

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6357807号
(P6357807)

(45) 発行日 平成30年7月18日(2018.7.18)

(24) 登録日 平成30年6月29日(2018.6.29)

(51) Int.Cl. F I
G 0 6 F 9 / 5 0 (2 0 0 6 . 0 1) G 0 6 F 9 / 5 0 1 5 0 E

請求項の数 9 (全 25 頁)

(21) 出願番号	特願2014-43333 (P2014-43333)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成26年3月5日(2014.3.5)		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2015-170054 (P2015-170054A)	(74) 代理人	100089118 弁理士 酒井 宏明
(43) 公開日	平成27年9月28日(2015.9.28)	(72) 発明者	上田 晴康 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	平成28年11月2日(2016.11.2)	(72) 発明者	松田 雄一 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	前田 高光 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

最終頁に続く

(54) 【発明の名称】 タスク割当プログラム、タスク実行プログラム、マスタサーバ、スレーブサーバおよびタスク割当方法

(57) 【特許請求の範囲】

【請求項1】

複数のMap処理タスクを含むMap処理および複数のReduce処理タスクを含むReduce処理を含むMapReduce処理の各タスクを複数のスレーブサーバに割当ててマスタサーバに、

前記複数のスレーブサーバ夫々に、前記複数のMap処理タスクを割り当て、

前記複数のMap処理タスクの全部が完了する前に、対象となる複数のMap処理タスクの実行結果を用いる前記複数のReduce処理タスク夫々が処理対象とするデータ量を、処理が完了した各Map処理タスクから送信される実行結果に含まれるデータ量に基づいて予測し、

前記複数のReduce処理タスクを前記複数のスレーブサーバ夫々に割当てて際に、予測された前記データ量が閾値を超えるReduce処理タスクである高負荷Reduce処理タスクを割当ててスレーブサーバに、前記高負荷Reduce処理タスクの分散処理を指示する、

処理を実行させることを特徴とするタスク割当プログラム。

【請求項2】

前記予測する処理は、全Map処理タスクにおける所定の割合の処理が終了した場合、または、最初のMap処理タスクが終了してから所定時間が経過した場合に、前記複数のReduce処理タスク夫々のデータ量を予測することを特徴とする請求項1に記載のタスク割当プログラム。

【請求項3】

複数のMap処理タスクを含むMap処理および複数のReduce処理タスクを含むReduce処理を含むMapReduce処理の各タスクを実行するスレーブサーバに、

10

20

前記各タスクの割当てを実行するマスタサーバからReduce処理タスクが割当てられた場合に、当該Reduce処理タスクが分散処理の実行指示が付された高負荷Reduce処理タスクであるか否かを判定し、

前記高負荷Reduce処理タスクであると判定された場合、前記高負荷Reduce処理タスクを分散処理で実行し、前記高負荷Reduce処理タスクではない通常のReduce処理タスクと判定された場合、前記通常のReduce処理タスクを分散させずに実行する

処理を実行させることを特徴とするタスク実行プログラム。

【請求項 4】

前記マスタサーバから割り当てられた前記Map処理タスクの実行が完了した場合に、完了した前記Map処理タスクの処理結果を用いるReduce処理タスクに関するデータ量を含む完了結果を前記マスタサーバに送信する処理を、前記スレーブサーバにさらに実行させることを特徴とする請求項 3 に記載のタスク実行プログラム。

10

【請求項 5】

前記実行する処理は、前記スレーブサーバが有するプロセッサの数、ディスクの数、予め指定された数の少なくとも 1 つを用いて、前記高負荷Reduce処理タスクをサブタスクに分割して、複数のプロセッサで並列実行させることを特徴とする請求項 3 または 4 に記載のタスク実行プログラム。

【請求項 6】

前記実行する処理は、前記高負荷Reduce処理タスクではない通常のReduce処理タスクの実行中に前記高負荷Reduce処理タスクが割当てられた場合、前記通常のReduce処理タスクの実行を中止して、前記高負荷Reduce処理タスクを前記サブタスクに分割して実行することを特徴とする請求項 5 に記載のタスク実行プログラム。

20

【請求項 7】

複数のMap処理タスクを含むMap処理および複数のReduce処理タスクを含むReduce処理を含むMapReduce処理の各タスクを複数のスレーブサーバに割当てるマスタサーバにおいて

前記複数のスレーブサーバ夫々に、前記複数のMap処理タスクを割り当てる割当部と、前記複数のMap処理タスクの全部が完了する前に、対象となる複数のMap処理タスクの実行結果を用いる前記複数のReduce処理タスク夫々が処理対象とするデータ量を、処理が完了した各Map処理タスクから送信される実行結果に含まれるデータ量に基づいて予測する予測部と、

30

前記複数のReduce処理タスクを前記複数のスレーブサーバ夫々に割当てる際に、予測された前記データ量が閾値を超えるReduce処理タスクである高負荷Reduce処理タスクを割当てるスレーブサーバに、前記高負荷Reduce処理タスクの分散処理を指示する送信部とを有することを特徴とするマスタサーバ。

【請求項 8】

複数のMap処理タスクを含むMap処理および複数のReduce処理タスクを含むReduce処理を含むMapReduce処理の各タスクを実行するスレーブサーバにおいて、

前記各タスクの割当てを実行するマスタサーバからReduce処理タスクが割当てられた場合に、当該Reduce処理タスクが分散処理の実行指示が付された高負荷Reduce処理タスクであるか否かを判定する判定部と、

40

前記高負荷Reduce処理タスクであると判定された場合、前記高負荷Reduce処理タスクを分散処理で実行し、前記高負荷Reduce処理タスクではない通常のReduce処理タスクと判定された場合、前記通常のReduce処理タスクを分散させずに実行する処理実行部とを有することを特徴とするスレーブサーバ。

【請求項 9】

複数のMap処理タスクを含むMap処理および複数のReduce処理タスクを含むReduce処理を含むMapReduce処理の各タスクを複数のスレーブサーバに割当てるマスタサーバが、

前記複数のスレーブサーバ夫々に、前記複数のMap処理タスクを割り当て、前記複数のMap処理タスクの全部が完了する前に、対象となる複数のMap処理タスクの実

50

行結果を用いる前記複数のReduce処理タスク夫々が処理対象とするデータ量を、処理が完了した各Map処理タスクから送信される実行結果に含まれるデータ量に基づいて予測し、

前記複数のReduce処理タスクを前記複数のスレーブサーバ夫々に割当てる際に、予測された前記データ量が閾値を超えるReduce処理タスクである高負荷Reduce処理タスクを割当てるスレーブサーバに、前記高負荷Reduce処理タスクの分散処理を指示する、

処理を実行することを特徴とするタスク割当方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、タスク割当プログラム、タスク実行プログラム、タスク割当装置、タスク実行装置およびタスク割当方法に関する。

10

【背景技術】

【0002】

クラウドコンピューティングの普及に伴い、クラウド上に保存される大量のデータを複数のサーバで分散して処理を実行する分散処理システムが利用されている。分散処理システムとしては、HDFS (Hadoop Distributed File System) とMapReduce処理とを基盤技術とするHadoop (登録商標) が知られている。

【0003】

HDFSは、複数のサーバにデータを分散格納するファイルシステムである。MapReduceは、HDFS上のデータをタスクと呼ばれる単位で分散処理する仕組みであり、Map処理、Shuffleソート処理、Reduce処理を実行する。

20

【0004】

MapReduceによる分散処理においては、マスタサーバが、ハッシュ関数等を用いて、複数のスレーブサーバに対してMap処理やReduce処理のタスクを割り当てるとともに、分割したデータを各スレーブサーバに送信する。そして、各スレーブサーバが、割り当てられたタスクを実行する。

【0005】

スレーブサーバに対するタスクの割り当ては、ハッシュ関数等を用いることにより、均等に行われる。その一方で、各Reduceタスクに対応する処理量は、Reduceタスクに対応するキー等に関連づけられたReduce対象のデータ量等により均等になるとは限らない。

30

【0006】

Reduceタスクに対応する処理量により、各スレーブサーバでの処理完了時間が異なることから、複数タスクからなるジョブ全体の完了が、最も処理が遅いスレーブサーバの処理完了に左右されることとなる。このため、全Reduceタスクをスレーブサーバに割り当てた後、各Reduceタスクに対応する処理量が均等となるように、データ量を調整する技術が知られている。

【先行技術文献】

【特許文献】

【0007】

【特許文献1】特開2012-118669号公報

40

【発明の概要】

【発明が解決しようとする課題】

【0008】

しかしながら、各Reduceタスクの処理量が不均一となることは、入力データやMap処理の結果など様々な影響により異なるので、上記技術のような調整処理を行うことが、ジョブ全体の完了時間を早くするものとは限らない。

【0009】

例えば、Reduceタスクの処理量を調整するには、全Map処理の完了を待ってから調整することになるので、各スレーブサーバでのReduce処理の実行開始が遅れることになり、却ってジョブ全体の処理時間が長くなる場合もある。

50

【0010】

1つの側面では、ジョブ全体の完了時間を短縮することができるタスク割当プログラム、タスク実行プログラム、タスク割当装置、タスク実行装置およびタスク割当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】

【0011】

第1の態様では、第1のサーバ装置に、複数の第2のサーバ装置夫々に第1の処理を割り当てる処理を実行させる。第1のサーバ装置に、前記複数の第2のサーバ装置夫々に割り当てる、前記第1の処理タスクの実行結果を用いて実行される第2の処理タスクに関連する前記第1の処理タスクの完了通知を受信した場合に、前記第2の処理タスクの処理量を見積もる処理を実行させる。第1のサーバ装置に、前記第2の処理タスクを割り当てる前記第2のサーバ装置に、見積もった前記処理量に関連した情報を送信する処理を実行させる。

