

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5594460号
(P5594460)

(45) 発行日 平成26年9月24日(2014.9.24)

(24) 登録日 平成26年8月15日(2014.8.15)

(51) Int.Cl. F 1
G 0 6 F 12/00 (2006.01) G 0 6 F 12/00 5 4 5 A
G 0 6 F 15/00 (2006.01) G 0 6 F 15/00 4 7 0

請求項の数 9 (全 16 頁)

<p>(21) 出願番号 特願2010-50042 (P2010-50042) (22) 出願日 平成22年3月8日(2010.3.8) (65) 公開番号 特開2011-186695 (P2011-186695A) (43) 公開日 平成23年9月22日(2011.9.22) 審査請求日 平成25年2月1日(2013.2.1)</p> <p>(出願人による申告)平成21年度、独立行政法人新エネルギー・産業技術総合開発機構、「グリーンネットワーク・システム技術研究開発プロジェクト(グリーンITプロジェクト)」委託研究、産業技術力強化法第19条の適用を受ける特許出願</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 (74) 代理人 100079005 弁理士 宇高 克己 (72) 発明者 上村 純平 東京都港区芝五丁目7番1号 日本電気株式会社内 (72) 発明者 柏木 岳彦 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 桜井 茂行</p>
---	---

最終頁に続く

(54) 【発明の名称】 送信情報制御装置、方法及びプログラム

(57) 【特許請求の範囲】

【請求項1】

分散データベースを構成する複数のストレージ処理部を備える分散データベースシステムにおいて用いられる送信情報制御装置であって、

データベース処理要求を受信し、該当するストレージ処理部に供給する処理要求供給手段と、

前記ストレージ処理部からの処理結果データを取得する処理結果取得手段と、

前記処理結果データを処理要求元へ送信する際のコストを、データの集約を行う場合と行わない場合についてそれぞれ計算するコスト計算手段と、

前記ストレージ処理部から取得した各処理結果データを集約するか否かを、データの集約を行う場合と行わない場合のコスト計算結果の比較結果に基づいて判定する判定手段と

10

、
 前記ストレージ処理部から取得した各処理結果データを集約すると判定した場合、当該各処理結果データを集約して処理要求元へ送信し、前記ストレージ処理部から取得した各処理結果データを集約しないと判定した場合、当該各処理結果データを集約せずに処理要求元へ送信する処理結果データ送信手段と、

を備えることを特徴とする送信情報制御装置。

【請求項2】

前記コスト計算手段は、前記処理結果データを処理要求元へ送信する際にかかる時間を前記コストとして、データの集約を行う場合と行わない場合についてそれぞれ計算し、

20

前記判定手段は、データの集約を行う場合のコストの方がデータの集約を行わない場合のコストより小さい場合にはデータの集約を行うと判定し、データの集約を行わない場合のコストの方がデータの集約を行う場合のコストより小さい場合にはデータの集約を行わないと判定する、

ことを特徴とする請求項 1 に記載の送信制御装置。

【請求項 3】

前記処理要求には、前記処理結果データの集約を必要とするか否かを示すポリシー情報が付与されており、

前記判定手段は、前記ポリシー情報に、前記処理結果データを集約してもよい旨が設定されている場合に、前記処理結果データを集約するか否かを前記コスト計算結果に基づいて判定する、

ことを特徴とする請求項 1 又は 2 に記載の送信情報制御装置。

【請求項 4】

分散データベースを構成する複数のストレージ処理部を備える分散データベースシステムにおける送信情報制御方法であって、

データベース処理要求を受信して、該当するストレージ処理部に供給し、

前記ストレージ処理部からの処理結果データを取得し、

前記処理結果データを処理要求元へ送信する際のコストを、データの集約を行う場合と行わない場合についてそれぞれ計算し、

前記ストレージ処理部から取得した処理結果データを集約するか否かを、データの集約を行う場合と行わない場合のコスト計算結果の比較結果に基づいて判定し、

前記ストレージ処理部から取得した処理結果データを集約すると判定した場合、当該処理結果データを集約して処理要求元に送信し、前記ストレージ処理部から取得した処理結果データを集約しないと判定した場合、当該処理結果データを集約せずに処理要求元に送信する、

ことを特徴とする送信情報制御方法。

【請求項 5】

前記処理結果データを処理要求元へ送信する際にかかる時間を前記コストとして、データの集約を行う場合と行わない場合についてそれぞれ計算し、

データの集約を行う場合のコストの方がデータの集約を行わない場合のコストより小さい場合にはデータの集約を行うと判定し、データの集約を行わない場合のコストの方がデータの集約を行う場合のコストより小さい場合にはデータの集約を行わないと判定する、

ことを特徴とする請求項 4 に記載の送信制御方法。

【請求項 6】

前記データベース処理要求は、前記処理結果データの集約を必要とするか否かを示すポリシー情報が処理要求元により付与されており、

前記コスト計算結果に基づく判定は、前記ポリシー情報に、前記処理結果データを集約してもよい旨が設定されている場合に行う、

ことを特徴とする請求項 4 又は 5 に記載の送信情報制御方法。

【請求項 7】

コンピュータを、

データベース処理要求を受信し、該当するストレージ処理部に供給する処理要求供給手段、

前記ストレージ処理部からの処理結果データを取得する処理結果取得手段、

前記処理結果データを処理要求元へ送信する際のコストを、データの集約を行う場合と行わない場合についてそれぞれ計算するコスト計算手段、

前記ストレージ処理部から取得した処理結果データを集約するか否かを、データの集約を行う場合と行わない場合のコスト計算結果の比較結果に基づいて判定する判定手段、

前記ストレージ処理部から取得した処理結果データを集約すると判定した場合、当該処理結果データを集約して処理要求元に送信し、前記ストレージ処理部から取得した処理結

10

20

30

40

50

果データを集約しないと判定した場合、当該処理結果データを集約せずに処理要求元に送信する処理結果データ送信手段、

として機能させるためのプログラム。

【請求項 8】

前記コスト計算手段は、前記処理結果データを処理要求元へ送信する際にかかる時間を前記コストとして、データの集約を行う場合と行わない場合についてそれぞれ計算し、

前記判定手段は、データの集約を行う場合のコストの方がデータの集約を行わない場合のコストより小さい場合にはデータの集約を行うと判定し、データの集約を行わない場合のコストの方がデータの集約を行う場合のコストより小さい場合にはデータの集約を行わないと判定する、

