

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第5984355号
(P5984355)

(45) 発行日 平成28年9月6日 (2016.9.6)

(24) 登録日 平成28年8月12日 (2016.8.12)

(51) Int.Cl.

F I

G O 6 F 12/00 (2006.01)

G O 6 F 17/30 (2006.01)

G O 6 F 12/00 5 1 3 J

G O 6 F 17/30 1 1 O C

G O 6 F 12/00 5 4 7 Q

請求項の数 30 (全 42 頁)

(21) 出願番号	特願2011-229280 (P2011-229280)	(73) 特許権者	391021710
(22) 出願日	平成23年10月18日 (2011.10.18)		株式会社インテック
(65) 公開番号	特開2012-108889 (P2012-108889A)		富山県富山市牛島新町5番5号
(43) 公開日	平成24年6月7日 (2012.6.7)	(74) 代理人	230104019
審査請求日	平成26年10月16日 (2014.10.16)		弁護士 大野 聖二
(31) 優先権主張番号	特願2010-237102 (P2010-237102)	(74) 代理人	100105038
(32) 優先日	平成22年10月22日 (2010.10.22)		弁理士 田中 久子
(33) 優先権主張国	日本国 (JP)	(74) 代理人	100131451
			弁理士 津田 理
		(72) 発明者	中川 郁夫
			東京都江東区新砂1-3-3 株式会社イ ンテック 先端技術研究所内
		審査官	打出 義尚

最終頁に続く

(54) 【発明の名称】 分散型データベースシステムおよび分散型データ処理システム

(57) 【特許請求の範囲】

【請求項1】

複数のデータベースノードと一つ以上のクライアントノードとを備え、前記複数のデータベースノードに複数のデータが分散してストアされる分散型データベースシステムであって、

前記クライアントノードは、
オブジェクト指向型プログラミングにおけるオブジェクトであるデータであって任意の構造を有することが可能なデータを、キー情報とペアにして、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信し、該データベースノードにストアさせる手段と、

オブジェクト指向型プログラミングにおけるオブジェクトであるプログラムによって前記複数のデータベースノードのいずれかにキー情報とペアでストアされている前記データに対する処理を行うプログラム情報を、該キー情報とペアにして、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信する手段とを備え、

前記プログラム情報は、前記プログラムのオブジェクトが属するクラスを示す情報と、該オブジェクトで用いられる変数であって任意の構造を有することが可能な変数の値の情報とを含むものであり、前記プログラムの定義である前記クラスの情報は、前記クライアントノードと前記データベースノードとで共有可能なものであり、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアでストアされている前記データのオブジェクトに対して、前記クライアントノードから送信された前記プログラム情報により定められる前記プログラムのオブジェクトを実行することにより、前記データベースノードで前記データに対する処理を行う手段と、

前記処理を行った結果の情報を、前記クライアントノードへ送信する手段とを備えることを特徴とする分散型データベースシステム。

【請求項 2】

前記データベースノードは、各クラスで実行されるメソッドの内容と用いられる変数の定義とを含むクラスライブラリに関する情報を記憶する手段をさらに備え、

前記データに対する処理は、前記クラスライブラリに関する情報を用いて前記プログラムのオブジェクトを実行することにより行われるものである、請求項 1 に記載の分散型データベースシステム。

【請求項 3】

前記クライアントノードは、前記クライアントノードから送信された前記プログラム情報により定められる前記プログラムのオブジェクトを前記データベースノードが実行するのに用いられる前記クラスの情報を作成する手段をさらに備える、請求項 1 又は 2 に記載の分散型データベースシステム。

【請求項 4】

前記クライアントノードは、前記プログラム情報を作成する手段をさらに備える、請求項 1 ~ 3 のいずれかに記載の分散型データベースシステム。

【請求項 5】

前記データに対する処理が、該データに新たなデータを追記する処理である場合、

前記クライアントノードから送信されるプログラム情報は、追記する処理の指示と、前記新たなデータとを含み、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データに前記クライアントノードから送信された新たなデータを加えて得られるデータを前記キー情報とペアでストアする処理を、前記データベースノードで行い、

前記追記の成功を示す通知を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものである、請求項 1 ~ 4 のいずれかに記載の分散型データベースシステム。

【請求項 6】

前記データに対する処理が、該データのうち少なくとも一部を新たなデータに書き換える処理である場合、前記クライアントノードから送信されるプログラム情報は、書き換える処理の指示と、書き換えられるべきデータの指示と、新たなデータとを含み、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中に前記クライアントノードから指示された書き換えられるべきデータがあればそれを新たなデータに置き換えて得られるデータを前記キー情報とペアでストアする処理を、前記データベースノードで行い、

前記書き換えの成功を示す通知を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものである、請求項 1 ~ 5 のいずれかに記載の分散型データベースシステム。

【請求項 7】

前記データに対する処理が、該データの一部を削除する処理である場合、

前記クライアントノードから送信されるプログラム情報は、削除する処理の指示と、削除すべきデータの指示とを含み、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中に前記クライアントノードから指示された削除すべきデータがあれば

10

20

30

40

50

それを削って得られるデータを前記キー情報とペアでストアする処理を、前記データベースノードで行い、

前記削除の成功を示す通知を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものである、請求項 1 ～ 6 のいずれかに記載の分散型データベースシステム。

【請求項 8】

前記データに対する処理が、該データに含まれるデータ要素を別のデータと比較する処理である場合、

前記クライアントノードから送信されるプログラム情報は、どのデータ要素と比較する処理であるかの指示と、比較対象となる前記別のデータとを含み、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中の前記クライアントノードから指示されたデータ要素と前記クライアントノードから送信された前記別のデータとを比較する処理を、前記データベースノードで行い、

前記比較により得られた結果を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものである、請求項 1 ～ 7 のいずれかに記載の分散型データベースシステム。

【請求項 9】

前記データに対する処理が、該データに含まれる一つ以上のデータ要素を用いて計算を行う処理である場合、

前記クライアントノードから送信されるプログラム情報は、どのデータ要素を用いて何の計算を行う処理であるかの指示を含み、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中の前記クライアントノードから指示されたデータ要素を用いて前記クライアントノードから指示された計算を行う処理を、前記データベースノードで行い、

前記計算により得られた結果を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものである、請求項 1 ～ 8 のいずれかに記載の分散型データベースシステム。

【請求項 10】

前記クライアントノードは、

複数のデータを対象にして一つの大処理を行うために、キー情報とプログラム情報のペアを複数準備し、各プログラム情報に、前記一つの大処理の前処理となる各データに対する処理の指示を含ませて、前記複数のデータベースノードのうち各キー情報に基づいて特定される一つ以上のデータベースノードへ送信させる手段と、

前記一つ以上のデータベースの各々から、前記各データに対する処理を行った結果の情報を受信し、受信した複数の情報を用いて、前記一つの大処理の結果を生成する手段とをさらに備えるものである、請求項 1 ～ 9 のいずれかに記載の分散型データベースシステム。

【請求項 11】

前記データに対する処理が、該データに含まれるデータ要素を前記クライアントノードが取得するという処理である場合、

前記クライアントノードから送信されるプログラム情報は、どのデータ要素を取得する処理であるかの指示を含み、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中の前記クライアントノードから指示されたデータ要素を抽出する処理を、前記データベースノードで行い、

前記抽出されたデータ要素を、前記処理を行った結果の情報として、前記クライアント

10

20

30

40

50

ノードへ送信するものである、請求項 1 ~ 10 のいずれかに記載の分散型データベースシステム。

【請求項 12】

前記データに対する処理が、該データを前記クライアントノードが取得するという処理である場合、

前記クライアントノードから送信されるプログラム情報は、取得する処理の指示を含み、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出す処理を行い、

読み出した前記データを、前記処理を行った結果の情報として、前記クライアントノードへ送信するものである、請求項 1 ~ 11 のいずれかに記載の分散型データベースシステム。

【請求項 13】

キー情報とデータのペアを、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードにストアさせる手段が、前記キー情報とプログラム情報のペアを、前記クライアントノードから前記データベースノードへ送信することによって実現される場合、

前記クライアントノードから送信されるプログラム情報は、ストアする処理の指示と、ストアすべきデータとを含み、

前記データベースノードは、

前記クライアントノードから送信されたキー情報とペアで、前記クライアントノードから送信されたプログラム情報に含まれるデータをストアする処理を行い、

前記ストアの成功を示す通知を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものである、請求項 1 ~ 12 のいずれかに記載の分散型データベースシステム。

【請求項 14】

前記クライアントノードから送信されたキー情報とデータのペアをストアしたデータベースノードは、前記複数のデータベースノードのうち自身の代替となり得るデータベースノードに、前記キー情報とデータのペアの複製をストアさせる手段をさらに備える、請求項 1 ~ 13 のいずれかに記載の分散型データベースシステム。

【請求項 15】

前記データに対する処理は、前記キー情報とペアでストアされるデータの更新を伴うものであるかデータを読み出すだけであるかが区別可能であり、

データを読み出すだけの処理である場合、前記クライアントノードが、前記キー情報とプログラム情報のペアを、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信する手段は、前記キー情報により指示されるデータベースノード及び当該データベースノードの代替となり得るデータベースノードのうちから選択したデータベースノードへ送信するものである、請求項 14 に記載の分散型データベースシステム。

【請求項 16】

前記複数のデータベースノードの各々は、自身の担当する値の範囲を有し、

前記キー情報に基づくデータベースノードの特定は、前記キー情報から算出されたハッシュ値を用いて、いずれのデータベースノードの担当する値の範囲に入るかを判断することにより行われる、請求項 1 ~ 15 のいずれかに記載の分散型データベースシステム。

【請求項 17】

前記データベースノードには、前記データが、バージョン情報が付されたオブジェクトとしてストアされており、

前記データに対する処理が、前記キー情報とペアでストアされるデータの更新を伴うものである場合、更新されたデータは、前記バージョン情報が更新されたオブジェクトとし

10

20

30

40

50

てストアされるものである、請求項 1 ~ 1 6 のいずれかに記載の分散型データベースシステム。

【請求項 1 8】

前記データベースノードは、あるクライアントノードから受信したプログラム情報に従って行っている処理が、前記データの更新を伴う処理である場合、該処理を行っている間に、別のクライアントノードから前記データの更新を指示するプログラム情報を受信したら、該別のクライアントノードに処理不能を通知するか、もしくは、該処理の終了を待って該別のクライアントノードから受け取った更新を指示するプログラムの処理を開始する手段をさらに備える、請求項 1 ~ 1 7 のいずれかに記載の分散型データベースシステム。

【請求項 1 9】

前記クライアントノードから送信されるプログラム情報は、前記プログラムのオブジェクトにより実行されるべきメソッドを示す情報と、該オブジェクトで用いられる変数に関する情報とを含み、

前記クライアントノードは、メソッド呼び出しに応じて、呼び出されたメソッドを定義するクラスの情報を参照することにより、前記プログラム情報を作成するものである、請求項 1 ~ 1 8 のいずれかに記載の分散型データベースシステム。

【請求項 2 0】

前記クライアントノードは、前記メソッド呼び出しの際に実行されるプロキシオブジェクトに、キーの情報を含ませることにより、前記キー情報と前記プログラム情報のペアを作成するものである、請求項 1 9 に記載の分散型データベースシステム。

【請求項 2 1】

前記分散型データベースシステムは、前記クライアントノードから送信されたプログラム情報に従って前記データベースノードが前記プログラムを実行するのに用いられるクラスの情報を記憶しているストレージ又はサーバをさらに備え、

前記複数のデータベースノードの各々は、前記クラスの情報を前記ストレージ又はサーバから取得する手段を備える、請求項 1 ~ 2 0 のいずれかに記載の分散型データベースシステム。

【請求項 2 2】

前記クライアントノードから送信されるキー情報とプログラム情報のペアは、前記データベースノードにおいて行うべき処理を一意に示す処理識別情報を伴うものであり、

前記データベースノードは、

前記処理識別情報に対応させて前記処理が実行済みであることを示す情報を記録する手段をさらに備え、

前記クライアントノードから受信した前記キー情報と前記プログラム情報のペアが伴う処理識別情報が実行済みとして記録されていない場合に、当該処理を実行させるものである、請求項 1 ~ 2 1 のいずれかに記載の分散型データベースシステム。

【請求項 2 3】

前記クライアントノードから送信されるキー情報とプログラム情報のペアは、前記データベースノードにおいて行うべき処理に係るユーザを特定するためのユーザ情報を伴うものであり、

前記ユーザ情報により特定されるユーザに対し、前記処理に関する通知を行う手段を備える、請求項 1 ~ 2 2 のいずれかに記載の分散型データベースシステム。

【請求項 2 4】

前記クライアントノードは、

前記データとペアにするキー情報を、前記データを保存するためのインデックスに基づいて生成する手段をさらに備え、

前記データベースノードは、

インデックスに対応してデータを保存する機能を有する共有ストレージに対し、前記クライアントノードから送信されたキー情報に基づいて求められるインデックスの識別情報に対応させて、前記クライアントノードから送信されたデータを書き出す手段を備えるこ

10

20

30

40

50

とを特徴とする請求項 1 ~ 2 3 のいずれかに記載の分散型データベースシステム。

【請求項 2 5】

前記クライアントノードは、

検索条件を含む検索要求を前記プログラム情報として、該プログラム情報とペアにするキー情報を、前記共有ストレージに保存されたデータを検索するためのインデックスに基づいて生成する手段をさらに備え、

前記データベースノードは、

前記クライアントノードから送信されたキー情報に基づいて求められるインデックスの識別情報に対応して保存されているデータを、前記共有ストレージから読み込み、前記プログラム情報により定められる前記プログラムのオブジェクトを実行することにより、前記検索条件に応じて検索されたデータを、前記クライアントノードへ送信する手段をさらに備える、請求項 2 4 に記載の分散型データベースシステム。

10

【請求項 2 6】

複数のデータベースノードと一つ以上のクライアントノードとを備え、前記複数のデータベースノードに複数のデータが分散してストアされる分散型データベースシステムで用いられるデータ処理方法であって、

オブジェクト指向型プログラミングにおけるオブジェクトであるデータであって任意の構造を有することが可能なデータが、キー情報とペアで、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ前記クライアントノードから送信され、該データベースノードにストアされており、

20

オブジェクト指向型プログラミングにおけるオブジェクトであるプログラムによって前記複数のデータベースノードのいずれかにキー情報とペアでストアされている前記データに対する処理を行うプログラム情報を、該キー情報とペアにして、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ前記クライアントノードから送信し、

前記プログラム情報は、前記プログラムのオブジェクトが属するクラスを示す情報と、該オブジェクトで用いられる変数であって任意の構造を有することが可能な変数の値の情報とを含むものであり、前記プログラムの定義である前記クラスの情報は、前記クライアントノードと前記データベースノードとで共有可能なものであり、

前記データベースノードにおいて、前記クライアントノードから送信されたキー情報とペアでストアされている前記データのオブジェクトに対して、前記クライアントノードから送信された前記プログラム情報により定められる前記プログラムのオブジェクトを実行することにより、前記データベースノードで前記データに対する処理を行い、

30

前記処理を行った結果の情報を、前記データベースノードから前記クライアントノードへ送信することを特徴とするデータ処理方法。

【請求項 2 7】

複数のデータベースノードに複数のデータが分散してストアされる分散型データベースシステムに接続されるクライアントノードであって、

オブジェクト指向型プログラミングにおけるオブジェクトであるデータであって任意の構造を有することが可能なデータを、キー情報とペアにして、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信し、該データベースノードにストアさせる手段と、

40

オブジェクト指向型プログラミングにおけるオブジェクトであるプログラムによって前記複数のデータベースノードのいずれかにキー情報とペアでストアされている前記データに対する処理を行うプログラム情報を、該キー情報とペアにして、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信する手段と、

前記プログラム情報は、前記プログラムのオブジェクトが属するクラスを示す情報と、該オブジェクトで用いられる変数であって任意の構造を有することが可能な変数の値の情報とを含むものであり、前記プログラムの定義である前記クラスの情報は、前記クライア

50

ントノードと前記データベースノードとで共有可能なものであり、

前記データベースノードにおいて、前記クライアントノードから送信されたキー情報とペアでストアされている前記データのオブジェクトに対して、前記クライアントノードから送信された前記プログラム情報により定められる前記プログラムのオブジェクトを実行することにより、前記データベースノードで前記データに対する処理を行ったときに、前記処理を行った結果の情報を、前記データベースノードから受信する手段と、を備えることを特徴とするクライアントノード。

【請求項 28】

複数のデータベースノードに複数のデータが分散してストアされる分散型データベースシステムを構成するデータベースノードであって、

前記分散型データベースシステムに接続される一つ以上のクライアントノードから、オブジェクト指向型プログラミングにおけるオブジェクトであるデータであって任意の構造を有することが可能なデータが、キー情報とペアで、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信され、