10

【発明の効果】

【0012】

1つの側面として、ジョブ全体の完了時間を短縮することができる。

【図面の簡単な説明】

【0013】

【図1】図1は、実施例1に係る分散処理システムの全体構成例を示す図である。

【図2】図2は、実施例1に係るマスタサーバの機能構成を示す機能ブロック図である。

20

【図3】図3は、ジョブリストDBに記憶される情報の例を示す図である。

【図4】図4は、タスクリストDBに記憶される情報の例を示す図である。

【図5】図5は、Map処理の完了通知の例を示す図である。

【図6】図6は、実施例1に係るスレーブサーバの機能構成を示す機能ブロック図である。

【図7】図7は、Map処理を説明する図である。

【図8】図8は、Shuffle処理を説明する図である。

【図9】図9は、Reduce処理を説明する図である。

【図10】図10は、Reduce処理タスクの処理量を予測してフラグを設定する処理を説明する図である。

30

【図11】図11は、実施例1に係るマスタサーバが実行する処理の流れを示すフローチャートである。

【図12】図12は、マスタサーバが実行する該当タスクの完了処理の流れを示すフローチャートである。

【図13】図13は、実施例1に係るスレーブサーバが実行する処理の流れを示すフローチャートである。

【図14】図14は、スレーブサーバが実行するReduce処理タスクの起動処理の流れを示すフローチャートである。

【図15】図15は、スレーブサーバが実行するReduce処理タスクの分割処理の流れを示すフローチャートである。

40

【図16】図16は、実施例2に係るReduce処理タスクの割当処理の流れを示すフローチャートである。

【図17】図17は、実施例3に係るReduce処理タスクの割当処理の流れを示すフローチャートである。

【図18】図18は、各サーバのハードウェア構成例を示す図である。

【発明を実施するための形態】

【0014】

以下に、本願の開示するタスク割当プログラム、タスク実行プログラム、タスク割当装置、タスク実行装置およびタスク割当方法の実施例を図面に基づいて詳細に説明する。なお、この実施例によりこの発明が限定されるものではない。なお、各実施例は、適宜組み

50

合わせることができる。

【実施例 1】

【0015】

[全体構成]

図 1 は、実施例 1 に係る分散処理システムの全体構成例を示す図である。図 1 に示すように、この分散処理システムは、入出力 D B (DataBase) サーバ 2、マスタサーバ 1 0、複数のスレーブサーバ 3 0 がネットワーク 1 を介して互いに通信可能に接続される。

【0016】

この分散処理システムでは、Hadoop (登録商標) などの分散処理フレームワークを使用した分散処理アプリケーションが各計算機で実行されており、データ基盤として HDFS などを使用する。

10

【0017】

入出力 D B サーバ 2 は、分散処理対象のデータのメタ情報等を記憶するデータベースサーバである。例えば、入出力 D B サーバ 2 が記憶するメタ情報は、どのデータがどのスレーブサーバ 3 0 に格納されているかを特定するのに使用される。

【0018】

マスタサーバ 1 0 は、分散処理システムを統括的に管理するサーバである。例えば、マスタサーバ 1 0 は、入出力 D B サーバ 2 が記憶するメタ情報から、どのデータがいずれのスレーブサーバ 3 0 に格納されているのかを特定する。また、マスタサーバ 1 0 は、各スレーブサーバ 3 0 に割当てたタスクやジョブなどを管理し、Map 処理や Reduce 処理などのタスクをスレーブサーバ 3 0 に割当てる。

20

【0019】

各スレーブサーバ 3 0 は、分散処理アプリケーションを実装し、Map 処理や Reduce 処理を実行して、HDFS で管理されるデータを分散処理するサーバである。例えば、スレーブサーバ 3 0 は、複数のプロセッサ、複数のディスクを有する。また、各スレーブサーバ 3 0 は、一意に識別される識別子が割り与えられる。

【0020】

このスレーブサーバ 3 0 は、入出力 D B サーバ 2 から取得したデータに対して、マスタサーバ 1 0 から割当てられた Map 処理のタスクを実行する。また、スレーブサーバ 3 0 は、各スレーブサーバの Map 処理結果を用いて、Shuffle ソート処理を実行し、マスタサーバ 1 0 から割り当てられた Reduce 処理のタスクを実行する。

30

【0021】

ここで、各処理について説明する。Map 処理は、ユーザが定義した Map 関数を実行する処理である。例えば、Map 処理は、入力データから中間結果として「Key、Value」のペアを出力する。Shuffle ソート処理は、Map 処理の結果を「Key」でソートし、同じ「Key」を有する「Key、Value」ペアをマージする。Reduce 処理は、ユーザが定義した Reduce 関数を実行する処理である。例えば、Reduce 処理は、Shuffle ソート処理の結果から、同じ「Key」の「Value」に対して重ね合わせ処理を実行して、新しい形式の「Key、Value」のペアを生成する。

【0022】

40

このような状態において、マスタサーバ 1 0 は、複数のスレーブサーバ 3 0 夫々に Map 処理を割り当てる。マスタサーバ 1 0 は、複数のスレーブサーバ 3 0 夫々に割り当てる、Map 処理タスクの実行結果を用いて実行される Reduce 処理タスクに関連する Map 処理タスクの完了通知を受信した場合に、Reduce 処理タスクの処理量を見積もる。マスタサーバ 1 0 は、Reduce 処理タスクを割り当てるスレーブサーバ 3 0 に、見積もった処理量に関連した情報を送信する。

【0023】

各スレーブサーバ 3 0 は、Map 処理タスクの実行結果を用いる Reduce 処理タスクをスレーブサーバ 3 0 に割り当てるマスタサーバ 1 0 から、Reduce 処理タスクに処理量に関連した情報を受信する。各スレーブサーバ 3 0 は、マスタサーバ 1 0 から割り当てられた Redu

50

ce処理タスクを実行する際に、関連した情報に応じて、Reduce処理タスクの処理方法を変更する。

【 0 0 2 4 】

このように、マスタサーバ10が、一部のMap処理タスクの終了結果から処理量が多いReduce処理タスクを検出してスレーブサーバ30に通知するので、スレーブサーバ30がそのReduce処理タスクを分割して並列処理でき、全処理の完了時間を短縮できる。

【 0 0 2 5 】

[マスタサーバの構成]

図2は、実施例1に係るマスタサーバの機能構成を示す機能ブロック図である。図2に示すように、マスタサーバ10は、通信制御部11と、記憶部12と、制御部13とを有する。

10

【 0 0 2 6 】

通信制御部11は、スレーブサーバ30などの他装置と通信を実行する処理部であり、例えばネットワークインタフェースカードなどである。例えば、通信制御部11は、各スレーブサーバ30に、Map処理タスクやReduce処理タスクを送信する。また、通信制御部11は、Map処理結果等を各スレーブサーバ30から受信する。

【 0 0 2 7 】

記憶部12は、ジョブリストDB12aとタスクリストDB12bとを有する記憶装置であり、例えばメモリやハードディスクなどである。また、記憶部12は、制御部13が実行するプログラムなどを記憶する。

20

【 0 0 2 8 】

ジョブリストDB12aは、分散処理対象のジョブ情報を記憶するデータベースである。図3は、ジョブリストDB12aに記憶される情報の例を示す図である。図3に示すように、ジョブリストDB12aは、「JobID、総Mapタスク数、総Reduceタスク数、Reduce割当許可」を対応付けて記憶する。

【 0 0 2 9 】

ここで記憶される「JobID」は、ジョブを識別する識別子である。「総Mapタスク数」は、ジョブに含まれるMap処理タスクの総数である。「総Reduceタスク数」は、ジョブに含まれるReduce処理タスクの総数である。「Reduce割当許可」は、Reduce処理タスクが割当可能な状態か否かを示し、可能な場合は「true」が設定され、可能ではない場合は「false」が設定され、ジョブの新規追加時も「false」が設定される。なお、「JobID、総Mapタスク数、総Reduceタスク数」は、管理者等によって設定更新される。

30

【 0 0 3 0 】

図3の例では、「JobID」が「1」のジョブは、4つのMap処理タスクと2つのReduce処理タスクで構成され、現在はまだ割当することができない状態であることを示す。同様に、「JobID」が「2」のジョブは、4つのMap処理タスクと2つのReduce処理タスクで構成され、現在はまだ割当することができない状態であることを示す。

【 0 0 3 1 】

タスクリストDB12bは、Map処理タスクやReduce処理タスクに関する情報を記憶するデータベースである。図4は、タスクリストDBに記憶される情報の例を示す図である。図4に示すように、タスクリストDB12bは、「JobID、TaskID、種別、Reduceの項番、データのあるスレーブID、状態、割り当てスレーブID、必要スロット数、処理データ量、フラグ」を記憶する。

40

【 0 0 3 2 】

ここで記憶される「JobID」は、ジョブを識別する識別子である。「TaskID」は、タスクを識別する識別子である。「種別」は、Map処理やReduce処理を示す情報である。「データのあるスレーブID」は、Map処理対象のデータを保持するスレーブサーバ30を識別する識別子であり、例えばホスト名などである。「状態」は、該当タスクが処理完了(Done)状態、実行中(Running)、割り当て前(Not assigned)のいずれであることを示す。「Reduceの項番」は、該当Reduceの実行順を示す。「割り当てスレーブID」は、タスクが

50

割当てられたスレーブサーバ30を識別する識別子であり、例えばホスト名などである。「必要スロット数」は、タスクを実行するのに使用するスロット数である。「処理データ量」は、該当Reduce処理タスクのデータ量である。「フラグ」は、該当Reduce処理タスクの処理方法の変更を指示するか否かを示し、変更を指示する場合は「true」が設定される。

【0033】

図4の場合、「jobID」が「1」であるジョブで、1スロットを用いるMap処理タスク「1_m_1」が「Node1」のスレーブサーバ30に割当てられる。そして、この「Node1」のスレーブサーバ30は、「Node1」のスレーブサーバ30と「Node2」のスレーブサーバ30とからデータを取得して、Map処理を実行し、実行が完了していることを示す。

10

【0034】

また、「jobID」が「1」であるジョブで、2番目に実行される1スロットを用いるReduce処理タスク「1_r_2」が「Node3」のスレーブサーバ30に割当てられる。また、Reduce処理タスク「1_r_2」のデータ量は、「25000」であり、フラグに「true」が設定されている。そして、この「Node3」のスレーブサーバ30は、Reduce処理タスクを分割して、並列に実行中であることを示す。

