1 の判定において登録不可だった場合、データ配置管理部 5 1 1 3 は、別のワーカーノード 5 上で動作するデータ配置管理部 5 1 1 3 に対して、データ登録可否の情報の要求を送信し、取得した情報に基づいて別のワーカーノード 5 のメモリを融通できるか、すなわち別のワーカーノード 5 でデータ登録可能かを判定する (S P 7 0 3)

【 0 1 1 6 】

ステップ S P 7 0 3 の判定において登録可能なワーカーノード 5 が 1 または複数存在した場合、データ配置管理部 5 1 1 3 は、その中の 1 つを配置先として選択し、スケジューリング管理部 2 1 1 1 に対して、これから登録するキャッシュデータを使用するタスクの再配置要求を送付する (S P 7 0 4)。

【 0 1 1 7 】

そして、選択されたワーカーノード 5 上で動作するローカルキャッシュ管理部 5 1 1 2 を介してデータ保管部 5 1 1 4 にキャッシュ登録要求を送付し (S P 7 0 5)、キャッシュ登録の処理を終了する。これにより、選択されたワーカーノード 5 のメモリがリモートメモリとなり、リモートメモリ上にキャッシュデータが配置される。

【 0 1 1 8 】

なお、配置先のワーカーノード 5 を選択する際、下記のように選択してもよい

- ・メモリの空き容量が最も多いものを選択
- ・CPU 負荷が最も低いものを選択
- ・ローカルタスク管理テーブル 5 1 2 1 に登録されているタスク数が最も少ないものを選択

また、1 つの選択では 1 つのワーカーノード 5 を決定できない場合、他の選択との組合せに基づいて選択してもよい。

【 0 1 1 9 】

ステップ S P 7 0 3 の判定において、登録可能なワーカーノード 5 が存在しなかった場合、データ配置管理部 5 1 1 3 は、同一ワーカーノード 5 上で動作するローカルキャッシュ管理部 5 1 1 2 を介し、ストレージ装置 7 に対してデータ配置要求を発行し (S P 7 0 6)、キャッシュ登録の処理を終了する。これにより、ストレージ装置 7 のディスク装置 7 1 上にデータを配置する。

【 0 1 2 0 】

以上で説明したように、ワーカーノード 5 におけるキャッシュ登録処理において、メモリ容量不足によりキャッシュデータをメモリ上には配置できない場合であっても、別のワーカーノード 5 上にキャッシュデータを配置可能となり、ストレージ装置 7 のディスク装置 7 1 へのアクセス数の低減を可能とし、通常時のタスク処理性能を向上させることができる。

10

20

30

40

50

【0121】

なお、以上の実施例1～3において、ダウンタイムに関するSLA要件を管理し、障害発生時の復旧処理において要件の厳しいデータから復旧処理を実行してもよい。また、キャッシュ対象のデータを圧縮及びまたは分割して管理するとしてもよい。

【符号の説明】

【0122】

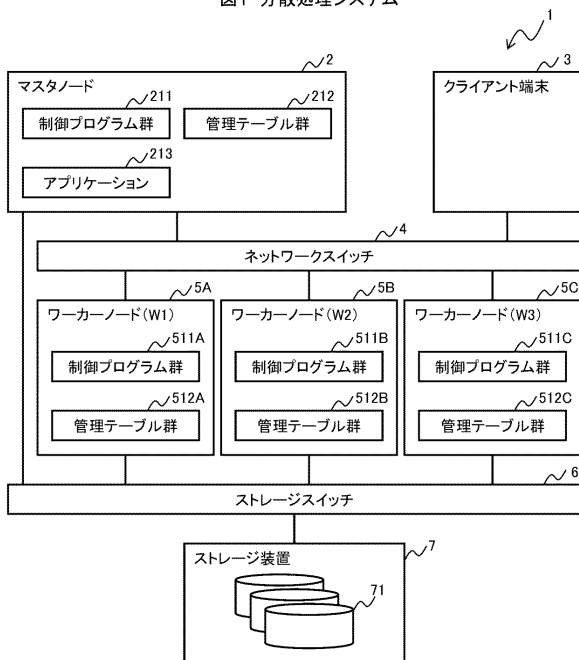
- 2 マスタノード
- 2 1 1 1 スケジュール管理部
- 2 1 1 2 タスク配置管理部
- 2 1 1 3 キャッシュ管理部
- 2 1 2 1 タスク管理テーブル
- 2 1 2 2 タスク配置管理テーブル
- 2 1 2 3 キャッシュ管理テーブル
- 5 ワーカーノード
- 5 1 1 1 タスク実行部
- 5 1 1 2 ローカルキャッシュ管理部
- 5 1 1 3 データ配置管理部
- 5 1 1 4、8 1 1 1 データ保管部
- 5 1 2 1 ローカルタスク管理テーブル
- 5 1 2 2 ローカルキャッシュ管理テーブル
- 5 1 2 3、8 1 2 1 コピーデータ管理テーブル
- 8 データノード

