

(19) 日本国特許庁 (JP)

## (12) 特 許 公 報 (B2)

(11) 特許番号

特許第5765416号  
(P5765416)

(45) 発行日 平成27年8月19日 (2015. 8. 19)

(24) 登録日 平成27年6月26日 (2015. 6. 26)

(51) Int. Cl.

F I

G 0 6 F 12/00 (2006. 01)

G 0 6 F 3/06 (2006. 01)

G 0 6 F 12/00 5 3 3 J

G 0 6 F 12/00 5 1 4 E

G 0 6 F 12/00 5 4 5 A

G 0 6 F 3/06 3 0 4 F

請求項の数 28 (全 46 頁)

(21) 出願番号 特願2013-503590 (P2013-503590)  
 (86) (22) 出願日 平成24年3月8日 (2012. 3. 8)  
 (86) 国際出願番号 PCT/JP2012/055917  
 (87) 国際公開番号 W02012/121316  
 (87) 国際公開日 平成24年9月13日 (2012. 9. 13)  
 審査請求日 平成25年9月6日 (2013. 9. 6)  
 (31) 優先権主張番号 特願2011-50151 (P2011-50151)  
 (32) 優先日 平成23年3月8日 (2011. 3. 8)  
 (33) 優先権主張国 日本国 (JP)

(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 100080816  
 弁理士 加藤 朝道  
 (72) 発明者 菅 真樹  
 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 鳥居 隆史  
 東京都港区芝五丁目7番1号 日本電気株  
 式会社内

審査官 池田 聡史

最終頁に続く

(54) 【発明の名称】 分散ストレージシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

それぞれがデータ格納部を備え、ネットワーク結合される複数のデータノードを備え、  
 データ複製先のデータノードが、前記データノード間で、論理的には同一であるが、物理  
 的には異なるデータ構造をそれぞれの前記データ格納部に保持する、少なくとも二つのデ  
 ータノードを含み、前記複製先のデータノードは、目的のデータ構造への変換を複製デー  
 タの受付とは非同期で行い、

データの識別子に対応して、前記データの格納先のデータノードと、データ構造の種類  
 を特定するデータ構造管理情報を記憶管理する構造情報管理手段を備え、

前記データノードは、

アクセス要求に基づき、データの更新処理を行う場合に、受け付けたデータを、一旦、  
 中間データ保持構造に保持して更新に対する応答を返すアクセス手段と、

設定された更新契機に応答して、前記中間データ保持構造に保持されるデータを、前記  
 データ構造管理情報で指定されたデータ構造に非同期で変換して前記データ格納部に格納  
 するデータ構造変換手段と、

を備えた、分散ストレージシステム。

【請求項 2】

前記構造情報管理手段は、前記中間データ保持構造に保持されるデータが目的のデータ  
 構造に変換して格納されるまでの時間情報である更新契機情報を、前記データ構造の種類  
 に対応させて保持する、請求項 1 記載の分散ストレージシステム。

## 【請求項 3】

前記構造情報管理手段は、前記データの識別子に対応して、1つ又は複数のデータ配置先のデータノード情報を特定するデータ配置特定情報をさらに記憶管理し、

前記更新処理に対応してアクセス先のデータノードを、前記データ構造管理情報と前記データ配置特定情報より選択するクライアント機能実現手段を備えた、請求項2記載の分散ストレージシステム。

## 【請求項 4】

予め定められたテーブル単位でデータの配置先のデータノード、配置先でのデータ構造、データ分割を可変に制御する手段を備えた請求項1記載の分散ストレージシステム。

## 【請求項 5】

データが配置されるデータノードを、コンシステント・ハッシングで求める手段を備えた、請求項1乃至4のいずれか1項に記載の分散ストレージシステム。

## 【請求項 6】

それぞれがデータ格納部を備え、ネットワーク結合される複数のデータノードを備え、データ複製先のデータノードが、前記データノード間で、論理的には同一であるが、物理的には異なるデータ構造をそれぞれの前記データ格納部に保持する、少なくとも二つのデータノードを含み、