【0035】

なお、JobID、TaskID、種別、Reduce項番については、ジョブリストDB12aに記憶される情報にしたがって生成される。データのあるスレーブIDは、入出力DBサーバ2が記憶するメタ情報等により特定することができる。状態は、タスクの割り当て状況やスレーブサーバ30からの処理結果等によって更新される。割り当てスレーブIDは、タスクを割当て時点で更新される。必要スロット数は、1タスクについて1スロットなどのように予め指定することができる。処理データ量は、Map処理の終了結果から予測することができる。フラグは、処理データ量が閾値を超えるか否かによって設定される。

20

【0036】

制御部13は、Map割当部14、予測部15、Reduce処理部16を有する処理部であり、例えばプロセッサなどの電子回路である。また、制御部13は、マスタサーバ10全体の処理を司る。

【0037】

Map割当部14は、各ジョブにおけるMap処理のタスクであるMap処理タスクを1つ以上スレーブサーバ30に割り当てる処理部である。具体的には、Map割当部14は、データのあるスレーブIDの情報等を用いて、各Map処理タスクをスレーブサーバ30に割り当てる。そして、Map割当部14は、図4に示した「割当スレーブID」や「状態」等を更新する。

30

【0038】

例えば、Map割当部14は、スレーブサーバ30等からMap処理タスクの割当要求を受信した場合に、タスクリストDB12bを参照して「状態」が「Not assigned」のMap処理タスクを特定する。続いて、Map割当部14は、割当要求を送信したスレーブサーバ30のIDが「データのあるスレーブID」に含まれるMap処理タスクがあればそのMap処理タスクを優先して選び、そのようなMap処理タスクがなければ任意の方法でMap処理タスクを選び、割当対象のMap処理タスクとする。その後、Map割当部14は、割当要求を送信したスレーブサーバ30のIDを、割当対象のMap処理タスクの「スレーブサーバID」に格納する。

40

【0039】

その後、Map割当部14は、特定した割当先のスレーブサーバ30に、TaskID、データのあるスレーブID、必要スロット数等を通知して、Map処理タスクを割当てる。また、Map割当部14は、割当てたMap処理タスクの「状態」を「Not assigned」から「Running」に更新する。

【0040】

予測部15は、Map処理タスクの実行結果を用いて、Reduce処理タスクの処理量を見積

50

もる処理部である。具体的には、予測部15は、スレーブサーバ30から通知されるMap処理の完了通知から、各Reduce処理タスクのデータ量を取得する。

【0041】

予測部15は、このようにして、所定数のMap処理の結果から取得されたReduce処理タスクのデータ量を加算して、Reduce処理タスクのデータ量を見積もる。そして、予測部15は、見積もったReduce処理タスクのデータ量をタスクリストDB12bの処理データ量に格納し、Reduce処理タスクのデータ量が所定値以上であればフラグにtrueを格納する。なお、予測部15は、完了通知を受信したMap処理タスクの「状態」を「Running」から「Done」に更新する。

【0042】

図5は、Map処理の完了通知の例を示す図である。図5に示す完了通知は、各スレーブサーバ30がマスタサーバ10に送信する完了通知である。図5に示すように、完了通知は、「通知種別、JobID、完了Map TaskID、Mapタスクを実行したスレーブID」から構成されるMap完了内容と、「通知種別、JobID、完了Map TaskID、Reduce TaskID、データ量」から構成されるMap完了内容を含む。

【0043】

ここで記憶される「通知種別」は、Map処理の完了通知かReduce情報かを示す情報であり、Map処理の完了通知の場合には「Map完了」が設定され、Reduce情報の場合には「Reduceデータ量」が設定される。「JobID」には、Map処理が属するジョブの識別子が設定される。「完了Map TaskID」は、完了したMap処理タスクを特定する識別子が設定される。「Mapタスクを実行したスレーブID」は、当該Map処理タスクを実行し、完了通知を送信したスレーブサーバの識別子が設定される。「Reduce TaskID」は、当該Map処理の実行結果からデータ量が判明したReduce処理タスクを特定する識別子が設定される。「データ量」は、当該Map処理の実行結果から判明したReduce処理タスクデータ量設定される。

【0044】

図5の例は、「JobID」が「13」のジョブにおける「13_m_5」のMap処理タスクの完了結果を示している。このMap処理タスク「13_m_5」は、スレーブサーバ「Node1」で実行されたタスクである。また、このMap処理タスク「13_m_5」によって、「JobID」が「13」のジョブにおけるReduce処理が「13_r_1」、「13_r_2」、「13_r_3」の3つあることが判明したことを示す。さらに、このMap処理タスク「13_m_5」によって、Reduce処理タスク「13_r_1」のデータ量が「1000」、Reduce処理タスク「13_r_2」のデータ量が「1200」、Reduce処理タスク「13_r_3」のデータ量が「8000」であることが判明したことを示す。

【0045】

予測部15は、このようにしてMap処理タスクの完了通知からReduce処理タスクのデータ量を取得して加算していく。そして、予測部15は、加算した結果が「10000」を超える場合に、フラグに「true」を設定する。

【0046】

ここで、予測部15がReduce処理の合計を見積もる契機を様々に設定することができる。つまり、どの程度のMap処理タスクが完了した時点で、処理データ量が閾値を超えるか否かを判定するのかを任意に設定することができる。

【0047】

例えば、予測部15は、全Map処理タスクのうち予め指定した割合のMap処理が完了した時点で判定することができる。また、予測部15は、最初のMap処理タスクが終了してから予め指定した時間が過ぎて時点で判定することができる。また、予測部15は、上記2つの時点のいずれか早い時点で判定することもできる。

【0048】

なお、最初のMap処理タスクについても、ランダムに指定することもできる。このように、Map処理タスクのタスク数等に基づいて予測タイミングを任意に変更することができるので、入力データによってカスタマイズすることができる。

10

20

30

40

50

【 0 0 4 9 】

Reduce割当部 1 6 は、スレーブサーバ 3 0 からReduce処理タスクの割当要求を受信した場合に、Reduce処理タスクを割当てる処理部である。具体的にはReduce割当部 1 6 は、振分キーに関するハッシュ関数等を用いて、各Reduce処理タスクをスレーブサーバ 3 0 に割り当てる。そして、Reduce割当部 1 6 は、図 4 に示した「割当スレーブID」や「状態」等を更新する。

【 0 0 5 0 】

例えば、Reduce割当部 1 6 は、スレーブサーバ 3 0 等からReduce処理タスクの割当要求を受信した場合に、タスクリスト D B 1 2 b を参照して「状態」が「Not assigned」のReduce処理タスクを特定する。続いて、Reduce割当部 1 6 は、ハッシュ関数等を用いて割当先のスレーブサーバを特定する。その後、Reduce割当部 1 6 は、特定した割当先のスレーブサーバ 3 0 のIDを、割当対象のReduce処理タスクの「スレーブサーバID」に格納する。

10

【 0 0 5 1 】

その後、Reduce割当部 1 6 は、特定した割当先のスレーブサーバ 3 0 に、TaskID、必要スロット数、処理データ量、フラグ等を通知して、Reduce処理タスクを割当てる。また、Reduce割当部 1 6 は、割当てたMap処理タスクの「状態」を「Not assigned」から「Running」に更新する。なお、Reduce割当部 1 6 は、Reduce処理タスクの完了通知を受信した場合、該当Reduce処理タスクの「状態」を「Running」から「Done」に更新する。

【 0 0 5 2 】

[スレーブサーバの構成]

図 6 は、実施例 1 に係るスレーブサーバの機能構成を示す機能ブロック図である。図 6 に示すように、スレーブサーバ 3 0 は、通信制御部 3 1 と、記憶部 3 2 と、制御部 3 3 とを有する。

20

【 0 0 5 3 】

通信制御部 3 1 は、マスタサーバ 1 0 や他のスレーブサーバ 3 0 などと通信を実行する処理部であり、例えばネットワークインタフェースカードなどである。例えば、通信制御部 3 1 は、マスタサーバ 1 0 から各種タスクの割当を受信し、各種タスクの完了通知を送信する。また、通信制御部 3 1 は、各種タスク処理の実行に伴って、該当するスレーブサーバ 3 0 から読み出されたデータを受信する。

30

【 0 0 5 4 】

記憶部 3 2 は、一時ファイル D B 3 2 a と入出力ファイル D B 3 2 b とを有する記憶装置であり、例えばメモリやハードディスクなどである。また、記憶部 3 2 は、制御部 3 3 が実行するプログラムなどを記憶する。

【 0 0 5 5 】

一時ファイル D B 3 2 a は、Map処理、Shuffle処理、Reduce処理等で生成される中間データ、他のスレーブサーバ 3 0 等から読み出されたデータや各処理部が処理を実行する際に使用するデータを一時的に記憶するデータベースである。入出力ファイル D B 3 2 b は、Map処理の入力およびReduce処理の出力を記憶するデータベースであり、入出力 D B サーバ 2 と連携したデータベースである。

40

【 0 0 5 6 】

制御部 3 3 は、Map処理部 3 4、Map結果送信部 3 5、Shuffle処理部 3 6、Reduce受信部 3 7、フラグ判定部 3 8、Reduce処理部 3 9 を有する処理部であり、例えばプロセッサなどの電子回路である。また、制御部 3 3 は、スレーブサーバ 3 0 全体の処理を司る。

【 0 0 5 7 】

Map処理部 3 4 は、Map処理タスクを実行する処理部である。具体的には、Map処理部 3 4 は、ハートビートなどを用いて、マスタサーバ 1 0 にMap処理タスクの割当を要求する。そして、Map処理部 3 4 は、マスタサーバ 1 0 から、「TaskID、データのあるスレーブID、必要スロット数」などを含むMap割当情報を受信する。

