

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6364727号
(P6364727)

(45) 発行日 平成30年8月1日(2018.8.1)

(24) 登録日 平成30年7月13日(2018.7.13)

(51) Int.Cl. F I
G 0 6 F 9/50 (2006.01) G 0 6 F 9/50 1 5 0 E

請求項の数 12 (全 22 頁)

<p>(21) 出願番号 特願2013-196635 (P2013-196635)</p> <p>(22) 出願日 平成25年9月24日 (2013.9.24)</p> <p>(65) 公開番号 特開2015-64636 (P2015-64636A)</p> <p>(43) 公開日 平成27年4月9日 (2015.4.9)</p> <p>審査請求日 平成28年8月16日 (2016.8.16)</p> <p>前置審査</p>	<p>(73) 特許権者 000004237 日本電気株式会社 東京都港区芝五丁目7番1号</p> <p>(74) 代理人 100109313 弁理士 机 昌彦</p> <p>(74) 代理人 100124154 弁理士 下坂 直樹</p> <p>(72) 発明者 安田 純一 東京都港区芝五丁目7番1号 日本電気株式会社内</p> <p>審査官 漆原 孝治</p>
------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

最終頁に続く

(54) 【発明の名称】 情報処理システム、分散処理方法、及び、プログラム

(57) 【特許請求の範囲】

【請求項1】

制御装置と複数の処理装置を備えた情報処理システムであって、
前記制御装置は、対象情報の開始または終了位置を示す所定のデリミタにより区切られた複数の対象情報を含む元データを、前記所定のデリミタの位置とは無関係に所定の大きさで複数のデータに分割し、当該複数のデータを、それぞれ、複数の処理装置へ送信し、前記複数の処理装置の各々は、
前記制御装置から、処理対象のデータを受信する、受信手段と、
前記制御装置から受信した前記データを関連データとして用いる可能性がある他の処理装置へ、当該データを送信する、送信手段と、
前記データと、他の処理装置から受信した、当該データの関連データと、を用いて、当該データに対する所定の処理を行う、処理手段と、
を含み、
前記処理手段は、前記所定の処理として、前記所定のデリミタを用いて、前記データと前記関連データとから前記対象情報を抽出し、当該対象情報に対する処理を行う、
情報処理システム。

【請求項2】

前記データの関連データは、前記複数のデータにおいて、当該データに隣接するデータであり、
 前記送信手段は、前記データに隣接するデータを処理対象とする他の処理装置へ、当該

データを送信する、
請求項 1 に記載の情報処理システム。

【請求項 3】

前記データの関連データは、前記複数のデータにおいて、当該データに隣接するデータ内の当該データに隣接する一部分であり、

前記送信手段は、前記データに隣接するデータを処理対象とする他の処理装置へ、当該データの内の、当該他の処理装置の処理対象のデータに隣接する一部分を送信する、
請求項 1 に記載の情報処理システム。

【請求項 4】

前記処理手段は、前記他の処理装置における障害が検出された場合に、当該他の処理装置から受信した前記関連データと、前記データと、を用いて、当該関連データに対する前記所定の処理を行う、

請求項 1 乃至 3 のいずれかに記載の情報処理システム。

【請求項 5】

前記所定の大きさは、前記対象情報の大きさをもとに決定される、
請求項 1 乃至 4 のいずれかに記載の情報処理システム。

【請求項 6】

前記制御装置は、前記処理装置の前記処理対象のデータを関連データとして用いる可能性がある他の処理装置の識別子を示す関連装置情報を、当該処理装置へ送信し、

前記処理装置の前記送信手段は、前記関連装置情報により示される他の処理装置へ、前記データを送信する、

請求項 1 乃至 5 のいずれかに記載の情報処理システム。

【請求項 7】

制御装置と複数の処理装置を備えた情報処理システムにおける分散処理方法であって、前記制御装置が、

対象情報の開始または終了位置を示す所定のデリミタにより区切られた複数の対象情報を含む元データを、前記所定のデリミタの位置とは無関係に所定の大きさで複数のデータに分割し、当該複数のデータを、それぞれ、複数の処理装置へ送信し、

前記複数の処理装置の各々が、

前記制御装置から、処理対象のデータを受信し、

前記制御装置から受信した前記データを関連データとして用いる可能性がある他の処理装置へ、当該データを送信し、

前記データと、他の処理装置から受信した、当該データの関連データと、を用いて、当該データに対する所定の処理を行い、

前記所定の処理として、前記所定のデリミタを用いて、前記データと前記関連データとから前記対象情報を抽出し、当該対象情報に対する処理を行う、

分散処理方法。

【請求項 8】

前記データの関連データは、前記複数のデータにおいて、当該データに隣接するデータであり、

前記複数の処理装置の各々は、前記データを送信する場合、前記データに隣接するデータを処理対象とする他の処理装置へ、当該データを送信する、

請求項 7 に記載の分散処理方法。

【請求項 9】

前記データの関連データは、前記複数のデータにおいて、当該データに隣接するデータ内の当該データに隣接する一部分であり、

前記複数の処理装置の各々は、前記データを送信する場合、前記データに隣接するデータを処理対象とする他の処理装置へ、当該データの内の、当該他の処理装置の処理対象のデータに隣接する一部分を送信する、

請求項 7 に記載の分散処理方法。

10

20

30

40

50

【請求項 10】

さらに、前記複数の処理装置の各々が、前記他の処理装置における障害が検出された場合に、当該他の処理装置から受信した前記関連データと、前記データと、を用いて、当該関連データに対する前記所定の処理を行う、
請求項 7 乃至 9 のいずれかに記載の分散処理方法。

【請求項 11】

前記所定の大きさは、前記対象情報の大きさをもとに決定される、
請求項 7 乃至 10 のいずれかに記載の分散処理方法。

【請求項 12】

さらに、前記制御装置が、前記処理装置の前記処理対象のデータを関連データとして用いる可能性がある他の処理装置の識別子を示す関連装置情報を、当該処理装置へ送信し、
前記複数の処理装置の各々は、前記データを送信する場合、前記関連装置情報により示される他の処理装置へ、前記データを送信する、
請求項 7 乃至 11 のいずれかに記載の分散処理方法。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、情報処理システム、分散処理方法、及び、プログラムに関し、特に、分割データを複数のノードで分散処理する情報処理システム、分散処理方法、及び、プログラムに関する。

【背景技術】**【0002】**

コンピュータのハードウェア、ソフトウェア、及び、ネットワークの高性能化に伴い、複数のコンピュータをネットワークで接続して、分散処理を行うことにより、高い処理性能を得る技術が開発されている。

【0003】

特に、近年では、分散処理技術の発展に伴い、大量データの高速分析が可能な分散並列処理基盤が提供され、大量データに対する傾向や知見の導出に適用されている。例えば、分散並列処理基盤としてよく知られている Hadoop は、顧客情報や行動履歴のマイニング、大量ログ情報からの傾向分析などに適用されている。

【0004】

分散並列処理基盤に大量データをインポートする技術が、例えば、非特許文献 1 に開示されている。非特許文献 1 のような技術において、大量データのインポートを高速に行う方法として、分散ストレージへの書き込みを複数のノードで並行に行う方法がある。図 16 は、分散並列処理基盤への大量データのインポートの方法の例を示す図である。図 16 の例では、データサーバが、大量データを含む元データから各データを抽出し、分散並列処理基盤上の複数のノードへ送信する。ここで、データサーバは、例えば、非特許文献 2 のような技術を用いて、元データにおけるレコード等のデリミタを検出し、各データを抽出する。各ノードは、各データに対する加工（例えば、型チェック、型変換等）や、分散ストレージへの書き込み等の処理を、並列して実行する。

【先行技術文献】**【非特許文献】****【0005】**

【非特許文献 1】 "Apache Sqoop"、The Apache Software Foundation、[online]、[平成25年8月13日検索]、インターネット URL : <http://sqoop.apache.org/>

【非特許文献 2】 "RFC4180 Common Format and MIME Type for Comma-Separated Values (CSV) Files"、Y. Shafranovich、[online]、[平成25年8月13日検索]、インターネット URL : <http://tools.ietf.org/html/rfc4180>

【発明の概要】**【発明が解決しようとする課題】**

10

20

30

40

50

【 0 0 0 6 】

上述の図 1 6 のような分散並列処理基盤へのインポートにおいて、データ間に関連性がある場合、各ノードがデータを処理するときに、他ノードの処理対象のデータ（関連データ）を必要とすることがある。この場合、ノードは、関連データを保持する他ノードを検索し、当該他ノードから関連データを取得する必要がある。特に、データやノードの数が多い場合、他ノードの検索や関連データの複製、転送に伴う、システムの負荷が増大する。

【 0 0 0 7 】

本発明の目的は、上述の課題を解決し、複数データを複数のノードで分散処理するシステムにおいて、システムの処理負荷を低減する情報処理システム、分散処理方法、及び、プログラムを提供することである。