オブジェクト指向型プログラミングにおけるオブジェクトであるプログラムによって前記複数のデータベースノードのいずれかにキー情報とペアでストアされている前記データに対する処理を行うプログラム情報が、該キー情報とペアで、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信され、

前記プログラム情報は、前記プログラムのオブジェクトが属するクラスを示す情報と、該オブジェクトで用いられる変数であって任意の構造を有することが可能な変数の値の情報とを含むものであり、前記プログラムの定義である前記クラスの情報は、前記クライアントノードと前記データベースノードとで共有可能なものであり、

前記データベースノードは、

前記クライアントノードから送信された前記キー情報と前記データのペアを受信してストアする手段と、

前記クライアントノードから送信された前記キー情報と前記プログラム情報のペアを受信する手段と、

前記クライアントノードから送信されたキー情報とペアでストアされている前記データのオブジェクトに対して、前記クライアントノードから送信された前記プログラム情報により定められる前記プログラムのオブジェクトを実行することにより、前記データベースノードで前記データに対する処理を行う手段と、

前記処理を行った結果の情報を、前記クライアントノードへ送信する手段と、を備えることを特徴とするデータベースノード。

【請求項 29】

コンピュータを、複数のデータベースノードに複数のデータが分散してストアされる分散型データベースシステムに接続されるクライアントノードとして機能させるためのプログラムであって、

前記コンピュータを、

オブジェクト指向型プログラミングにおけるオブジェクトであるデータであって任意の構造を有することが可能なデータを、キー情報とペアにして、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信し、該データベースノードにストアさせる手段、

オブジェクト指向型プログラミングにおけるオブジェクトであるプログラムによって前記複数のデータベースノードのいずれかにキー情報とペアでストアされている前記データに対する処理を行うプログラム情報であって、前記プログラムのオブジェクトが属するクラスを示す情報と、該オブジェクトで用いられる変数であって任意の構造を有することが可能な変数の値の情報とを含むプログラム情報を、該キー情報とペアにして、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信する手段、

前記データベースノードにおいて、前記クライアントノードから送信されたキー情報と

10

20

30

40

50

ペアでストアされている前記データのオブジェクトに対して、前記クライアントノードから送信された前記プログラム情報により定められる前記プログラムのオブジェクトを、該プログラムの定義である前記クラスの情報であって前記クライアントノードと共有可能な情報を用いて、実行することにより、前記データベースノードで前記データに対する処理を行ったときに、前記処理を行った結果の情報を、前記データベースノードから受信する手段、

として機能させることを特徴とするプログラム。

【請求項 30】

コンピュータを、複数のデータベースノードに複数のデータが分散してストアされる分散型データベースシステムを構成するデータベースノードとして機能させるプログラムであって、

前記分散型データベースシステムに接続される一つ以上のクライアントノードから、オブジェクト指向型プログラミングにおけるオブジェクトであるデータであって任意の構造を有することが可能なデータが、キー情報とペアで、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信され、

オブジェクト指向型プログラミングにおけるオブジェクトであるプログラムによって前記複数のデータベースノードのいずれかにキー情報とペアでストアされている前記データに対する処理を行うプログラム情報であって、前記プログラムのオブジェクトが属するクラスを示す情報と、該オブジェクトで用いられる変数であって任意の構造を有することが可能な変数の値の情報とを含むプログラム情報が、該キー情報とペアで、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信され、

前記プログラムは、前記コンピュータを、

前記クライアントノードから送信された前記キー情報と前記データのペアを受信してストアする手段、

前記クライアントノードから送信された前記キー情報と前記プログラム情報のペアを受信する手段、

前記クライアントノードから送信されたキー情報とペアでストアされている前記データのオブジェクトに対して、前記クライアントノードから送信された前記プログラム情報により定められる前記プログラムのオブジェクトを、該プログラムの定義である前記クラスの情報であって前記クライアントノードと共有可能な情報を用いて、実行することにより、前記データベースノードで前記データに対する処理を行う手段、

前記処理を行った結果の情報を、前記クライアントノードへ送信する手段、

として機能させることを特徴とするプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、分散型データベースシステムを構成するデータベースノードにストアされているデータを処理する技術に関する。

【背景技術】

【0002】

分散型データベースシステムは、複数のデータベースサーバによって構成され、複数のデータベースサーバに複数のデータが分散してストアされる。従来、分散型データベースシステムとして、KVS（キー・バリュー・ストア）が知られており、KVSを利用した技術には、例えば、ROMA、kumofs、Flareなどがある。KVSでは、各データベースサーバに、キーとバリューのペアがストアされる。バリューは、文字列やバイト列等のデータであり、キーは、バリュー（データ）を特定するための情報である。

【0003】

KVSの長所の一つは、スケーラビリティが高いことにある。したがって、KVSを利

10

20

30

40

50

用してシステムを運用するときに、当初は少ないサーバ台数で低コストで運用しながら、システムが大きくなる（システムで取り扱うデータ量が増える）につれて、後からサーバ台数を増やしてスケールアウトするという運用が可能である。そのため、KVSでは、膨大なデータが多数のデータベースサーバに適切に分散してストアされる（分散効率を高くする）ことができるように運用されることが多い。

【0004】

このような従来のシステム（KVS）においてデータの処理をする場合には、クライアントが、キーを手がかりにして、処理対象のデータ（バリュー）がストアされているデータベースサーバを特定し、特定されたデータベースサーバへクライアントがキーを送ってデータを要求すると、データベースサーバからクライアントにキーに対応してストアされているデータが送られて、クライアントでその取得したデータに対して処理を実行する。その後、クライアントでデータの処理が完了したら、再び、キーを基にデータベースサーバを特定し、特定されたデータベースサーバへ、クライアントが、処理後のデータをキーとともに送り、ストアさせる。そうすると、データベースサーバに、キーとペアでストアされていたデータが、処理後のデータに置き換わる。このようにして、複数のデータベースサーバに分散してストアされているデータの処理が行われる。

10

【先行技術文献】

【非特許文献】

【0005】

【非特許文献1】中田 敦、外12名、「クラウド大全 第2版 サービス詳細から基盤技術まで」、日経BP社、平成22年7月12日

20

【発明の概要】

【発明が解決しようとする課題】

【0006】

従来のシステム（KVS）では、データの処理をする際には、処理のたびに、処理されるべきデータをデータベースサーバからクライアントへ送り、処理が完了したら、その処理後のデータをクライアントからデータベースサーバへ再び送る必要がある。このとき、データベースサーバからクライアントが取得するデータは、取得要求に記入されているキーに対応してストアされているデータ（バリュー）の全体となる（部分的に取得することはできない）ため、一つのキーに対応させて、大容量のデータがストアされている場合には、1回の処理あたりのデータ通信量が大きくなる。このようなデータ通信量の問題を考慮すると、従来のシステムでは、各キーに対応させてストアできるデータ量には、自ずと制限が加えられることになる。

30

【0007】

そこで、従来のシステム（KVS）においては、各キーに対応させてストアするデータ量を小さくする運用が行われていた。しかし、そうすると、一つの処理に必要なデータが複数のキーにまたがってしまい、キーが複数になれば、適度な分散効果により対応するデータがストアされているデータサーバも複数にまたがることになるから、一つの処理を実行するために、複数のデータベースサーバから複数のデータを別々に取得することが必要になる。そのため、一つの処理あたりのデータ通信回数が多くなってしまう。

40

【0008】

また、従来のシステム（KVS）では、データの処理をする毎に、キーに対応してストアされているデータそのものが、データベースサーバからクライアントへ送られるので、秘匿性の高いデータを扱うのには適しておらず、セキュリティ性を担保するためには、暗号化通信などの何らかの措置を講じる必要があった。

【0009】

本発明は、上記の課題に鑑みてなされたもので、例えば、データベースノードにストアできるデータ量に、従来のシステム（KVS）のような制限が課せられずに、一つのキーに対応させて大容量のデータをストアすることができるとともに、一つの処理に必要なデータの全部が一つのデータベースノードに存在する可能性を高くして、通信回数を減らす

50

ことができ、また、秘匿性の高いデータを扱うのにも適した分散型データベースシステ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】

【0010】

本発明の原理に従う一つの例に係る分散型データベースシステムは、複数のデータベースノードと一つ以上のクライアントノードとを備え、前記複数のデータベースノードに複数のデータが分散してストアされる分散型データベースシステムであって、前記クライアントノードは、キー情報とデータのペアを、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信し、該データベースノードにストアさせる手段と、前記複数のデータベースノードのいずれかにキー情報とペアでストアされているデータに対する処理を行うために、該キー情報とプログラム情報のペアを、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信する手段とを備え、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアでストアされているデータに対して、前記クライアントノードから送信されたプログラム情報に従ってプログラムを実行することにより、前記データベースノード内で前記データに対する処理を行う手段と、前記処理を行った結果の情報を、前記クライアントノードへ送信する手段とを備えている。

10

【0011】

上記の構成では、キー情報により特定されるデータベースノードに、そのキー情報とペアをなすデータが送信されてストアされる。つまり、あるデータをストアするデータベースノードは、そのデータとペアをなすキー情報により特定できるようになっている。そして、クライアントノードから、データがストアされたデータベースノード（キー情報により特定されるデータベースノード）に、そのキー情報とペアをなすプログラム情報が送信されると、プログラム情報に従ってプログラムが実行され、データベースノード内でデータに対する処理が行われる。そして、その処理を行った結果の情報が、データベースからクライアントノードへ送信される。このように、データに対する処理がデータベースノード内で行われるので、データの処理が行われるときに、データベースノードからクライアントノードへ、処理されるべきデータが送信されることがない。この場合、データベースノードからクライアントノードへ送信されるのは、結果の情報であり、これは、処理されるデータよりはるかにデータ量が少ない。

20

30

【0012】

そのため、上記の構成では、一つのデータベースノードに大容量のデータがストアされている場合であっても、データに対する処理を1回行ったときにデータベースノードとクライアントノードとの間で通信されるデータ量（1回の処理あたりのデータ通信量）は、プログラム情報と処理を行った結果の情報という小さなデータ量だけで済む。したがって、データベースノードにストアできるデータ量に、従来のシステム（KVS）のような制限が課せられることがない。そして、上記の構成では、一つのデータベースノードに大容量のデータをストアすることができるので、一つの処理に必要なデータが、一つのデータベースノードにまとめて存在する可能性を高くすることができる。したがって、データベースノードがクライアントノードとの間で通信を行う回数（一つの処理あたりのデータ通信回数）が少なく（例えば、1回で）済む。また、上記の構成では、データそのものがデータベースノードからクライアントノードへ送信されることがない（結果の情報だけが送信される）ので、秘匿性の高いデータを扱うのに適している。

40

【0013】

上述した分散型データベースシステムにおいて、前記データベースノードで実行されるプログラムは、オブジェクト指向型プログラミングにおけるオブジェクトであり、前記クライアントノードから送信されるプログラム情報は、前記オブジェクトが属するクラスを示す情報と、該オブジェクトで用いられる変数の値の情報とを含んでもよい。

【0014】

データベースノードで実行されるプログラムが、オブジェクト指向型プログラミングに

50

おけるオブジェクトであり、クライアントノードから送信されるプログラム情報には、オブジェクトが属するクラスを示す情報と、そのオブジェクトで用いられる変数の値の情報が含まれるので、データベースノードでプログラム（オブジェクト）の実行が可能になる。このように、オブジェクト指向型プログラミングを利用することにより、データベースノード内でのデータの処理（プログラムの実行）が実現できる。

【 0 0 1 5 】

上述した分散型データベースシステムにおいて、前記データベースノードは、各クラスで実行されるメソッドの内容と用いられる変数の定義とを含むクラスライブラリに関する情報を記憶する手段をさらに備え、前記データに対する処理は、前記クラスライブラリに関する情報を用いて前記プログラムを実行することにより行われるものであってもよい。

10

【 0 0 1 6 】

データベースノードに記憶されるクラスライブラリに関する情報には、例えば、各クラスについて、クラスで実行されるメソッドの内容と用いられる変数の定義が含まれているので、クライアントノードからプログラム情報（クラスを示す情報と変数の値の情報が含まれる）が送信されると、このクラスライブラリに関する情報を用いて、データベースノードでプログラム（オブジェクト）の実行が可能になる。このように、オブジェクト指向型プログラミングを利用することにより、データベースノード内でのデータの処理（プログラムの実行）が実現できる。

【 0 0 1 7 】

上記の場合には、そのクラス情報を必要とするノードの各々に、クラス情報（オリジナルまたはそのコピー）を登録してもよい。また、別の例として、共有ストレージや共有ウェブサーバを利用してよく、その場合、まず、本システムの利用者（例えば、データベースノードにストアされているデータを処理するサービスを提供する者）が、サービス開始前に、共有ファイルサーバ（または共有ウェブサーバ）にクラス情報を登録し、各データベースノードに対しては、保存場所（またはURL）を指示し、各データベースノードが、サービス開始時に、記憶された保存場所（またはURL）により示される共有ファイルサーバ（または共有ウェブサーバ）からクラス情報を取得するようにしてもよい。

20

【 0 0 1 8 】

上述した分散型データベースシステムにおいて、前記クライアントノードからキー情報とペアで送信され、前記データベースノードにストアされるデータが、オブジェクト指向型プログラミングにおけるオブジェクトで表現され、任意の構造を有することが可能なものであってもよい。

30

【 0 0 1 9 】

データベースノードにストアされるデータが、オブジェクト指向型プログラミングにおけるオブジェクトであるので、オブジェクト指向型プログラミングを利用することにより、データベースノード内でのデータの処理が実現できる。また、データベースノードにストアされるデータが、任意の構造を有することが可能なものであるので、クライアントノードから送信されるプログラム情報が含む変数の値の情報を、変数が任意の構造を有するものとして記述することができ、該プログラム情報に従って実行されるオブジェクトで用いられる変数を、記述された構造に沿って指定することができるようになる。例えば、データの構造が、複数のデータオブジェクト内キーと各キーに対応するバリューである場合、クライアントノードから送信されるプログラム情報によって、処理対象のデータの構造内に存在するあるキー（データオブジェクト内キー）に対応してストアされているバリューに対して、選択的に処理を行うことが可能になる。なお、ここでのデータ構造は、クライアントノードだけでなくデータベースノードでも扱うことができるものであり、データベースノードでプログラムを実行する際に、文字列や数値だけではなく、ある一定以上の構造やサイズを持ったデータに対する処理を一括で行うことが可能になる。

40

【 0 0 2 0 】

上述した分散型データベースシステムにおいて、前記データに対する処理が、該データに新たなデータを追記する処理である場合、前記クライアントノードから送信されるプロ

50

グラム情報は、追記する処理の指示と、前記新たなデータとを含み、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データに前記クライアントノードから送信された新たなデータを加えて得られるデータを前記キー情報とペアでストアする処理を、前記データベースノード内で行い、前記追記の成功を示す通知を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものであってもよい。

【0021】

例えば、あるデータベースノードに「a b c」というデータがストアされており、そのデータに「x y z」を追記する場合には、クライアントノードから、追記する処理の指示と「x y z」のデータを含んだプログラム情報が、「a b c」のデータのキー情報とともに送信される。データベースノードでは、キー情報に基づいて「a b c」のデータを読み出して、それに「x y z」のデータを加える。このようにして得られた新たなデータ「a b c x y z」がキー情報とともにデータベースノードにストアされる。この追記の処理が完了すると、データベースノードからクライアントノードへ通知（追記の成功を示す通知）が送信される。このようにして、データベースノードにストアされたデータに新たなデータを追記する処理を、そのデータベースノード内で実現することができる。

10

【0022】

上述したの分散型データベースシステムにおいて、前記データに対する処理が、該データのうち少なくとも一部を新たなデータに書き換える処理である場合、前記クライアントノードから送信されるプログラム情報は、書き換える処理の指示と、書き換えられるべきデータの指示と、新たなデータとを含み、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中に前記クライアントノードから指示された書き換えられるべきデータがあればそれを新たなデータに置き換えて得られるデータを前記キー情報とペアでストアする処理を、前記データベースノード内で行い、前記書き換えの成功を示す通知を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものであってもよい。

20

【0023】

例えば、あるデータベースノードに「a b c」というデータがストアされており、そのデータの2文字目を「b」から「x」に書き換える場合には、クライアントノードから、書き換える処理の指示と2文字目を示すデータと「x」のデータを含んだプログラム情報が、「a b c」のデータのキー情報とともに送信される。データベースノードでは、キー情報に基づいて「a b c」のデータを読み出して、2文字目のデータを「x」に書き換える。このようにして得られた新たなデータ「a x c」がキー情報とともにデータベースノードにストアされる。この書き換えの処理が完了すると、データベースノードからクライアントノードへ通知（書き換えの成功を示す通知）が送信される。このようにして、データベースノードにストアされたデータのうち少なくとも一部を新たなデータに書き換える処理を、そのデータベースノード内で実現することができる。