格納対象のデータを識別する識別子であるテーブル識別子に対応させて、複製を特定するレプリカ識別子と、前記レプリカ識別子に対応したデータ構造の種類を特定するデータ構造情報と、指定されたデータ構造に変換して格納されるまでの時間情報である更新契機情報と、を、前記データ構造の種類の数に対応させて備えたデータ構造管理情報と、

前記テーブル識別子に対応して、前記レプリカ識別子と、前記レプリカ識別子に対応した1つ又は複数のデータ配置先のデータノード情報と、を備えたデータ配置特定情報と、

を記憶管理する構造情報保持部を有する構造情報管理装置と、

前記データ構造管理情報と前記データ配置特定情報とを参照して、更新処理及び参照処理のアクセス先を特定するデータアクセス部を備えたクライアント機能実現部と、

それぞれが前記データ格納部を備え、前記構造情報管理装置と前記クライアント機能実現部とに接続される複数の前記データノードと、

を備え、

前記データノードは、

前記クライアント機能実現部からのアクセス要求に基づき、更新処理を行う場合に、中間データ保持構造にデータを保持して前記クライアント機能実現部に応答を返すデータ管理・処理部と、

前記データ構造管理情報を参照し、指定された更新契機に応答して、前記中間データ保持構造に保持されるデータを、前記データ構造管理情報で指定されたデータ構造に変換する処理を行うデータ構造変換部と、

を備えることを特徴とする、分散ストレージシステム。

## 【請求項 7】

予め定められたテーブル単位でデータの配置先のデータノード、配置先でのデータ構造、データ分割を可変に制御する手段を備えた請求項6記載の分散ストレージシステム。

## 【請求項 8】

前記中間データ保持構造は、指定された目的のデータ構造としてデータが前記データ格納部に格納されるまでの間、前記データを保持する、請求項6記載の分散ストレージシステム。

## 【請求項 9】

前記クライアント機能実現部が、前記更新処理又は前記参照処理の内容に応じてアクセス先のデータノードを、前記データ構造管理情報と前記データ配置特定情報より選択する、請求項6記載の分散ストレージシステム。

## 【請求項 10】

前記クライアント機能実現部は、前記構造情報管理装置の前記構造情報保持部に保持さ

10

20

30

40

50

れている前記データ配置特定情報、又は、前記構造情報保持部に保持される情報をキャッシュする構造情報キャッシュ保持部に保持されているデータ配置特定情報を取得し、データ配置先のデータノードに対して、アクセス命令を発行する、請求項6記載の分散ストレージシステム。

【請求項11】

前記データノードは、アクセス受付部、アクセス処理部、データ構造変換部を備え、

前記データノードの前記データ格納部は、構造別データ格納部を備え、

前記アクセス受付部は、前記クライアント機能実現部からの更新要求を受け付け、前記データ配置特定情報においてレプリカ識別子に対応して指定されているデータノードに対して更新要求を転送し、

10

前記データノードの前記アクセス処理部は、受け取った更新要求の処理を行い、前記データ構造管理情報の情報を参照して更新処理を実行し、その際、前記データ構造管理情報の情報から、前記データノードに対する前記更新契機が零の場合、更新データを、前記データ構造管理情報に指定されるデータ構造に変換して前記構造別データ格納部を更新し、

前記更新契機が零でない場合、前記中間データ保持構造に、一旦、更新データを書き込み、処理完了を応答し、

前記アクセス受付部は、前記アクセス処理部からの完了通知と、レプリカ先のデータノードの完了通知を受けると、前記クライアント機能実現部に対して応答し、

前記データ構造変換部は、前記中間データ保持構造のデータを、前記データ構造管理情報に指定されているデータ構造に変換し変換先の前記構造別データ格納部に格納する、請求項6又は10記載の分散ストレージシステム。

20

【請求項12】

前記クライアント機能実現部は、参照系アクセスの場合、データノードに対して行われるデータアクセスに適しているデータ構造を選択し、レプリカ識別子を特定した後、アクセスすべきデータノードを算出し、選択されたデータノードに対してアクセス要求を発行し前記データノードからアクセス処理結果を受け取る、請求項6記載の分散ストレージシステム。