10

20

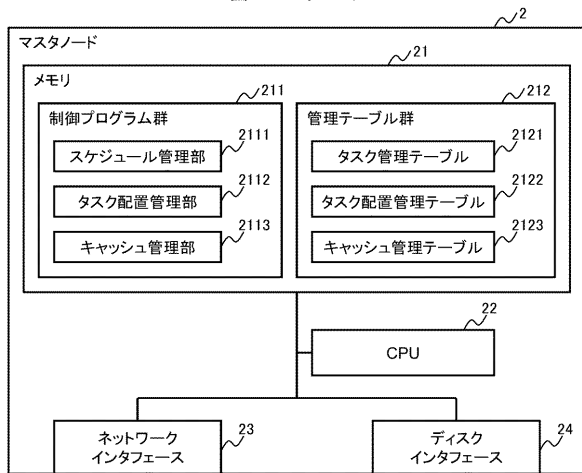
【図1】

図1 分散処理システム



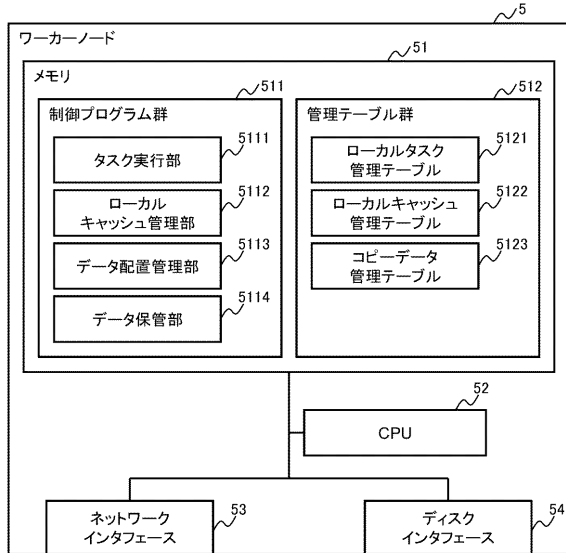
【図2】

図2 マスタノード



【 図 3 】

図3 ワーカーノード



【 図 4 】

図4 タスク管理テーブル

2121

タスク識別子	処理	入力データ	出力データ	キャッシュ要求	状態	優先度
T1	load	file1	D1	×	処理済	高
T2	filter	D1	D2	○	処理中	高
T3	count	D2	D3	×	処理中	中
T4	filter	D2	D4	○	割り当て済	低
T5	map	D4	D5	×	未割り当て	低