ことを特徴とする請求項 7 に記載のプログラム。

【請求項 9】

前記処理要求には、前記処理結果データの集約を必要とするか否かを示すポリシー情報が付与されており、

前記判定手段は、前記ポリシー情報に、前記処理結果データを集約してもよい旨が設定されている場合に、前記処理結果データを集約するか否かを前記コスト計算結果に基づいて判定する、

ことを特徴とする請求項 7 又は 8 に記載のプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、分散データベースシステムにおいて、データベース処理の結果データの送信処理を制御する技術に関する。

【背景技術】

【0002】

データベースが保持するデータに対するクエリ処理は、多段階のデータ操作を行うことで達成される。一般的に多段階のデータ操作を行う処理方法として、ある処理が終わってから次の処理を開始するバッチ型の処理方法と、ある処理が終わった部分に対して次の処理を行うパイプライン型の処理の処理方法がある。分散データベースにおいていくつかの処理を分散して行う場合においても、ノード間で、パイプライン型の処理を行う方法とバッチ型の処理を行う方法が考えられる。

【0003】

パイプライン型処理は、前処理の一部の結果を用いて次の処理をすぐに開始できるため、次処理開始までの待ち時間が少なく、並列に処理を行うことができる。また、バッチ型処理は、処理結果を集約して次のノードに転送するため、総転送データサイズを小さくして短い時間で転送を行うことができる。

【0004】

例えば、特許文献 1 には、分散データベースにおいて、リンク情報を用いて階層テーブル同士を結合することで入れ子節を含む X Query 処理の高速化を可能とするデータベース処理装置が開示されている。

【先行技術文献】

【特許文献】

【0005】

【特許文献 1】特開 2009 - 211154 号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

特許文献 1 の分散データベースシステムでは、クエリの処理結果データを各 DB サーバから要求元サーバに送信する際、ネットワークやノードの負荷等のような動的に変化するシステム状況を全く考慮していない。このため、システム状況に応じて最適なデータ送信

10

20

30

40

50

処理方法を選択してデータベース処理の高速化を図ることができなかった。

【0007】

本発明は、上記問題点に鑑みてなされたもので、システム状況に適したデータ送信処理方法を選択することにより、分散データベースシステムにおける処理を高速化することができる送信情報制御装置、方法及びプログラム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】

【0008】

上記課題を解決する本発明は、分散データベースを構成する複数のストレージ処理部を備える分散データベースシステムにおいて用いられる送信情報制御装置であって、データベース処理要求を受信し、該当するストレージ処理部に供給する処理要求供給手段と、前記ストレージ処理部からの処理結果データを取得する処理結果取得手段と、前記処理結果データを処理要求元へ送信する際のコストを、データの集約を行う場合と行わない場合についてそれぞれ計算するコスト計算手段と、前記ストレージ処理部から取得した各処理結果データを集約するか否かを、データの集約を行う場合と行わない場合のコスト計算結果の比較結果に基づいて判定する判定手段と、前記ストレージ処理部から取得した各処理結果データを集約すると判定した場合、当該各処理結果データを集約して処理要求元に送信し、前記ストレージ処理部から取得した各処理結果データを集約しないと判定した場合、当該各処理結果データを集約せずに処理要求元に送信する処理結果データ送信手段と、を備えることを特徴とする送信情報制御装置である。

【0009】

上記課題を解決する本発明は、分散データベースを構成する複数のストレージ処理部を備える分散データベースシステムにおける送信情報制御方法であって、データベース処理要求を受信して、該当するストレージ処理部に供給し、前記ストレージ処理部からの処理結果データを取得し、前記処理結果データを処理要求元へ送信する際のコストを、データの集約を行う場合と行わない場合についてそれぞれ計算し、前記ストレージ処理部から取得した処理結果データを集約するか否かを、データの集約を行う場合と行わない場合のコスト計算結果の比較結果に基づいて判定し、前記ストレージ処理部から取得した処理結果データを集約すると判定した場合、当該処理結果データを集約して処理要求元に送信し、前記ストレージ処理部から取得した処理結果データを集約しないと判定した場合、当該処理結果データを集約せずに処理要求元に送信する、ことを特徴とする送信情報制御方法である。

【0010】

上記課題を解決する本発明は、コンピュータを、データベース処理要求を受信し、該当するストレージ処理部に供給する処理要求供給手段、前記ストレージ処理部からの処理結果データを取得する処理結果取得手段、前記処理結果データを処理要求元へ送信する際のコストを、データの集約を行う場合と行わない場合についてそれぞれ計算するコスト計算手段、前記ストレージ処理部から取得した処理結果データを集約するか否かを、データの集約を行う場合と行わない場合のコスト計算結果の比較結果に基づいて判定する判定手段、前記ストレージ処理部から取得した処理結果データを集約すると判定した場合、当該処理結果データを集約して処理要求元に送信し、前記ストレージ処理部から取得した処理結果データを集約しないと判定した場合、当該処理結果データを集約せずに処理要求元に送信する処理結果データ送信手段、として機能させるためのプログラムである。