10

【課題を解決するための手段】

【 0 0 0 8 】

本発明の情報処理システムは、複数のデータの内の処理対象のデータを関連データとして用いる可能性がある他の処理装置へ、当該データを送信する、送信手段と、前記データと、他の処理装置から受信した、当該データの関連データと、を用いて、当該データに対する所定の処理を行う、処理手段と、を含む処理装置を備える。

【 0 0 0 9 】

本発明の分散処理方法は、処理装置において、複数のデータの内の処理対象のデータを関連データとして用いる可能性がある他の処理装置へ、当該データを送信し、前記処理装置において、前記データと、他の処理装置から受信した、当該データの関連データと、を用いて、当該データに対する所定の処理を行う。

20

【 0 0 1 0 】

本発明のプログラムは、処理装置用のコンピュータを、複数のデータの内の処理対象のデータを関連データとして用いる可能性がある他の処理装置へ、当該データを送信する、送信手段と、前記データと、他の処理装置から受信した、当該データの関連データと、を用いて、当該データに対する所定の処理を行う、処理手段と、して機能させる。

【発明の効果】

【 0 0 1 1 】

本発明の効果は、複数データを複数のノードで分散処理するシステムにおいて、システムの処理負荷を低減できることである。

30

【図面の簡単な説明】

【 0 0 1 2 】

【図 1】本発明の第 1 の実施の形態の特徴的な構成を示すブロック図である。

【図 2】本発明の第 1 の実施の形態における、分散処理システム 1 の構成を示すブロック図である。

【図 3】本発明の実施の形態において、データサーバ 1 0 0、及び、ノード 2 0 0 がコンピュータにより実現される場合の、分散処理システム 1 の構成を示すブロック図である。

【図 4】本発明の第 1 の実施の形態における、元データ 5 0 0 のインポート処理を示すフローチャートである。

40

【図 5】本発明の第 1 の実施の形態における、分散並列処理基盤への元データ 5 0 0 のインポートを示す図である。

【図 6】本発明の第 1 の実施の形態における、元データ 5 0 0、分割データ 5 1 0、及び、メタデータ 5 2 0 の例を示す図である。

【図 7】本発明の第 1 の実施の形態における、サーバ設定情報 1 6 1 の例を示す図である。

【図 8】本発明の第 1 の実施の形態における、転送計画 1 3 1 の例を示す図である。

【図 9】本発明の第 1 の実施の形態における、ノード設定情報 2 5 1 の例を示す図である。

【図 1 0】本発明の第 1 の実施の形態における、対象情報の抽出、及び、加工の例を示す

50

図である。

【図 1 1】本発明の第 2 の実施の形態における、分散並列処理基盤への元データ 5 0 0 のインポートを示す図である。

【図 1 2】本発明の第 2 の実施の形態における、対象情報の抽出、及び、加工の例を示す図である。

【図 1 3】本発明の第 3 の実施の形態における、分散処理システム 1 の構成を示すブロック図である。

【図 1 4】本発明の第 3 の実施の形態における、引継ぎ処理を示すフローチャートである。

【図 1 5】本発明の第 3 の実施の形態における、引継ぎ処理における対象情報の抽出、及び、加工の例を示す図である。 10

【図 1 6】分散並列処理基盤への大量データのインポートの方法の例を示す図である。

【発明を実施するための形態】

【0 0 1 3】

(第 1 の実施の形態)

本発明の第 1 の実施の形態について説明する。

【0 0 1 4】

はじめに、本発明の第 1 の実施の形態における、分散並列処理基盤への元データ 5 0 0 のインポートについて説明する。

【0 0 1 5】

図 5 は、本発明の第 1 の実施の形態における、分散並列処理基盤への元データ 5 0 0 のインポートを示す図である。 20

【0 0 1 6】

本発明の第 1 の実施の形態においては、データサーバ 1 0 0 に保存されている元データ 5 0 0 は、例えば、データベースや、ログファイルであり、複数の対象情報を含む。ここで、対象情報は、データベースにおけるレコードや、ログにおけるログレコード等、マイニングや分析が行われる処理単位である。

【0 0 1 7】

データサーバ 1 0 0 は、元データ 5 0 0 を所定長の分割データ（または、単に、データ）5 1 0 に分割し、複数のノード 2 0 0 に送信する。そして、各ノード 2 0 0 は、データサーバ 1 0 0 から受信した分割データ 5 1 0（処理対象の分割データ 5 1 0）に対して、対象情報の抽出、型チェック、変換、及び、複数のノード 2 0 0 上に構築される分散ストレージへの書き込み等の所定の処理を行う。 30

【0 0 1 8】

各ノード 2 0 0 は、処理対象の分割データ 5 1 0 に、抽出しようとする対象情報の一部しか含まれていない場合、当該処理対象の分割データ 5 1 0 に隣接する分割データ 5 1 0（隣接分割データ）のレプリカ（複製）を用いて対象情報を抽出する。本発明の第 1 の実施の形態においては、分割データ 5 1 0 の隣接分割データのレプリカを、分割データ 5 1 0 の関連データと呼ぶ。各ノード 2 0 0 は、データサーバ 1 0 0 から分割データ 5 1 0 を受信したときに、当該分割データ 5 1 0 を関連データとして用いる（当該分割データ 5 1 0 の隣接分割データを処理対象とする）他のノード 2 0 0 に、当該分割データ 5 1 0 のレプリカを生成する。 40

【0 0 1 9】

次に、本発明の第 1 の実施の形態における、分散処理システム 1 の構成を説明する。

【0 0 2 0】

図 2 は、本発明の第 1 の実施の形態における、分散処理システム 1 の構成を示すブロック図である。図 2 を参照すると、本発明の第 1 の実施の形態における分散処理システム 1 は、データサーバ（または、制御装置）1 0 0、及び、分散並列処理基盤上の複数のノード（または、処理装置）2 0 0 を含む。

【0 0 2 1】

分散処理システム 1 は、本発明の情報処理システムの一実施形態である。

【0022】

データサーバ 100、及び、複数のノード 200 は、ネットワーク等により、互いに通信可能に接続される。図 2 の例では、データサーバ 100、及び、ノード 200 「N1」、 「N2」、... が接続されている。ここで、「」内は、ノード 200 の識別子を示す。以下、後述する他の識別子についても、同様の表現を用いる。

【0023】

データサーバ 100 は、データ記憶部 110、データ取得部 120、転送計画部 130、分割部 140、分割データ送信部 150、及び、サーバ設定記憶部 160 を含む。

【0024】

データ記憶部 110 は、元データ 500 を記憶する。

【0025】

図 6 は、本発明の第 1 の実施の形態における、元データ 500、分割データ 510、及び、メタデータ 520 の例を示す図である。

【0026】

本発明の第 1 の実施の形態では、図 6 に示すように、元データ 500 のデータ形式は、XML (eXtensible Markup Language) 形式である。また、元データ 500 は、対象情報として、イベント識別子 (イベント ID) で識別されるイベント情報を含む。各対象情報は、<event></event> を開始ポイント/終了ポイントとするデリミタにより抽出される。

【0027】

データ取得部 120 は、データ記憶部 110 から元データ 500 を取得する。

【0028】

サーバ設定記憶部 160 は、データサーバ 100 の処理に係る情報である、サーバ設定情報 161 を記憶する。サーバ設定情報 161 は、例えば、管理者等により、予め設定される。

【0029】

図 7 は、本発明の第 1 の実施の形態における、サーバ設定情報 161 の例を示す図である。図 7 の例では、サーバ設定情報 161 は、送信先ノード群、送信先決定方法、送信並列度、及び、分割データサイズを含む。

【0030】

ここで、送信先ノード群は、分割データ 510 の送信先の候補であるノード 200 の識別子を示す。送信先決定方法は、送信先ノード群に含まれるノード 200 の中から、分割データ 510 の送信先を決定する方法を示す。送信並列度は、送達確認を待たずに、並列に送信可能な分割データ 510 の数を示す。分割データサイズは、分割データ 510 の大きさを示す。

【0031】

転送計画部 130 は、サーバ設定情報 161 に従って、分割データ 510 のノード 200 への送信に係る情報である、転送計画 131 を生成する。

【0032】

図 8 は、本発明の第 1 の実施の形態における、転送計画 131 の例を示す図である。図 8 の例では、転送計画 131 は、分割データ ID ごとに、送信先ノード ID、及び、メタデータ (または、関連装置情報) 520 を含む。

【0033】

ここで、分割データ ID は、分割データ 510 の識別子を示す。送信先ノード ID は、分割データ 510 の送信先であるノード 200 の識別子を示す。

【0034】

メタデータ 520 は、分割データ 510 とともにノード 200 に送信される情報である。メタデータ 520 は、分割データ ID、レプリカ生成先ノード ID (前、または、後)、及び、関連データ ID (前、または、後) を含む。レプリカ生成先ノード ID (前、ま

10

20

30

40

50

たは、後)は、分割データ510のレプリカの生成先(送信先)であるノード200の識別子を示す。レプリカ生成先ノードID(前)は、分割データ510の前方の隣接分割データを処理対象とするノード200の識別子である。レプリカ生成先ノードID(後)は、分割データ510の後方の隣接分割データを処理対象とするノード200の識別子である。関連データID(前)は、分割データ510の前方の隣接分割データの識別子を示す。関連データID(後)は、分割データ510の後方の隣接分割データの識別子を示す。