【 図 5 】

図5 タスク配置管理テーブル

2122

タスク識別子	ワーカーノード識別子
T1	W2
T2	W1
T5	W2

【 図 6 】

図6 キャッシュ管理テーブル

2123

データ識別子	ワーカーノード識別子
D1	W2
D2	W1
D5	W2

【 図 7 】

図7 ローカルタスク管理テーブル

5121

タスク識別子	処理	入力データ	出力データ	キャッシュ要求	状態	優先度
T2	filter	D1	D2	○	処理中	高

【 図 8 】

図8 ローカルキャッシュ管理テーブル

5122	51221
データ識別子	
D2	

【 図 9 】

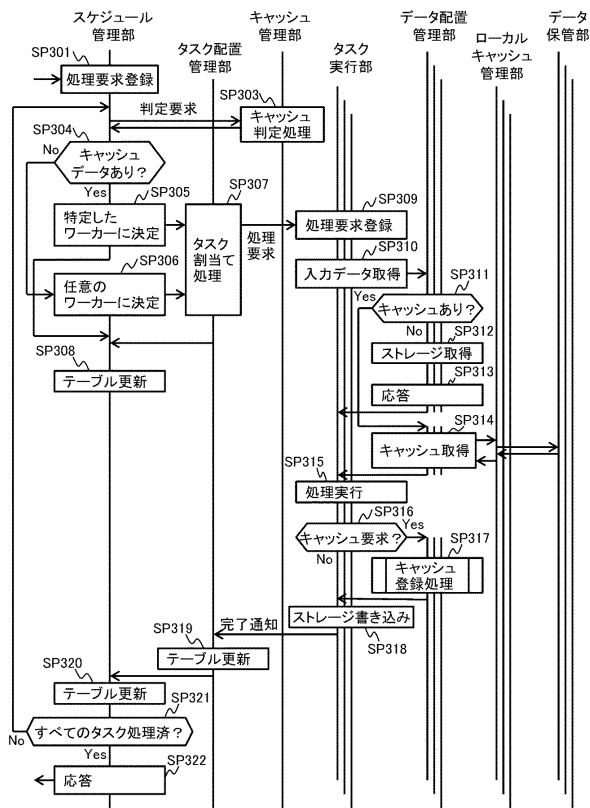
図9 コピーデータ管理テーブル

5123

51231	51232	51233	51234
データ識別子	コピーデータ識別子	ワーカーノード識別子	状態
D2	M2	W1	M
	C2	W2	S
D4	M4	W2	M
	C4	W3	S