【発明の効果】

【0011】

本発明によれば、システム状況に適したデータ送信処理方法を選択することにより、分散データベースシステムにおける処理を高速化することができる。

【図面の簡単な説明】

【0012】

【図1】図1は本発明の実施形態に係るルータ（送信制御装置）を用いる分散データベース管理システムの機能ブロック図である。

10

20

30

40

50

【図2】図2は分散データベースを構成するデータベース・テーブルT B Lの一例を概略的に示す図である。

【図3】図3は本発明の実施形態に係るルータの機能ブロック図である。

【図4】図4はストレージ処理部の機能ブロック図である。

【図5】図5は集約実行判定処理を説明するためのフローチャートである。

【図6】図6は分散データベース管理システムの動作を説明するための図である。

【図7】図7はデータの集約処理を説明するための図である。

【発明を実施するための形態】

【0013】

以下、本発明の実施形態について図面を参照して説明する。

10

【0014】

図1は、本発明の実施形態に係る送信制御装置（ルータ）を用いる分散データベース管理システム10の全体構成を概略的に示す機能ブロック図である。この分散データベース管理システム10は、ロードバランサ11と、クエリサーバ20A、20B、20Cと、データサーバ22₁～22_Nと、管理サーバ30とを備え、これはLAN（Local Area Network）等のネットワークにより接続されている。

【0015】

データサーバ22₁～22_Nには、それぞれ分散データベースを構成する部分データベースが格納されている。分散データベース管理システム10は、この分散データベースに対するデータ操作を行う。分散データベースは、少なくとも1つのテーブル構造を有し、部分データベースは、このテーブル構造のサブセット（部分集合）を構成する。

20

【0016】

図2は、分散データベースを構成するデータベース・テーブルT B Lの一例を概略的に示す図である。このデータベース・テーブルT B Lは、複数のタプル（行）と、列方向に定義されたカラム（属性フィールド）A₁、A₂、・・・、A_pとを有する。タプルとカラムとの交差領域で定まる領域にはデータが格納される。図2に示すように、このデータベース・テーブルT B Lを行方向に分割（水平分割）することで複数の部分集合T G₁、T G₂、・・・、T G_Nを構成することができる。このような部分集合T G₁、T G₂、・・・、T G_Nを、部分データベースのテーブルとしてそれぞれデータサーバ22₁～22_Nに格納することができる。

30

【0017】

なお、データベース・テーブルT B Lをカラム方向に分割（垂直分割）することで複数の部分データベース・テーブルを構成してもよいし、あるいは、水平分割と垂直分割との組み合わせにより複数の部分データベース・テーブルを構成してもよい。

【0018】

図1に例示されるように、通信網NWには、分散データベース管理システム10とクライアント端末T1とが接続されている。これら分散データベース管理システム10とクライアント端末T1だけでなく、多数のクライアント端末が通信網NWに接続される。ネットワークNWは、例えばインターネット等の広域ネットワークとLANを少なくとも含む。

40

【0019】

クライアント端末T1は、分散データベース管理システム10が有するデータベースについてSQLやXQuery（XML Query Language: XML問い合わせ言語）等のデータベース言語（データ操作言語）で記述されたクエリを生成し、このクエリを分散データベース管理システム10に宛てて送信する機能を有する。クエリには、分散データベースに対してデータの検索、挿入、更新あるいは削除等のデータ操作を規定するデータベース言語が記述されている。

【0020】

ロードバランサ11は、通信網NWを介してクライアント端末T1から送信されたクエリをデータ処理要求として受信し、このクエリをクエリサーバ20A～20Cに均等に振

50

り分けて処理不可を分散する機能を有する。ロードバランサ 11 は、例えばラウンドロビン方式に従ってクエリサーバ 20A ~ 20C のいずれかを選択してもよい。

【0021】

クエリサーバ 20A ~ 20C は、それぞれクエリ解析部 21A ~ 21C を備えている。クエリ解析部 21A ~ 21C は、ロードバランサ 11 により振り分けられたクエリを解析し最適化する機能を有している。クエリ解析部 21A ~ 21C は、受信したクエリを、ストレージ処理部 25 のデータ構造に基づいて最適化された解析ツリー形式のクエリに変換する。このとき、クエリを例えば抽象構文木 (AST: Abstract Syntax Tree) 形式のクエリに変換することができる。

【0022】

データサーバ 22₁ ~ 22_N の各々は、ルータ (送信情報制御装置) 24 と複数のストレージ処理部 25₁ ~ 25_M とを有する。データサーバ 22₁ ~ 22_N は、LAN などの有線伝送路又は無線伝送路を介して相互に接続されている。データサーバ 22_i 内のルータ 24 は、ストレージ処理部 25₁ ~ 25_M のうちの任意のストレージ処理部間のデータ転送を制御する機能と、他のデータサーバ 22_j 内のルータ 24 との間でデータ通信を行う機能を有する。

【0023】

図 3 は、ルータ 24 の構成を概略的に示す機能ブロック図である。ルータ 24 は、ルーティングテーブル 241 と、クエリ配送部 242 と、資源情報テーブル 243 と、回答受付部 244 と、集約実行部 245 と、データ送信部 246 とを備える。

【0024】

ルーティングテーブル 241 には、ストレージ処理部 25 と、ストレージ処理部 25 にそれぞれ格納されているデータベース・テーブルとの対応関係が設定されている。

【0025】

クエリ配送部 242 は、ルーティングテーブル 241 や管理サーバ 30 を参照して、受信したクエリの配送先となるストレージ処理部 25 や他のルータ 24 を決定し、クエリを配送する。クエリを他のルータ 24 に配送する場合には、クエリに付与されている静的集約実行ポリシーが「集約する」を示し且つルータ 24 自身の負荷状況が閾値以下である場合、そのクエリ情報中の静的集約実行ポリシーを「集約してもよい」に変更する。

【0026】

クエリに付与される静的集約実行ポリシーは、要求元が実行結果データをどのような形式で受け取りたいか (データの集約を必要とするか否か) を示しており、「集約する」 (集約を必要とする) と「集約してもよい」 (集約を必要としない) の 2 つの値のいずれかをとり。静的集約実行ポリシーは、クエリサーバ 20A ~ 20C のクエリ解析部 21A ~ 21C と、後述するストレージ処理部 25₁ ~ 25_M のクエリ解析部 252 により、クエリに付与される。