【0035】

分割部140は、転送計画131に従って、元データ500を分割データ510に分割する。

【0036】

分割データ送信部150は、転送計画131に従って、分割データ510とメタデータ520とをノード200に送信する。分割データ送信部150は、分割データ510に対するACKをノード200から受信することにより、ノード200との間で分割データ510の送達確認を行ってもよい。

【0037】

ノード200は、分割データ受信部210、分割データ送信部(または、単に送信部)220、処理部230、分割データ記憶部240、及び、ノード設定記憶部250を含む。

【0038】

分割データ受信部210は、データサーバ100から分割データ510とメタデータ520とを受信する。なお、分割データ受信部210は、分割データ510に対するACKをデータサーバ100へ返信することにより、データサーバ100との間で分割データ510の送達確認を行ってもよい。この場合、分割データ受信部210は、分割データ510のレプリカが他のノード200に生成されたときに、データサーバ100へACKを返信する。

【0039】

分割データ送信部220は、データサーバ100から分割データ510を受信した場合に、メタデータ520に従って、他のノード200に分割データ510のレプリカを生成する。本発明の第1の実施の形態では、ノード200の分割データ記憶部240が、他のノード200からも書き込み可能であると仮定する。分割データ送信部220は、メタデータ520のレプリカ生成先ノードID(前、及び、後)で示されるノード200の分割データ記憶部240に、分割データ510を書き込むことにより、レプリカを生成する。

【0040】

なお、分割データ送信部220は、レプリカ生成先ノードID(前、及び、後)で示されるノード200の関連データ受信部(図示せず)に分割データ510を送信し、関連データ受信部が分割データ記憶部240に当該分割データ510を書き込むことにより、レプリカを生成してもよい。

【0041】

ノード設定記憶部250は、ノード200の処理に係る情報である、ノード設定情報251を記憶する。ノード設定情報251は、例えば、管理者等により、予め設定される。

【0042】

図9は、本発明の第1の実施の形態における、ノード設定情報251の例を示す図である。ノード設定情報251は、処理定義を含む。

【0043】

ここで、処理定義は、抽出された対象情報に対して行うべき、加工処理(型チェック、型変換等)の処理内容を示す。図9の例では、処理定義において、XML形式からCSV形式への変換が定義されている。

【0044】

分割データ記憶部240は、分割データ受信部210がデータサーバ100から受信した分割データ510とメタデータ520、及び、他のノード200により生成された分割

10

20

30

40

50

データ 510 のレプリカを記憶する。

【0045】

処理部 230 は、メタデータ 520、及び、ノード設定情報 251 に従って、分割データ 510 に対する所定の処理（対象情報の抽出、加工、及び、分散ストレージへの書き込み）を行う。処理部 230 は、分割データ 510 に、抽出しようとする対象情報の一部しか含まれていない場合、分割データ 510、及び、当該分割データ 510 の隣接分割データのレプリカから、対象情報を抽出する。

【0046】

なお、データサーバ 100、及び、ノード 200 は、それぞれ、CPU (Central Processing Unit) とプログラムを記憶した記憶媒体を含み、プログラムに基づく制御によって動作するコンピュータであってもよい。また、データサーバ 100 における、データ記憶部 110、及び、サーバ設定記憶部 160 は、それぞれ個別の記憶媒体（例えば、メモリ、ハードディスク等）でも、1つの記憶媒体によって構成されてもよい。同様に、ノード 200 における、分割データ記憶部 240、及び、ノード設定記憶部 250 は、それぞれ個別の記憶媒体（例えば、メモリ、ハードディスク等）でも、1つの記憶媒体によって構成されてもよい。

【0047】

図 3 は、本発明の実施の形態において、データサーバ 100、及び、ノード 200 がコンピュータにより実現される場合の、分散処理システム 1 の構成を示すブロック図である。

【0048】

図 3 を参照すると、データサーバ 100 は、CPU 101、記憶媒体 102、及び、通信部 103 を含む。CPU 101 は、データ取得部 120、転送計画部 130、分割部 140、及び、分割データ送信部 150 の機能を実現するためのコンピュータプログラムを実行する。記憶媒体 102 は、データ記憶部 110、及び、サーバ設定記憶部 160 のデータを記憶する。通信部 103 は、ノード 200 に分割データ 510 を送信する。

【0049】

ノード 200 は、CPU 201、記憶媒体 202、及び、通信部 203 を含む。CPU 201 は、分割データ受信部 210、分割データ送信部 220、及び、処理部 230 の機能を実現するためのコンピュータプログラムを実行する。記憶媒体 202 は、分割データ記憶部 240、及び、ノード設定記憶部 250 のデータを記憶する。通信部 203 は、データサーバ 100 から分割データ 510 を受信する。また、通信部 203 は、他のノード 200 から隣接分割データのレプリカを受信、他のノード 200 へ分割データ 510 のレプリカを送信してもよい。

【0050】

次に、本発明の第 1 の実施の形態の動作について説明する。

【0051】

ここでは、図 7 のサーバ設定情報 161、図 9 のノード設定情報 251 が、それぞれ、サーバ設定記憶部 160、ノード設定記憶部 250 に記憶されていると仮定する。

【0052】

図 4 は、本発明の第 1 の実施の形態における、元データ 500 のインポート処理を示すフローチャートである。

【0053】

はじめに、データサーバ 100 のデータ取得部 120 は、データ記憶部 110 から元データ 500 を取得する（ステップ S101）。

【0054】

例えば、データ取得部 120 は、図 6 の元データ 500 を取得する。

【0055】

次に、転送計画部 130 は、転送計画 131 を生成する（ステップ S102）。ここで、転送計画部 130 は、元データ 500 をサーバ設定情報 161 の分割データサイズで分

10

20

30

40

50

割し、各分割データ510に対して、分割データIDを付与する。そして、転送計画部130は、サーバ設定情報161の送信先決定方法に従って、送信先ノード群に含まれるノード200の中から各分割データ510の送信先を決定する。さらに、転送計画部130は、各分割データ510のメタデータ520におけるレプリカ生成先ノードID(前)に、当該分割データ510のレプリカを関連データ(後)として用いる(当該分割データ510の前方の隣接分割データを処理対象とする)他のノード200の識別子を設定する。また、転送計画部130は、各分割データ510のメタデータ520におけるレプリカ生成先ノードID(後)に、当該分割データ510のレプリカを関連データ(前)として用いる(当該分割データ510の後方の隣接分割データを処理対象とする)他のノード200の識別子を設定する。

10

【0056】

例えば、転送計画部130は、図6の元データ500を、図7のサーバ設定情報161における分割データサイズに従って分割した場合の各分割データ510に対して、図8に示すように、分割データID「D1」、「D2」、...を付与する。転送計画部130は、図8に示すように、図7のサーバ設定情報161の送信先決定方法(ラウンドロビン)に従って、分割データ510「D1」、「D2」、...の送信先を、それぞれ、ノード200「N1」、「N2」、...に決定する。また、転送計画部130は、図8に示すように、分割データ510「D1」のメタデータ520におけるレプリカ生成先ノードID(後)に、分割データ510「D1」のレプリカを使用する(隣接分割データ「D2」を処理対象とする)ノード200「N2」を、また、関連データID(後)に後方の隣接分割データ「D2」を設定する。また、転送計画部130は、分割データ510「D2」のメタデータ520におけるレプリカ生成先ノードID(前)に、分割データ510「D2」のレプリカを使用する(隣接分割データ「D1」を処理対象とする)ノード200「N1」を、レプリカ生成先ノードID(後)に、分割データ510「D2」のレプリカを使用する(隣接分割データ「D3」を処理対象とする)ノード200「N3」を、さらに、関連データID(前)に前方の隣接分割データ「D1」を、関連データID(後)に後方の隣接分割データ「D3」を、それぞれ設定する。

20

【0057】

分割部140は、転送計画131に含まれる分割データIDの先頭から順番に、分割データIDを1つ選択する(ステップS103)。

30

【0058】

分割部140は、元データ500から、選択した分割データIDに対応する分割データ510を生成する(ステップS104)。

【0059】

分割データ送信部150は、生成した分割データ510と、転送計画131に含まれる当該分割データ510に対応するメタデータ520とを、転送計画131における当該分割データ510に対応する送信先ノードIDのノード200に送信する(ステップS105)。分割データ送信部150は、ノード200から送信した分割データ510に対するACKを受信した場合、当該分割データ510を送信済みにする。

【0060】

分割部140、分割データ送信部150は、転送計画131に含まれる全ての分割データIDについて、ステップS103~S105を繰り返す(ステップS106)。

40

【0061】

なお、分割部140、分割データ送信部150は、サーバ設定情報161における送信並列度に応じて、ステップS103~S105を、複数の分割データ510に対して、送達確認を待たずに並列に実施してもよい。

【0062】

例えば、図7のサーバ設定情報161の送信並列度は3であるため、分割部140は、図8の転送計画131をもとに、図6に示すように、元データ500から分割データ510「D1」、「D2」、「D3」を生成する。そして、分割データ送信部150は、図6