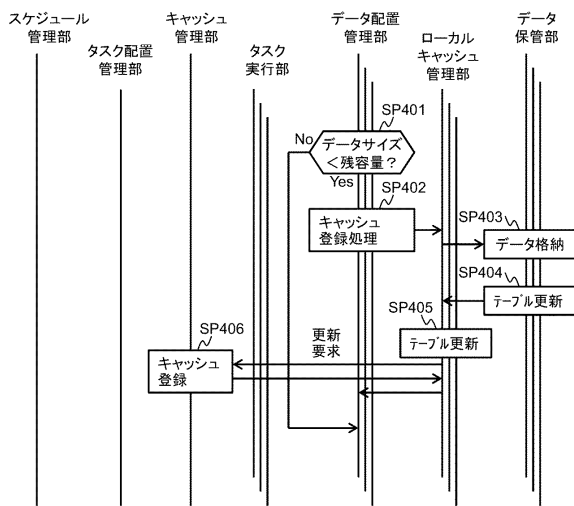
【 図 10 】

図10 タスク登録処理フロー



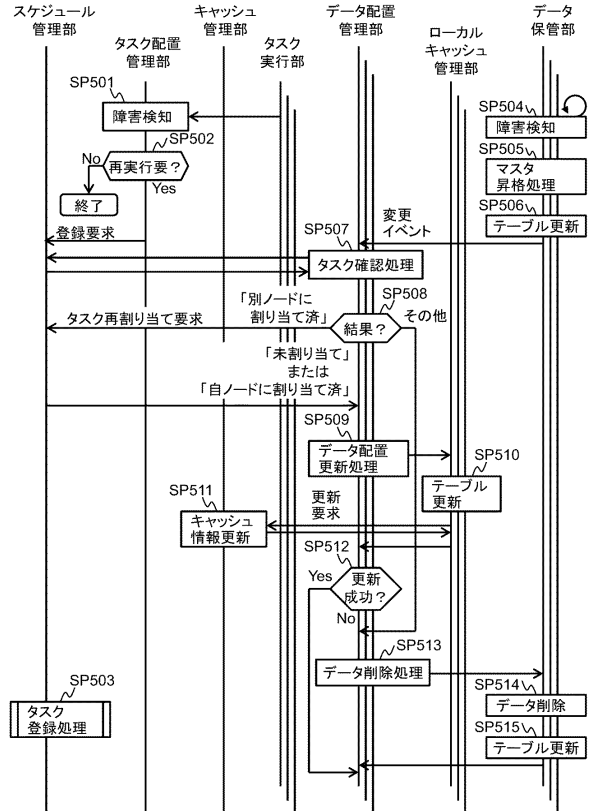
【 図 11 】

図11 キャッシュ登録処理フロー



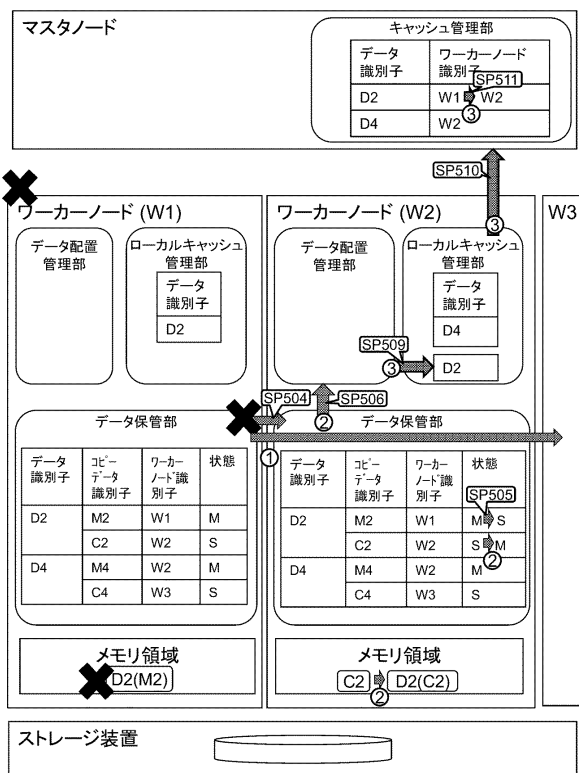
【 図 12 】

図12 障害タスク復旧処理フロー



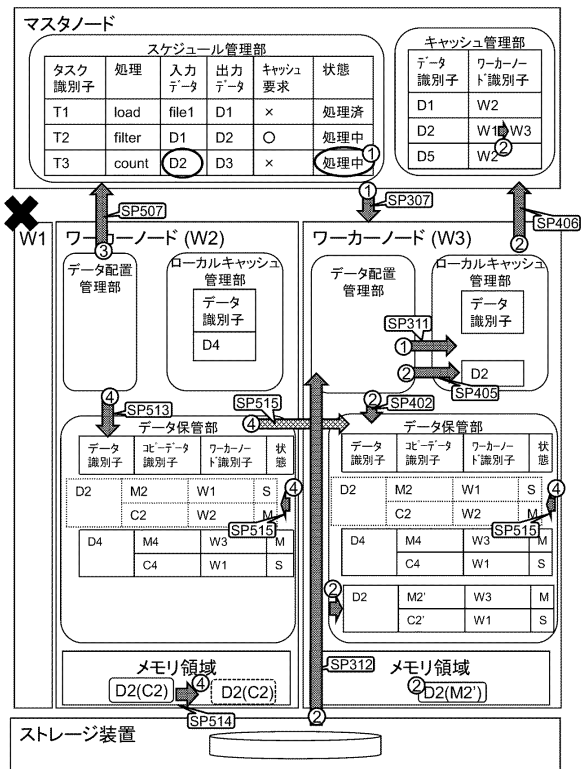
【図13】

図13 データ配置更新処理



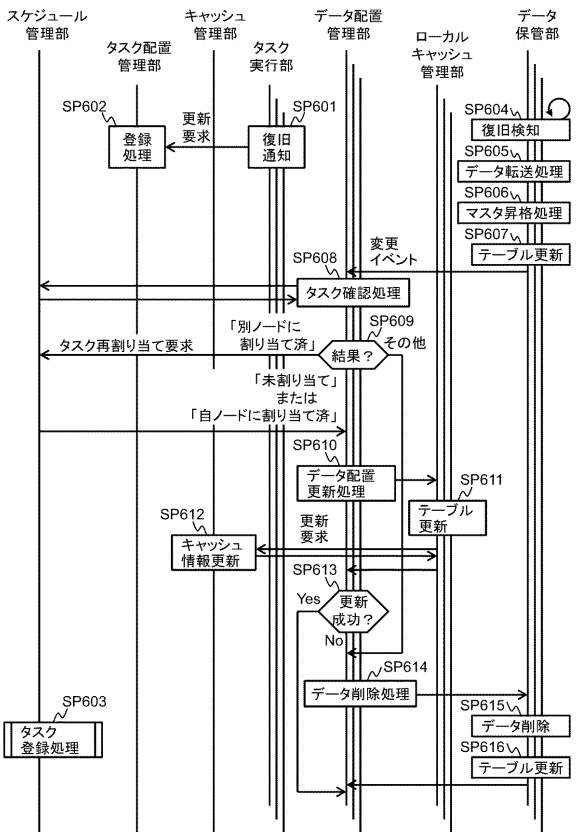
【図14】

図14 データ削除処理



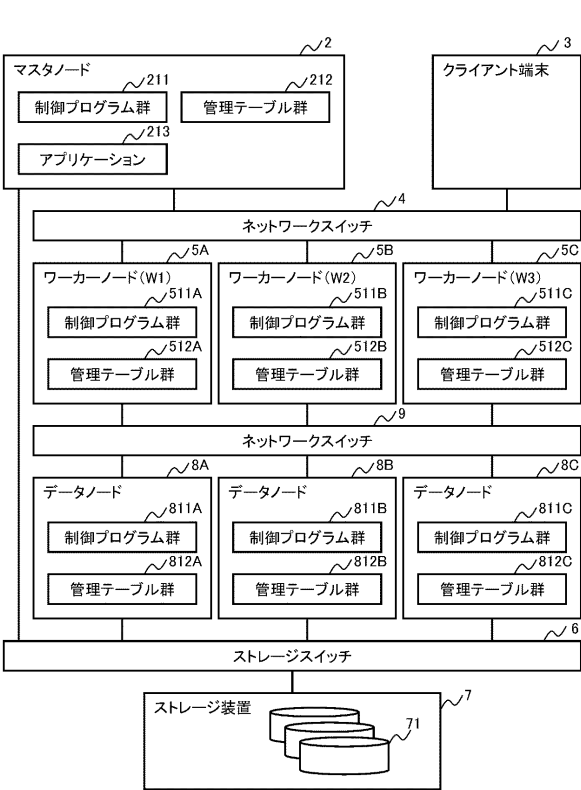
【図15】

図15 ワーカーノード復旧処理フロー



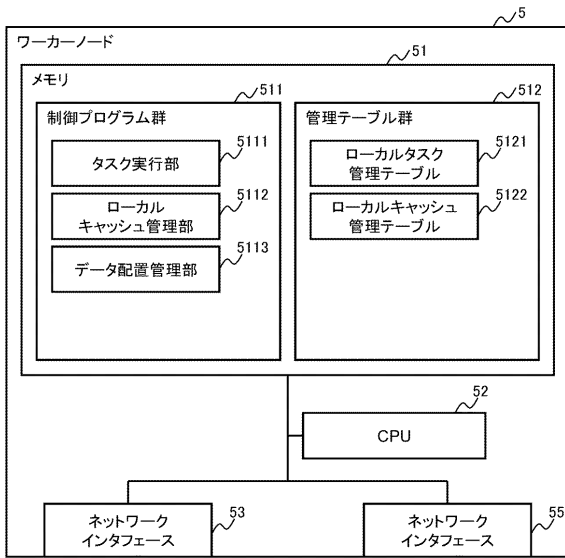