【0027】

資源情報テーブル 243 には、システム資源に関する資源情報が設定されている。資源情報は、クエリサーバ 20 や他のルータ 24 とのネットワーク速度、ノード間の負荷、CPU 未使用率 * CPU クロック数、他のノードの CPU 情報 (使用率、クロック数) 等の情報を含む。これらの資源情報のうち、動的に変化する情報 (ネットワーク負荷や各ノードの負荷等) についてはルータ 24 自身が定期的に取得して更新設定する。

【0028】

回答受付部 244 は、ストレージ処理部 25 や他のルータ 24 から返ってきた回答データ (クエリ結果) を一時的に保存する。

【0029】

集約実行部 245 は、クエリに付与されている静的集約実行ポリシーが「集約する」の場合は、回答受付部 244 にあるデータを集約する。また、静的集約実行ポリシーが「集約してもよい」の場合は、データの集約を行う場合と行わない場合で、資源情報に基づくコスト比較を行い、比較結果に基づいてデータの集約を行うか否かを判定する。データの集約

10

20

30

40

50

を行う場合のコストの方が小さい場合には集約すると判定し、回答受付部 2 4 4 に保存されているクエリ結果のデータを集約する。また、データの集約を行わない場合のコストの方が小さい場合には集約しないと判定し、データの集約は行わない。

【 0 0 3 0 】

データ送信部 2 4 6 は、集約実行部 2 4 5 が集約しないと判定した各クエリ結果を直ちにクエリの発行元へ送信する。また、集約すると判定されたクエリ結果については、集約後にクエリの発行元へ送信する。

【 0 0 3 1 】

図 4 は、ストレージ処理部 2 5 の構成を概略的に示す機能ブロック図である。ストレージ処理部 2 5 は、キュー部 2 5 0 と、データ操作部 2 5 1 と、ストレージ装置 2 5 5 を備える。データ操作部 2 5 1 は、クエリ解析部 2 5 2 と、トランザクション実行部 2 5 3 と、内部クエリ発行部 2 5 4 を含む。ストレージ装置 2 5 5 は、複数のストレージを搭載しており、これらストレージを制御するコントローラや入出力ポートを有している。

10

【 0 0 3 2 】

キュー部 2 5 0 は、ルータ 2 4 から順次入力された複数のクエリを一時的に保存する機能を有し、先に入力され保持されたクエリを優先的にデータ操作部 2 5 1 に供給する。

【 0 0 3 3 】

データ操作部 2 5 1 のクエリ解析部 2 5 2 は、キュー部 2 5 0 から供給されたクエリを解析し、実行プランを生成する。トランザクション実行部 2 5 3 は、この実行プランに従ったトランザクションを実行する。

20

【 0 0 3 4 】

トランザクション実行部 2 5 3 は、トランザクション実行のために必要なデータセットがストレージ装置 2 5 5 内の部分データベースに格納されていないとき、内部クエリ発行部 2 5 4 に対して当該データセットのデータ取得要求を発する。そして、内部クエリ発行部 2 5 4 により取得されたデータセットを用いてトランザクションを実行する。

【 0 0 3 5 】

内部クエリ発行部 2 5 4 は、トランザクション実行部 2 5 3 からのデータ取得要求に応じて、内部クエリを生成し、生成した内部クエリを含むデータ転送要求をルータ 2 4 に対して発し、これに対応するデータセットを取得する。

【 0 0 3 6 】

30

管理サーバ 3 0 は、分散データベースを構成する複数の部分データベースと、データサーバ 2 2₁ ~ 2 2_N との対応関係を示す管理テーブル 3 0 を有している。クエリサーバ 2 0 A ~ 2 0 C のいずれかが、受信クエリの解析結果を管理サーバ 3 0 に転送すると、管理サーバ 3 0 は、その解析結果に基づいて管理テーブル 3 0 T を参照してデータサーバ 2 2₁ ~ 2 2_N の中からクエリの供給先を決定し、この結果をそのクエリサーバ 2 0 に通知する。通知を受けたクエリサーバ 2 0 は、この通知内容に従って、データサーバ 2 2₁ ~ 2 2_N のうちの単数又は複数のデータサーバに、変換後のクエリを送信する。

【 0 0 3 7 】

次に本実施形態に係るルータ（送信情報制御装置）2 4 による集約実行判定処理について図 5 のフローチャートを参照して説明する。なお、本処理は、ルータ 2 4 がクエリを受け付けた後から、ストレージ処理部 2 5 或いは他のルータ 2 4 からクエリ結果が返ってくるまでの間の任意のタイミングで開始される。

40

【 0 0 3 8 】

まず、集約実行部 2 4 5 は、受信したクエリに付与されている静的集約実行ポリシーが「集約する」であるかを判別する（ステップ S 1）。

【 0 0 3 9 】

ステップ S 1 において静的集約実行ポリシーが「集約する」である場合には、集約を行うことを決定して本処理を終了する（ステップ S 2）。

【 0 0 4 0 】

また、ステップ S 1 において静的集約実行ポリシーが「集約してもよい」の場合、集約実

50

行部 2 4 5 は、クエリ結果を集約する場合と集約しない場合のコストを計算する（ステップ S 3）。コスト計算に用いられる算出式は、例えば資源情報を用いてコストを計算するものを採用してもよく、式の内容は任意に設定できる。資源情報とは、本システムにおけるネットワークや各ノードの負荷情報（転送速度、転送時間、利用率、負荷値等）を含む。また、コスト計算式は、任意のタイミングで更新できる。

【 0 0 4 1 】

コスト計算式の一例として、（集約する場合のコスト）＝（データが全て揃うまでの時間）＋（集約に要する時間）＋（集約後データの転送時間）＋（集約後データの解凍処理時間）、（集約しない場合のコスト）＝ \quad ＋（データの転送時間）、というような、時間に基づいた式を用いてもよい。