50

に示すように、分割データ510「D1」、「D2」、「D3」に、図8の転送計画131における対応するメタデータ520を付与し、それぞれ、ノード200「N1」、「N2」、「N3」に送信する。

【0063】

次に、ノード200の分割データ受信部210は、データサーバ100から分割データ510とメタデータ520とを受信する(ステップS201)。分割データ受信部210は、受信した分割データ510とメタデータ520とを、分割データ記憶部240に保存する。

【0064】

例えば、ノード200「N1」、「N2」、「N3」の分割データ受信部210は、それぞれ、図6に示すような分割データ510「D1」、「D2」、「D3」とメタデータ520とを受信する。

10

【0065】

分割データ送信部220は、メタデータ520のレプリカ生成先ノードID(前、及び、後)で示されるノード200の分割データ記憶部240に、分割データ510のレプリカを生成する(ステップS202)。分割データ受信部210は、分割データ510のレプリカが他のノード200に生成された時点で、データサーバ100へ当該分割データ510に対するACKを返信する。

【0066】

例えば、ノード200「N1」の分割データ送信部220は、図6における分割データ510「D1」のメタデータ520に従って、図5に示すように、分割データ510「D1」のレプリカをノード200「N2」に生成する。同様に、ノード200「N2」の分割データ送信部220は、分割データ510「D2」のレプリカを、ノード200「N1」、「N3」に生成する。

20

【0067】

次に、処理部230は、分割データ記憶部240から分割データ510を取得し、当該分割データ510から対象情報が抽出可能かどうか判定する(ステップS203)。ここで、処理部230は、対象情報の開始ポイント/終了ポイントのデリミタを検出することにより、対象情報が抽出可能かどうか判定する。処理部230は、分割データ510に、開始ポイントのデリミタと当該開始ポイントに対応する終了ポイントのデリミタとが含まれている場合、対象情報を抽出可能と判断する。また、処理部230は、分割データ510に、開始ポイントのデリミタが含まれているが、当該開始ポイントに対応する終了ポイントのデリミタが含まれていない場合、対象情報を抽出不可と判断する。

30

【0068】

ステップS204で、対象情報の抽出可能の場合(ステップS203/Y)、処理部230は、分割データ510から対象情報を抽出する(ステップS205)。

【0069】

ステップS204で、対象情報の抽出不可の場合(ステップS203/N)、処理部230は、分割データ記憶部240から、メタデータ520の関連データID(後)で示される、分割データ510の後方の隣接分割データのレプリカを取得する。

40

【0070】

処理部230は、分割データ510と隣接分割データのレプリカとから対象情報が抽出可能かどうか判定する(ステップS204)。ここで、処理部230は、隣接分割データのレプリカに、分割データ510に含まれる開始ポイントに対応する終了ポイントのデリミタが含まれている場合、対象情報を抽出可能と判断する。

【0071】

ステップS204で、対象情報の抽出可能の場合(ステップS204/Y)、処理部230は、分割データ510と隣接分割データのレプリカとから対象情報を抽出する(ステップS206)。

【0072】

50

図10は、本発明の第1の実施の形態における、対象情報の抽出、及び、加工の例を示す図である。

【0073】

例えば、図10に示すように、ノード200「N1」において、分割データ510「D1」には、イベント情報「E1」の開始ポイントのデリミタ<event>は含まれるが、終了ポイントのデリミタ</event>は含まれない。また、隣接分割データ「D2」のレプリカに、終了ポイントのデリミタ</event>が含まれる。従って、ノード200「N1」の処理部230は、図10に示すように、分割データ510「D1」と隣接分割データ「D2」のレプリカとから、イベント情報「E1」を抽出する。

【0074】

同様に、図10に示すように、ノード200「N2」において、分割データ510「D2」には、イベント情報「E2」の開始ポイントのデリミタ<event>は含まれるが、終了ポイントのデリミタ</event>は含まれない。また、隣接分割データ「D3」のレプリカに、終了ポイントのデリミタ</event>が含まれる。従って、ノード200「N2」の処理部230は、図10に示すように、分割データ510「D2」と隣接分割データ「D3」のレプリカとから、イベント情報「E2」を抽出する。

【0075】

処理部230は、抽出された対象情報に対して、ノード設定情報251の処理定義で示される加工処理を実行する(ステップS207)。

【0076】

例えば、ノード200「N1」、「N2」の処理部230は、図9のノード設定情報251における処理定義に従って、それぞれ、図10に示すように、イベント情報「E1」、「E2」をXML形式からCSV形式に変換する。

【0077】

処理部230は、加工処理された対象情報を、分散ストレージへ書き込む(ステップS208)。

【0078】

例えば、ノード200「N1」、「N2」の処理部230は、それぞれ、図10に示す、CSV形式のイベント情報「E1」、「E2」を、分散ストレージへ書き込む。

【0079】

以上により、本発明の第1の実施の形態の動作が完了する。

【0080】

なお、本発明の第1の実施の形態では、処理部230は、分割データ510に開始ポイントのデリミタが含まれている対象情報を抽出している。しかしながら、処理部230は、分割データ510に終了ポイントのデリミタが含まれている対象情報を抽出してもよい。この場合、分割データ510に終了ポイントに対応する開始ポイントのデリミタが含まれていなければ、処理部230は、当該分割データ510と前方の隣接分割データのレプリカとを用いて、対象情報を抽出する。

【0081】

また、各ノード200の処理部230は、例えば、複数のノード200における、全ての分割データ510に対する所定の処理が完了した時点で、分割データ記憶部240に記憶されている分割データ510や隣接分割データを削除してもよい。

【0082】

また、本発明の第1の実施の形態では、元データ500のデータ形式として、XML形式を用いたが、データ形式は、CSV(comma-separated values)形式やJSON(Java(登録商標)Script Object Notation)形式、ログファイル等、XML形式以外の形式であってもよい。データ形式が、JSON形式の場合は、XML形式の場合と同様に、対象情報を囲むタグを対象情報の開始ポイント/終了ポイントを表すデリミタとして用いることができる。また、データ形式がCSV形式の場合は改行コード、ログファイルの場合は日時を、対象情報の開始ポイント/終了ポイントを表すデリミタとして用いることができる

10

20

30

40

50

。

【0083】

また、本発明の第1の実施の形態では、各ノード200が分割データ510に対する所定の処理として、対象情報の抽出、加工、及び、分散ストレージへの書き込みを行っているが、分散ストレージへの書き込みは行わなくてもよい。また、所定の処理は、これらの処理とは異なる他の処理でもよい。

【0084】

また、データサーバ100は、分割データ510の圧縮や暗号化を行って、各ノード200に送信してもよい。この場合、各ノード200は、圧縮された分割データ510のレプリカを、他のノード200に生成してもよい。これにより、レプリカの生成に係るノード200間の通信量や、メモリ使用量を低減できる。

10

【0085】

また、データサーバ100は、分割データサイズを動的に変更してもよい。この場合、データサーバ100は、例えば、各ノード200で抽出された対象情報の平均サイズ等をもとに、分割データサイズを決定する。また、この場合、エラー時のログレコード等、異常な大きさの対象情報を除外して、分割データサイズを決定してもよい。

【0086】

また、本発明の第1の実施の形態では、各ノード200は、分割データ510の関連データとして、分割データ510の前または後に隣接する、1つの分割データ510のレプリカを用いているが、分割データ510の前または後に隣接する、連続する2つ以上の分割データ510のレプリカを用いてもよい。これにより、対象情報が大きい場合でも、各ノード200で対象情報を抽出できる。

20

【0087】

また、分割データ510の関連データは、例えば、分割データ510にリンクにより関連づけられた他の分割データ510等、分割データ510に対する所定の処理において利用される他の分割データ510あれば、元データ500上で隣接する分割データ510以外の分割データ510でもよい。

【0088】

また、本発明の第1の実施の形態では、各ノード200は、メタデータ520におけるレプリカ送信先ノードIDに従って、分割データ510のレプリカを他のノード200に生成したが、例えば、データサーバ100からノード200への分割データ510の送信がラウンドロビンに行われる場合等、ノード200が、処理対象の分割データ510を関連データとして用いる他のノード200が認識できる場合は、メタデータ520を用いることなく、当該他のノード200に分割データ510のレプリカを生成してもよい。

30

【0089】

次に、本発明の第1の実施の形態の特徴的な構成を説明する。図1は、本明の第1の実施の形態の特徴的な構成を示すブロック図である。

【0090】

分散処理システム（情報処理システム）1は、ノード（処理装置）200を含む。ノード200は、分割データ送信部（送信部）220と処理部230とを含む。分割データ送信部220は、複数の分割データ（データ）510の内の処理対象の分割データ510を関連データとして用いる可能性がある他のノード200へ、当該分割データ510を送信する。処理部230は、分割データ510と、他のノード200から受信した、当該分割データ510の関連データと、を用いて、当該分割データ510に対する所定の処理を行う。