30

【0024】

上述した分散型データベースシステムにおいて、前記データに対する処理が、該データの一部を削除する処理である場合、前記クライアントノードから送信されるプログラム情報は、削除する処理の指示と、削除すべきデータの指示とを含み、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中に前記クライアントノードから指示された削除すべきデータがあればそれを削って得られるデータを前記キー情報とペアでストアする処理を、前記データベースノード内で行い、前記削除の成功を示す通知を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものであってもよい。

40

【0025】

例えば、あるデータベースノードに「a b c」というデータがストアされており、そのデータの2文字目の「b」を削除する場合には、クライアントノードから、削除する処理の指示と2文字目を示すデータを含んだプログラム情報が、「a b c」のデータのキー情

50

報とともに送信される。データベースノードでは、キー情報に基づいて「a b c」のデータを読み出して、2文字目の「b」のデータを削除する。このようにして得られた新たなデータ「a c」がキー情報とともにデータベースノードにストアされる。この削除の処理が完了すると、データベースノードからクライアントノードへ通知（削除の成功を示す通知）が送信される。このようにして、データベースノードにストアされたデータの一部を削除する処理を、そのデータベースノード内で実現することができる。

【0026】

上述した分散型データベースシステムにおいて、前記データに対する処理が、該データに含まれるデータ要素を別のデータと比較する処理である場合、前記クライアントノードから送信されるプログラム情報は、どのデータ要素を比較する処理であるかの指示と、比較対象となる前記別のデータとを含み、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中の前記クライアントノードから指示されたデータ要素と前記クライアントノードから送信された前記別のデータとを比較する処理を、前記データベースノード内で行い、前記比較により得られた結果を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものであってもよい。

【0027】

例えば、あるデータベースノードに、データ（キー情報とペアをなすデータ）に含まれるデータ要素として「パスワード」がストアされており、クライアントノードの利用者が入力した「パスワード」との比較を行って、クライアントノードの利用者のユーザ認証を行う場合を考える。この場合には、クライアントノードから、「パスワード」のデータ要素を比較する処理であることの指示と、比較対象となる「パスワード」（クライアントノードの利用者が入力した「パスワード」）を含んだプログラム情報がキー情報とともに送信される。データベースノードでは、キー情報に基づいてデータを読み出し、そのデータ中に含まれる「パスワード」と比較対象となる「パスワード」を比較する。この比較の処理が完了すると、比較によって得られた結果（認証OKまたは認証NG）がデータベースノードからクライアントノードへ送信される。このようにして、データベースノードにストアされたデータに含まれるデータ要素を別のデータと比較する処理を、そのデータベースノード内で実現することができる。

【0028】

上述した分散型データベースシステムにおいて、前記データに対する処理が、該データに含まれる一つ以上のデータ要素を用いて計算を行う処理である場合、前記クライアントノードから送信されるプログラム情報は、どのデータ要素を用いて何の計算を行う処理であるかの指示を含み、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中の前記クライアントノードから指示されたデータ要素を用いて前記クライアントノードから指示された計算を行う処理を、前記データベースノード内で行い、前記計算により得られた結果を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものであってもよい。

【0029】

例えば、あるデータベースノードに、データ（キー情報とペアをなすデータ）に含まれる一つ以上のデータ要素として「過去のテストの点数」が幾つかストアされており、それらのデータ要素を用いた計算として「過去のテストの平均点の算出」を行う場合を考える。この場合には、クライアントノードから、「過去のテストの点数」のデータ要素を用いて「過去のテストの平均点の算出」を行う処理であることの指示を含んだプログラム情報がキー情報とともに送信される。データベースノードでは、キー情報に基づいてデータを読み出し、そのデータ中の「過去のテストの点数」を用いて「過去のテストの平均点の算出」を行う。この計算の処理が完了すると、計算によって得られた結果（算出された平均点）がデータベースノードからクライアントノードへ送信される。このようにして、データベースノードにストアされたデータに含まれる一つ以上のデータ要素を用いて計算を行う処理を、そのデータベース内で実現することができる。

【0030】

上述した分散型データベースシステムにおいて、前記クライアントノードは、複数のデータを対象にして一つの大処理を行うために、キー情報とプログラム情報のペアを複数準備し、各プログラム情報に、前記一つの大処理の前処理となる各データに対する処理の指示を含ませて、前記複数のデータベースノードのうち各キー情報に基づいて特定される一つ以上のデータベースノードへ送信させる手段と、前記一つ以上のデータベースの各々から、前記各データに対する処理を行った結果の情報を受信し、受信した複数の情報を用いて、前記一つの大処理の結果を生成する手段とをさらに備えるものであってもよい。

【0031】

例えば、ある通信事業者のメールサービスの利用者のメールのデータが、複数のデータベースノードに分散してストアされている場合に、全利用者の全メール数の平均を計算したい場合を考える。この場合、複数のデータベースノードに分散してストアされている複数のデータを対象として「全利用者の全メール数の平均の計算」という大処理を行うことになる。この場合には、クライアントノードから、大処理の前処理として「（キー情報で指定される）利用者のメール数の平均の計算」という処理の指示を含んだ各プログラム情報が各キー情報とともに各データベースノードにそれぞれ送信される。それぞれのデータベースノードでは、自身にストアされている各利用者についてメール数の平均の計算の処理が行われ、その処理を行った結果の情報（各利用者のメール数の平均）が、クライアントノードへ送信される。クライアントノードでは、各データベースノードから受信した情報を用いて、全利用者の全メール数の平均の計算を行うことができる。

【0032】

上述した分散型データベースシステムにおいて、前記データに対する処理が、該データに含まれるデータ要素を前記クライアントノードが取得するという処理である場合、前記クライアントノードから送信されるプログラム情報は、どのデータ要素を取得する処理であるかの指示を含み、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出し、該データ中の前記クライアントノードから指示されたデータ要素を抽出する処理を、前記データベースノード内で行い、前記抽出されたデータ要素を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものであってもよい。

【0033】

例えば、あるデータベースノードに、データ（キー情報とペアをなすデータ）として「過去10年分のテストの点数」がストアされており、そのデータに含まれるデータ要素として「最近1年分のテストの点数」を取得する場合には、クライアントノードから、「最近1年分のテストの点数」というデータ要素を取得する処理であることの指示を含んだプログラム情報がキー情報とともに送信される。データベースノードでは、キー情報に基づいて「過去10年分のテストの点数」のデータを読み出し、そのデータ中の「最近1年分のテストの点数」のデータ要素を抽出する。そして、抽出されたデータ要素（最近1年分のテストの点数）が、クライアントノードへ送信される。このようにして、データベースノードにストアされたデータに含まれるデータ要素（データの一部）を、クライアントノードが取得する処理を実現できる。

【0034】

上述した分散型データベースシステムにおいて、前記データに対する処理が、該データを前記クライアントノードが取得するという処理である場合、前記クライアントノードから送信されるプログラム情報は、取得する処理の指示を含み、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアでストアされているデータを読み出す処理を行い、読み出した前記データを、前記処理を行った結果の情報として、前記クライアントノードへ送信するものであってもよい。

【0035】

あるデータベースノードにストアされているデータを、クライアントノードが取得する場合、クライアントノードから、データを取得する処理の指示を含んだプログラム情報が

キー情報とともに送信される。データベースノードでは、キー情報に基づいてデータを読み出して、クライアントノードへ送信する。このようにして、データベースノードにストアされているデータを、クライアントノードが取得する処理を実現できる。

【0036】

上述した分散型データベースシステムにおいて、キー情報とデータのペアを、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードにストアさせる手段が、前記キー情報とプログラム情報のペアを、前記クライアントノードから前記データベースノードへ送信することによって実現される場合、前記クライアントノードから送信されるプログラム情報は、ストアする処理の指示と、ストアすべきデータとを含み、前記データベースノードは、前記クライアントノードから送信されたキー情報とペアで、前記クライアントノードから送信されたプログラム情報に含まれるデータをストアする処理を行い、前記ストアの成功を示す通知を、前記処理を行った結果の情報として、前記クライアントノードへ送信するものであってもよい。

10

【0037】

あるデータベースノードにデータをストアさせる場合、クライアントノードから、データをストアする処理の指示を含んだプログラム情報がキー情報とともに送信される。データベースノードでは、キー情報とペアでデータがストアされる。このストアの処理が完了すると、データベースノードからクライアントノードへ通知（ストアの成功を示す通知）が送信される。このようにして、データベースノードにデータをストアさせる処理を実現できる。

20

【0038】

上述した分散型データベースシステムにおいて、前記クライアントノードから送信されたキー情報とデータのペアをストアしたデータベースノードは、前記複数のデータベースノードのうち自身の代替となり得るデータベースノードに、前記キー情報とデータのペアの複製をストアさせる手段をさらに備えてもよい。

【0039】

キー情報とデータのペアが、あるデータベースノード（キー情報により特定されるデータベースノード）にストアされると、そのキー情報とデータのペアの複製が、他のデータベースノード（代替となり得るデータベースノード）にもストアされる。あるデータベースノードに障害が発生した場合には、他のデータベースノードにストアされている複製を用いるようにデータベースノードの担当範囲を変更することにより、容易にシステムを復旧することができる。

30

【0040】

例えば、あるデータベースノードAにキー情報とデータのペアがストアされており、その代替となり得るデータベースノードBにキー情報とデータのペアの複製がストアされていたとする。データベースノードAに障害が発生した場合、データベースノードBの担当範囲に、データベースノードAの担当範囲を含めるようにする。そうすると、データベースノードBが、データベースノードAの代わりをするようになる。このようにして、データベースノードの担当範囲を変更することにより、システムを復旧することができる。

【0041】

または、あるデータベースノードに障害が発生した場合に、他のデータベースノードにストアされている複製を用いてデータベースノードの置き換え（追加・削除）をすることにより、容易にシステムを復旧することもできる。

40

【0042】

例えば、あるデータベースノードAにキー情報とデータのペアがストアされており、その代替となり得るデータベースノードBにキー情報とデータのペアの複製がストアされていたとする。データベースノードAに障害が発生した場合、新しいデータベースノードXが追加される。データベースノードXには、データベースノードBからキー情報とデータのペアの複製がコピーされてストアされる。このデータベースノードXは、データベースノードAと同じキー情報によって特定されるようにされる。このようにして、新しいデー

50

データベースノードXを古いデータベースノードAに置き換えて、システムを復旧してもよい。

【0043】

上述した分散型データベースシステムにおいて、前記データに対する処理は、前記キー情報とペアでストアされるデータの更新を伴うものであるかデータを読み出すだけであるかが区別可能であり、データを読み出すだけの処理である場合、前記クライアントノードが、前記キー情報とプログラム情報のペアを、前記複数のデータベースノードのうち前記キー情報に基づいて特定されるデータベースノードへ送信する手段は、前記キー情報により指示されるデータベースノード及び当該データベースノードの代替となり得るデータベースノードのうちから選択したデータベースノードへ送信するものであってもよい。

10

【0044】

データを読み出すだけの処理は、キー情報により特定されるデータベースノードとその代替となり得るデータベースのうちから選択したデータベースノードにより実行可能である。すなわち、読み出しのみの処理を分散することが可能である。

【0045】

例えば、あるデータベースノードAにキー情報とデータのペアがストアされており、その代替となり得るデータベースノードBにキー情報とデータのペアの複製がストアされていた場合、データベースノードAとBから選択されたデータベースノード（データベースノードAとBのいずれか一方または両方）から、データを読み出す処理を実行することができる。したがって、状況に応じて適切なデータベースノードからデータを読み出すことができる。例えば、データを読み出す処理を実行しようとしたときにデータベースノードAの負荷が高かった場合には、代わりにデータベースノードBからデータを読み出すことができる。別の例として、データベースノードAとBの負荷が低い場合には、両方へデータ読み出し要求を送って応答の早い方を使うことにより、データを読み出す処理の高速化が可能である。

20

【0046】

上述した分散型データベースシステムにおいて、前記複数のデータベースノードの各々は、自身の担当する値の範囲を有し、前記キー情報に基づくデータベースノードの特定は、前記キー情報から算出されたハッシュ値を用いて、いずれのデータベースノードの担当する値の範囲に入るかを判断することにより行われてもよい。

30

【0047】

キー情報に基づくデータベースノードの特定は、キー情報から算出されたハッシュ値を用いて行われる。この場合、キー情報から算出されたハッシュ値の一部（例えば、ハッシュ値の先頭の二文字など）を用いて、データベースノードを特定してもよい。例えば、キー情報から算出されたハッシュ値が「39021650ae9e43ecb53f66bf865d9730135f43b1」であったとする。このハッシュ値の先頭の二文字は「39」である。一方、各データベースノードには、それぞれストアされるべきデータの範囲（担当範囲）を設定しておく。例えば、ノードAの担当範囲が「08～34」、ノードBの担当範囲が「35～48」・・・に設定されていたとする。このような場合、キー情報のハッシュ値の一部「39」からノードBを特定することができる。

40

【0048】

上述した分散型データベースシステムにおいて、前記データベースノードには、前記データが、バージョン情報が付されたオブジェクトとしてストアされており、前記データに対する処理が、前記キー情報とペアでストアされるデータの更新を伴うものである場合、更新されたデータは、前記バージョン情報が更新されたオブジェクトとしてストアされるものであってもよい。

【0049】

これにより、データベースノードにストアされるデータのバージョン管理が可能になる。例えば、あるデータベースノードAにキー情報とデータのペアがストアされており、その代替となり得るデータベースノードBにキー情報とデータのペアの複製がストアされて

50

いたとする。データベースノードAのデータが更新され、その更新がまだデータベースノードBに反映されない間、データベースノードAのデータのバージョン情報は更新されるが、データベースノードBのデータのバージョンは古いままである。したがって、このとき、例えば、データベースノードAとBの両方からデータを読み出す処理が実行されると、データベースノードBのデータが古いバージョンのデータ（更新前のデータ）であることが分かる。

【0050】

また別の例として、データの冗長化のために、あるキー情報とデータのペアが、二つのデータベースノード（DBノードAとDBノードB）に保存されており、一方のデータベースノード（DBノードA）が障害により一時的にダウンし、その後に、そのデータベースノード（DBノードA）が復旧してサービスを提供できるようになったとする。この場合、該当するキー情報とデータのペアを、再度、もう一方のデータベースノード（DBノードB）からコピーする必要があるが、その際、二つのデータベースノード（DBノードAとDBノードB）の間で、該当のキー情報に対するバージョン番号を比較し、更新が必要な場合にだけ、データの転送（コピー）を行うようにすることができる。このようにして、データベースノード間でのデータ交換の効率化を図ることが可能になる。

【0051】

上述した分散型データベースシステムにおいて、前記データベースノードは、あるクライアントノードから受信したプログラム情報に従って行っている処理が、前記データの更新を伴う処理である場合、該処理を行っている間に、別のクライアントノードから前記データの更新を指示するプログラム情報を受信したら、該別のクライアントノードに処理不能を通知するか、もしくは、該処理の終了を待って該別のクライアントノードから受け取った更新を指示するプログラムの処理を開始する手段をさらに備えてもよい。

【0052】

例えば、データベースノードが、あるクライアントノードAから受信したプログラム情報に従ってデータの更新を実行しているときに、別のクライアントノードBからデータの更新を指示するプログラム情報を受信した場合には、その別のクライアントノードBには、処理不能が通知される、あるいは、クライアントノードAからのデータ更新の処理の終了を待って、クライアントノードBからのデータ更新の処理を開始する。このようにして、複数のクライアントノードによるデータ更新のコンフリクトを制御することができる。

【0053】

上述した分散型データベースシステムにおいて、前記クライアントノードは、前記プログラム情報を作成する手段をさらに備えてもよい。

【0054】

クライアントノードの利用者は、自分でオリジナルのプログラムを作成でき、そのプログラムを用いて、データベースノードにストアされているデータを処理することができる。したがって、様々なデータ処理を実行させることが可能になる。

【0055】

上述した分散型データベースシステムにおいて、前記データベースノードで実行されるプログラムは、オブジェクト指向型プログラミングにおけるオブジェクトであり、前記クライアントノードは、前記クライアントノードから送信されたプログラム情報に従って前記データベースノードが前記プログラムを実行するのに用いられるクラスの情報を作成する手段をさらに備えてもよい。