【 0 0 5 8 】

50

その後、Map処理部34は、受信したMap割当情報にしたがって、「データのあるスレーブID」で特定されるスレーブサーバ30が処理を行っているスレーブサーバであれば、入出力DB32bからデータを取得し、そうでなければ、「データのあるスレーブID」で特定されるスレーブサーバ30からデータを取得して一時ファイルDB32a等に保存し、「必要スロット数」で指定されるスロット数を用いてMap処理を実行する。そして、Map処理部34は、Map処理結果を一時ファイルDB32a等に格納する。ここで、生成されるMap処理結果は、例えば図5に示すように、Reduce処理のタスクIDやデータ量などが含まれる。

【0059】

Map結果送信部35は、Map処理部34が実行したMap処理の結果をマスタサーバ10に送信する処理部である。例えば、Map結果送信部35は、Map処理部34からMap処理が終了したことが通知されると、一時ファイルDB32aなどからMap処理結果の一部を読み出す。そして、Map結果送信部35は、図5に示した完了通知を生成してマスタサーバ10に送信する。

10

【0060】

Shuffle処理部36は、Map処理の結果を「Key」でソートし、同じ「Key」を有する「Key、Value」ペアをマージして、Reduce処理の処理対象を生成する処理部である。具体的には、Shuffle処理部36は、マスタサーバ10からMap処理が終了したことを通知されると、当該Map処理が属するジョブのReduce処理を実行する準備として、各スレーブサーバ30から該当するMap処理結果を取得する。そして、Shuffle処理部36は、Map処理の結果を予め指定された「Key」でソートし、同じ「Key」を有する処理結果をマージして、一時ファイルDB32aに格納する。

20

【0061】

例えば、Shuffle処理部36は、「JobID」が「1」のMap処理タスクである「1_m_1、1_m_2、1_m_3、1_m_4」が終了したこと、つまり、「JobID」が「1」のReduce処理タスクの実行開始をマスタサーバ10から受信する。すると、Shuffle処理部36は、Node1、Node2、Node3、Node4からMap処理結果を取得する。続いて、Shuffle処理部36は、Map処理結果のソートおよびマージを実行し、その結果を一時ファイルDB32a等に格納する。

【0062】

Reduce受信部37は、マスタサーバ10から割当てられたReduce処理タスクを受信する処理部である。例えば、Reduce受信部37は、「JobID、TaskID、必要スロット数、処理データ量、フラグ」などから構成されるReduce処理タスクの情報を受信する。そして、Reduce受信部37は、受信した情報を一時ファイルDB32a等に格納する。

30

【0063】

フラグ判定部38は、マスタサーバ10から割当てられたReduce処理タスクにフラグが設定されているか否かを判定する処理部である。具体的には、フラグ判定部38は、Reduce受信部37が一時ファイルDB32a等に格納したReduce処理タスクの情報を参照して、フラグが設定されているか否かを判定する。そして、フラグ判定部38は、判定結果をReduce処理部39に通知する。

40

【0064】

例えば、フラグ判定部38は、Reduce処理タスクの情報が「JobID=2、TaskID=2_r_1、必要スロット数=1、処理データ量=24000、フラグ=true」であった場合、「フラグ=true」であることから、フラグが設定されていると判定する。なお、フラグ判定部38は、「フラグ」に「true」が設定されていない場合、フラグが設定されていないと判定する。

【0065】

Reduce処理部39は、フラグ判定部38による判定結果に基づいて、Reduce処理タスクの処理方法を変更して、Reduce処理タスクを実行する処理部である。具体的には、Reduce処理部39は、割当てられたReduce処理タスクにフラグが設定されている場合には、割当

50

てられたReduce処理タスクを分散処理する。一方、Reduce処理部39は、割当てられたReduce処理タスクにフラグが設定されていない場合には、割当てられたReduce処理タスクを分散処理せずに実行する。そして、Reduce処理部39は、Reduce処理タスクの処理結果を入出力ファイルDB32b等に格納する。

【0066】

例えば、Reduce処理部39は、スレーブサーバ30が有するプロセッサの数、ディスクの数、予め指定された数の少なくとも1つを用いて、割り当てられたReduce処理タスクをサブタスクに分割して、複数のプロセッサで並列実行する。

【0067】

一例を挙げると、Reduce処理部39は、プロセッサの数またはディスクの数が4つである場合、Reduce処理タスクを4つのサブタスクに分割し、4つプロセッサを用いて各サブタスクを並列に実行する。

【0068】

(Map処理の説明)

ここで、スレーブサーバ30が実行するMap処理について説明する。図7は、Map処理を説明する図である。図7に示すように、各スレーブサーバ30は、入力データとして「Hello Apple!」と「Apple is red」を受信し、それぞれの入力データに対してMap処理を実行して、「Key、Value」のペアを出力する。

【0069】

図7の例では、スレーブサーバ30は、「Hello Apple!」に対してMap処理を実行して、入力データの各要素の数を計数し、要素を「Key」、計数結果を「Value」とする「Key、Value」のペアを出力する。具体的には、スレーブサーバ30は、入力データ「Hello Apple!」から「Hello、1」、「Apple、1」、「!、1」を生成する。同様に、スレーブサーバ30は、入力データ「Apple is red」から「Apple、1」、「is、1」、「red、1」を生成する。

【0070】

(Shuffle処理)

次に、スレーブサーバ30が実行するShuffle処理について説明する。図8は、Shuffle処理を説明する図である。図8に示すように、各スレーブサーバ30は、各スレーブサーバからMap処理結果を取得してShuffle処理を実行する。

【0071】

図8の例では、スレーブサーバ(A)、(B)、(C)・・・が同じジョブ(例えば、JobID=20)に属するMap処理タスクを実行し、スレーブサーバ(D)と(Z)とが、JobID=20に属するReduce処理タスクを実行する。

【0072】

例えば、スレーブサーバ(A)がMap処理1を実行して「Apple、1」、「is、3」を生成し、スレーブサーバ(B)がMap処理2を実行して「Apple、2」、「Hello、4」を生成し、スレーブサーバ(C)がMap処理3を実行して「Hello、3」、「red、5」を生成する。スレーブサーバ(X)がMap処理1000を実行して「Hello、1000」、「is、1002」を生成する。

【0073】

続いて、スレーブサーバ(D)およびスレーブサーバ(Z)は、割当てられたReduce処理タスクで使用する各スレーブサーバのMap処理結果を取得して、ソートおよびマージを実行する。具体的には、スレーブサーバ(D)には、「Apple」と「Hello」についてのReduce処理タスクが割当てられて、スレーブサーバ(Z)には、「is」と「red」についてのReduce処理タスクが割当てられたとする。

【0074】

この場合、スレーブサーバ(D)は、スレーブサーバ(A)からMap処理1の結果「Apple、1」を取得し、スレーブサーバ(B)からMap処理2の結果「Apple、2」および「Hello、4」を取得する。また、スレーブサーバ(D)は、スレーブサーバ(C)からMap処理

10

20

30

40

50

3の結果「Hello、3」を取得し、スレーブサーバ(X)からMap処理1000の結果「Hello、1000」を取得する。そして、スレーブサーバ(D)は、これらの結果をソートおよびマージして、「Apple、[1,2]」および「Hello、[3,4,1000]」を生成する。

【0075】

同様に、スレーブサーバ(Z)は、スレーブサーバ(A)からMap処理1の結果「is、3」を取得し、スレーブサーバ(C)からMap処理3の結果「red、5」を取得し、スレーブサーバ(X)からMap処理1000の結果「is、1002」を取得する。そして、スレーブサーバ(Z)は、これらの結果をソートおよびマージして、「is、[3,1002]」および「red、[5]」を生成する。

【0076】

(Reduce処理)

次に、スレーブサーバ30が実行するReduce処理について説明する。図9は、Reduce処理を説明する図である。図9に示すように、各スレーブサーバ30は、各スレーブサーバのMap処理結果から生成したShuffle結果を用いて、Reduce処理を実行する。具体的には、Shuffle処理の説明と同様、スレーブサーバ(D)には、「Apple」と「Hello」についてのReduce処理タスクが割当てられて、スレーブサーバ(Z)には、「is」と「red」についてのReduce処理タスクが割当てられたとする。

【0077】

この例では、スレーブサーバ(D)は、Shuffle処理の結果である「Apple、[1,2]」および「Hello、[3,4,1000]」から、Reduce処理結果として「Apple、3」および「Hello、1007」を生成する。同様に、スレーブサーバ(Z)は、Shuffle処理の結果である「is、[3,1002]」および「red、[5]」から、Reduce処理結果として「is、1005」および「red、5」を生成する。

【0078】

(Reduce処理タスクのフラグ設定)

次に、マスタサーバ10がMap処理結果からReduce処理タスクにフラグを設定する例を説明する。図10は、Reduce処理タスクの処理量を予測してフラグを設定する処理を説明する図である。この図10は、マスタサーバ10が保持するタスクリストを示している。

【0079】

図10示すタスクリストのうち「JobID=1」についてはReduce処理の割当てが完了しており、Reduce処理が既に実行されていることを示す。このような状態で、「JobID=2」におけるMap処理タスク「2_m_1、2_m_2、2_m_3、2_m_4」のうち「2_m_1」および「2_m_2」が完了したとする(S1)。

【0080】

そして、マスタサーバ10は、Map処理タスク「2_m_1」を実行したNode1からReduce処理のデータ量を含むMap完了通知1を受信し、Map処理タスク「2_m_2」を実行したNode2からReduce処理のデータ量を含むMap完了通知2を受信する(S2)。

【0081】

続いて、マスタサーバ10は、受信したMap完了通知1とMap完了通知2とから、「JobID=2」におけるReduce処理タスク「2_r_1」と「2_r_2」の処理データ量をそれぞれ「24000」と「13000」と見積もる(S3)。

【0082】

その後、マスタサーバ10は、処理データ量が閾値「20000」を超えるReduce処理タスク「2_r_1」に対して、フラグ「true」を設定する(S4)。そして、マスタサーバ10は、Reduce処理タスク「2_r_1」と「2_r_2」の割当先をハッシュ関数で決定し、Reduce処理タスク「2_r_1」が割当てられたNodeに、Reduce処理タスクとともにフラグ「true」を送信する(S5)。