【請求項13】

前記クライアント機能実現部が、前記データノード内に配設されている、請求項6記載の分散ストレージシステム。

30

【請求項14】

前記クライアント機能実現部が、前記構造情報保持部に保持される情報をキャッシュする構造情報キャッシュ保持部を備えた請求項13記載の分散ストレージシステム。

【請求項15】

前記クライアント機能実現部の前記構造情報キャッシュ保持部の構造情報と、前記構造情報管理装置の前記構造情報保持部に保持される構造情報を同期させる構造情報同期部を備えた請求項14記載の分散ストレージシステム。

【請求項16】

前記データ構造管理情報が、データを複数のデータノードに分割して格納する分割数であるパーティション数をレプリカ識別子に対応して備え、

40

前記データ配置特定情報は、前記データ構造管理情報においてパーティション数が2以上に対応するレプリカ識別子に対応した配置ノードとして、複数のデータノードを含み、

アクセス要求を受けた前記データノードのアクセス受付部は、パーティショニングされたデータの配置先が複数のデータノードにまたがる場合に、前記複数のデータノードを構成する他のデータノードのアクセス処理部にアクセス要求を発行する、請求項6記載の分散ストレージシステム。

【請求項17】

アクセス要求を受けた前記データノードの前記データ構造変換部は、前記更新契機が零のとき、他のデータノードの前記データ構造変換部に対してアクセス要求を発行する、請求項6又は11記載の分散ストレージシステム。

50

## 【請求項 18】

アクセス要求の履歴を記録する履歴記録部と、  
前記履歴記録部の履歴情報を用いてデータ構造の変換を行うか否かを判定する変更判定部と、  
を備えた請求項 6 記載の分散ストレージシステム。

## 【請求項 19】

前記変更判定部は、データ構造の変換が必要と判定した場合、前記構造情報管理装置の構造情報変更部に変換要求を出力し、  
前記構造情報管理装置の前記構造情報変更部は、前記構造情報保持部の情報を変更し、  
前記データノードの前記データ構造変換部に変換要求を出力し、  
前記データノードの前記データ構造変換部は前記データノードの前記データ格納部に保持されるデータ構造の変換を行う、請求項 18 記載の分散ストレージシステム。

10

## 【請求項 20】

それぞれがデータ格納部を備え、ネットワーク結合される複数のデータノードを備えたシステムでの分散ストレージ方法であって、  
データ複製先のデータノードの少なくとも二つのデータノードが、前記データノード間で、論理的には同一であるが、物理的には異なるデータ構造をそれぞれの前記データ格納部に保持し、前記複製先のデータノードは、目的のデータ構造への変換を複製データの受付とは非同期で行い、

データの識別子に対応して、前記データの格納先のデータノードと、データ構造の種類を特定するデータ構造管理情報を構造情報管理手段で記憶管理し、

20

前記データノードは、

アクセス要求に基づき、データの更新処理を行う場合に、受け付けたデータを、一旦、中間データ保持構造に保持して、更新に対する応答を返し、

データの更新契機にตอบสนองして、前記中間データ保持構造に保持されるデータを、前記データ構造管理情報で指定されたデータ構造に非同期で変換し前記データ格納部に格納する、分散ストレージ方法。

## 【請求項 21】

前記構造情報管理手段は、前記中間データ保持構造に保持されるデータが目的のデータ構造に変換して格納されるまでの時間情報である更新契機情報を、前記データ構造の種類に対応させて保持する、請求項 20 記載の分散ストレージ方法。

30

## 【請求項 22】

前記構造情報管理手段では、前記データの識別子に対応して、1つ又は複数のデータ配置先のデータノード情報を特定するデータ配置特定情報をさらに記憶管理し、

前記更新処理に対応してアクセス先のデータノードを、前記データ構造管理情報と前記データ配置特定情報より選択する、請求項 21 記載の分散ストレージ方法。

## 【請求項 23】

予め定められたテーブル単位でデータの配置先のデータノード、配置先でのデータ構造、データ分割を可変に制御する、請求項 20 記載の分散ストレージ方法。

## 【請求項 24】

データが配置されるデータノードをコンシステント・ハッシングで求める、請求項 20 乃至 23 のいずれか 1 項に記載の分散ストレージ方法。

40

## 【請求項 25】

それぞれがデータ格納部を備え、ネットワーク結合される複数のデータノードを備えたシステムでの分散ストレージ方法であって、

データ複製先のデータノードの少なくとも二つのデータノードが、前記データノード間で、論理的には同一であるが、物理的には異なるデータ構造をそれぞれの前記データ格納部に保持し、