【0056】

データベースノードで実行されるプログラムが、オブジェクト指向型プログラミングにおけるオブジェクトである場合に、クライアントノードの利用者は、自分でオリジナルのクラスを作成でき、そのクラスで定義されたオブジェクトを用いて、データベースノードにストアされているデータを処理することができる。したがって、様々なデータ処理を実行させることが可能になる。

【0057】

10

20

30

40

50

データベースノードで実行されるプログラムが、オブジェクト指向型プログラミングにおけるオブジェクトである場合に、上記の分散型データベースシステムにおいて、前記クライアントノードから送信されるプログラム情報は、前記オブジェクトにより実行されるべきメソッドを示す情報と、該オブジェクトで用いられる変数に関する情報とを含み、前記クライアントノードは、メソッド呼び出しに応じて、呼び出されたメソッドを定義するクラスの情報を参照することにより、前記プログラム情報を作成するようにしてもよい。

【0058】

これにより、クライアントノードの利用者は、ネットワークの向こう側に存在するデータベースノードにおいて処理が実行されるにもかかわらず、あたかも手元のクライアントノードにおいて処理が実行されるかのような感覚で、メソッド呼び出しを行うプログラムを開発することが可能になる。

10

【0059】

上記の構成において、前記クライアントノードは、前記メソッド呼び出しの際に実行されるプロキシオブジェクト（擬似的なオブジェクトであり、例えばJava（登録商標）言語ではProxy技術で利用可能）に、キーの情報を含ませることにより、前記キー情報と前記プログラム情報のペアを作成するようにしてもよい。

【0060】

これにより、プロキシオブジェクトが、クライアントノードにおけるメソッド呼び出しに基づいて、データベースノードにおいて実行されるべきメソッドを示す情報を生成するとともに、これをキー情報とペアにしてデータベースノードへ送信することが可能になり、クライアントノードの利用者は、キー情報とプログラム情報のペアを渡す先が、ネットワークの向こう側に存在するデータベースノードであることを意識しなくても、プログラムを開発することが可能になる。

20

【0061】

上述した分散型データベースシステムにおいて、前記データベースノードで実行されるプログラムは、オブジェクト指向型プログラミングにおけるオブジェクトであり、前記分散型データベースシステムは、前記クライアントノードから送信されたプログラム情報に従って前記データベースノードが前記プログラムを実行するのに用いられるクラスの情報を記憶しているストレージ又はサーバをさらに備え、前記複数のデータベースノードの各々は、前記クラスの情報を前記ストレージ又はサーバから取得する手段を備えてもよい。

30

【0062】

データベースノードで実行されるプログラムが、オブジェクト指向型プログラミングにおけるオブジェクトである場合に、クラスの情報は、ストレージ又はサーバに記憶されており、各データベースノードは、これらのストレージ又はサーバから、クラスの情報を取得して利用することができる。このようにして、クラスの情報を、複数のデータベースノードで共有することができる。

【0063】

上述した分散型データベースシステムにおいて、前記クライアントノードから送信されるキー情報とプログラム情報のペアを、前記データベースノードにおいて行うべき処理を一意に示す処理識別情報を伴うものとし、前記データベースノードは、前記処理識別情報に対応させて前記処理が実行済みであることを示す情報を記録する手段をさらに備え、前記クライアントノードから受信した前記キー情報と前記プログラム情報のペアが伴う処理識別情報が実行済みとして記録されていない場合に、当該処理を実行させるようにしてもよい。

40

【0064】

これにより、例えば、クライアントノードからデータベースノードへ依頼した処理を行った結果の情報が、何らかの事情でクライアントノードに届かず、クライアントノードが再依頼する場合に、以前の依頼をデータベースノードが受信しておらず処理を行っていなかったのか、それとも処理後のデータベースノードからクライアントノードへの通信に障害があったのかを区別して、重複処理をすることなくリトライを扱うことが可能になる。

50

【 0 0 6 5 】

上述した分散型データベースシステムにおいて、前記クライアントノードから送信されるキー情報とプログラム情報のペアを、前記データベースノードにおいて行うべき処理に係るユーザを特定するためのユーザ情報を伴うものとし、前記ユーザ情報により特定されるユーザに対し、前記処理に関する通知を行う手段を備えるようにしてもよい。

【 0 0 6 6 】

これにより、例えば、データベースノードあるいは複数のデータベースノードを管理するサーバにおいて、ユーザに対する認証を行ってから依頼された処理に進むようにしたり、どのユーザがどのような処理を依頼したかを把握して、課金、ログの管理、レポート等のサービスを提供したりすることが可能になる。

10

【 0 0 6 7 】

上述した分散型データベースシステムの発明は、システム全体の方法の発明としても、汎用のコンピュータシステムを本システムとして動作させるためのプログラム（又はそのプログラムを記録した記録媒体）の発明としても、クライアントノード装置の発明としても、データベースノード装置の発明としても、汎用のコンピュータを本システムのクライアントノード又はデータベースノードとして動作させるためのプログラム（又はそのプログラムを記録した記録媒体）の発明としても、クライアントノードにおいて実行される方法の発明としても、データベースノードにおいて実行される方法の発明としても、勿論成立するものである。

【 0 0 6 8 】

20

また、本発明の原理は、上述した分散型データベースシステムのように、複数のデータベースノードの各々に、自身の担当するキー情報とペアをなすデータが分散してストアされる構成に適用することができるほか、共有ストレージをバックに有する複数のデータ処理ノードが、共有ストレージに保存されるデータの書き込みや読み出しを分散して行う構成に適用することもできる。

【 0 0 6 9 】

そのような本発明の原理に従う一つの例に係る分散型データ処理システムは、インデックス（データを一意に特定するための情報）に対応してデータを保存する機能を有する共有ストレージに対し分散して処理を行う分散型データ処理システムであって、前記分散型データ処理システムは、複数のデータ処理ノードと一つ以上のクライアントノードとを備え、前記クライアントノードは、データを保存するためのインデックスに基づいてキー情報を生成する手段と、前記キー情報と前記データのペアを、前記複数のデータ処理ノードのうち前記キー情報に基づいて特定されるデータ処理ノードへ、送信する手段とを備え、前記データ処理ノードは、前記共有ストレージに対し、前記クライアントノードから送信されたキー情報に基づいて求められるインデックスの識別情報に対応させて、前記クライアントノードから送信されたデータを書き出す手段を備える。

30

【 0 0 7 0 】

上記の構成によれば、共有ストレージにおいてインデックスに対応してデータが保存されるところ、そのインデックスに基づいてキー情報が生成され、このキー情報により特定されるデータ処理ノードが、共有ストレージへデータを書き込む処理を行うことになる。したがって、各々のインデックスのデータ領域に対するアクセスは、そのインデックスを担当するデータ処理ノードから生じるようになり、それ以外のデータ処理ノードからは生じないから、アクセスの競合が生じにくく、処理分散の効率を向上することが可能になる。

40

【 0 0 7 1 】

上記の分散型データ処理システムにおいて、前記クライアントノードは、前記共有ストレージに保存されたデータを検索するためのインデックスに基づいてキー情報を生成する手段と、前記キー情報と検索要求のペアを、前記複数のデータ処理ノードのうち前記キー情報に基づいて特定されるデータ処理ノードへ、送信する手段とをさらに備え、前記データ処理ノードは、前記クライアントノードから送信されたキー情報に基づいて求められる

50

インデックスの識別情報に対応して保存されているデータを、前記共有ストレージから読み出し、前記クライアントノードへ送信する手段をさらに備えるようにしてもよい。

【0072】

これにより、共有ストレージからデータを読み出す処理についても、各々のインデックスのデータ領域に対するアクセスが、そのインデックスを担当するデータ処理ノードから生じるようになり、アクセスの競合を生じにくくして、処理分散の効率をさらに向上することが可能になる。また、キャッシュ（読み込んだデータを一時的に高速なメモリに保存しておくこと）によるアクセス性能の向上が可能になる。

【0073】

上記の構成において、前記クライアントノードから送信される検索要求を、検索条件を含むプログラム情報とし、前記データ処理ノードは、前記プログラム情報に従ってプログラムを実行することにより、前記共有ストレージから読み出すべきデータを検索する処理を行うようにしてもよい。

【0074】

上記の構成では、例えば、大量のログを長時間にわたって継続的に記録するような場合でも、各クライアントノードにおけるログの受付及び送信（分散化）と、各データ処理ノードにおけるインデックス処理及び共有ストレージへのアクセスとが、パイプラインで行え、クライアントノードやデータ処理ノードの台数を増やすことにより性能を上げるスケールアウトの実現が可能になる。さらに、そのようなログ記録からの検索を様々な検索条件に応じて行うことも可能になる。

【0075】

上述した分散型データ処理システムの発明は、システム全体の方法の発明としても、汎用のコンピュータシステムを本システムとして動作させるためのプログラム（又はそのプログラムを記録した記録媒体）の発明としても、クライアントノード装置の発明としても、データ処理ノード装置の発明としても、汎用のコンピュータを本システムのクライアントノード又はデータ処理ノードとして動作させるためのプログラム（又はそのプログラムを記録した記録媒体）の発明としても、クライアントノードにおいて実行される方法の発明としても、データ処理ノードにおいて実行される方法の発明としても、勿論成立するものである。

【0076】

また、上述した分散型データベースシステムのデータベースノードにおいて採用可能な各種構成は、上記の分散型データ処理システムのデータ処理ノードにおいても適宜採用できるものである。

【発明の効果】

【0077】

本発明によれば、例えば、一つのキーに対応させて大容量のデータをストアすることができるとともに、一つの処理に必要なデータの全部が一つのデータベースノードに存在する可能性を高くして、通信回数を減らすことができる。また、本発明は、秘匿性の高いデータを扱うのにも適している。

【図面の簡単な説明】

【0078】

【図1】本発明の実施の形態における分散型データベースシステムの一例（本システム）を説明するための図

【図2】本システムのクライアントノードの構成の一例を示すブロック図

【図3】本システムのデータベースノードの構成の一例を示すブロック図

【図4】本システムの概要（データの処理と分散）を説明するための図

【図5】本システムにおけるクラス情報の共有の一例を説明するための図

【図6】本システムの特徴（プログラミング環境）を説明するための図

【図7】本システムの特徴（データベース技術）を説明するための図

【図8】本システムの適用例（個人情報扱うシステム）を説明する図

10

20

30

40

50

【図 9】本システムの適用例（オンラインショップの購買履歴を扱うシステム）を説明する図

【図 10】本システムの適用例（商品在庫管理システム）を説明する図

【図 11】本システムの適用例（健康管理システム）を説明する図

【図 12】本システムの適用例（通信事業者のメールサービス）を説明する図

【図 13】本システムの適用例（通信事業者のメールサービス）を説明する図

【図 14】本システムの適用例（オンラインショップのサービスの一部）の詳細を説明する図

【図 15】図 14 の適用例と同じサービスを別の仕組みで実現した適用例を説明するための図

10

【図 16】図 15 の仕組みを採用したシステムの一例を説明するための図

【図 17】図 15 の仕組みの利点の一例を説明するための図

【図 18】本システムの利点の一例を説明するための図

【図 19】本システムの利点の一例を説明するための図

【図 20】本発明の実施の形態における分散型データ処理システムの一例を説明するための図

【図 21】図 20 のシステムの適用例（通信ログの処理）を説明する図

【発明を実施するための形態】

【0079】

以下、本発明の実施の形態に係る分散型データベースシステムについて、例示のために、図面を用いて説明する。本実施の形態では、データストレージサービス等に用いられる分散型データベースシステムの場合を例示する。

20

【0080】

本発明の実施の形態における分散型データベースシステムの一例（本システム）の構成を、図面を参照して説明する。図 1 は、本システムの構成を示す説明図である。図 1 に示すように、本システム 1 は、複数のデータが分散してストアされる複数のデータベースノード 2 と、本システム 1 の利用者が用いるクライアントノード 3 で構成されている。本システム 1 によって提供されるサービスは、複数のエンドユーザが用いるユーザ端末 4 から利用できるようになっている。なお、図 1 では、説明の便宜のため、データベースノード 2 やクライアントノード 3 やユーザ端末 4 の台数が少なく例示されているが、データベースノード 2 やクライアントノード 3 やユーザ端末 4 の台数がこれに限定されないのは勿論である。

30

【0081】

図 2 は、本システムのクライアントノード 3 の構成を示すブロック図である。図 2 に示すように、クライアントノード 3 は、データ送受信部 5 と、データ処理部 6 と、プログラム作成部 7 を備えている。データ送受信部 5 は、データベースノード 2 との間でデータ通信を行う機能を備えている。データ処理部 6 は、データベースノード 2 と間で通信されたデータを処理する機能を備えている。プログラム作成部 7 は、本システム 1 で利用される種々のプログラムを作成する機能を備えている。

【0082】

40

図 3 は、本システムのデータベースノード 2 の構成を示すブロック図である。図 3 に示すように、データベースノード 2 は、データ送受信部 8 と、データ処理部 9 と、データ格納部 10 と、クラスライブラリ格納部 11 と、コンフリクト制御部 12 と、バージョン管理部 13 を備えている。データ送受信部 8 は、クライアントノード 3 との間でデータ通信を行う機能を備えている。データ処理部 9 は、クライアントノード 3 と間で通信されたデータを処理する機能を備えている。また、データ送受信部 8 は、他のデータベースノード 2 との間でデータ通信を行う機能も備えており、データ処理部 9 は、他のデータベースノード 2 と間で通信されたデータを処理する機能も備えている。

【0083】

データ格納部 10 には、クライアントノード 3 から受信したデータ（キー情報とデータ

50

のペアなど)が格納される。また、クラスライブラリ格納部11には、本システムで利用するプログラムのクラス(オブジェクト指向型プログラミングにおけるオブジェクトが属するクラス)の情報が格納される。例えば、各クラスで実行されるメソッドの内容や、各クラスで用いられる変数の定義などが格納される。コンフリクト制御部12は、データ格納部に格納されているデータを更新するときのコンフリクトを制御する機能を備えている。バージョン管理部13は、データ格納部に格納されているデータのバージョンを管理する機能を備えている。

【0084】

コンフリクト制御は、複数のクライアントノードが、同一のキーに対して同時にデータ更新の処理を要求してきた場合に実行される。例えば、あるクライアントノードAからデータの更新の要求を処理している間に、別のクライアントノードBからデータの更新の要求があった場合には、その別のクライアントノードBには、処理不能が通知されるか、あるいは、クライアントノードAからのデータ更新の処理の終了を待って、クライアントノードBからのデータ更新の処理が開始される。

【0085】

次に、本システムの概要について説明する。本システムは、スケール型の大規模分散データベースシステムであり、データベースノードの台数は、例えば数百台以上である。図4は、本システムの概要を示した説明図である。図4では、説明の便宜上、8台のデータベースノード(DBノードA~H)のみが例示されている。図4に示すように、各データベースノードには、それぞれID(例えば、16進の数字)が割り当てられている。これらのデータベースノードは、割り当てられたIDの番号順に「リング状」になるように、仮想的に並べられている。つまり、最初のデータベースノード(例えば、DBノードA)からIDの番号順に一巡すると、最後のデータベースノード(DBノードH)の次は、再び最初のデータベースノード(DBノードA)に戻るようになっている。

【0086】

また、この割り当てられるIDの大小によって「向き」が決められており、あるデータベースノード(例えば、DBノードA)よりIDの番号が次に大きいデータベースノード(DBノードB)は、そのデータベースノード(DBノードA)の「successor」と呼ばれる。また、あるデータベースノード(例えば、DBノードB)よりIDの番号が次に小さいデータベースノード(DBノードA)は、そのデータベースノード(DBノードB)の「predecessor」と呼ばれる。

【0087】

本システムでは、データベースノードの配置、範囲の割り当て、バックアップなどの仕組みに、DHT(Distributed Hash Table)やConsistent Hashと呼ばれる技術と同様の技術を用いている。図4の例で説明すると、例えば、あるデータベースノード(例えば、DBノードB)の担当範囲は、そのデータベースノードのID(DBノードBのID:35)から次のデータベースノードのID(DBノードCのID:49)までに設定される。そして、クライアントノードから送られたキーとデータのペアは、そのキーのハッシュ値を計算したときの最初の2文字が担当範囲に含まれるデータベースノードにストアされる。例えば、キー情報から算出されたハッシュ値が「39021650ae9e43ecb53f66bf865d9730135f43b1」であると、このハッシュ値の先頭の二文字は「39」であり、DBノードBにストアされることになる。