40

【0091】

次に、本発明の第1の実施の形態の効果の説明する。

【0092】

本発明の第1の実施の形態によれば、複数データを複数のノード200で分散処理するシステムにおいて、システムの処理負荷を低減できる。その理由は、各ノード200の分

50

分割データ送信部 220 が、複数の分割データ 510 の内の処理対象の分割データ 510 を関連データとして用いる可能性がある他のノード 200 へ、当該分割データ 510 を送信し、処理部 230 が、処理対象の分割データ 510 と、他のノード 200 から受信した、当該分割データ 510 の関連データと、を用いて、当該分割データ 510 に対する所定の処理を行うためである。これにより、各ノード 200 は、分割データ 510 の関連データを保持するノード 200 を検索する必要はなく、ノード 200 における処理負荷が低減される。

【0093】

また、本発明の第 1 の実施の形態によれば、データサーバ 100 の処理負荷も低減できる。その理由は、データサーバ 100 が、元データ 500 を所定の大きさに分割し、ノード 200 が、処理対象の分割データ 510 と関連データとから、対象情報を抽出するためである。これにより、データサーバ 100 は、元データ 500 からデリミタを検出して対象情報を抽出する必要はなく、データサーバ 100 における処理負荷が低減される。また、これにより、各ノード 200 で対象情報の抽出が、分散して、並列に行われるため、システムの処理速度が向上する。

10

【0094】

(第 2 の実施の形態)

次に、本発明の第 2 の実施の形態について説明する。

【0095】

本発明の第 2 の実施の形態においては、分割データ 510 の全部のレプリカを生成する代わりに、分割データ 510 の一部のレプリカを生成する点において、本発明の第 1 の実施の形態と異なる。

20

【0096】

次に、本発明の第 2 の実施の形態における、分散並列処理基盤への元データ 500 のインポートについて説明する。

【0097】

図 11 は、本発明の第 2 の実施の形態における、分散並列処理基盤への元データ 500 のインポートを示す図である。

【0098】

各ノード 200 は、データサーバ 100 から受信した分割データ 510 (処理対象の分割データ 510) に、抽出しようとする対象情報の一部しか含まれていない場合、受信した分割データ 510 の隣接分割データの一部(前半分、または、後半分)のレプリカを用いて対象情報を抽出する。本発明の第 2 の実施の形態においては、分割データ 510 の隣接分割データの一部のレプリカを、関連データと呼ぶ。各ノード 200 は、データサーバ 100 から分割データ 510 を受信したときに、当該分割データ 510 の一部(前半分、または、後半分)を関連データとして用いる(当該分割データ 510 の隣接分割データを処理対象とする)他のノード 200 に、当該分割データ 510 の一部(前半分、または、後半分)のレプリカを生成する。

30

【0099】

次に、本発明の第 2 の実施の形態における、分散処理システム 1 の構成を説明する。

40

【0100】

本発明の第 2 の実施の形態における分散処理システム 1 の構成は、本発明の第 1 の実施の形態(図 2)と同様となる。

【0101】

ノード 200 の分割データ送信部 220 は、データサーバ 100 から分割データ 510 を受信した場合に、メタデータ 520 に従って、他のノード 200 に分割データ 510 の一部(前半分、または、後半分)のレプリカを生成する。

【0102】

処理部 230 は、分割データ 510 に、抽出しようとする対象情報の一部しか含まれていない場合、分割データ 510、及び、当該分割データ 510 に係る隣接分割データの

50

部のレプリカから、対象情報を抽出する。

【0103】

次に、本発明の第2の実施の形態の動作について説明する。

【0104】

本発明の第2の実施の形態における、データサーバ100、及び、ノード200の処理を示すフローチャートは、本発明の第1の実施の形態(図4)と同様となる。

【0105】

図4のステップS202において、分割データ送信部220は、メタデータ520のレプリカ生成先ノードID(前)で示されるノード200の分割データ記憶部240に、分割データ510の前半分のレプリカを生成する。同様に、分割データ送信部220は、メタデータ520のレプリカ生成先ノードID(後)で示されるノード200の分割データ記憶部240に、分割データ510の後半分のレプリカを生成する。

10

【0106】

例えば、ノード200「N1」の分割データ送信部220は、図6における分割データ510「D1」のメタデータ520に従って、図11に示すように、分割データ510「D1」の後半分のレプリカをノード200「N2」に生成する。同様に、ノード200「N2」の分割データ送信部220は、分割データ510「D2」の前半分のレプリカをノード200「N1」に、後半分のレプリカをノード200「N3」に、それぞれ生成する。

【0107】

図4のステップS206において、処理部230は、分割データ510と隣接分割データの一部のレプリカとから対象情報を抽出する。

20

【0108】

図12は、本発明の第2の実施の形態における、対象情報の抽出、及び、加工の例を示す図である。

【0109】

例えば、ノード200「N1」の処理部230は、図12に示すように、分割データ510「D1」と隣接分割データ「D2」の前半分のレプリカとから、イベント情報「E1」を抽出する。同様に、ノード200「N2」の処理部230は、図12に示すように、分割データ510「D2」と隣接分割データ「D3」の前半分のレプリカとから、イベント情報「E2」を抽出する。

30

【0110】

以上により、本発明の第2の実施の形態の動作が完了する。

【0111】

なお、本発明の第2の実施の形態では、各ノード200は、他のノード200に、分割データ510の前半分または後半分のレプリカを生成しているが、他のノード200の処理対象である分割データ510に隣接する部分を含めば、レプリカの大きさは、半分より大きくても、小さくてもよい。

【0112】

次に、本発明の第2の実施の形態の効果を説明する。

40

【0113】

本発明の第2の実施の形態によれば、本発明の第1の実施の形態に比べて、分割データ510のレプリカの生成に係るコストを低減し、システムの処理速度をより高速化できる。その理由は、各ノード200が、他のノード200に、当該分割データ510の一部のレプリカを生成するためである。特に、分割データサイズと対象情報の大きさが近い場合は、分割データ510から対象情報の全てが含まれていない場合でも、隣接分割データの一部があれば、分割データ510と隣接分割データとから、対象情報を抽出できる可能性が高いため、上述の効果をえられる。

【0114】

(第3の実施の形態)

50

次に、本発明の第3の実施の形態について説明する。

【0115】

本発明の第3の実施の形態においては、ノード200において障害が発生した場合に、他のノード200が所定の処理を引き継ぐ点において、本発明の第1の実施の形態と異なる。

【0116】

次に、本発明の第3の実施の形態における、分散処理システム1の構成を説明する。

【0117】

図13は、本発明の第3の実施の形態における、分散処理システム1の構成を示すブロック図である。

10

【0118】

図13を参照すると、本発明の第3の実施の形態における分散処理システム1のデータサーバ100は、本発明の第3の実施の形態のデータサーバ100の構成に加えて、障害監視部170、及び、引継制御部180を含む。

【0119】

障害監視部170は、ノード200における障害を検出する。

【0120】

引継制御部180は、ノード200における障害が検出された場合に、当該ノード200の所定の処理を引き継ぐべきノード200（引き継ぎ先ノード200）を決定し、引き継ぎ指示を送信する。

20

【0121】

ノード200の処理部230は、データサーバ100から受信した分割データ510（処理対象の分割データ510）の隣接分割データのレプリカと、当該処理対象の分割データ510とを用いて、当該隣接分割データに対する所定の処理を行う（障害が検出されたノード200が実行すべき所定の処理を引き継ぐ）。

【0122】

次に、本発明の第3の実施の形態の動作について説明する。

【0123】

本発明の第3の実施の形態における元データ500のインポート処理は、本発明の第1の実施の形態と同様となる。

30

【0124】

図14は、本発明の第3の実施の形態における、引き継ぎ処理を示すフローチャートである。

【0125】

ここでは、既に、インポート処理において、データサーバ100からノード200への分割データ510の送信、及び、ノード200間での分割データ510のレプリカの生成が行われており、各ノード200が所定の処理（対象情報の抽出、加工、及び、分散ストレージへの書き込み）を実行中であると仮定する。

【0126】

はじめに、データサーバ100の障害監視部170は、ノード200の障害を検出する（ステップS301）。ここで、障害監視部170は、例えば、各ノード200との間で、死活確認のメッセージを送受信することにより、障害を検出する。

40

【0127】

例えば、障害監視部170は、図5におけるノード200「N1」の障害を検出する。

【0128】

引継制御部180は、引き継ぎ先ノード200を決定する（ステップS302）。ここで、引継制御部180は、転送計画131のメタデータ520を参照し、障害が検出されたノード200の処理対象の分割データ510に対するレプリカ送信先ノードID（後）で示されるノード200を、引き継ぎ先ノード200に決定する。

【0129】

50

例えば、引継制御部 180 は、図 8 の転送計画 131 のメタデータ 520 を参照し、ノード 200 「N1」の処理対象の分割データ 510 「D1」のレプリカ送信先であるノード 200 「N2」を、引き継ぎ先ノード 200 に決定する。

【0130】

引継制御部 180 は、引き継ぎ先ノード 200 に、引き継ぎ指示を送信する（ステップ S303）。ここで、引き継ぎ指示は、引継ぐべき分割データ 510 の分割データ ID、及び、当該分割データ 510 に対する関連データ ID（後）を含む。