10

【 0 0 4 2 】

（データが全て揃うまでの時間）は見積もり時間でよく、前回のクエリの実測値等を用いて求めてもよい。また、 \quad は定数である。また、（集約に要する時間）は、データの見積量（バイト）と集約処理の速度（バイト／秒）で求めてもよい。データの見積量は前回のクエリ処理時の実測値などを利用して求めても良い。（集約後データの転送時間）は、集約後データの見積量（バイト）と、資源情報テーブル 2 4 3 に保持している回線速度（バイト／秒）から求めてもよい。（集約後データの解凍処理時間）は、集約後データの見積量（バイト）と解凍処理の速度（バイト／秒）から求めてもよい。（データの転送時間）は、集約前データの見積量（バイト）と資源情報テーブル 2 4 3 に保持されている回線速度（バイト／秒）から求めてもよい。

20

【 0 0 4 3 】

なお、集約処理や解凍処理の速度、どの程度データが小さくなるかを示す集約効率は、集約アルゴリズムの静的な性質として予め取得され、ルータ 2 4 の記憶部に記憶されていてもよい。

【 0 0 4 4 】

コスト計算式の他の例として、（集約する場合のコスト）＝ \quad ×（ルータ自身の CPU 利用率）、（集約しない場合のコスト）＝ \quad ×（結果送信先の CPU 利用率）、というような、CPU 資源等の利用率に基づいた式を用いてもよい。 \quad 、 \quad は定数である。

【 0 0 4 5 】

次に、集約実行部 2 4 5 は、クエリ結果を集約する場合のコストと集約しない場合で計算したコストを比較し（ステップ S 4）、集約する場合のコストの方が集約しない場合のコストよりも小さい場合には（ステップ S 4：YES）、集約を行うことを決定して本処理を終了する（ステップ S 2）。また、集約しない場合のコストの方が集約する場合のコストよりも小さい場合には（ステップ S 4：NO）、集約を行わないことを決定して本処理を終了する（ステップ S 5）。

30

【 0 0 4 6 】

次に、分散データベース管理システム 1 0 の動作を示す通信シーケンスの一例を、図 6 を参照して説明する。

【 0 0 4 7 】

まず、クエリサーバ 2 0 A が、ロードバランサ 1 1 を介してクライアント端末 T 1 からクエリを受信する。

40

【 0 0 4 8 】

クエリサーバ 2 0 A のクエリ解析部 2 1 A は、受信したクエリを解析し、この解析結果に基づいて、受信クエリを、ストレージ処理部 2 5 のデータベース構造に基づいて最適化された解析ツリー形式のクエリに変換する（L 1）。次に、クエリ解析部 2 1 A は、クエリの解析結果に基づいて、クエリを送信すべきデータサーバ 2 2 i を決定し（L 2）、データサーバ i にクエリを送信する。

【 0 0 4 9 】

データサーバ 2 2 i では、ルータ 2 4 i のクエリ配送部 2 4 2 が、ルーティングテーブル 2 4 1 を参照し、受信クエリの配送先をストレージ処理部 2 5 a に決定し（L 3）、ス

50

ストレージ処理部 2 5 a に配送する。

【 0 0 5 0 】

ストレージ処理部 2 5 a のデータ操作部 2 5 1 は、クエリを解析し、最適化された実行プランを生成する (L 4)。クエリサーバ 2 0 A のクエリ解析部 2 1 A が、ストレージ処理部 2 5 が管理する部分データベースの構造に合わせてクエリの最適化を既に行っている場合には、データ操作部 2 5 1 は、クエリの最適化を行う必要はない。

【 0 0 5 1 】

その後、ストレージ処理部 2 5 a のトランザクション実行部 2 5 3 が、実行プランに従ってトランザクションを実行していく。処理の途中で、トランザクション実行部 2 5 3 が、トランザクション実行のために必要なデータセットがストレージ装置 2 5 5 内の部分データベースに格納されていないと判定したとする。この場合、トランザクション実行部 2 5 3 は、内部クエリ発行部 2 5 4 に対して、必要なデータセットのデータ取得要求を発する。

10

【 0 0 5 2 】

例えば、トランザクション実行部 2 5 3 が、選択操作 (特定の条件に合致するタプルを抽出し、これら抽出されたタプルから新たなテーブルを生成するためのデータ操作) や結合操作 (ジョイン操作 : 複数のカラムを結合して新たなテーブルを生成するためのデータ操作) を実行しようとしたが、自己のストレージ処理部 2 5 で管理する部分テーブルに、その選択操作や結合操作に必要なタプルやカラムが存在しない場合、これらのデータセットのデータ取得要求を内部クエリ発行部 2 5 4 に発する。

20

【 0 0 5 3 】

このデータ取得要求に応じて、内部クエリ発行部 2 5 4 は、内部クエリを生成し (L 5)、この内部クエリを含むデータ転送要求をルータ 2 4 i に送信する。ここでは、内部クエリ発行部 2 5 4 は、内部クエリ内の静的集約実行ポリシーを「集約する」に設定したとする。

【 0 0 5 4 】

ルータ 2 4 i は、管理サーバ 3 0 やルーティングテーブル 2 4 1 を参照することで、内部クエリの配送先をストレージ処理部 2 5 b、データサーバ 2 2 j、2 2 k と決定し (L 6)、内部クエリを配送する。ここで、ルータ 2 4 i は、静的集約実行ポリシーを「集約する」から「集約してもよい」に変更したとする。

30

【 0 0 5 5 】

ストレージ処理部 2 5 b は、転送された内部クエリを解析して最適化する (L 7)。そして、データ操作を実行し (L 8)、得られたデータセットを内部クエリ結果としてルータ 2 4 i に送信する。

【 0 0 5 6 】

また、データサーバ 2 2 j のルータ 2 4 j は、ルーティングテーブル 2 4 1 を参照して配送先を決定し (L 9)、内部クエリをストレージ処理部 2 5 c に配送する。ストレージ処理部 2 5 c は、内部クエリを解析して最適化する (L 1 0)。そして、データ操作を実行し (L 1 1)、得られたデータセットを内部クエリ結果としてルータ 2 4 j に送信する。この例では、ルータ 2 4 j は、上述した集約実行判定処理を実行し、集約を行わないことを決定したとする。ルータ 2 4 j は、内部クエリ結果を受信すると、それをそのままルータ 2 4 i に送信する。