【0088】

また、本システムでは、データの多重度設定により、データの「複製」の数を指定することができ、同時に複数のデータベースノードに障害が発生しても、それに対応することができる。より具体的に説明すると、本システムでは、あるIDを持つデータベースノード(例えば、DBノードB)が担当すべき情報を、そのデータベースノード(DBノードB)が持つほかに、その情報のコピー(複製)を、そのデータベースノードのpredecessorとなるデータベースノード(DBノードA)が持っている。このように二つのデータベースノードが情報を持つ状態は、多重度が「2」に設定されているともいえる。多重度は

、全データで共通である。多重度を大きくするには、predecessorの向きに、コピーを持つデータベースノードの数を増やせばよい。なお、データのコピーは、データベースノード間で行われる。つまり、あるデータベースノード（例えばDBノードB）でデータの保存、更新、削除が行われた場合、その情報は、そのデータベースノードのpredecessorとなるデータベースノード（DBノードA）にも伝えられ、predecessorとなるデータベースノード（DBノードA）でも、同様に、データの保存、更新、削除が行われる。ただし、データは同じものを複製するか、または差分の変換によってデータの複製を更新するだけであるため、複製されたデータのキーは変わらない。したがって、predecessorとなるデータベースノード（DBノードA）が持っている「複製」は、本来、そのデータベースノード（DBノードA）の担当範囲外のデータである。そのため、そのデータベースノード（DBノードA）のデータには、そのキーに対応するプログラム情報が、クライアントノードから直接届くことがない。そのキーに対応するプログラム情報（job）は、本来のデータベースノード（DBノードB）のみに届き、そのデータベースノード（DBノードA）のデータの更新は、本来のデータベースノード（DBノードB）から伝えられる情報のみによって行われる。

10

【0089】

なお、この「複製」は、データの並列読み出しに利用することができる。そのため、本システムでは、プログラム情報（job）に「READ（読み込みのみ）」か「EDIT（更新を伴う）」かの区別ができるようにフラグが付けられている。「READ」のフラグが付けられた処理（読み込みのみの処理）については、多重度の数までpredecessorの方向に処理を分散することも可能である。なお、書き込みの処理は「一貫性」を保つため、本来の担当範囲を受け持っているデータベースノードでしか行われない。

20

【0090】

また、本システムでは、動的なデータベースノードの追加・削除が容易である。データベースノードを追加するときには、新しいIDを持ったデータベースノード（例えばDBノードI）を準備して、そのデータベースノードが担当しなければならない範囲のデータを、現在その範囲のデータを担当しているデータベースノード（例えばDBノードH）からコピーし、その後、新しいデータベースノード（DBノードI）が「リング」に追加される。この場合、データの複製を終えてから、データベースノードを追加することにより、ダウンタイムを最小限にできる。

30

【0091】

一方、データベースノードを削除するときには、削除対象のデータベースノード（例えばDBノードH）を削除する。それ以降は、そのデータベースノード（DBノードH）が担当していた範囲のデータに対するプログラム情報（job）は、そのデータベースノードのpredecessorとなるデータベースノード（DBノードG）に届くことになる。このデータベースノード（DBノードG）は、削除対象のデータベースノード（DBノードH）が担当していた範囲のデータの複製を持っているので、サービスの提供が継続される。この場合、predecessorとなるデータベースノードがデータの複製を持っていることを前提に、データベースノードを削除することにより、ダウンタイムを最小限にできる。

【0092】

40

また、本システムは、KVSと同等もしくはそれ以上のスケールと拡張性を有しており、さらに「分散プログラミング」のインターフェースもサポートしている。本システムでは、データベースノードにデータが送られて、データの書き込みなどが行われるだけでなく、データベースノードにプログラム情報が送られて、そのデータに対する「処理」が行われる。この場合、データは「オブジェクト」としてデータベースノードにストアされており、データベースノード上のオブジェクトは「Java（登録商標）Object」（Java（登録商標）で扱うことができる任意のオブジェクト）として表現できる。そして、データベースノード上で実行するプログラムも「Java（登録商標）のクラス情報」として記述できるため、通常のJava（登録商標）のプログラミングのような直感的で容易なツールによる「分散処理」のプログラム開発が可能になり、Java（登録商標）Class Libraryによる「分

50

散開発環境」を構築することができる。

【 0 0 9 3 】

本システムでは、本システムの利用者が指定したプログラムをデータベースノード上で実行できるようにするため、クライアントノードとデータベースノードで、プログラムの定義である「クラス情報」を共有することができる。共有の仕方には、3つのパターンがある。第1のパターンは、そのクラス情報を必要とするノードの各々に、クラス情報（オリジナルまたはそのコピー）を登録するパターンである。第2のパターンは、共有ストレージを使うパターンである。第3のパターンは、共有のウェブサーバを使うパターンである。

【 0 0 9 4 】

図5を参照して、クラス情報の共有について更に説明する。共有ストレージを使う場合には、利用者が、共有ファイルサーバにクラス情報を登録し、各データベースノードに、保存場所を教えるから、サービスを開始させる。そうすると、各データベースノードが、共有ファイルサーバからクラス情報を取得する。共有のウェブサーバを使う場合には、利用者が、共有ウェブサーバにクラス情報を登録し、各データベースノードに、URLを教えるから、サービスを開始させる。そうすると、各データベースノードは、そのURL（共有ウェブサーバ）からクラス情報をダウンロードする。

【 0 0 9 5 】

次に、図6と図7を参照して、本システムの特徴について説明する。図6に示すように、本システムでは、クライアントノード側で処理内容（job）を定義して、そのプログラム（job）をデータベースノードへ送り、データベースノード側で分散して処理を実行し、クライアントノードは、その結果だけを取得する。一方、従来のシステム（KVS）では、クライアント側で処理を行うため、データベースサーバからデータを取得（get）し、その処理が終わったらデータをデータベースサーバへ送信（put）する。

【 0 0 9 6 】

したがって、本システムは、従来のシステム（KVS）に比べて、データベースノードとの通信の回数が少なく済み、データベースノードから送られるデータ量も小さく、しかも、データベースノードからデータが外に出ないため秘匿性の高いデータの取り扱いに適している、という特徴を有している。

【 0 0 9 7 】

図7には、本システムと従来のシステム（KVS）との違いが概念的に示されている。なお、図7では、ノードやサーバで実行される処理が「略四角形」の図形で示されており、処理の対象となるデータが「円形」の図形で示されている。図7に示すように、従来のシステム（KVS）では、処理そのものはクライアント側で実行される。そのため、従来のシステムでは、クライアントの処理能力に制限があり、大規模な処理を実行するのは困難である。それに対して、本システムでは、処理そのものをデータベースノードに分散させることができる（理論上、データベースノードの数だけ分散することができる）。したがって、本システムでは、クライアントノードの処理能力の制限を受けることなく（複数のデータベースノードの処理能力を利用して）、大規模な処理を実行することが可能である。

【 0 0 9 8 】

本システムは、データベースノードへの処理の実装が容易であり、処理の分散が可能である点で、従来のシステム（KVS）より優れている。したがって、本システムは、扱うデータ量が多い（例えば、数千万件～数百億件）サービスで大規模なデータ処理を行うのに適している。例えば、携帯事業者などがメールを処理するためのシステムでは、1億人分のユーザのメール（ユーザ1人あたり数千件とすると、数千億件のメール）が、1万台のデータベースサーバに分散してストアされる。この場合、データベースサーバ1台あたり、1万人分のユーザのメール（ユーザ1人あたり数千件とすると、数千万件のメール）をストアし、処理するように、適度にデータ及びその処理が分散されることになる。また、本システムは、スケーラビリティ（拡張性）が高いため、データの増加率が高いサービ

10

20

30

40

50

スに適している。例えば、取り扱うデータ量が増えたときには、データベースノードの台数を追加するだけで対応可能である。また、本システムは、データがデータベースノードの外部に出ることがないので、秘匿性の高いデータを扱うサービスに適している。また、各データがいくつかの要素から構成される場合、単純なデータ取得だけでなく、データの構造に従った分析や処理が可能である。

【 0 0 9 9 】

以下、図面を参照しながら、本システムの具体的な適用例について説明する。図 8 には、個人情報扱うシステムへの適用例が示されている。図 8 に示すように、本システムでは、キー情報 (key) として識別番号 (例えば、0 1 2 3 4 5) を用い、そのキー情報とペアで各種のデータ (パスワード「h o g e」, 年齢「2 5」, 氏名「山田太郎」, 郵便番号「0 1 2 - 3 4 5 6」, 住所「東京都千代田区」など) がストアされている。

【 0 1 0 0 】

図 8 の「本システム」の「処理 1」に示すように、本システムでは、「パスワードチェック」の処理を行うことができる。この場合、「ユーザ名とパスワード」が入力されると、「認証の成否」が出力として得られる。

【 0 1 0 1 】

ここで、本システムの job の定義について、図 8 の「本システム」の「処理 1」を例にして、詳しく説明する。図 8 の「処理 1」の例において、第 1 行目は、プログラムの内容を記述するためのクラスの定義を示しており、この例では「Authenticate」というクラスが定義されている。第 2 行目は、このプログラムで利用する変数の定義を示しており、この例では、「pass」という String 型の変数が定義されている。第 3 ~ 5 行目は、プログラムをオブジェクトとして生成するための方法 (メソッド) を定義しており、この例では、文字列を引数とし、内部変数「pass」に、該引数を保持するよう指示している。なお、このメソッドは、Java (登録商標) プログラミングで「constructor」と呼ばれているものである。第 6 ~ 9 行目は、データベースノードで実行されるメソッドを定義しており、この例では「Boolean applyTo」というメソッドが定義されている。また、「Object o」は、データベースノード上に保存されたデータ (オブジェクト) を引数として受け取ることの意味している。第 7 行目は、引数として受け取った (データベースノード上の) オブジェクトに対して行う処理の内容を記述しており、この例では「入力されたパスワードと保存されているパスワードを比較する」処理が示されている。ここでは、データベースノードに保持されたデータ (オブジェクト) を「Map (Java (登録商標) プログラムで定義されるクラスのひとつ)」として扱うものとする。また、このプログラムを利用するアプリケーションでは、データベースノードに保持されたデータ (オブジェクト) は「Map」というクラスに属するものとして扱い、その中で「password」というキーに対応する値として、ユーザのパスワード情報が保存されているものとする。このように、図 8 の「処理 1」の例では、データベースノードに保存されているオブジェクト内の「password」というキーに対応する値と、本システムの利用者がパスワードとして入力した内部変数の「pass」とを比較して、一致するかどうかを判断し、「Boolean」の真偽値 (一致したか否かの値) を返す処理が行われる。

【 0 1 0 2 】

また、図 8 の「本システム」の「処理 2」に示すように、本システムでは、「所定の住所と年齢に該当するかの判断」の処理を行うことができる。この場合、「識別番号と住所と年齢の範囲」が入力されると、「該当するか否か」が出力として得られる。この例では、「GetCampain」というクラスが用いられ、「applyTo」という String 型のメソッドが定義されている。そして、データベースノードに保存されているオブジェクトの「address (住所)」というキーに対応する値が「東京都」で始まるものであり、かつ、「age (年齢)」というキーに対応する値が 2 0 以上であれば「選挙に行こう」という文字列を返す処理が行われる。

【 0 1 0 3 】

図 9 には、オンラインショップの購買履歴を扱うシステムへの適用例が示されている。

図 9 に示すように、本システムでは、キー情報 (key) として識別番号 (例えば、0 1 2 3 4 5) を用い、そのキー情報とペアで各種のデータ (パスワード「h o g e」, 氏名「山田太郎」, 購買履歴 1「ビール」, 購買履歴 2「ウイスキー」など) がストアされている。

【 0 1 0 4 】

図 9 の「本システム」の「処理 1」でも、図 8 と同様の「パスワードチェック」の処理が行われる。ここでは、詳しい説明を省略するが、この例でも、「Authenticate」というクラスが用いられ、「ユーザ名とパスワード」が入力されると、「認証の成否」が出力として得られる。

【 0 1 0 5 】

また、図 9 の「本システム」の「処理 2」に示すように、本システムでは、「過去に酒類を買ったことがあるかの判断」の処理を行うことができる。この場合、「ユーザ ID」が入力され、過去に酒類の購買履歴があれば「酒類のキャンペーン」の情報が出力として返される。この例では、「GetCampain」というクラスが用いられ、「applyTo」という String 型のメソッドが定義されている。そして、データベースノードに保存されているオブジェクトの「history」というキーに購買履歴に関する情報が配列で格納されており、該配列から、過去の購買履歴を検索し、購買履歴の中に「ビール」の文字列がある場合には、「ビールのキャンペーン」という文字列を返し、そうでなければ、「キャンペーンなし」という文字列を返す処理が行われる。このように、データベースノードに保存されるデータは、処理に応じた構造を持つデータをオブジェクトとして持つことができる (本例では Map の中に、文字列や数値だけでなく、さらに配列を格納している)。

【 0 1 0 6 】

図 1 0 には、商品在庫管理システムへの適用例が示されている。図 1 0 に示すように、本システムでは、キー情報 (key) として商品 ID (例えば、0 1 2 3 4 5) を用い、そのキー情報とペアで各種のデータ (品名「チョコレート」, 価格「1 5 0」, 2 0 1 0 年 4 月の販売数「1 2 0 0」, 2 0 1 0 年 5 月の販売数「1 5 0 0」など) がストアされている。なお、この例では、格納するデータを表現するための「SellInfo」という独自のクラスを定義している。このように、データベースノードに格納するデータは、オブジェクト指向型プログラミングで定義できる、任意の独自クラスのオブジェクトを持つ事もできる。

【 0 1 0 7 】

図 1 0 の「本システム」の「処理 1」に示すように、本システムでは、「現在の在庫数」を求める処理を行うことができる。この場合、「商品 ID」が入力されると、「在庫数」が出力として得られる。この例では、「GetRemains」というクラスが用いられ、「applyTo」という Integer 型のメソッドが定義されている。そして、データベースノードに保存されているオブジェクトを、独自クラス「SellInfo」のオブジェクトとみなして、該オブジェクトの「remains (在庫数)」という変数に格納された値を返す処理が行われる。

【 0 1 0 8 】

また、図 1 1 の「本システム」の「処理 2」に示すように、本システムでは、「月次の平均販売数の計算」の処理を行うことができる。この場合、「商品 ID」が入力されると、「月次販売数の平均値」が出力として得られる。この例では、「GetSellAverage」というクラスが用いられ、「applyTo」という Double 型のメソッドが定義されている。そして、データベースノードに保存されているオブジェクトを、独自クラス「SellInfo」のオブジェクトと見なして、該オブジェクトの「sold (販売履歴)」に格納されている「sell-2010-01 (2 0 1 0 年 1 月の販売数)」、・・・、「sell-2010-12 (2 0 1 0 年 1 2 月の販売数)」というキーに対応する値を合計し、その平均値を返す処理が行われる。

【 0 1 0 9 】

図 1 1 には、健康管理システムへの適用例が示されている。図 1 1 に示すように、本システムでは、キー情報 (key) として個人 ID (例えば、0 1 2 3 4 5) を用い、そのキー情報とペアで各種のデータ (2 0 1 0 年 9 月 1 日の取得カロリー「1 6 0 0」, 2 0 1

10

20

30

40

50

0 年 9 月 2 日の取得カロリー「1750」、体重「65kg」など)がストアされている。

【0110】

図11の「本システム」の「処理1」に示すように、本システムでは、「現在の体重を求める」処理を行うことができる。この場合、「ユーザID」が入力されると、「現在の体重」が出力として得られる。この例では、「GetWeight」というクラスが用いられ、「applyTo」というDouble型のメソッドが定義されている。そして、データベースノードに保存されているオブジェクトの「weight(体重)」というキーに対応する値を返す処理が行われる。

【0111】

また、図11の「本システム」の「処理2」に示すように、本システムでは、「毎日の平均取得カロリーの計算」の処理を行うことができる。この場合、「ユーザID」が入力されると、「取得カロリーの平均値」が出力として得られる。この例では、「GetAverage」というクラスが用いられ、「applyTo」というDouble型のメソッドが定義されている。そして、データベースノードに保存されているオブジェクトの「start(開始日。この例では2010年9月1日)の取得カロリー」、「2010年9月2日の取得カロリー」、・・・というキーに対応する値を合計し、その平均値を返す処理が行われる。

【0112】

図12と図13には、通信事業者のメールサービスへの適用例が示されている。図12は、従来のシステム(KVS)ではデータが大きすぎるのに対して、本システムでは効率的に処理を行うことができる例である。図12に示すように、本システムでは、キー情報(key)として個人ID(例えば、03a9e0)を用い、そのキー情報とペアで各種のデータ(最大保存メール数「950」、第1番目のメールの本文、第2番目のメールの本文など)がストアされている。