【図16】

図16 分散処理システム



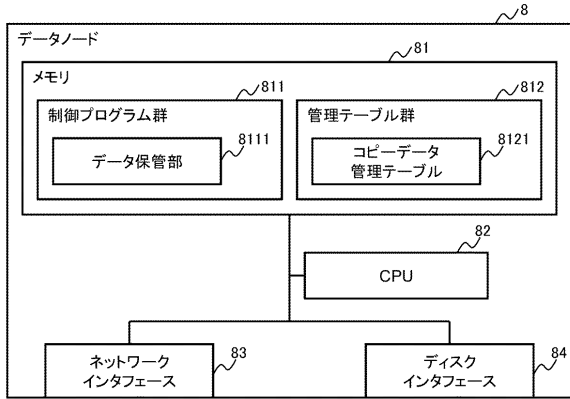
【図17】

図17 ワーカーノード



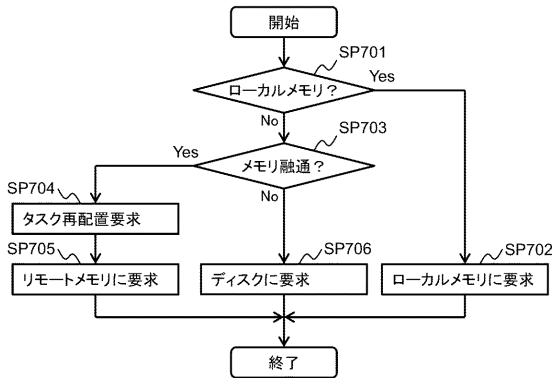
【図18】

図18 データノード



【図19】

図19 キャッシュ登録処理フロー



フロントページの続き

審査官 田中 幸雄

(56)参考文献 特開2012-73975(JP, A)

BU, Yingyi et al., HaLoop: Efficient Iterative Data Processing on Large Clusters, Proceedings of the VLDB Endowment, 2010年 9月, Vol. 3, Issue 1-2, Pages285-296

(58)調査した分野(Int.Cl., DB名)

G06F 9/50

G06F 11/16

G06F 11/20