40

【 0 0 5 7 】

また、データサーバ 2 2 k のルータ 2 4 k は、ルーティングテーブル 2 4 1 を参照して配送先を決定し (L 1 2)、内部クエリをストレージ処理部 2 5 d、ストレージ処理部 2 5 e に配送する。ストレージ処理部 2 5 d、ストレージ処理部 2 5 e は、内部クエリを解析して最適化する (L 1 3、L 1 4)。そして、データ操作を実行し (L 1 5、L 1 6)、得られたデータセットを内部クエリ結果としてルータ 2 4 k に送信する。この例では、ルータ 2 4 k は、上述した集約実行判定処理を実行し、集約を行うことを決定したとする。ルータ 2 4 k は、受信した内部クエリ結果の集約を行い (L 1 7)、それをルータ 2 4

50

i に送信する。集約処理は、ストレージ処理部 25 のデータ構造を考慮した集約処理や、zip 等のデータの可逆圧縮アルゴリズムによる集約処理を含んでもよい。

【0058】

次に、ルータ 24 i は、ストレージ処理部 25 b、データサーバ 22 j、22 k からそれぞれ得られた内部クエリ結果を集約し (L18)、ストレージ処理部 25 a に送信する。

【0059】

ストレージ処理部 25 a のトランザクション処理部 25 3 は、内部クエリ発行部 25 4 により取得されたデータセットを用いてデータ操作を実行し (L19)、その実行結果 (クエリ結果) をルータ 24 i に送信する。ルータ 24 i は、受信したクエリ結果をクエリサーバ 20 A に送信する。クエリサーバ 20 A は、受信したクエリ結果をクライアント端末 T1 に送信する。クライアント端末 T1 は、受信したクエリ結果を表示する (L20)。

【0060】

データの集約処理について、上述した処理シーケンスにおいてルータ 24 k がクエリ結果の集約を行う場面 (図 6 : L17) を例に、図 7 を参照して説明する。図示されるように、この例では、ストレージ処理部 25 が、キーの一覧と、キーの示す値を保持する辞書からなるデータ構造を有するテーブルデータを保持する。図 7 では、ストレージ処理部 25 d が、テーブルデータを、キーの一覧である RTa と、列ごとの辞書である Ca1、Ca2 からなるデータセット DSa として保持している。RTa において、TID は行番号を示す記号であり、Col1Ref は列 1 のキーを示す記号で、Col2Ref は列 2 のキーを示す記号である。Ca1 は列 1 の辞書であり、Ca2 は列 2 の辞書である。これらの辞書は、キー値と実際のデータ値の対応関係を保持している。

【0061】

また、ストレージ処理部 25 e は、データセット DSb を保持している。データサーバ 22 k のストレージ処理部 25 d、25 e が、データセット DSa、DSb をルータ 24 k に送信すると、ルータ 24 k は、これらのデータセット DSa、DSb を集約して、新たなテーブル RTd、Cd1、Cd2 を生成し、これらのテーブルを含むデータセット DSd をルータ 24 i に送信する。

【0062】

なお、図 7 に示されるように、テーブル Ca1 とテーブル Cb1 では、同一の実体データ値 "AA" に対して異なるキー (参照識別子) CRV11、CRV12 が使用されている。また、テーブル Ca2 とテーブル Cb2 では、同一の実体データ値 "AD" に対して異なるキー CRV21、CRV22 が使用されている。この場合、ルータ 24 k は、同一の実体データ値 "AA" に対して一意のキー CRV11 を割り当て、同一の実体データ値 "AD" に対して一意のキー CRV21 を割り当て、新たなテーブル RTd、Cd1、Cd2 を生成する。これにより、参照識別子の不整合を解消することができ、辞書のデータサイズを小さくすることができる。

【0063】

以上のように、本実施の形態によれば、データサーバ 22 のルータ 24 が、クエリ結果データを転送する際に、データの集約を行うか行わないかをコスト計算結果に基づいて動的に決定するため、ネットワークの混雑の回避、データ取得の待機時間の削減、負荷の分散等を効率的に行うことができ、高速なクエリ処理を実現することができる。

【0064】

なお、上記実施形態では、ルータ 24 は、集約実行判定処理においてコスト計算を行っているがこれに限定されず、コスト計算は他の所定のタイミングで予め行って、計算結果を記憶しておき、集約実行判定処理において、この保存された計算結果を参照するようにしてもよい。

【0065】

上述した本発明の実施形態に係るルータ 24 は、データサーバ 22 の CPU (Central

10

20

30

40

50

Processing Unit) が記憶部に格納された動作プログラム等を読み出して実行することにより実現されてもよく、また、ハードウェアで構成されてもよい。上述した実施の形態の一部の機能のみをコンピュータプログラムにより実現することもできる。

【0066】

以上、好ましい実施の形態をあげて本発明を説明したが、本発明は必ずしも上記実施の形態に限定されるものではなく、その技術的思想の範囲内において様々に変形し実施することが出来る。

【0067】

上記の実施形態の一部又は全部は、以下の付記のようにも記載されうるが、以下には限られない。

【0068】

(付記1)

分散データベースを構成する複数のストレージ処理部を備える分散データベースシステムにおいて用いられる送信情報制御装置であって、

データベース処理要求を受信し、該当するストレージ処理部に供給する処理要求供給手段と、

前記ストレージ処理部からの処理結果データを取得する処理結果取得手段と、

前記処理結果データを処理要求元へ送信する際のコストの計算を行うコスト計算手段と

、
前記ストレージ処理部から取得した処理結果データを集約するか否かを、コスト計算結果に基づいて判定する判定手段と、

前記ストレージ処理部から取得した処理結果データを集約すると判定した場合、当該処理結果データを集約して処理要求元に送信し、前記ストレージ処理部から取得した処理結果データを集約しないと判定した場合、当該処理結果データを集約せずに処理要求元に送信する処理結果データ送信手段と、

を備えることを特徴とする送信情報制御装置。

【0069】

(付記2)