【0131】

例えば、引継制御部 180 は、分割データ ID 「D1」、関連データ ID（後）「D2」を含む引き継ぎ指示をノード 200 「N2」に送信する。

10

【0132】

ノード 200 の処理部 230 は、引き継ぎ指示を受信する（ステップ S401）。

【0133】

次に、処理部 230 は、分割データ記憶部 240 から、引き継ぎ指示の分割データ ID で示される分割データ 510 のレプリカ、すなわち、処理対象の分割データ 510 の前方の隣接分割データのレプリカを取得する。処理部 230 は、当該隣接分割データのレプリカから対象情報が抽出可能かどうか判定する（ステップ S402）。ここで、処理部 230 は、隣接分割データのレプリカに、開始ポイントのデリミタと当該開始ポイントに対応する終了ポイントのデリミタとが含まれている場合、対象情報を抽出可能と判断する。また、処理部 230 は、隣接分割データのレプリカに、開始ポイントのデリミタが含まれているが、当該開始ポイントに対応する終了ポイントのデリミタが含まれていない場合、対象情報を抽出不可と判断する。

20

【0134】

ステップ S402 で、対象情報の抽出可能の場合（ステップ S402 / Y）、処理部 230 は、当該隣接分割データのレプリカから対象情報を抽出する（ステップ S404）。

【0135】

ステップ S402 で、対象情報の抽出不可の場合（ステップ S402 / N）、処理部 230 は、分割データ記憶部 240 から、引き継ぎ指示の関連データ ID（後）で示される分割データ 510、すなわち、処理対象の分割データ 510 を取得する。

【0136】

処理部 230 は、隣接分割データのレプリカと処理対象の分割データ 510 とから対象情報が抽出可能かどうか判定する（ステップ S403）。ここで、処理部 230 は、処理対象の分割データ 510 に、隣接分割データのレプリカに含まれる開始ポイントに対応する終了ポイントのデリミタが含まれている場合、対象情報を抽出可能と判断する。

30

【0137】

ステップ S403 で、対象情報の抽出可能の場合（ステップ S403 / Y）、処理部 230 は、隣接分割データのレプリカと処理対象の分割データ 510 とから対象情報を抽出する（ステップ S405）。

【0138】

図 15 は、本発明の第 3 の実施の形態における、引き継ぎ処理における対象情報の抽出、及び、加工の例を示す図である。

40

【0139】

例えば、ノード 200 「N2」の処理部 230 は、図 10 に示すように、隣接分割データ「D1」のレプリカと分割データ 510 「D2」とから、イベント情報「E1」を抽出する。

【0140】

以降、処理部 230 は、ステップ S207、S208 と同様に、抽出された対象情報に対して加工処理を実行し、分散ストレージへ書き込む（ステップ S406、S407）。

【0141】

以上により、本発明の第 3 の実施の形態の動作が完了する。

50

【 0 1 4 2 】

なお、本発明の第3の実施の形態では、データサーバ100の障害監視部170が、ノード200の障害を検出し、引継制御部180が、引き継ぎ先ノード200に引き継ぎ指示を送信しているが、各ノード200が、引き継ぎ対象のノード200の障害を検出し、当該ノード200の所定の処理を引き継いでもよい。この場合、各ノード200は、例えば、メタデータ520のレプリカ生成先ノードID(前)で示されたノード200の障害を検出した場合に、関連データID(前)で示される、処理対象の分割データ510の前方の隣接分割データのレプリカと当該処理対象の分割データ510とから、隣接分割データに対する所定の処理を実行する。

【 0 1 4 3 】

また、データサーバ100が、ノード200の障害に代わって、ノード200における処理対象の分割データ510のロスト(紛失)を検出し、引き継ぎ先ノード200が、分割データ510をロストしたノード200の所定の処理を引き継いでもよい。

【 0 1 4 4 】

次に、本発明の第3の実施の形態の効果の説明する。

【 0 1 4 5 】

本発明の第3の実施の形態によれば、複数のノード200の内のいずれかに障害や分割データ510のロストが発生した場合でも、所定の処理を継続できる。その理由は、ノード200において障害や分割データ510のロストが発生した場合に、他のノード200が、障害や分割データ510のロストが発生したノード200から受信した、処理対象の分割データ510の隣接分割データのレプリカと、当該処理対象の分割データ510とを用いて、当該隣接分割データに対する所定の処理を引き継ぐためである。これにより、ノード200において障害や分割データ510のロストが発生した場合に、データサーバ100が引き継ぎ先ノードに対して、分割データ510の再送を行うことなく、引き継ぎ処理が実行でき、データサーバ100の負荷を低減できるとともに、引き継ぎ処理が高速化される。また、メタデータ520に、分割データ510のレプリカ送信先や、分割データ510の隣接分割データに係る情報が含まれているため、データサーバ100は、メタデータ520を参照することにより、引き継ぎ先ノードの決定や、引き継ぎ指示を容易に行うことができる。

【 0 1 4 6 】

以上、実施形態を参照して本願発明を説明したが、本願発明は上記実施形態に限定されるものではない。本願発明の構成や詳細には、本願発明のスコープ内で当業者が理解し得る様々な変更をすることができる。

【 符号の説明 】

【 0 1 4 7 】

- 1 分散処理システム
- 100 データサーバ
- 101 CPU
- 102 記憶媒体
- 103 通信部
- 110 データ記憶部
- 120 データ取得部
- 130 転送計画部
- 131 転送計画
- 140 分割部
- 150 分割データ送信部
- 160 サーバ設定記憶部
- 161 サーバ設定情報
- 200 ノード
- 201 CPU

10

20

30

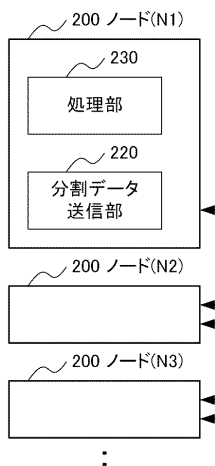
40

50

- 2 0 2 記憶媒体
- 2 0 3 通信部
- 2 1 0 分割データ受信部
- 2 2 0 分割データ送信部
- 2 3 0 処理部
- 2 4 0 分割データ記憶部
- 2 5 0 ノード設定記憶部
- 2 5 1 ノード設定情報
- 5 0 0 元データ
- 5 1 0 分割データ
- 5 2 0 メタデータ