【0083】

このように、マスタサーバ10は、ジョブにおける全てのMap処理タスクが終了する前に、一部のMap処理タスクが完了した時点で、Reduce処理タスクの処理データ量を見積り

10

20

30

40

50

、見積もった結果に応じてフラグを設定する。

【 0 0 8 4 】

[マスタサーバの処理]

図 1 1 は、実施例 1 に係るマスタサーバが実行する処理の流れを示すフローチャートである。図 1 1 に示すように、マスタサーバ 1 0 は、管理者等から登録されるジョブ登録の情報にしたがって、ジョブリスト DB 1 2 a にジョブリストやタスクリスト DB 1 2 b にタスクリストを追加する (S 1 0 1) 。

【 0 0 8 5 】

その後、マスタサーバ 1 0 は、スレーブサーバ 3 0 からハートビートなどの通知を受信するまで待機し (S 1 0 2) 、通知を受信した場合に、当該通知がタスク要求かタスクの完了通知かを判定する (S 1 0 3) 。

10

【 0 0 8 6 】

そして、マスタサーバ 1 0 のMap割当部 1 4 は、受信された通知がタスク要求であると判定された場合 (S 1 0 3 : タスク要求) 、ハッシュ関数等を用いてタスク割当を実施する (S 1 0 4) 。その後、Map割当部 1 4 は、通知要求元のスレーブサーバ 3 0 に、割当てたタスク情報を応答する (S 1 0 5) 。ここで、タスク情報には、タスクに属するジョブに関するジョブリストの該当行の一行分の情報およびタスクリストの該当行の一行分の情報などが含まれる。

【 0 0 8 7 】

一方、マスタサーバ 1 0 の制御部 1 3 は、受信された通知がタスクの完了通知であると判定された場合 (S 1 0 3 : 完了通知) 、該当タスクの完了処理を実行する (S 1 0 6) 。その後、制御部 1 3 は、ジョブの全タスクが完了した場合 (S 1 0 7 : Y e s) 、 S 1 0 1 に戻って以降の処理を繰り返す。一方、制御部 1 3 は、ジョブの全タスクが完了していない場合 (S 1 0 7 : N o) 、 S 1 0 2 に戻って以降の処理を繰り返す。

20

【 0 0 8 8 】

(タスクの完了処理)

図 1 2 は、マスタサーバが実行する該当タスクの完了処理の流れを示すフローチャートである。この処理は、図 1 1 の S 1 0 6 で実行される処理である。

【 0 0 8 9 】

図 1 2 に示すように、マスタサーバ 1 0 の予測部 1 5 は、受信された完了通知がMap処理タスクの完了通知であると判定された場合 (S 2 0 1 : Map) 、終了したMap処理の完了通知からReduce処理タスクの処理データ量を加算する (S 2 0 2) 。

30

【 0 0 9 0 】

例えば、予測部 1 5 は、受信された完了通知のヘッダにMap処理タスクを示す識別子が付与されていた場合に、Map処理タスクの完了通知を判定し、当該完了通知に含まれるReduceデータ量に基づいて、図 4 に示す該当Reduce処理のデータ量を計算する。

【 0 0 9 1 】

その後、予測部 1 5 は、所定のMap処理タスクが完了し、フラグ判定タイミングであると判定すると (S 2 0 3 : Y e s) 、図 4 を参照し、処理データ量が所定値を超えるReduce処理タスクが存在するか否かを判定する (S 2 0 4) 。

40

【 0 0 9 2 】

そして、予測部 1 5 は、処理 (転送) データ量が所定値を超えるReduce処理タスクにフラグ「true」を設定して、フラグの付いたReduce処理タスクの割当を実行する (S 2 0 5) 。例えば、予測部 1 5 は、ハッシュ値等を用いて、フラグ「true」を設定したReduce処理タスクの割当先を決定し、決定した割当先のスレーブサーバ 3 0 にReduce処理タスクを送信する。

【 0 0 9 3 】

続いて、予測部 1 5 は、完了通知を受信したMap処理タスクの「状態」を「Done」に変更し (S 2 0 6) 、該当Map処理タスクの完了をMap処理タスクの完了通知領域に登録する (S 2 0 7) 。

50

【 0 0 9 4 】

一方、マスタサーバ10のReduce割当部16は、受信された完了通知がMap処理タスクではなく、Reduce処理タスクの完了通知であると判定された場合（S201:Reduce）、完了通知を受信したReduce処理タスクの「状態」を「Done」に変更する（S208）。

【 0 0 9 5 】

[スレーブサーバの処理]

図13は、実施例1に係るスレーブサーバが実行する処理の流れを示すフローチャートである。図13に示すように、スレーブサーバ30は、マスタサーバ10に、ハートビートでタスク要求を送信する（S301）。

【 0 0 9 6 】

続いて、スレーブサーバ30は、タスク要求の応答としてジョブ情報とタスク情報を取得し（S302）、取得したタスク情報がMap処理タスクの情報か否かを判定する（S303）。

【 0 0 9 7 】

そして、スレーブサーバ30のMap処理部34は、取得したタスク情報がMap処理タスクの情報であると判定された場合（S303:Map）、入力データを読み込み（S304）、Map処理タスクを起動する（S305）。

【 0 0 9 8 】

例えば、Map処理部34は、取得したMap処理タスクの情報における「データのあるスレーブID」で特定されるスレーブサーバ30から入力データを取得し、取得したMap処理タスクの情報によって割当てられたMap処理タスクを起動する。

【 0 0 9 9 】

その後、Map処理部34は、一時ファイルDB32aに、Reduce処理タスクごとに分けて処理結果を保存し（S306）、Map処理タスクが終了するまで待機する（S307）。そして、Map結果送信部35は、ハートビート等でMap処理タスクの完了とReduce向けデータ量をマスタサーバ10に送信する（S308）。

【 0 1 0 0 】

一方、スレーブサーバ30のReduce受信部37が取得したタスク情報がReduce処理タスクの情報であると判定した場合（S303:Reduce）、Shuffle処理部36が、各スレーブサーバ30からMap処理結果を取得して、Shuffle処理を実行する（S309）。

【 0 1 0 1 】

Reduce処理部39は、Reduce受信部37が取得したReduce処理タスクを実行して（S310）、タスクが終了するまで待機し（S311）、タスクが完了すると、ハートビート等で完了通知をマスタサーバ10に送信する（S312）。

【 0 1 0 2 】

(Reduce処理タスクの起動処理)

図14は、スレーブサーバが実行するReduce処理タスクの起動処理の流れを示すフローチャートである。図14に示すように、フラグ判定部38は、Reduce受信部37が受信したReduce処理タスクにフラグ「true」が付加されているか否かを判定する（S401）。

【 0 1 0 3 】

そして、フラグ判定部38がReduce処理タスクにフラグ「true」が付加されていると判定した場合（S401:Yes）、Reduce処理部39は、Reduce処理タスクの入力を分割する（S402）。

【 0 1 0 4 】

続いて、Reduce処理部39は、各分割された入力でS403からS405までループ処理する。具体的には、Reduce処理部39は、各分割されたReduce処理タスクの入力に関して、Reduce処理タスクのサブタスクを起動する（S404）。

【 0 1 0 5 】

そして、Reduce処理部39は、Reduce処理タスクの全サブタスクが完了するまで待機し（S406）、全サブタスクが完了すると、処理を終了する。

10

20

30

40

50

【 0 1 0 6 】

一方、S 4 0 1において、フラグ判定部 3 8がReduce処理タスクにフラグ「true」が付け加されていないと判定した場合（S 4 0 1：No）、Reduce処理部 3 9は、受信されたReduce処理タスクをそのまま実行する（S 4 0 7）。そして、Reduce処理部 3 9は、Reduce処理タスクが完了すると、処理を終了する。

【 0 1 0 7 】

（Reduce処理タスクの分割処理）

図 1 5は、スレーブサーバが実行するReduce処理タスクの分割処理の流れを示すフローチャートである。図 1 5に示すように、スレーブサーバ 3 0のReduce処理部 3 9は、自サーバ内におけるスレーブの受け入れ可能スロット数を「S」と設定する（S 5 0 1）。 10

【 0 1 0 8 】

そして、Reduce処理部 3 9は、変数「i」がS - 1となるまで順にS 5 0 2からS 5 0 8までを実行するループ処理する。具体的には、Reduce処理部 3 9は、「i × Reduce処理タスクの入力の全レコード数 / S」を「開始位置」に設定する（S 5 0 3）。その後、Reduce処理部 3 9は、Reduce処理タスクの入力「開始位置」が「Key」ではない間、S 5 0 4からS 5 0 6の処理を繰り返すループ処理を実行する。

【 0 1 0 9 】

具体的には、Reduce処理部 3 9は、Reduce処理タスクの入力「開始位置」が「Key」になるまで、「開始位置」をインクリメントする（S 5 0 5）。そして、Reduce処理部 3 9は、Reduce処理タスクの入力「開始位置」が「Key」になると、分割された入力の開始位置「i」に、S 5 0 4からS 5 0 6で算出された「開始位置」を代入する（S 5 0 7）。その後、Reduce処理部 3 9は、S 5 0 2以降のループ処理を実行する。 20

【 0 1 1 0 】

上述したように、マスタサーバ 1 0は、一部のMap処理タスクの実行後に、データ量の集中する可能性が高いReduce処理タスクを検知する。そして、マスタサーバ 1 0は、検出したReduce処理タスクをスレーブサーバ 3 0で処理させる際に、スレーブサーバ 3 0内で並列実行させることができる。

【 0 1 1 1 】

このように、マスタサーバ 1 0は、タスク割り当ての際には知ることができないReduceタスクの処理量の不均一性に関する情報を、各スレーブサーバ 3 0に伝えて処理方法を変更させることができる。したがって、該当スレーブサーバのみReduce処理タスクを優先させることにより、ジョブ全体の完了時間を短縮することができる。 30