前記コスト計算手段は、ネットワーク負荷情報とノードの負荷情報の少なくとも一方を含む、システム資源に関する資源情報に基づいてコスト計算を行う、

ことを特徴とする付記1に記載の送信情報制御装置。

【0070】

(付記3)

前記処理要求には、前記処理結果データの集約を必要とするか否かを示すポリシー情報が付与されており、

前記判定手段は、前記ポリシー情報に、前記処理結果データの集約を必要としない旨が設定されている場合に、前記処理結果データを集約するか否かを前記コスト計算結果に基づいて判定する、

ことを特徴とする付記1又は2に記載の送信情報制御装置。

【0071】

(付記4)

分散データベースを構成する複数のストレージ処理部を備える分散データベースシステムにおける送信情報制御方法であって、

データベース処理要求を受信して、該当するストレージ処理部に供給し、

前記ストレージ処理部からの処理結果データを取得し、

前記処理結果データを処理要求元へ送信する際のコストの計算を行い、

前記ストレージ処理部から取得した処理結果データを集約するか否かを、コスト計算結果に基づいて判定し、

前記ストレージ処理部から取得した処理結果データを集約すると判定した場合、当該処理結果データを集約して処理要求元に送信し、前記ストレージ処理部から取得した処理結

10

20

30

40

50

果データを集約しないと判定した場合、当該処理結果データを集約せずに処理要求元に送信する、

ことを特徴とする送信情報制御方法。

【0072】

(付記5)

ネットワーク負荷情報とノードの負荷情報の少なくとも一方を含む、システム資源に関する資源情報に基づいて前記コストの計算を行う、

ことを特徴とする付記4に記載の送信情報制御方法。

【0073】

(付記6)

前記データベース処理要求は、前記処理結果データの集約を必要とするか否かを示すポリシー情報が処理要求元により付与されており、

前記コスト計算結果に基づく判定は、前記ポリシー情報に、前記処理結果データの集約を必要としない旨が設定されている場合に行う、

ことを特徴とする付記4又は5に記載の送信情報制御方法。

【0074】

(付記7)

コンピュータを、

データベース処理要求を受信し、該当するストレージ処理部に供給する処理要求供給手段、

前記ストレージ処理部からの処理結果データを取得する処理結果取得手段、

前記処理結果データを処理要求元へ送信する際のコストの計算を行うコスト計算手段、

前記ストレージ処理部から取得した処理結果データを集約するか否かを、コスト計算結果に基づいて判定する判定手段、

前記ストレージ処理部から取得した処理結果データを集約すると判定した場合、当該処理結果データを集約して処理要求元に送信し、前記ストレージ処理部から取得した処理結果データを集約しないと判定した場合、当該処理結果データを集約せずに処理要求元に送信する処理結果データ送信手段、

として機能させるためのプログラム。

【0075】

(付記8)

前記コスト計算手段は、ネットワーク負荷情報とノードの負荷情報の少なくとも一方を含む、システム資源に関する資源情報に基づいてコスト計算を行う、

ことを特徴とする付記7に記載のプログラム。

【0076】

(付記9)

前記処理要求には、前記処理結果データの集約を必要とするか否かを示すポリシー情報が付与されており、

前記判定手段は、前記ポリシー情報に、前記処理結果データの集約を必要としない旨が設定されている場合に、前記処理結果データを集約するか否かを前記コスト計算結果に基づいて判定する、

ことを特徴とする付記7又は8に記載のプログラム。

【符号の説明】

【0077】

- 10 分散データベース管理システム
- 11 ロードバランサ
- 20A ~ 20C クエリサーバ
- 21A ~ 21C クエリ解析部
- 22₁ ~ 22_N データサーバ
- 24 ルータ(送信情報制御装置)

10

20

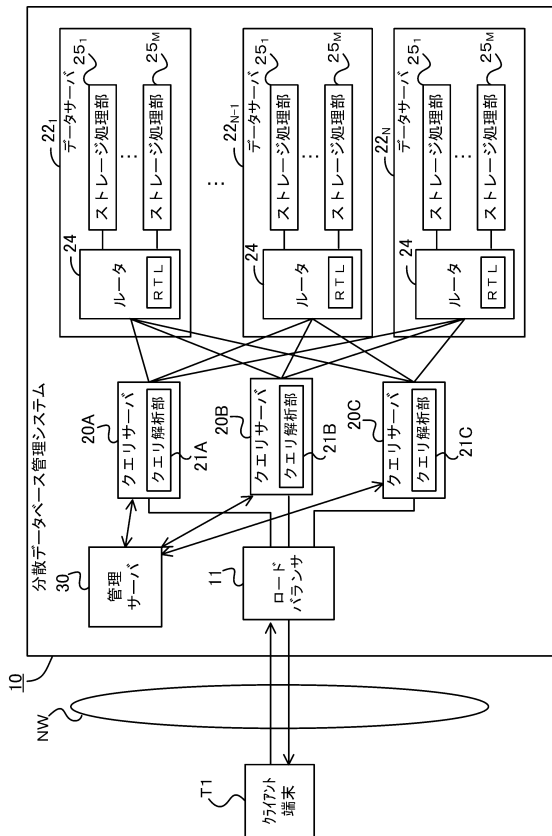
30

40

50

- 2 5₁ ~ 2 5_M ストレージ処理部
- 3 0 管理サーバ
- 2 4 1 ルーティングテーブル
- 2 4 2 クエリ配送部
- 2 4 3 資源情報テーブル
- 2 4 4 回答受付部
- 2 4 5 集約実行部
- 2 4 6 データ送信部
- 2 5 0 キュー部
- 2 5 1 データ操作部
- 2 5 2 クエリ解析部
- 2 5 3 トランザクション実行部
- 2 5 4 内部クエリ発行部
- 2 5 5 ストレージ装置