【0113】

図12の「本システム」の「処理1」に示すように、本システムでは、「最新のメールを取得する」処理を行うことができる。この場合、「ユーザID」が入力されると、「最新メールの本文」が出力として得られる。この例では、「GetNewMail」というクラスが用いられ、「applyTo」というString型のメソッドが定義されている。そして、データベースノードに保存されているオブジェクトの「maxId(最新メール)」というキーに対応するメールの本文を返す処理が行われる。

【0114】

また、図12の「本システム」の「処理2」に示すように、本システムでは、「保存メールの数を求める」処理を行うことができる。この場合、「ユーザID」が入力されると、「削除されていないメールの数」が出力として得られる。この例では、「GetCount」というクラスが用いられ、「applyTo」というInteger型のメソッドが定義されている。そして、データベースノードに保存されている「第1番目のメール」から「最新メール」までの全メールについて、本文が「null」でないメール(削除されていないメール)の数をカウントし、カウントした数(削除されていないメール数)を返す処理が行われる。

【0115】

図13は、複数のキーにまたがる処理の例である。図12と同様、本システムでは、キー情報(key)として個人ID(例えば、03a9e0)を用い、そのキー情報とペアで各種のデータ(最大保存メール数「950」、第1番目のメールの本文、第2番目のメールの本文など)がストアされている。

【0116】

図13の「本システム」の「処理3」に示すように、本システムでは、複数の「ユーザID」にまたがる場合であっても、「保有メール数の平均値を算出する」処理を行うことができる。この場合、「ユーザID」の集合が入力されると、「保有メール数の平均値」が出力として得られる。この例では、「GetCount」というクラスが用いられ、「applyTo」というInteger型のメソッドが定義されている。そして、それぞれの「ユーザID」に

10

20

30

40

50

について、データベースノードに保存されている「第1番目のメール」から「最新メール」までの全メールについて、本文が「null」でないメール（削除されていないメール）の数をカウントし、カウントした数（削除されていないメール数）を返す処理が、各データベースノードにおいて行われる。その後、それぞれの「ユーザID」について得られた「メール数」を合計し、その平均値を求める処理が、クライアントノードにおいて行われる。

【0117】

以下には、本実施の形態に係る分散型データベースシステムを利用するアプリケーションの提供者が、プログラム開発やオペレーション（運用管理）をより簡易に行えるようにするためのフレームワークを含む本システムについて、図14～17を参照しながら説明する。

10

【0118】

本フレームワークにより、例えば、Android端末や、Java（登録商標）（Apple、OSGI、他、任意のJava（登録商標）仮想マシンを含む）が稼働するコンピュータ等のユーザ端末から、クラウド（本システム）上のオブジェクトに簡単にアクセスし、クラウド上のデータオブジェクトをあたかもユーザ端末の手元にあるオブジェクトであるかのように扱えるようにするサービスを、容易に提供することが可能になる。

【0119】

このようなサービスを提供するためのアプリケーションプログラムは、既存技術では、クライアント（ユーザ端末）/サーバ（本システムではクライアントノード）/データベース（本システムではデータベースノード）の三層で開発することになるため、クライアントとサーバ側との間の通信に加え、サーバとデータベースとの間の通信をも意識して、アプリケーションプログラムを開発する必要がある。

20

【0120】

これに対し、本フレームワーク（例えばJava（登録商標）のようなオブジェクト指向型プログラミングに適用される）では、サーバとデータベースとの間の通信や、データが保存される場所、処理が行われる場所等を意識することなく、あたかもローカルに動くアプリケーションプログラムであるかのように開発することができる。

【0121】

また、既存技術では、サーバ側（サーバ/データベース）の管理運用に手間がかかるのに対し、本フレームワークでは、サーバ側の一切の運用管理（物理的なサーバの冗長化やサイジング（スケールの調整）等を含む）を事業者任せにすることが可能になる。

30

【0122】

本フレームワークを含む本システムでは、クライアントノードにおいて、必要に応じて、キーに対応するデータベースノード上のオブジェクトと通信するためのプロキシ（Proxy）オブジェクトを保持する。プロキシオブジェクト自体は、Java（登録商標）の技術として知られているものであるが、本フレームワークでは、プロキシオブジェクト（クライアント側オブジェクト）に、本システムにおけるキーの情報を保持するようにする。そして、後述するメソッド呼び出し時に、このプロキシオブジェクトを用いる。

【0123】

具体的には、クライアントノードにおいて、上記のプロキシオブジェクトに対するJava（登録商標）のメソッド呼び出し（method call）をフック（hook）し、自動的に「リモートプロシージャコール」に変換することにより、データベースノード上のオブジェクトに対するメソッド呼び出しを実行する。

40

【0124】

プロキシによるフックを行うと、Java（登録商標）のメソッド呼び出し時に、呼び出したメソッドの本来のメソッド名と引数に関する情報が取得できる。これを利用して、取得した情報をもとに、該当するメソッドを呼び出すジョブ（job）を自動生成する仕組みを、本システムに含ませておく。

【0125】

本フレームワークにより自動生成されたジョブは、呼び出されるメソッド名と引数に関

50

する情報を含むものであり、本システムによって、上記キーの値により特定されるデータベースノードへ配信され、分散実行される。

【0126】

図16は、上述した本フレームワーク（自動ジョブ生成技術と呼ぶことがある）を含む本システムを説明するための図である。

【0127】

クライアントノードは、プロキシオブジェクトを保持しており、プロキシオブジェクトには、メソッドが呼び出された時に実行されるフック用のメソッドが含まれており、メソッドが呼び出された時に、本来のメソッド名と引数に関する情報（例えば引数のクラス及び値）が渡される。これにより、フック用のメソッドの実行結果が、本来のメソッドの実行結果になる。

10

【0128】

このプロキシオブジェクトに、データベースノード上のオブジェクトを特定するキー情報が保持される。そして、上記フック用のメソッドが呼び出された時に、ジョブが自動生成され、キーのハッシュ値等に基づいて特定されるデータベースノードへ配信される。

【0129】

データベースノードへ配信されるのは、キー情報とジョブであり、ジョブには、メソッド名と、引数に関する情報（引数のクラスと値）と、上記メソッドを実行して結果を返すという処理内容とが含まれる。

【0130】

20

このジョブを受け取ったデータベースノードは、キーに対応する対象オブジェクトをデータベースからロードし、その対象オブジェクト内のデータに対して、ジョブにより規定される処理を実行し、実行結果をクライアントノードへ返す。このとき、処理によってデータに変更があった場合、データベースノードは、変更後の対象オブジェクトをデータベースにストアする。

【0131】

以上に説明した自動ジョブ生成技術を採用した本システムの具体的な適用例として、図15の買い物カートの例を説明する。本例では、オンラインショッピングにおいて買い物カートに入れられたみかんの値段と数から全体の金額を計算する。

【0132】

30

まず、図15と同じ適用例を、本フレームワークを含まない本システムにより実現する場合について、図14を用いて説明しておく。図14では、図8～13に示したのと同様の手法で、みかんの値段と数をセットする処理のためのジョブ（「OrangeSet」）と、みかんの合計金額を計算する処理のためのジョブ（「OrangeTotal」）とが、定義されている。

【0133】

図8～13では、定義されたジョブをキーにより特定されるデータベースノードへ配信して処理を実行させるためのプログラムの記述を省略していたが、図14には、そのプログラムが「Application例1」で記述されている。

【0134】

40

「Application例1」におけるプログラムは、まず、「jobcast.store」という本システムへの処理の登録を行うためのメソッドを使って、処理の対象としたいキーと実行したい処理（ここでは「orange」）とを登録する。そして、みかんの値段と数をセットする処理を行うために、対応するジョブ（「OrangeSet」）を生成し、「jobcast.execute」という本システムにおけるジョブの配信を行うためのメソッドを使って、生成したジョブをキーにより特定されるデータベースノードへ配信し、処理を実行させる。さらに、みかんの合計金額を計算する処理を行うために、対応するジョブ（「OrangeTotal」）を生成し、「jobcast.execute」という本システムにおけるジョブの配信を行うためのメソッドを使って、生成したジョブをキーにより特定されるデータベースノードへ配信し、処理を実行させる。

50

【 0 1 3 5 】

つまり、図 8 ～ 1 4 に例示したスキームでは、アプリケーションプログラムの開発者は、各々の処理のためのジョブを定義するプログラムと、そのジョブを生成、配信、実行するためのプログラムとを、本システムに特有の方式で記述することになる。

【 0 1 3 6 】

これに対し、図 1 5 に例示するスキームでは、手元にあるコンピュータ単体で処理を行う（リモートで分散して処理を行うのではない）場合と、ほぼ同様の方式で本システムを利用するアプリケーションプログラムを記述することができる。これが、自動ジョブ生成技術を採用した本システムの利点であり、図 1 7 に具体的に例示されている。

【 0 1 3 7 】

図 1 7 は、Java（登録商標）のプログラミングの例であり、単体での処理の場合、その処理（ここでは「Orange」）のインターフェースを定義した上で、「OrangeImpl」として処理内容を定義する。そして、定義された処理を実行するための「Application例 2」が、「Orange」オブジェクトを生成し、「OrangeImpl」に定義された各メソッドを呼び出すように、記述される。

【 0 1 3 8 】

この単体での処理のためのプログラムと、図 1 5 の自動ジョブ生成技術を使った分散処理のためのプログラムとを比較すると、インターフェースの定義と処理内容の定義は共通であり、定義された処理を実行するための「Application例 2」におけるメソッド呼び出しも共通であり、「Application例 2」におけるオブジェクト生成を、プロキシを用いながらキーを設定するものに変えるだけでよいことが分かる。したがって、開発者は、本システムを利用したアプリケーションプログラムを、従前に慣れ親しんだ方式で記述することが可能になる。

【 0 1 3 9 】

ここで、図 1 5 に戻り、本システムで自動生成を実現するための仕組みについて、より詳細に説明する。

【 0 1 4 0 】

Java（登録商標）では、プロキシという技術を使うことによって、あるインターフェース（本例では「Orange」）で定義されるメソッド（*1）を持つ仮想的なオブジェクト（*2）を生成して、上記（*1）のメソッドが呼ばれたときに、メソッド名と引数（メソッド呼び出しの際のパラメータ）の情報を受け取る仮想的なメソッドを定義することができる。すなわち、プロキシで仮想的なオブジェクトを作っておくことで、それ以降、メソッド呼び出しの際に、当該メソッド呼び出し時のメソッド名と引数の情報を、『特定のメソッド』（*3）に渡すことができる。

【 0 1 4 1 】

さらに、本フレームワークでは、上記（*2）のオブジェクトにキー情報を持たせておいて、上記（*3）の『特定のメソッド』において、呼び出された本来のメソッド名（*4）とメソッドが呼び出された際の引数の情報（*5）とからジョブ（*6）を生成し、キーで特定されるデータベースノードへジョブを送信し、その結果（送信先のデータベースノードにおいてキーで特定されるデータオブジェクトに対してジョブを実行することにより得られた結果）を受け取る、ということを行う。すなわち、上記（*3）の『特定のメソッド』というものが、ジョブの生成と、ジョブの配信と、結果の受信とを自動的に行うことになる。

【 0 1 4 2 】

上記（*3）の『特定のメソッド』が生成するジョブは、メソッド呼び出しに応じて実行すべき処理内容が特定されたオブジェクトであるが、この生成を実現するために、図 1 5 に「jobの定義例」として記載されているように、一般化したジョブのクラス（「class GenericJob」）が予め定義されている。

【 0 1 4 3 】

個々に生成される上記（*6）のジョブは、一般化したジョブのクラスの一つのオブジェ

10

20

30

40

50

クト（インスタンス）であり、上記（*4）のメソッド名と、上記（*5）の引数の情報とを持つ。そして、ジョブの実行時には、これらメソッド名と引数の情報から、データベースノード上に保存されている対象オブジェクトに定義されている本来のメソッド情報を探してきて、対象オブジェクトのメソッドを呼び出すことで、本来のメソッドの実行がデータベースノード上で行われることになる。

【0144】

上述した仕組みにより、分散処理を考えない場合と同様に「obj.setOrange(10, 100)」というメソッド呼び出しを行うと、事前に「getProxy」で生成されているプロキシオブジェクトの、『特定のメソッド』が呼び出され、本来呼び出されたメソッド名（「setOrange」）と、メソッドが呼び出されたときの引数情報（10(int), 100(int)）とをもとにjob 10
インスタンスが生成される。このjobの中身は、メソッド名である setOrangeと、引数のタイプ情報であるint, intという情報と、引数の値である10, 100という情報を持つ。

【0145】

そして、上記『特定のメソッド』は、プロキシオブジェクトを生成したときに渡したキー情報をもとに、対象となるデータベースノードを選択し、そのデータベースノードへ上記のjob インスタンスを送付して、そのデータベースノードがjob を実行した結果を、ネットワークを介して受け取る。受け取った結果が、最終的に、上記「obj.setOrange(10, 100)」の実行結果として扱われる。

【0146】

なお、データベースノードがjob を実行する際には、job オブジェクトの applyTo が 20
実行されるため、その中では、メソッド名（setOrange）及び引数の情報（int, int型の10, 100という情報）から、対象のオブジェクトの該当メソッドを探し出して実行する。ここで、保存されているのはOrangeImplクラスのオブジェクトであるので、OrangeImplのsetOrange(10, 100)が呼び出され、その結果が返ることになる。

【0147】

「obj.getTotal()」というメソッド呼び出しについても、同様に、事前に「getProxy」で生成されているプロキシオブジェクトの、『特定のメソッド』が呼び出され、本来呼び出されたメソッド名（「getTotal」）をもとにjobインスタンスが生成される。このjobの中身は、メソッド名である getTotalという情報を持つ。

【0148】

そして、上記『特定のメソッド』は、プロキシオブジェクトを生成したときに渡したキー情報をもとに、対象となるデータベースノードを選択し、そのデータベースノードへ上記のjob インスタンスを送付して、そのデータベースノードがjob を実行した結果を、ネットワークを介して受け取る。受け取った結果が、最終的に、上記「obj.getTotal()」 30
の実行結果として扱われる。

【0149】

なお、データベースノードがjob を実行する際には、job オブジェクトの applyTo が
実行されるため、その中では、メソッド名（getTotal）から、対象のオブジェクトの該当メソッドを探し出して実行する。ここで、保存されているのはOrangeImplクラスのオブジェクトであるので、OrangeImplのgetTotal()が呼び出され、その結果が返ることになる。 40

【0150】

図15の仕組みにより自動生成されるジョブは、図14で生成される各処理のジョブと機能的には同等のものになる。つまり、図15の「obj.setOrange(10, 100)」に対応して生成されるジョブは、「class GenericJob」の一つのオブジェクト(インスタンス)であり、メソッド名がorangeSet、パラメータ情報が(int, int)型の(10, 100)であるジョブになるため、図14の「OrangeSet」クラスのオブジェクト(job1 = new OrangeSet(10, 100))とほぼ同等である。また、図15の「obj.getTotal()」に対応して生成されるジョブは、「class GenericJob」の一つのオブジェクト(インスタンス)であり、メソッド名がgetTotal、パラメータ情報が無しであるジョブになるため、図14の「OrangeTotal」 50
クラスのオブジェクト(job2 = new OrangeTotal())とほぼ同等である。

【0151】

以上、図14の適用例に対して、図15の実現手法が可能であることを説明したが、図8～13の適用例に対しても、同様に、自動ジョブ生成技術を用いた実現手法を採用することが可能なことは、勿論である。このように自動ジョブ生成技術を用いれば、本システムを利用する大規模な並列処理プログラムの開発が可能となり、そのように開発したプログラムを実行することにより、自動的にデータの処理や保存がクラウド上で分散して行われることになる。

【0152】

上述したジョブ自動生成を行う本システムにおいて、データベースノード上のデータを「読み込む」処理と「変更を伴う」処理とを区別するために、以下のような拡張機能を付加してもよい。

【0153】

ここで、「変更を伴う」処理とは、値に変更があったら保存しなおすことが必要な処理であり、複数同時に実行できないように、サーバ側で排他ロックを行う。「読み込み」処理は、値の変更を意識しないため、保存しなおしはスキップされる処理であり、複数同時に実行できるため、サーバ側では共有ロックを行う。

【0154】

具体的には、例えば、Java（登録商標）のアノテーション（annotation）と呼ばれる技術で、「読み込み」か「変更を伴う」をマーク付けしておく。マークは、メソッド毎に定義することが可能である。そして、ジョブを自動生成する際に、マークの有無で実際に生成するジョブを区別し、例えば、デフォルト（マークがない場合）は「変更を伴う」処理として実行すればよい。

【0155】

以上に詳述した本システムにおいて、各ジョブにオペレーションIDを付与することにより、以下のように処理の確実性を担保することが可能になる。なお、オペレーションIDを付与するジョブは、自動生成されるジョブに限らず、本システムで使われるジョブのすべてに付与することが可能である。