【 0 1 1 2 】

また、マスタサーバ 1 0は、全Map処理タスクの所定の割合が終了した場合や最初のMap処理タスクが終了してから所定時間が経過した場合に、Reduce処理タスクの処理量を見積もることができる。この結果、マスタサーバ 1 0は、どの程度のMap処理タスクが終了したところで、フラグ付けの判断を行うかを任意の方式で決定することができるので、ジョブに適した方式を採用することができ、汎用性を高めることができる。

【 0 1 1 3 】

スレーブサーバ 3 0は、処理量の多いReduce処理タスクについては、フラグが通知されるので、当該Reduce処理タスクを分割して実行することができ、処理時間を短縮することができる。また、スレーブサーバ 3 0は、Reduce処理タスクを分割する際に、プロセッサの数、ディスクの数、予め指定された数の少なくとも1つを用いて分割するので、自サーバの処理性能にあわせて分割することができる。 40

【実施例 2】

【 0 1 1 4 】

ところで、マスタサーバ 1 0は、フラグ付きのReduce処理タスクを割り当てる先のスレーブサーバ 3 0に、他のReduce処理タスクが既に割り当てられている場合に、フラグ付きのReduce処理タスクを優先させることができる。

【 0 1 1 5 】

そこで、実施例2では、フラグ付きのReduce処理タスクを優先させる例について説明する。図16は、実施例2に係るReduce処理タスクの割当処理の流れを示すフローチャートである。

【0116】

図16に示すように、マスタサーバ10のReduce割当部16は、タスクリストDB12b中のReduce処理タスクを処理データ量が多い順に並び替える(S601)。続いて、Reduce割当部16は、タスクリストDB12b中のReduce処理タスクを処理データ量の多い順に、スレーブサーバ30の数だけ選択し、「優先タスクリストP」と設定する(S602)。そして、Reduce割当部16は、「未割当優先タスク数」に「0」を設定する(S603)。

10

【0117】

その後、Reduce割当部16は、「優先タスクリストP」の数のReduce処理タスクについて、S604からS608までをループ処理する。具体的には、Reduce割当部16は、処理対象のReduce処理タスクが既に割当済みである場合(S605:Yes)、当該タスクと同じスレーブサーバ30に割当てられた他のReduce処理タスクを中止する(S606)。一方、Reduce割当部16は、処理対象のReduce処理タスクが割当済みではない場合(S605:No)、「未割当優先タスク数」をインクリメントする(S607)。

【0118】

S605からS608までのループ処理が終了すると、Reduce割当部16は、「未割当優先タスク数」が0より大きいかな否かを判定する(S609)。そして、Reduce割当部16は、「未割当優先タスク数」が0である場合(S609:No)、処理を終了する。

20

【0119】

一方、Reduce割当部16は、「未割当優先タスク数」が0より大きい場合(S609:Yes)、全スレーブサーバ30に対してS610からS615までをループ処理する。

【0120】

具体的には、Reduce割当部16は、該当スレーブサーバ30に優先タスクリストのいずれかのReduce処理タスクが割当済みである場合には(S611:Yes)、当該スレーブサーバ30に対するループ処理を終了し(S615)、次のスレーブサーバ30に対してループ処理を実行する。

【0121】

一方、Reduce割当部16は、該当スレーブサーバ30に、優先タスクリストにおけるいずれのReduce処理タスクも割当てられていない場合には(S611:No)、対象のスレーブサーバ30に割当てられた他のReduce処理タスクを中止する(S612)。続いて、Reduce割当部16は、該当する優先タスクを対象のスレーブサーバ30に割り与え、「未割当優先タスク数」を1減算する(S613)。

30

【0122】

その後、Reduce割当部16は、1減算後の「未割当優先タスク数」が0である場合には(S614:Yes)、処理を終了する。一方、Reduce割当部16は、1減算後の「未割当優先タスク数」が0より大きい場合には(S614:No)、次のスレーブサーバ30についてS610以降のループ処理を実行する。

40

【0123】

このように、マスタサーバ10は、フラグ付きのReduce処理タスクを割当てる先のスレーブサーバ30に、他のReduce処理タスクが既に割当てられている場合に、フラグ付きのReduce処理タスクを優先させることができる。このため、マスタサーバ10は、後からフラグが付いたReduce処理タスクを優先させて処理させることができるので、Reduce処理タスクの見積り順序に依存せずに、ジョブ全体の完了時間を短縮することができる。

【0124】

一般には実行中のタスクを中止することで無駄になる計算時間、また送信済みのデータを使わなくなるために無駄になる通信時間やデータの読み書き時間が発生する。しかし、いくつかのMap処理タスクが終了してフラグが立った段階で、フラグの経ったReduce処理

50

タスクと同じスレーブサーバ30に割り当て済みのReduce処理タスクを中止するようにすることで、一部のデータ転送の無駄とそのデータの処理の無駄だけに抑えることができる。また、クリティカルパス上にあってもっとも早く開始したいフラグの立ったReduce処理タスクの開始を従来通り素早く開始できるため全体としての処理時間が短くなる。

【実施例3】

【0125】

ところで、マスタサーバ10は、Reduce処理タスクのフラグが設定されてすぐに割当を実施せずに、いくつかのMap処理タスクが終了して割当てすることもできる。そこで、実施例3では、すぐにReduce処理タスクの割当を実行せずに、いくつかのMap処理タスクが終了してからReduce処理タスクの割当を実行する例を説明する。

10

【0126】

図17は、実施例3に係るReduce処理タスクの割当処理の流れを示すフローチャートである。図17に示すように、マスタサーバ10のMap割当部14は、タスク要求してきたスレーブサーバ30がMap処理タスクの受け入れが可能である場合(S701:Yes)、未割り当てのローカルMap処理タスクが存在するか否かを判定する(S702)。ここで、ローカルMap処理タスクとは、タスク要求してきたスレーブサーバ30がタスクリスト中の「データがあるスレーブID」の列に含まれているタスクを指す。

【0127】

そして、Map割当部14は、未割り当てのローカルMap処理タスクが存在すると判定した場合(S702:Yes)、タスク要求してきたスレーブサーバ30に、当該未割り当てのローカルMap処理タスクを割当てる(S703)。その後、Map割当部14は、S701以降の処理を繰り返す。

20

【0128】

一方、Map割当部14は、未割り当てのローカルMap処理タスクが存在せず(S702:No)、未割り当てのMap処理タスクが存在すると判定した場合(S704:Yes)、タスク要求してきたスレーブサーバに、未割り当てのMap処理タスクを割当てる(S703)。

【0129】

そして、Map処理タスクの受け入れが可能ではなく(S701:No)、または、未割り当てのMap処理タスクが存在しない場合(S704:No)、S705を実行する。具体的には、Reduce割当部16は、タスク要求してきたスレーブサーバ30がReduce処理タスクの受け入れが可能であるか否かを判定する。ここで、Reduce処理タスクの受け入れが可能であるとは、「スレーブサーバ30の処理可能なReduce処理タスク数>スレーブサーバ30に割り当て済みReduce処理タスクの必要スロット数の合計」である状態を指す。

30

【0130】

Reduce割当部16は、タスク要求してきたスレーブサーバ30がReduce処理タスクの受け入れが可能であると判定した場合(S705:Yes)、S706を実行する。具体的には、Reduce割当部16は、ジョブリストDB12aのReduce割当許可がtrueか否かを参照して、Reduce処理タスクの割当が許可されているか否かを判定する。

【0131】

そして、Reduce割当部16は、Reduce処理タスクの割当が許可されていると判定した場合(S706:Yes)、未割り当てのReduce処理タスクが存在するか否かを判定する(S707)。

40

【0132】

ここで、Reduce割当部16は、未割り当てのReduce処理タスクが存在すると判定した場合(S707:Yes)、タスク要求してきたスレーブサーバ30に、未割り当てのReduce処理タスクを割り当てる(S708)。その後、Reduce割当部16は、S701以降の処理を繰り返す。

【0133】

一方、Reduce処理タスクの受け入れが可能ではない場合(S705:No)、Reduce処

50

理タスクの割当が許可されていない場合 (S 7 0 6 : N o)、未割り当てのReduce処理タスクが存在しない場合 (S 7 0 7 : N o)、Reduce割当部 1 6 は、処理を終了する。

【 0 1 3 4 】

このように、マスタサーバ 1 0 は、Reduce処理タスクのフラグが設定されてすぐに割当を実施せずに、いくつかのMap処理タスクが終了して割当てすることもできる。この方式では、shuffle処理のデータ転送の開始が遅くなる場合もあるが、多くのMap処理タスクのデータ量は均等であり、最初に割り当てられたMapタスクは一斉に終わることが期待される。このため、遅れの影響は限定的である一方、中止されるReduce処理タスクはないので、データ転送とその処理に関する無駄は発生しない。

【 0 1 3 5 】

したがって、若干の余分な処理時間がかかる場合があるが、ボトルネックとなるデータ量の多いReduce処理タスクの処理時間を短縮することで、ジョブ全体の処理時間を短縮できる。

【実施例 4】

【 0 1 3 6 】

さて、これまで本発明の実施例について説明したが、本発明は上述した実施例以外にも、種々の異なる形態にて実施されてよいものである。そこで、以下に異なる実施例を説明する。

【 0 1 3 7 】

(Reduce処理タスクのフラグ判断)

マスタサーバ 1 0 は、所定値以上の処理データ量であるReduce処理タスクをフラグ対象として決定する例を説明したが、これに限定されず、様々な判断基準で決定することができる。

【 0 1 3 8 】

例えば、マスタサーバ 1 0 は、Reduce処理タスクのデータ量の平均を「 m 」とした場合、例えば 2 などのあらかじめ定められた係数 k に基づき、「 $k \times m$ 」よりも多くのデータ量を有するReduce処理タスクにフラグを付けることもできる。

【 0 1 3 9 】

また、マスタサーバ 1 0 は、Reduce処理タスクのデータ量の平均を「 m 」、分散を「 σ 」とした場合に、例えば 3 などのあらかじめ定められた係数 k に基づき、「 $m + k \times \sigma$ 」よりも多くのデータ量を有するReduce処理タスクにフラグを付けることもできる。