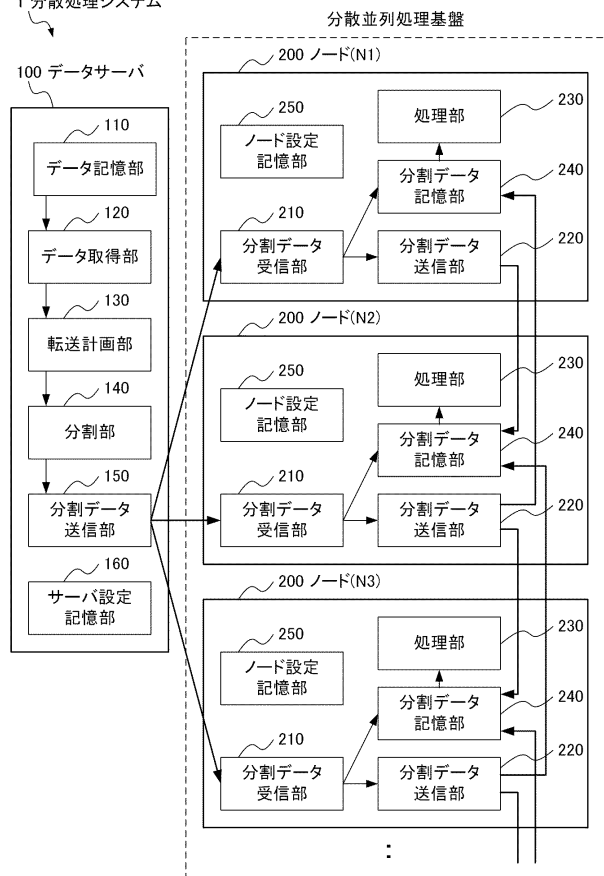
【図 1】

1 分散処理システム

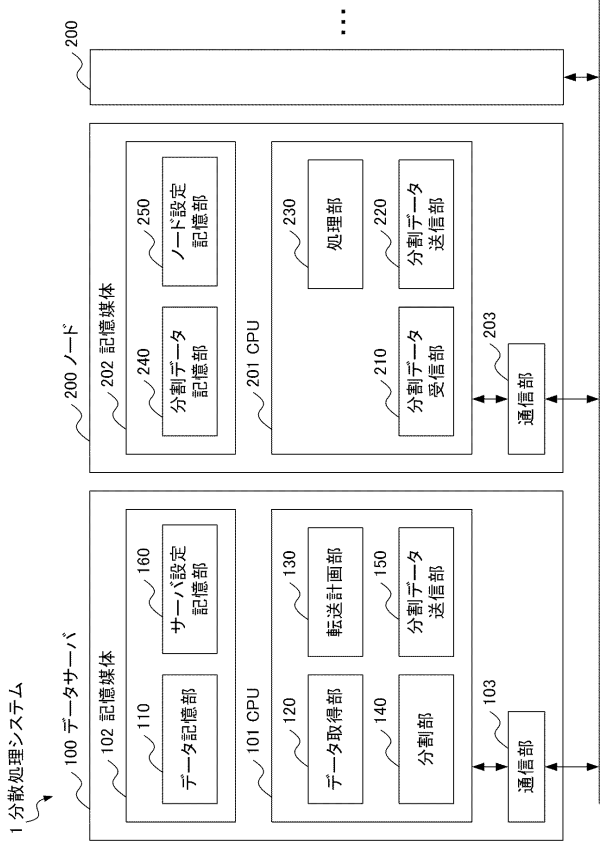


【図 2】

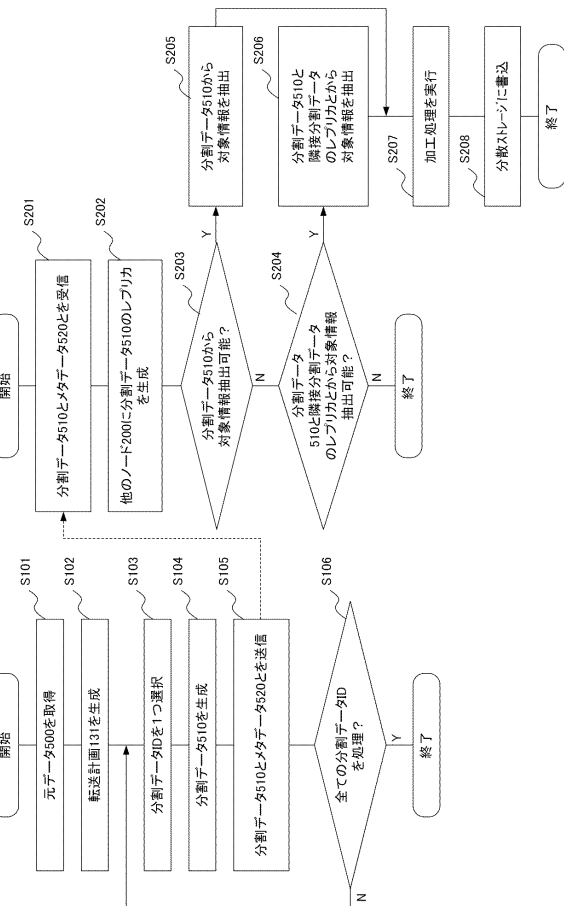
1 分散処理システム



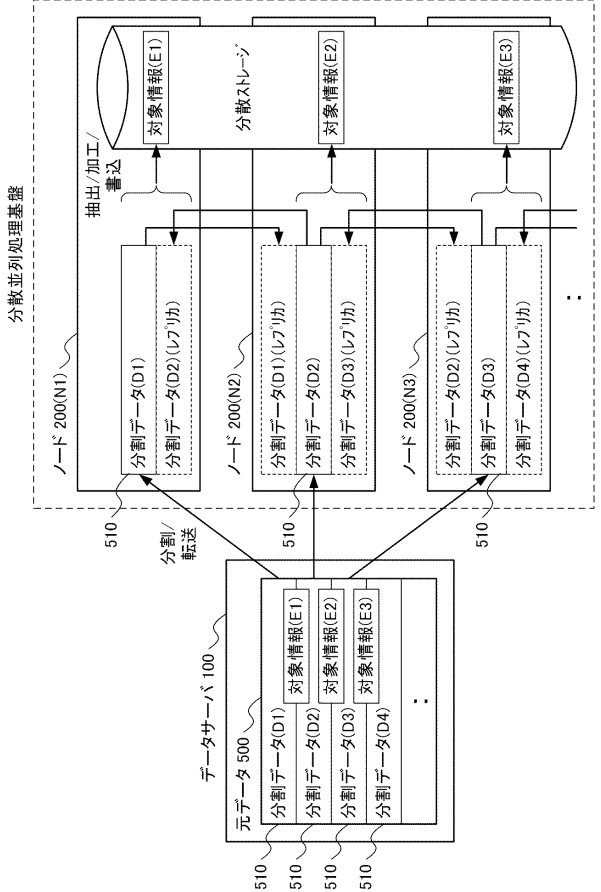
【図3】



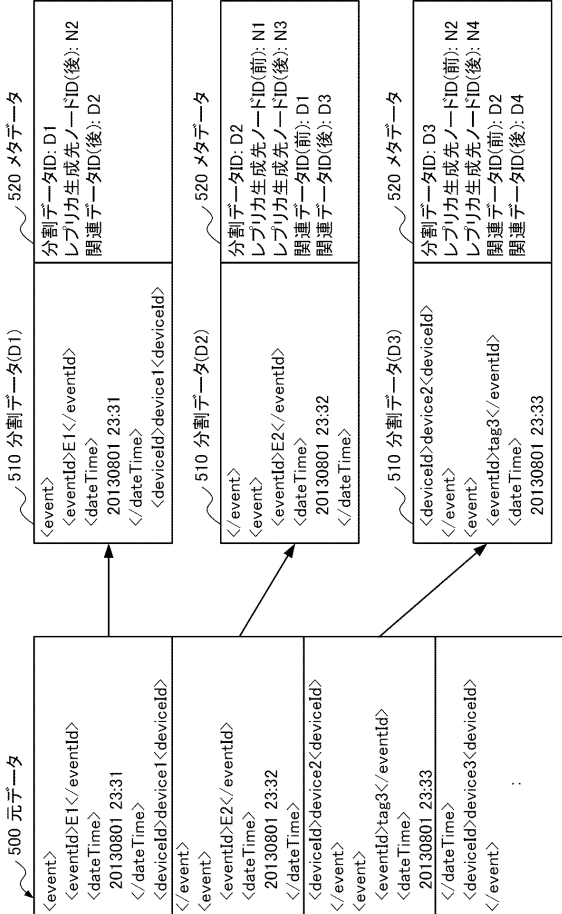
【図4】



【図5】



【図6】



【図7】

161 サーバ設定情報

送信先ノード群	N1, N2, N3, ...
送信先決定方法	ラウンドロビン
送信並列度	3
分割データサイズ	M

【図8】

131 転送計画

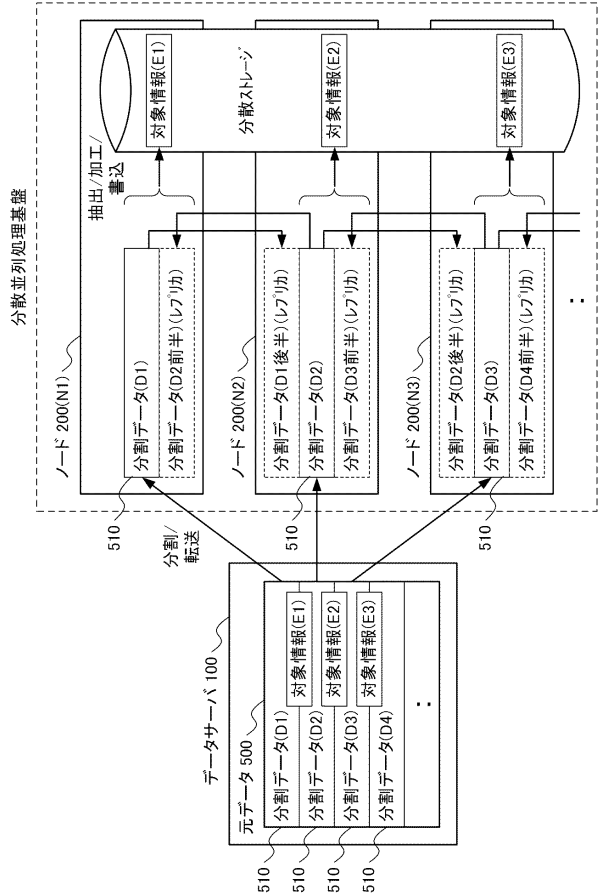
分割データID	送信先ノードID	メタデータ
D1	N1	分割データID: D1 レプリカ生成先ノードID(後): N2 関連データID(後): D2
D2	N2	分割データID: D2 レプリカ生成先ノードID(前): N1 レプリカ生成先ノードID(後): N3 関連データID(前): D1 関連データID(後): D3
D3	N3	分割データID: D3 レプリカ生成先ノードID(前): N2 レプリカ生成先ノードID(後): N4 関連データID(前): D2 関連データID(後): D4
⋮	⋮	⋮