【0156】

まず、クライアントノードにおいて、各ジョブに一意的なオペレーションID（例えばUUIDを利用）を付与してから、データベースノードへ配信する。そして、データベースノードにおいて、あるジョブによる処理が「実行済み」の場合には、そのジョブに付与されているオペレーションIDに対応させて結果を記録しておく。

【0157】

これにより、ネットワークもしくはデータベースノードで障害が発生し、処理の実行結果がクライアントノードに返ってこない場合に、クライアントノードからリトライをすると、ジョブを受信したデータベースノードは、受信されたオペレーションIDと記録とを照らし合わせることにより、そのジョブが以前に受信されて「実行済み」のもの（サーバ側でジョブを実行した後、結果がクライアントに届くまでの間に障害があったもの）なのか、そうでない（初めて実行する）ものなのかを判断することができる。

【0158】

そして、データベースノードは、「実行済み」とであると判断すれば、そのオペレーションIDに対応させて記録してあった結果を、クライアントノードに返し、「実行済み」とないと判断すれば、改めて処理を実行して、その結果をクライアントノードに返す。

【0159】

図18と図19には、本システムの利点がまとめられている。図18に示すように、本システムは、1回のデータ転送量が小さく、複数要素にまたがる処理でも、通信回数が少ない。また、処理の結果だけを返すため、秘匿性のあるデータの取り扱いに適している。さらに、図19に示すように、本システムでは、単一の項目へのアクセスをするときの通信量が少なく、複数の項目へアクセスするときの通信量も少ない。そして、バックエンド処理（データベースノード側での処理）のため、秘匿性のあるデータがネットワークに流

10

20

30

40

50

出ることがない。なお、図18と図19における「従来」のシステムは、KVSで細かくデータを分割して保存するようにしたシステムであり、「参考例」のシステムは、KVSで大きなデータを単位として保存するようにしたシステムである。

【0160】

以上には、本発明の実施の形態に係るシステムとして、分散型データベースシステムを例にとって詳述したが、本発明の原理に従う別のシステムとして、図20に例示するような分散型データ処理システム（これも「本システム」と呼ばれる）を構成することも可能である。

【0161】

例えば、大量のログを、長時間にわたって継続的に、高性能かつスケーラブルに処理、保存したい場合に、分散型データベースシステムの各データベースノードにデータをストアするのもよいが、それとは別のバックエンドのストレージにデータを保存するようにし、そのストレージに対する処理を分散して行う分散型データ処理システムを構築することも効果的である。

10

【0162】

図20の分散型データ処理システムは、デバイスからログのレコード情報を入力するクライアントノードと、並列分散的にレコードの保存／検索の処理を行う複数のデータ処理ノードとから構成される。

【0163】

図20におけるクライアントノードは、レコードを受付部として機能するのに加えて、分散型データベースシステムのクライアントノードに対応する機能を有する。具体的には、受け付けたレコード情報からキーを生成し、キーにより特定されるデータ処理ノードへジョブを配信することにより、処理を分散化する。

20

【0164】

図20における各データ処理ノードは、データの保存先が各ノードではなくバックエンドに設けられた共有ストレージである点が異なるが、分散型データベースシステムのデータベースノードに対応する機能を有し、クライアントノードから送信されてきたジョブを実行することにより、レコードを処理する。

【0165】

上記のデータ処理ノードは、ログのレコード情報にインテックス処理（indexing）を行い、インデックスに応じて、ストレージに書き込む。例えば、クライアントノードから受信したレコードを、その都度、インデックスに対応するデータ領域に追加書き込みしていくことにより、ログが蓄積される。よって、バッチ処理ではなくリアルタイムで、継続的に（24時間365日の稼働で）ログの処理が可能である。

30

【0166】

このような本システムは、通信事業者の通信ログや、企業のネット利用ログや、機械のセンサーログ等の処理、保存を行う応用の他、商店・販売店・オンラインショップの販売／購買記録等、大量の記録を行う応用に、適している。

【0167】

当該応用においては、デバイスが多数存在して、各々のデバイスにおいてレコードが任意のタイミングで生成され、生成されたレコードを各デバイスがネットワーク経由でクライアントノード（デバイスに対し、フロントエンドサーバとして働く）へ入力する。多数のデバイス中に、異なる種類のレコードを生成するデバイスが混在しても構わない。

40

【0168】

クライアントノードは、すべて同じ、ステートレスな処理を行うため、多数のクライアントノードを並列で動作させることができ、各デバイスは、どのクライアントノードへレコードを送信してもよい。

【0169】

データ処理ノードは、図1～19を参照して説明したデータベースノードと同様に、キーをもとに複数のデータ処理ノードで処理を分散するものであるため、スケールアウト（

50

台数を増やすことで全体性能が向上する)を実現可能である。

【0170】

以上のように、図20のシステムにおいては、ログの生成(デバイス)、ログの受付(クライアントノード)、インデックス処理(データ処理ノード)、書き込み(ストレージ)の各段階が直列に接続され、続々とパイプライン的に処理を進めることができ、かつ、各段階が多数の装置の並列分散処理で実行されるため、大量のログを高速に処理する高性能化が可能である。

【0171】

図20におけるストレージは、共有ストレージであればよいが、本発明者が提案した特願2010-237828のストレージ技術を用いれば、ストレージの書き込み/読み出しにおいてもスケールアウトと高性能化を実現することが可能になる。つまり、追加書き込み(蓄積)を行うデータ領域を、上記先行出願のストレージ上に設けるようにすれば、ストレージ容量もスケラブルに拡張可能となる。

10

【0172】

図21には、図20のシステムの適用例として、通信ログの処理を示す。例えば、デバイスとしての通信機器が、「どのアドレスから、どのURLに対してアクセスがあった」というレコード情報を生成し、任意に選択したクライアントノードへ送信する。

【0173】

クライアントノードは、デバイスから受け取ったレコード情報に基づいて、例えば「送信元アドレス」をキーにし、後述するジョブを生成して、キーに応じて特定されるデータ処理ノードへ送信する。

20

【0174】

データ処理ノードは、クライアントノードから受け取ったジョブを実行する。例えば、キーである「送信元アドレス」をインデックス(データベースでいう検索キーに相当するものとしてもよい)として扱い、キー毎にレコード情報を追加書き込み(蓄積)する。

【0175】

例えば、レコード(ログの内容)が、先頭から「通信時刻(8バイト);送信元アドレス(4バイト);送信先アドレス(4バイト);送信元ポート(2バイト);送信先ポート(2バイト);アクセスするURL(文字列)」のように並んでおり、インデックスとすべき要素が送信元アドレスの場合、クライアントノードは、レコードの先頭から8バイト目から4バイト分をキーとする。そして、クライアントノードは、キーのハッシュ値を計算し、計算値に従って、ジョブの配信先となるデータ処理ノードを特定する。

30

【0176】

上記のキーとジョブを受信したデータ処理ノードは、インデックスとすべき要素(ここでは送信元アドレス)であるキーに対応するデータ領域に、受け取ったジョブに含まれているレコード情報を追加書き込み(蓄積)する。例えば、受け取ったキーに基づいて、出力先のファイル名を決定し、レコード情報を共有ストレージ内に保存されるファイルとして出力する。ファイル名は、例えば、キーと受信した日時とから「key-date-time.log」のように定めればよい。

【0177】

以上では、レコード情報の保存について説明したが、クライアントノードからデータ処理ノードへ、レコード情報を検索するためのジョブを送信することにより、共有ストレージからの読み出しの処理も分散化することが可能である。この場合のジョブは、キーに対応するデータ領域からレコード情報を読み出して、レコード一覧を作成し、このレコード一覧を結果としてクライアントノードへ返すものである。

40

【0178】

図21に例示したジョブでは、検索条件としてインデックス(ここでは送信元アドレス)だけが指定されており、データ処理ノードは、インデックスの示すデータ領域から読み出されたレコード情報をすべて結果に含めて返すことになるが、クライアントノードから送信するジョブに詳細な検索条件(例えば、時刻、送信先アドレス、ポート番号、それら

50

の任意の組合せ等)を設定することにより、データ処理ノードが、インデックスの示すデータ領域から読みだされたレコード情報のうち検索条件に合致するものを抽出して、抽出したレコードの一覧をクライアントノードへ返すことも可能である。

【0179】

本システムでは、インデックスとすべき要素をキーとしてジョブが配信されるため、アクセス対象となるデータ領域によって、処理を行うデータノードが分散化されることになり、複数のデータ処理ノードから同じデータ領域がアクセスされる競合が生じる可能性を低減することができる。

【0180】

さらに、本システムでは、そのインデックスについて保存(ストレージへの書き込みを含む)の処理を行う同じデータ処理ノードが、検索(ストレージからの読み出しを含む)の処理も行うことになる。

10

【0181】

ここで、一つのインデックスについて保存や検索の処理を行うデータ処理ノードが、複数存在して、その都度選択されるようにしてもよい。そのために、例えば、送信元アドレス(インデックスとすべき要素)と時刻情報(その他の要素)とを組み合わせたものをキーとして、データ処理ノードの特定を行ってもよい。この場合、データ処理ノードでは、キー情報全体をインデックスとするのではなく、キー情報からインデックスとなる部分を抜き出してインデックスとすることになる。

【0182】

20

なお、図1～19を参照して説明した分散型データベースシステムにおける様々な仕組みや工夫は、図20～21を参照して説明した分散型データ処理システムにも適用することが可能である。

【0183】

以上、本発明の実施形態について説明したが、上述の実施形態を本発明の範囲内で当業者が種々に変形、応用して実施できることは勿論である。

【産業上の利用可能性】

【0184】

本発明の原理に従う分散型データベースシステムあるいは分散型データ処理システムによれば、例えば、様々なデータストレージサービスを提供することができ、有用である。

30

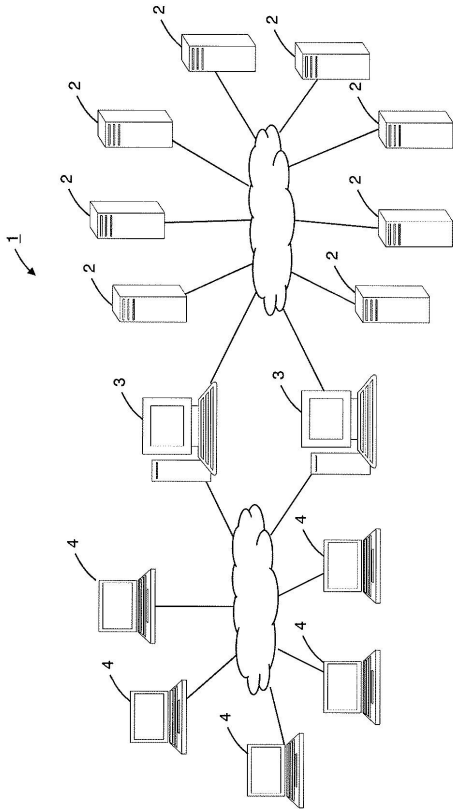
【符号の説明】

【0185】

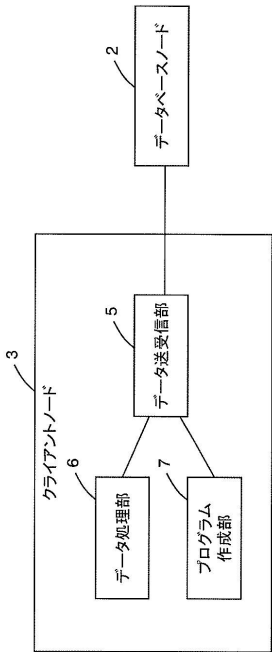
- 1 分散型データベースシステム(本システム)
- 2 データベースノード
- 3 クライアントノード
- 4 ユーザ端末
- 5 データ送受信部
- 6 データ処理部
- 7 プログラム作成部
- 8 データ送受信部
- 9 データ処理部
- 10 データ格納部
- 11 クラスライブラリ格納部
- 12 コンフリクト制御部
- 13 バージョン管理部

40

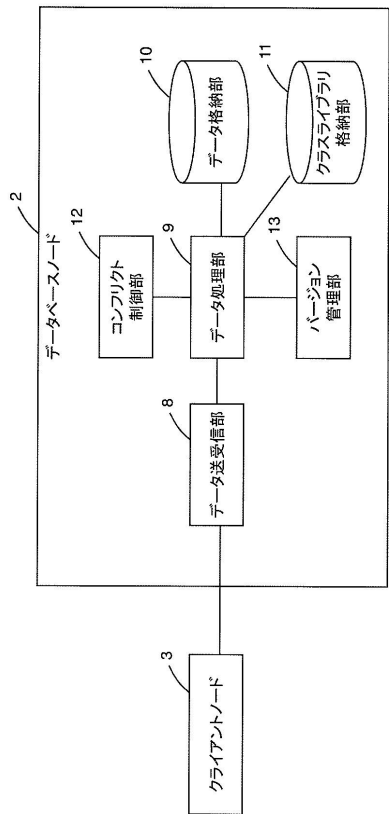
【図 1】



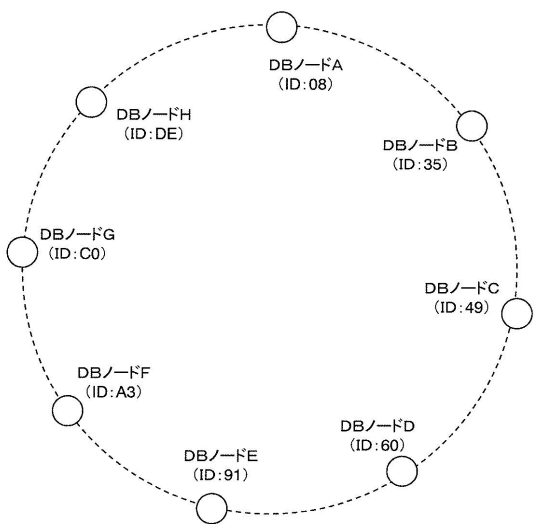
【図 2】



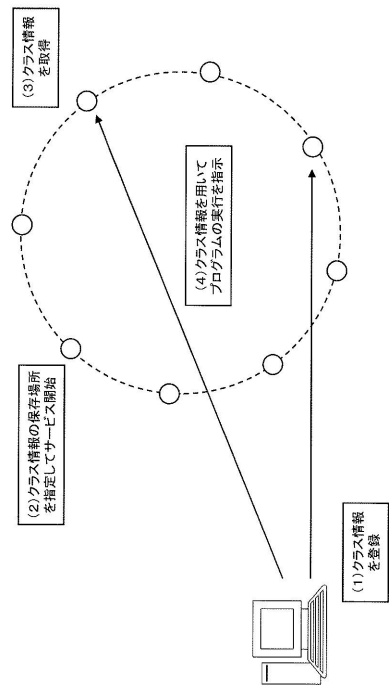
【図 3】



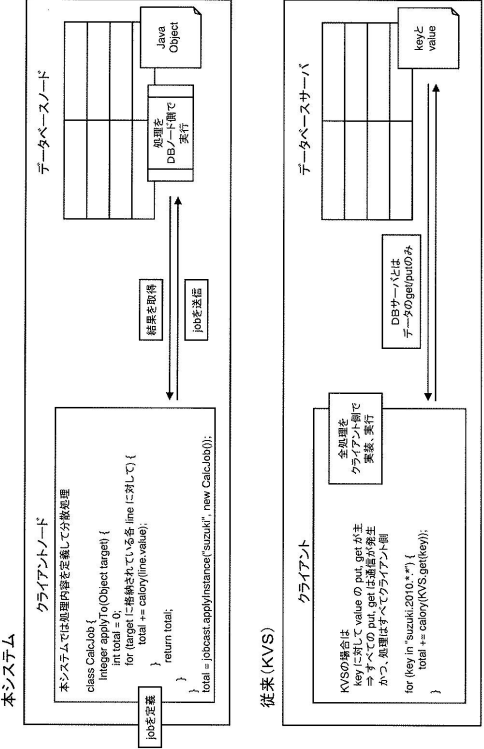
【図 4】



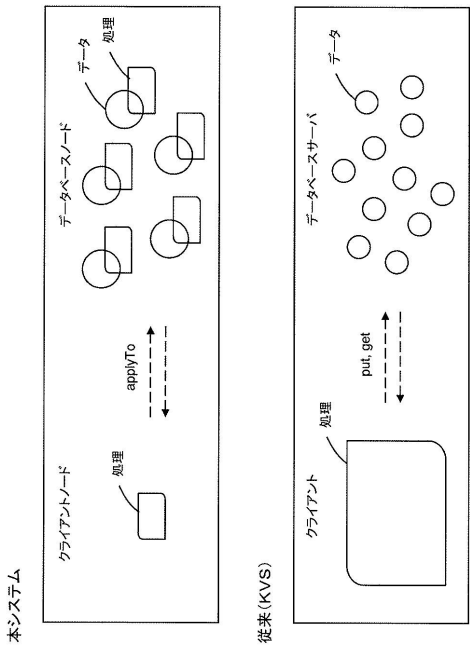
【図 5】



【図 6】



【図 7】



【図 8】

<p>個人情報を使うシステム</p>	<p>処理 1 「パスワードチェック」</p> <p>概要: ユーザーがデータにアクセスする際の認証 入力: ユーザーID、パスワードが入力されたパスワード 出力: accepted or denied</p> <pre>class Authenticate { String pass; Authenticate(String in) { pass = in; } Boolean applyTo(Object o) { return pass.equals(((Map)o).get("password")); } }</pre> <p>処理 2 「住所と年齢に該当するかどうかの判断」</p> <p>概要: 都内で20歳以上は選挙の呼びかけ、など 入力: 選挙番号、住所、年齢の範囲 出力: true or false</p> <pre>class GetCampaign { String applyTo(Object o) { if (((Map)o).get("address").startsWith("東京都")) && (((Map)o).get("age") >= 20) { return "選挙に行こう"; } return ""; } }</pre>
--------------------	---