【図 1】

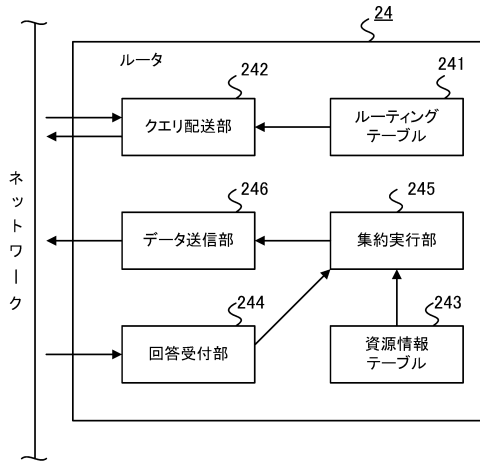


【図 2】

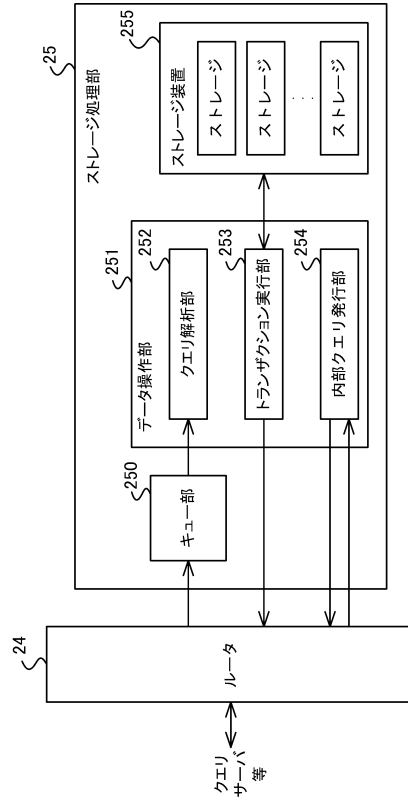
データベース・テーブルTBL

	A ₁	A ₂	A ₃	A _{p-1}	A _p
TG ₁					
TG ₂					
...					
TG _N					

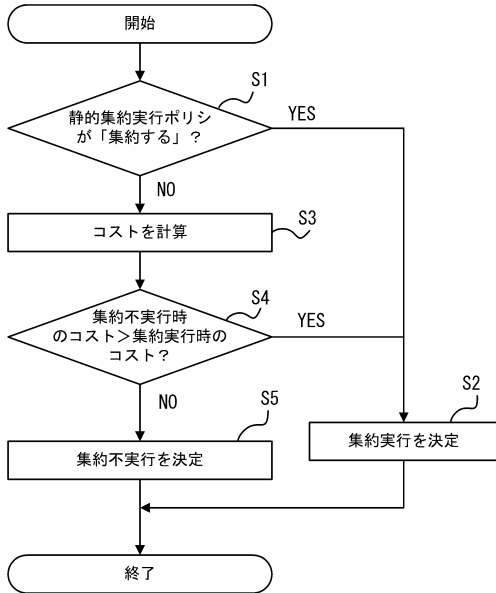
【図3】



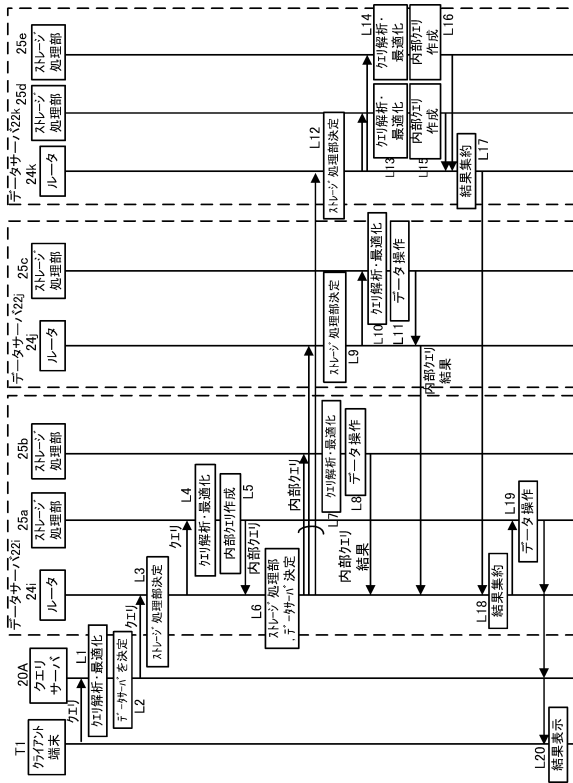
【図4】



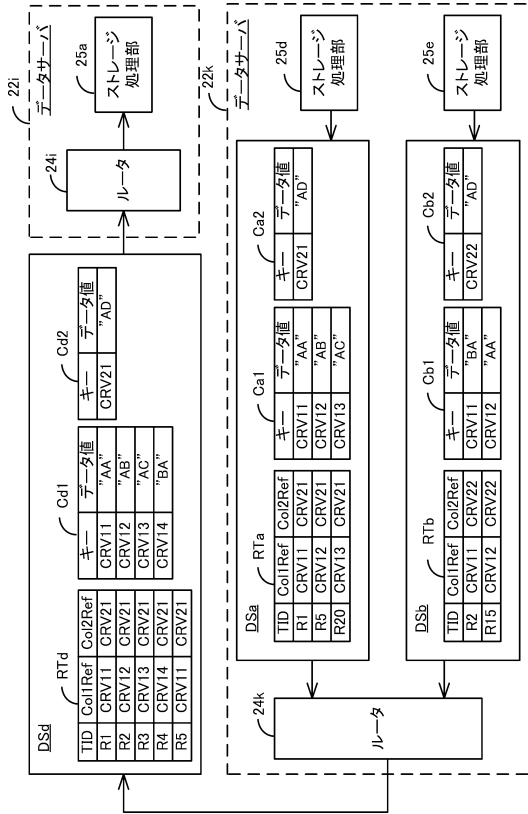
【図5】



【図6】



【 図 7 】



フロントページの続き

(56)参考文献 特開2009-205486(JP,A)
国際公開第2005/001700(WO,A1)
特開平07-141399(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F 12/00
G06F 15/00