【図9】

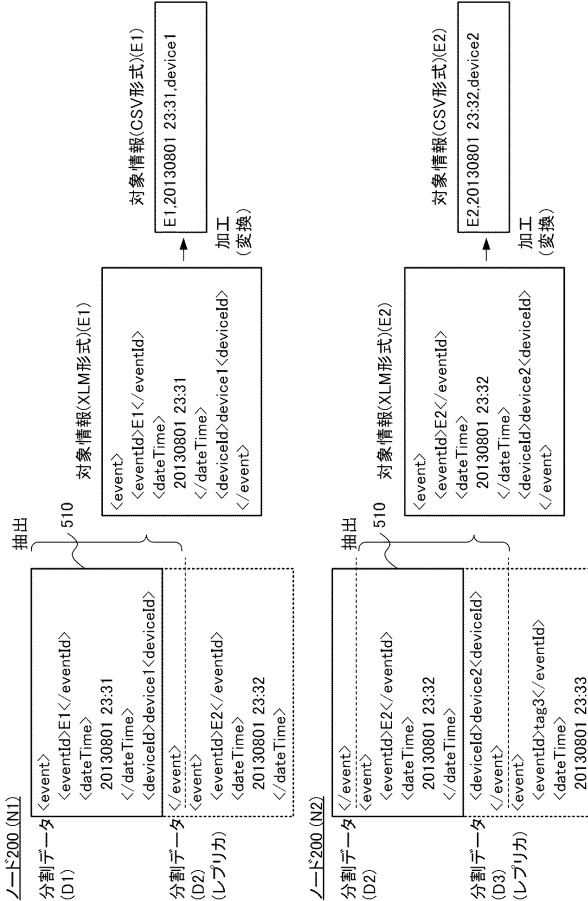
251 ノード設定情報

処理定義	XML形式 → CSV形式 XML: <event>単位 CSV: eventId, dateTime, deviceId
------	---------------------------------------------------------------------

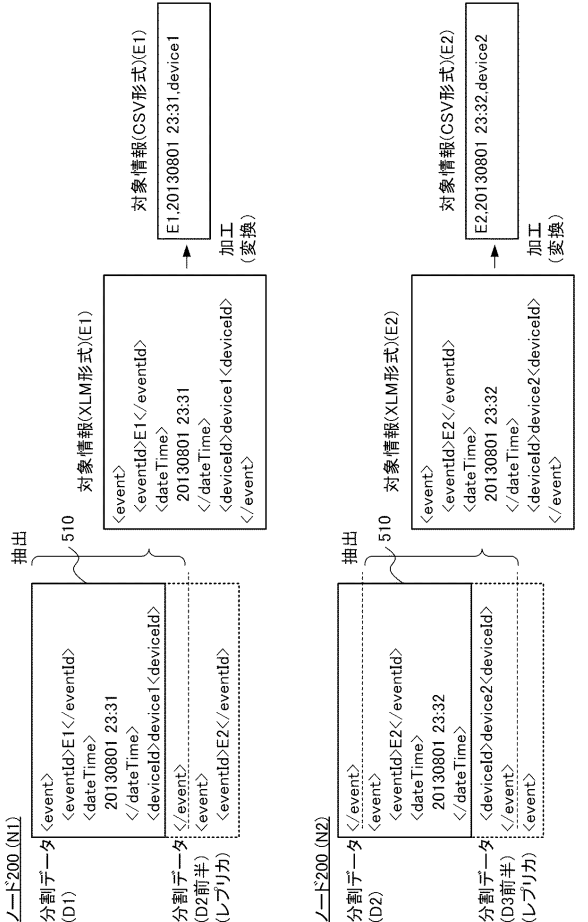
【図11】



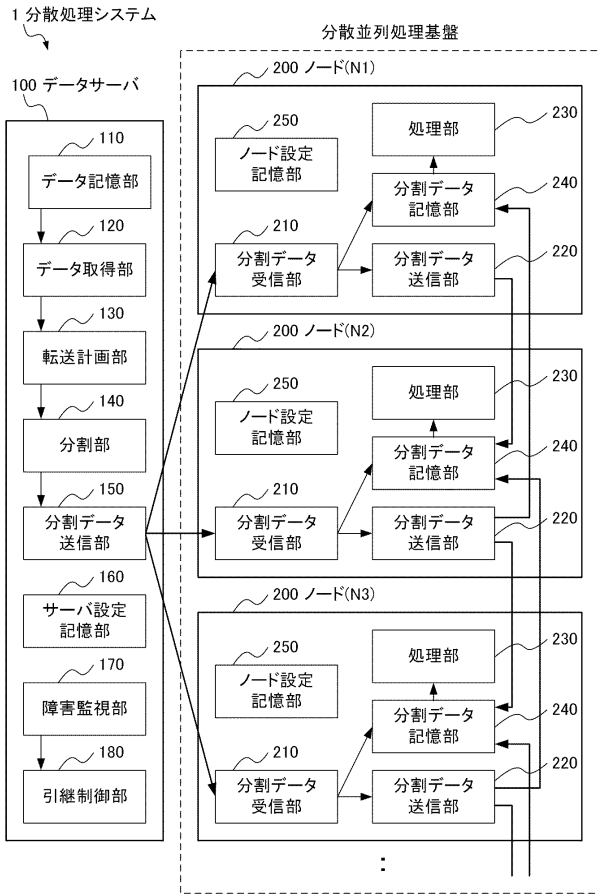
【図10】



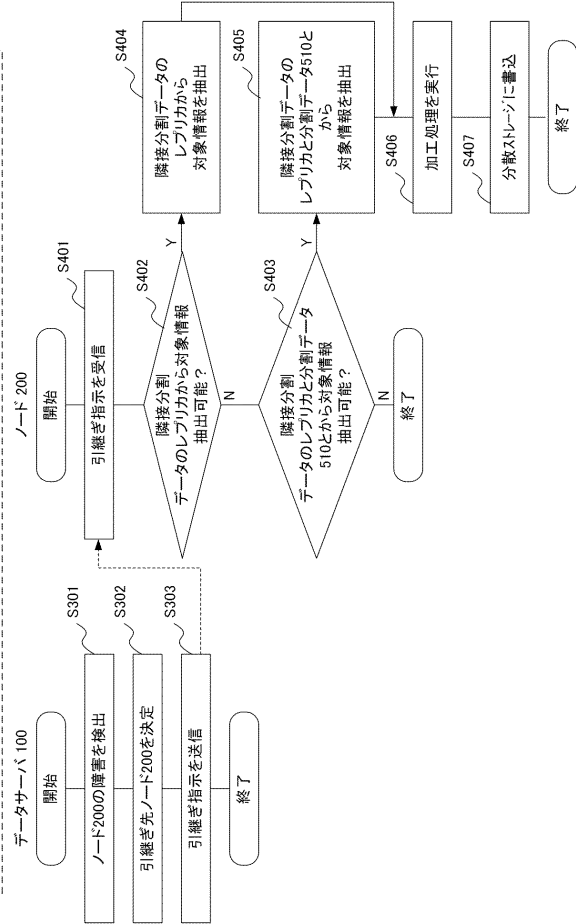
【図12】



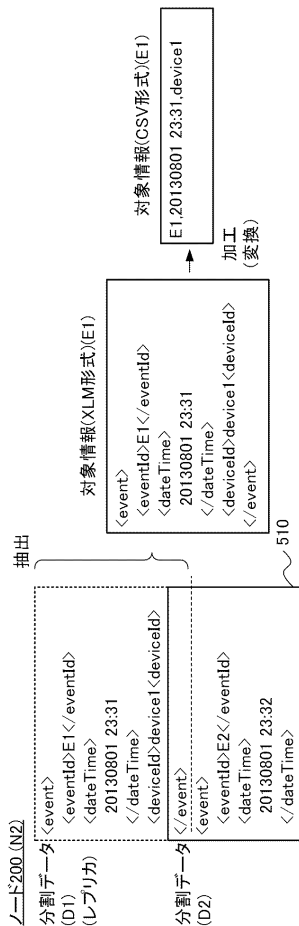
【図13】



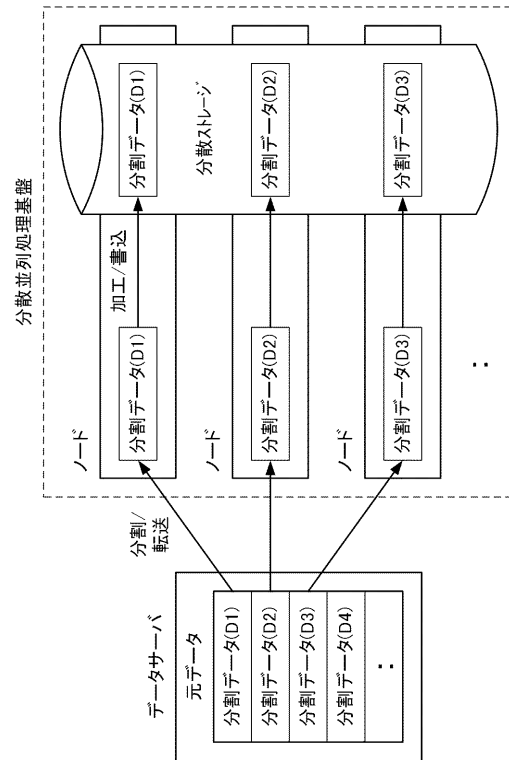
【図14】



【図15】



【図16】



フロントページの続き

(56)参考文献 特開2006-005524(JP,A)
特開2006-303952(JP,A)
特開2011-198122(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F 9/50