健康管理システム

【図 1 1】

	処理 1 「現在の体重」 概要: ユーザーの現在体重を取得 入力: ユーザID 出力: 現在体重	処理 2 「今日の平均取得カロリー」 概要: ユーザーの取得カロリーの平均を計算 入力: ユーザID 出力: 取得カロリーの平均値
本システム KEY=ユーザID, VALUE=要素+値のテーブル 012345 = { 2010-09-01 = 1600 2010-09-02 = 1750 weight = 65kg } ※各情報は Map オブジェクト内に格納する	処理 1 の定義例 class GetWeight { Double applyTo(Objects o) { return ((Map)o).get("weight"); } } ※各情報は Map オブジェクト内に格納する	処理 2 の定義例 class GetAverage { Double applyTo(Objects o) { t = ((Map)o).get("start"); while ((v = ((Map)o).get(t)) != null) { total += v; count++; t++; } } }

オンラインショップの購買履歴を扱うシステム

【図 9】

	処理 1 「パスワードチェック」 概要: ユーザーがデータにアクセスする際の認証 入力: ユーザID 出力: accepted or denied	処理 2 「過去の注文履歴をまとめたリストがある」 概要: 履歴からキャンペーン情報を提供 入力: ユーザID 出力: true or false
本システム KEY=ユーザID, VALUE = { 各要素と値のテーブル } 012345 = { password = hoge, name = 山田太郎, history = [[ユーザID, ユニバーサル, ...]] } ※各情報は Map オブジェクト内に格納する history は Map 内にもさらに配列情報として格納する	処理 1 の定義例 class Authenticate { String pass; Authenticate(String in) { pass = in; } Boolean applyTo(Objects o) { return pass.equals(((Map)o).get("password")); } }	処理 2 の定義例 class GetCampaign { String applyTo(Objects o) { String[] array = ((Map)o).get("history"); for (i = 0; i < array.length; i++) { if (array[i].equals("ユーザID")) { return エラーのキャンペーン; } } return "キャンペーンなし"; }

通信事業者のメールサービス

【図 1 2】

	処理 1 「最新のメールを取得」 概要: メールボックスから最新メールを取得 入力: ユーザID 出力: 最新メールの本文	処理 2 「未読メールの数」 概要: 削除されていないメールの数をカウント 入力: ユーザID 出力: 削除されていないメールの数
本システム KEY=ユーザID, VALUE=メッセージ番号+本文のテーブル 034560 = { msgid = 456, 2 = <mail-1の本文> 2 = <mail-2の本文> 950 = <mail-950の本文> } ※各情報は Map オブジェクト内に格納する	処理 1 の定義例 class GetNewMail { String applyTo(Objects o) { return ((Map)o).get("msgid"); } } ※各情報は Map オブジェクト内に格納する	処理 2 の定義例 class GetCount { Integer applyTo(Objects o) { tail = ((Map)o).get("msgid"); for (int i = 1; i <= tail; i++) { if (((Map)o).get(i)) != null) count++; } return count; }

商品在庫管理システム

【図 1 0】

処理 1 「現在の在庫数」 概要: 指定の商品の在庫数を確認する 入力: 商品ID 出力: 在庫数	処理 2 「月次の平均販売数」 概要: 指定の商品の毎月の売れ行きを確認する 入力: 商品ID 出力: 月次販売数の平均値
本システム KEY = 商品ID, VALUE = { 各要素と値のテーブル } 012345 = { item = アイコネット, price = 150, sell-2010-04 = 1200, sell-2010-05 = 1500, ... remains = 200 } ただし、上記は以下の独自クラスに格納する。 class SellInfo { String item; // item 名 int price; // 値段 int remains; // 在庫数 Map sold; // 販売数のテーブル }	処理 1 の定義例 class GetRemains { Integer applyTo(Objects o) { return ((SellInfo)o).remains; } } 処理 2 の定義例 class GetSellAverage { Double applyTo(Objects o) { for (t = 1; t <= 12; t++) { v = ((SellInfo)o).sold.get("sell-2010-" + t); total += v; count++; } } }

【 図 1 3 】

通信事業者のメールサービス

本システム	<pre> KEY* 個人ID VALUE=メッセージ番号+本文のテーブル 03a9d0 = { mail = 950 maxId = 950 1 = <mail>の本文> 2 = <mail>の本文> ... 950 = <mail>950の本文> } </pre>
図 3 「所有メールアドレスの平均値」	<p>個人ユーザが保持しているメールアドレスの平均を計算 入力ユーザIDの値を 出力: 所有メールアドレスの平均値</p> <pre> // 0300 定義欄 class GetCount{ Integer applyToObjects o { int tail = (Map)o.get("mail"); for (int i = 1; i <= tail; i++) { if (((Map)o.get(i)) != null) count++; } return count; } } </pre> <p>注記: ジョブ (プログラム = 処理を表すクラス) を用いて以下を要約 for (Userid : a set of user id) { count += applyTo(Userid, new GetCount()); } users++; return count / users;</p>

【 図 1 5 】

みかんの値段と数から、全体の金額を計算する例

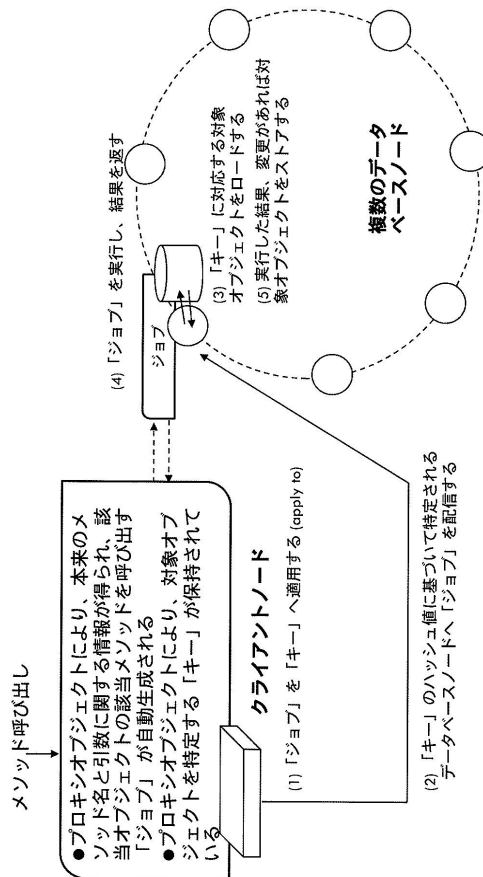
データ構造の定義 - クラスの定義	メンバの呼び出し方に関して インスタンスを生成	型宣言 1 「みかんの値段と数を扱う」 例案：みかんの値段と数を扱う 出力：みかんの値段と数 出力：なし	処理 2 「みかんの値段を計算」 例案：みかんの値段と数を計算 出力：みかんの値段と数 出力：みかんの値段と数
データ構造の定義 - クラスの定義	メンバの呼び出し方に関して インスタンスを生成	型宣言 1 「みかんの値段と数を扱う」 例案：みかんの値段と数を扱う 出力：みかんの値段と数 出力：なし	処理 2 「みかんの値段を計算」 例案：みかんの値段と数を計算 出力：みかんの値段と数 出力：みかんの値段と数
データ構造の定義 - クラスの定義	メンバの呼び出し方に関して インスタンスを生成	型宣言 1 「みかんの値段と数を扱う」 例案：みかんの値段と数を扱う 出力：みかんの値段と数 出力：なし	処理 2 「みかんの値段を計算」 例案：みかんの値段と数を計算 出力：みかんの値段と数 出力：みかんの値段と数

【 図 1 4 】

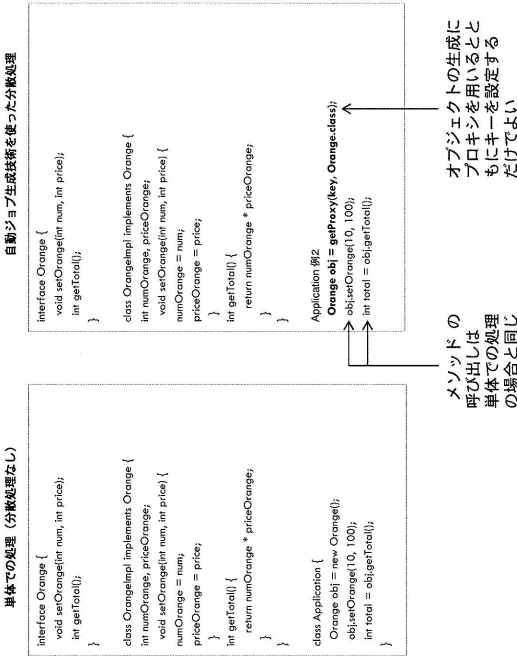
みかんの値段と数から、全体の金額を計算する例

<p>本システム</p> <pre> KEY: 買い物カートのID VALUE: { かんかの値段と数 } 012345 = { numOrange = 10 priceOrange = 100 } </pre> <p>※データベースモードでは、各種情報をこのJava クラスのオブジェクトとして保存</p> <pre> class OrangeData { int num, price; } </pre>	<p>処理 1 「かんかの値段と数をセット」 概要: みかんの値段と数をセット 入力: みかんの値段と数 出力: なし</p>	<p>処理 2 「合計金額を計算」 概要: みかんの合計金額を計算 入力: なし 出力: 合計金額</p>	<p>処理 3 「合計金額を計算」 概要: みかんの合計金額を計算 入力: なし 出力: 合計金額</p>
<pre> class OrangeSet implements Job { int num, price; OrangeSet(int n, int p) { num = n; price = p; } boolean applyTo(Orange o) { OrangeData odata = new OrangeData(); odata.num = num; odata.price = price; } } </pre>	<p>Application 例 1</p> <pre> jobdata.store(key, orange); Job job1 = new OrangeSet(10, 100); jobdata.execute(key, job1); Job job2 = new OrangeTotal(); int total = jobdata.execute(key, job2); </pre>	<p>Application 例 2</p> <pre> class OrangeTotal implements Job { int num, price; OrangeTotal() { num = 0; price = 0; } boolean applyTo(Orange o) { num += o.num; price += o.price; } } </pre>	<p>Application 例 3</p> <pre> class OrangeTotal implements Job { int num, price; OrangeTotal() { num = 0; price = 0; } boolean applyTo(Orange o) { num += o.num; price += o.price; } } </pre>
<p>※データベースモードでは、各種情報をこのJava クラスのオブジェクトとして保存</p> <pre> class OrangeData { int num, price; } </pre>	<p>Application 例 1</p> <pre> jobdata.store(key, orange); Job job1 = new OrangeSet(10, 100); jobdata.execute(key, job1); Job job2 = new OrangeTotal(); int total = jobdata.execute(key, job2); </pre>	<p>Application 例 2</p> <pre> class OrangeTotal implements Job { int num, price; OrangeTotal() { num = 0; price = 0; } boolean applyTo(Orange o) { num += o.num; price += o.price; } } </pre>	<p>Application 例 3</p> <pre> class OrangeTotal implements Job { int num, price; OrangeTotal() { num = 0; price = 0; } boolean applyTo(Orange o) { num += o.num; price += o.price; } } </pre>

【 図 1 6 】



【図 17】



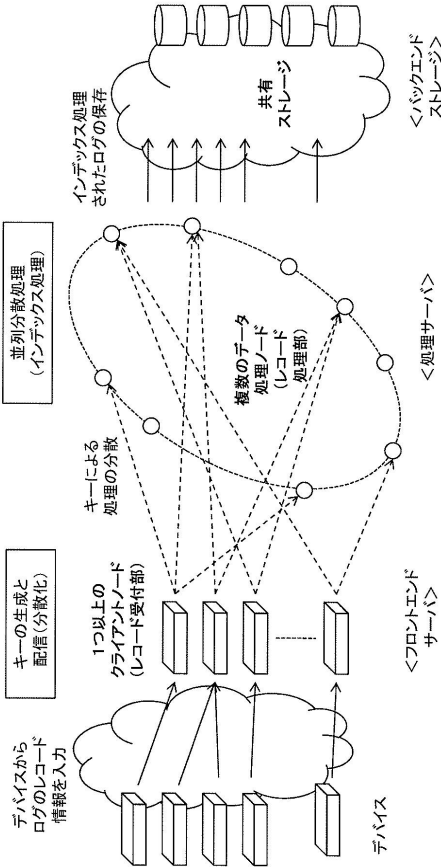
【図 18】

	1回のデータ転送量	複数要素にまたがる処理	秘匿性のあるデータ
本システム	◎（転送量＝小）	◎（通信回数＝少）	◎（結果だけ返す）
従来	◎（転送量＝小）	△（通信回数＝多）	×
参考例	△（転送量＝大）	○（通信回数＝少）	×

【図 19】

	単一の項目 へのアクセス	複数の項目 へのアクセス	データの秘匿性
本システム	◎ 通信量：少ない	◎ 通信回数：少ない	◎ バックエンド処理のため秘匿性のあるデータをネットワークに流す必要なし
従来	◎ 通信量：少ない	△ 通信回数：多い	×
参考例	△ 通信量：多い	○ 通信回数：少ない	△ フロントエンド処理のため、秘匿性のあるデータをネットワークに流す必要あり

【図 20】



ログの高速な保存

レコードの高速な保存	レコード構造の定義＝クラスの定義	高度「レコード情報をファイルに保存 入力：レコード情報 出力なし」	高度「レコード情報の検索」 概要：送信元アドレスからレコード情報を検索 入力：送信元アドレス 出力：レコード一覧
レコードの構造例 (以下はログの例) 通信時間 = 1234.56 送信元アドレス = 10.1.2.3 送信元アドレス = 10.10.10.10 送信元ポート = 2000 送信元ポート = 1000	本システム (送信元アドレスをキーとする例) KEY = 10.1.2.3 VALUE = レコードの情報。 として、データを生成 10.1.2.3 = { 通信時間 = 1234.56 送信元アドレス = 10.1.2.3 送信元アドレス = 10.10.10.10 送信元ポート = 2000 送信元ポート = 1000 } → クライアントノード(レコード受け側)が 「キー」により特定されるデータ処理ノード (レコード処理部)へ 「ジョブ」を配信して実行させる	高度「レコード情報の検索」 概要：送信元アドレスからレコード情報を検索 入力：送信元アドレス 出力なし	高度「レコード情報の検索」 概要：送信元アドレスからレコード情報を検索 入力：送信元アドレス 出力：レコード一覧

フロントページの続き

- (56)参考文献 高槻芳, 特選フリーソフト ROMA Rubyで動くキーバリュ型データストア, 日経Linux, 日本, 日経BP社, 2010年 4月 8日, 第12巻, 第5号, pp. 92 - 94
まつもとゆきひろ, まつもとゆきひろ 技術を斬る, 日経Linux, 日本, 日経BP社, 2009年10月 8日, 第11巻, 第11号, pp. 150 - 156
西澤無我, ROMAとそのプラグイン機構, Software Design, 日本, (株)技術評論社, 2010年 2月18日, No. 232, pp. 36 - 41
前橋孝広, 大規模データ処理のための分散システムの実装とその応用, 電子情報通信学会論文誌, 日本, 社団法人電子情報通信学会, 2010年 7月 1日, 第J93-D巻, 第7号, pp. 1072 - 1081

(58)調査した分野(Int.Cl., DB名)

G 0 6 F 1 2 / 0 0
G 0 6 F 1 7 / 3 0