【 0 1 4 0 】

また、マスタサーバ 1 0 は、スレーブタスク 3 0 の同時実行可能タスク数を「 s 」、最大のReduce処理タスクのデータ量を「 d 」とした場合に、サブタスクのデータ量「 $d / s = d'$ 」とすると、「 d' 」よりも多くのデータ量を持つ全てのReduce処理タスクにフラグを付けることもできる。

【 0 1 4 1 】

(システム)

また、本実施例において説明した各処理のうち、自動的におこなわれるものとして説明した処理の全部または一部を手動的におこなうこともできる。あるいは、手動的におこなわれるものとして説明した処理の全部または一部を公知の方法で自動的におこなうこともできる。この他、上記文書中や図面中で示した処理手順、制御手順、具体的名称、各種のデータやパラメータを含む情報については、特記する場合を除いて任意に変更することができる。

【 0 1 4 2 】

また、図示した各装置の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散や統合の具体的形態は図示のものに限られない。つまり、その全部または一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的または物理的に分散・統合して構成することができる。さらに、各装置にて行なわれる各処理機能は、その全部または任意の一部が、CPUおよび当該

10

20

30

40

50

C P Uにて解析実行されるプログラムにて実現され、あるいは、ワイヤードロジックによるハードウェアとして実現され得る。

【 0 1 4 3 】

(ハードウェア)

次に、各サーバのハードウェア構成例を説明するが、各サーバは同様の構成を有するので、ここでは一例を説明する。図 1 8 は、サーバのハードウェア構成例を示す図である。図 1 7 に示すように、サーバ 1 0 0 は、通信インタフェース 1 0 1、メモリ 1 0 2、複数の H D D (ハードディスクドライブ) 1 0 3、プロセッサ装置 1 0 4 を有する。

【 0 1 4 4 】

通信インタフェース 1 0 1 は、図 2 や図 6 に示した通信制御部に該当し、例えばネットワークインタフェースカードなどである。複数の H D D 1 0 3 は、図 2 や図 6 に示した機能を動作させるプログラムや D B 等を記憶する。

10

【 0 1 4 5 】

プロセッサ装置 1 0 4 が有する複数の C P U 1 0 5 は、図 2 や図 6 に示した各処理部と同様の処理を実行するプログラムを H D D 1 0 3 等から読み出してメモリ 1 0 2 に展開することで、図 2 や図 6 等で説明した各機能を実行するプロセスを動作させる。すなわち、このプロセスは、マスタサーバ 1 0 が有する Map 割当部 1 4、予測部 1 5、Reduce 処理部 1 6 と同様の機能を実行する。また、このプロセスは、スレーブサーバ 3 0 が有する Map 処理部 3 4、Map 結果送信部 3 5、Shuffle 処理部 3 6、Reduce 処理部 3 7、フラグ判定部 3 8、Reduce 処理部 3 9 と同様の機能を実行する。

20

【 0 1 4 6 】

このようにサーバ 1 0 0 は、プログラムを読み出して実行することで、タスク割当方法またはタスク実行方法を実行する情報処理装置として動作する。また、サーバ 1 0 0 は、媒体読取装置によって記録媒体から上記プログラムを読み出し、読み出された上記プログラムを実行することで上記した実施例と同様の機能を実現することもできる。なお、この他の実施例でいうプログラムは、サーバ 1 0 0 によって実行されることに限定されるものではない。例えば、他のコンピュータまたはサーバがプログラムを実行する場合や、これらが協働してプログラムを実行するような場合にも、本発明を同様に適用することができる。

【 符号の説明 】

30

【 0 1 4 7 】

- 1 0 マスタサーバ
- 1 1 通信制御部
- 1 2 記憶部
- 1 2 a ジョブリスト D B
- 1 2 b タスクリスト D B
- 1 3 制御部
- 1 4 Map 割当部
- 1 5 予測部
- 1 6 Reduce 割当部
- 3 0 スレーブサーバ
- 3 1 通信制御部
- 3 2 記憶部
- 3 2 a 一時ファイル D B
- 3 2 b 入出力ファイル D B
- 3 3 制御部
- 3 4 Map 処理部
- 3 5 Map 結果送信部
- 3 6 Shuffle 処理部
- 3 7 Reduce 受信部

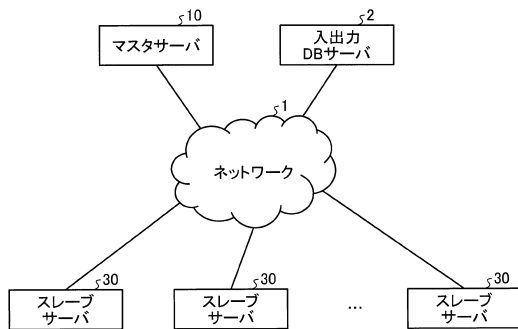
40

50

- 3 8 フラグ判定部
- 3 9 Reduce処理部

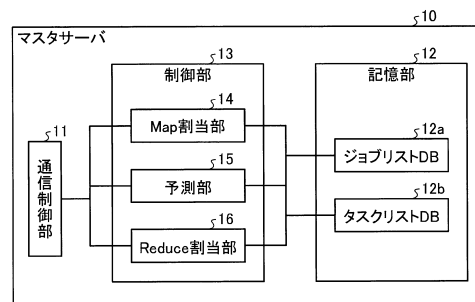
【図 1】

実施例1に係る分散処理システムの全体構成例を示す図



【図 2】

実施例1に係るマスターサーバの機能構成を示す機能ブロック図



【図 3】

ジョブリストDBに記憶される情報の例を示す図

JobID	総Mapタスク数	総Reduceタスク数	Reduce割当許可
1	4	2	false
2	4	2	false

【図4】

タスクリストDBに記憶される情報の例を示す図

JobID	TaskID	種別	Reduceの順番	データのあるスレーブID	状態	割り当てスレーブID	必要スロット数	処理データ量	フラグ
1	1_m_1	Map	-	Node1, Node2	Done	Node1	1		
1	1_m_2	Map	-	Node3, Node4	Done	Node4	1		
1	1_m_3	Map	-	Node1, Node3	Running	Node1	1		
1	1_m_4	Map	-	Node2, Node4	Running	Node2	1		
1	1_r_1	Reduce	1	-	Running	Node2	1	12000	
1	1_r_2	Reduce	2	-	Running	Node3	1	25000	true
2	2_m_1	Map	-	Node3, Node4	Done	Node1	1		
2	2_m_2	Map	-	Node1, Node3	Done	Node2	1		
2	2_m_3	Map	-	Node2, Node4	Running	Node3	1		
2	2_m_4	Map	-	Node1, Node4	Running	Node4	1		
2	2_r_1	Reduce	1	-	Not assigned	-	1	24000	true
2	2_r_2	Reduce	2	-	Not assigned	-	1	13000	

【図5】

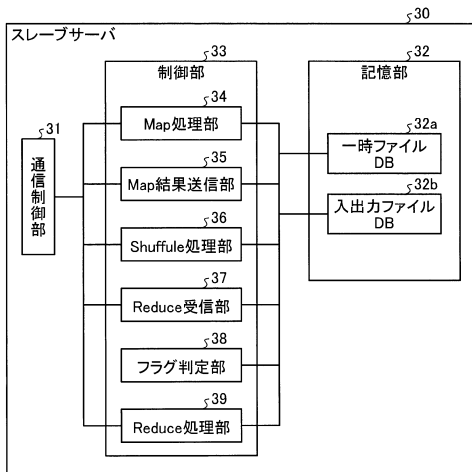
Map処理の完了通知の例を示す図

通知種別	JobID	完了Map TaskID	Mapタスクを実行したスレーブID
Map完了	13	13_m_5	Node1

通知種別	JobID	完了Map TaskID	Reduce TaskID	データ量
Reduceデータ量	13	13_m_5	13_r_1	1000
Reduceデータ量	13	13_m_5	13_r_2	1200
Reduceデータ量	13	13_m_5	13_r_3	8000

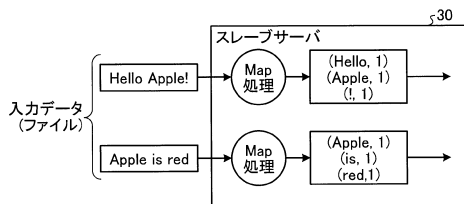
【図6】

実施例1に係るスレーブサーバの機能構成を示す機能ブロック図



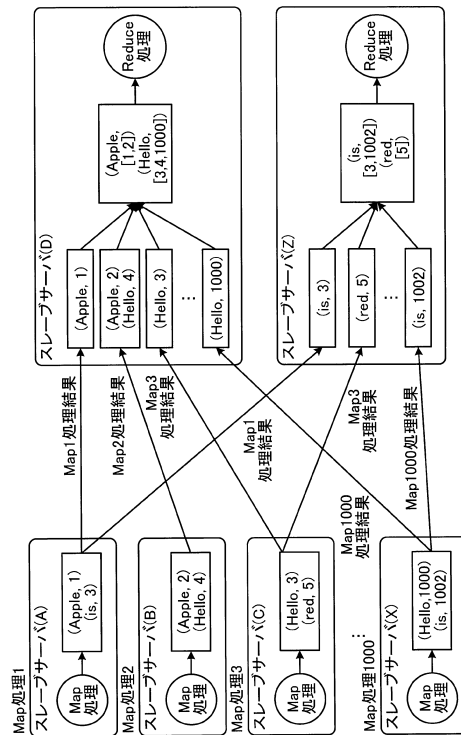
【図7】

Map処理を説明する図

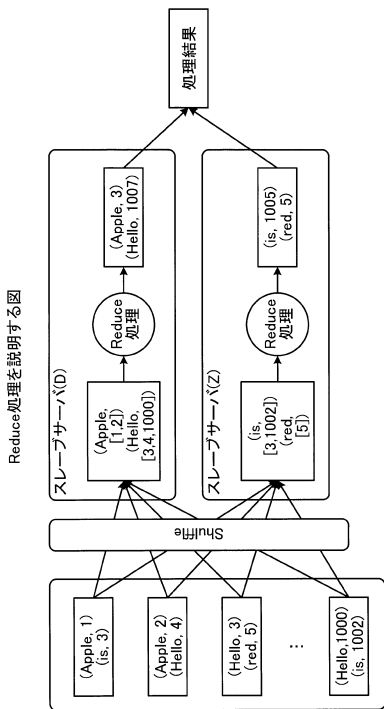


【図8】

Shuffle処理を説明する図



【図9】



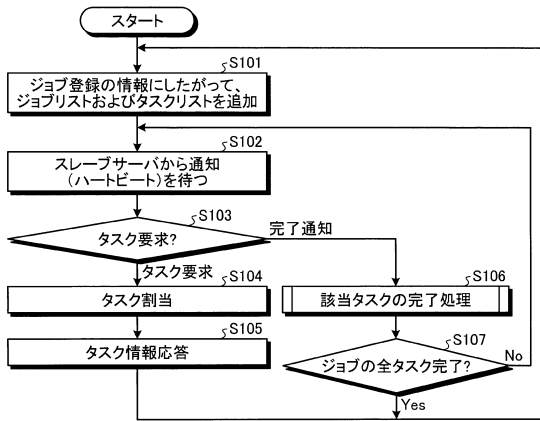
【図10】

Reduce処理タスクの処理量を予測してフラグを設定する処理を説明する図

JobID	TaskID	種別	Reduceの項番	データのあるスレーブID	状態	割り当てスレーブID	必要スロット数	処理データ量	フラグ
1	1_m_1	Map	-	Node1, Node2	Done	Node1	1		
1	1_m_2	Map	-	Node3, Node4	Done	Node4	1		
1	1_m_3	Map	-	Node1, Node3	Running	Node1	1		
1	1_m_4	Map	-	Node2, Node4	Running	Node2	1		
1	1_r_1	Reduce	1	-	Running	S1	1	12000	
1	1_r_2	Reduce	2	-	Running	Node3	1	25000	true
2	2_m_1	Map	-	Node3, Node4	Done	Node1	1		
2	2_m_2	Map	-	Node1, Node3	Done	Node2	1		
2	2_m_3	Map	-	Node2, Node4	Running	Node3	1		
2	2_m_4	Map	-	Node1, Node4	Running	Node4	1		
2	2_r_1	Reduce	1	-	Not assigned	-	1	24000	S3
2	2_r_2	Reduce	2	-	Not assigned	-	1	13000	S4, S5

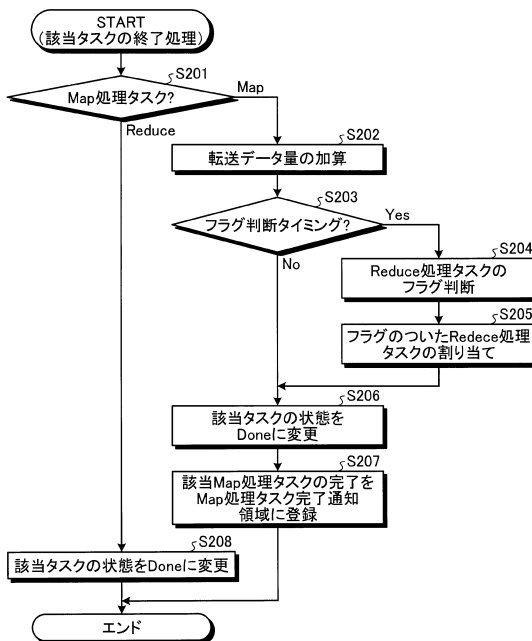
【図11】

実施例1に係るマスタサーバが実行する処理の流れを示すフローチャート



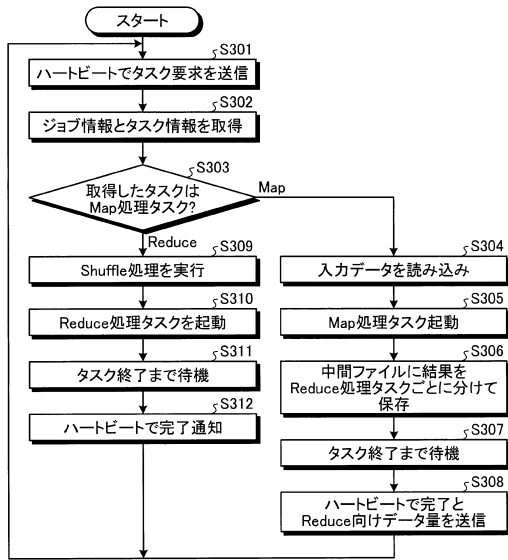
【図12】

マスタサーバが実行する該当タスクの完了処理の流れを示すフローチャート



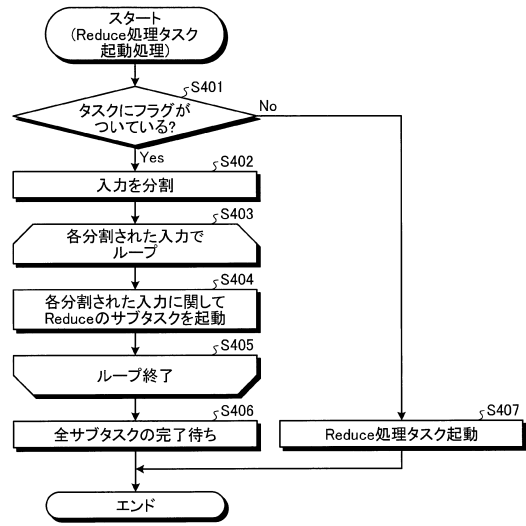
【図13】

実施例1に係るスレーブサーバが実行する処理の流れを示すフローチャート



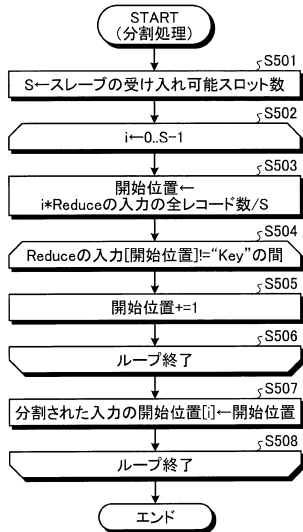
【図14】

スレーブサーバが実行するReduce処理タスクの起動処理の流れを示すフローチャート



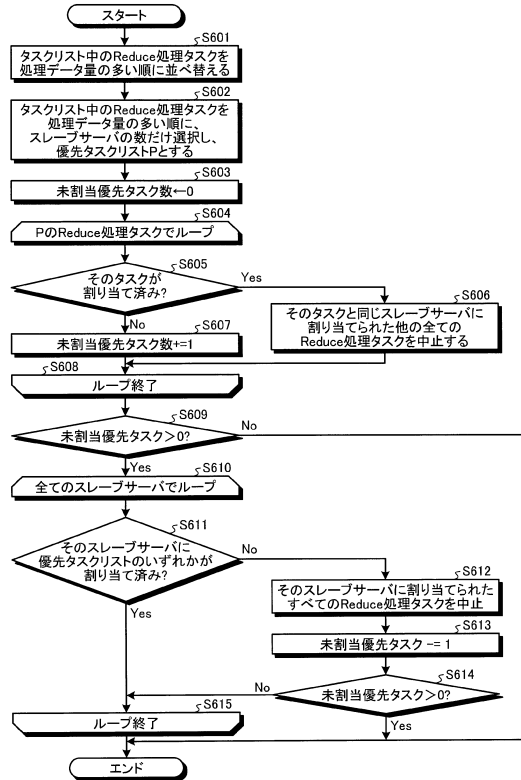
【図15】

スレーブサーバが実行するReduce処理タスクの分割処理の流れを示すフローチャート



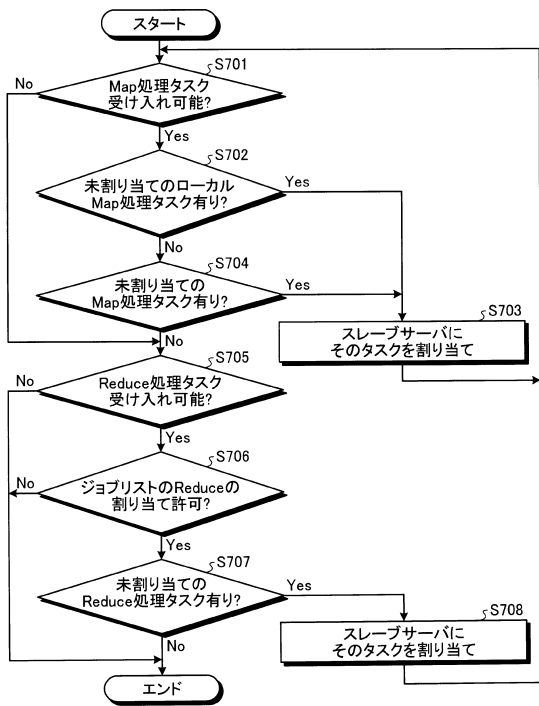
【図16】

実施例2に係るReduce処理タスクの割当処理の流れを示すフローチャート



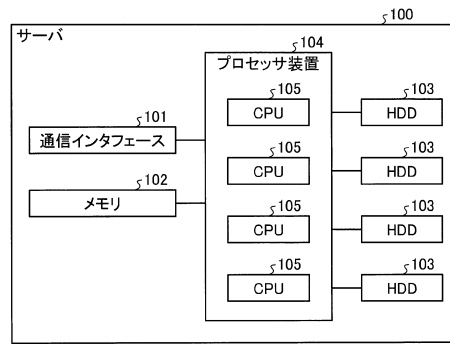
【図17】

実施例3に係るReduce処理タスクの割当処理の流れを示すフローチャート



【図18】

各サーバのハードウェア構成例を示す図



フロントページの続き

審査官 井上 宏一

- (56)参考文献 特開2012-118669(JP,A)
国際公開第2013/024597(WO,A1)
特開2013-254280(JP,A)
国際公開第2014/020735(WO,A1)

- (58)調査した分野(Int.Cl., DB名)
G06F 9/455 - 9/54