データ更新時に行われるデータの複製において、前記複製先のデータノードでは、更新対象のデータを、それぞれ、指定された目的のデータ構造に変換して前記データ格納部に

50

格納し、その際、前記データノードは、更新対象のデータを一旦、中間構造を保持して前記更新に対する応答を返し、更新要求とは非同期で、目的のデータ構造に変換して格納し

、  
格納対象のデータを識別する識別子であるテーブル識別子に対応させて、複製を特定するレプリカ識別子と、前記レプリカ識別子に対応したデータ構造の種類を特定するデータ構造情報と、指定されたデータ構造に変換して格納されるまでの時間情報である更新契機情報と、を前記データ構造の種類の数に対応させて備えたデータ構造管理情報と、

前記テーブル識別子に対応して、前記レプリカ識別子と、前記レプリカ識別子に対応した1つ又は複数のデータ配置先のデータノード情報と、を備えたデータ配置特定情報と、を含む構造情報を構造情報管理部で記憶管理し、

クライアント側では、前記データ構造管理情報と前記データ配置特定情報を参照して、更新処理及び参照処理のアクセス先を特定し、

前記データノードは、

前記クライアント側からのアクセス要求に基き、更新処理を行う場合に、中間データ保持構造にデータを保持して前記クライアントに応答を返し、

前記データ構造管理情報を参照し、指定された更新契機に応じて、前記中間データ保持構造から指定されたデータ構造に変換する、

ことを特徴とする分散ストレージ方法。

【請求項26】

前記データ構造管理情報が、データを複数のデータノードに分割して格納する分割数であるパーティション数を、レプリカ識別子に対応して備え、

前記データ配置特定情報は、前記データ構造管理情報においてパーティション数が2以上に対応するレプリカ識別子に対応した配置ノードとして、複数のデータノードを含み、

アクセス要求を受けた前記データノードでは、パーティショニングされたデータの配置先が複数のデータノードにまたがる場合に、前記複数のデータノードを構成する他のデータノードに対してアクセス要求を発行する、請求項25記載の分散ストレージ方法。

【請求項27】

アクセス要求の履歴を記録する履歴記録部での履歴情報を用いて、データ構造の変換を行うか否かを判定し、変換が必要な場合、前記構造情報を変換し、さらに前記データノードのデータ構造を変換する、請求項25記載の分散ストレージ方法。

【請求項28】

データ格納部を備え、他のデータノードとネットワーク結合され、複数のデータノードが分散ストレージシステムを構成し、

更新対象のデータを複数のデータノードに複製する場合、前記データに関して、少なくとも一つの他のデータノードとの間で、論理的には同一であるが、物理的には異なるデータ構造を前記データ格納部に保持し、目的のデータ構造への変換を複製データの受付とは非同期で行うデータノード装置であって、

アクセス要求に基づき、データの更新処理を行う場合に、受け付けたデータを、一旦、中間データ保持構造に保持して更新に対する応答を返すアクセス手段と、

データの識別子に対応して、前記データの格納先のデータノードと、データ構造の種類を特定し、前記データの更新契機を設定するデータ構造管理情報を記憶管理する構造情報管理装置からの前記データ構造管理情報に基づき、設定された更新契機に応答して、前記中間データ保持構造に保持されるデータを、前記データ構造管理情報で指定されたデータ構造に非同期で変換して前記データ格納部に格納するデータ構造変換手段と、

を備えたデータノード装置。

【発明の詳細な説明】

【技術分野】

【0001】

(関連出願についての記載)

本発明は、日本国特許出願：特願2011-050151号(2011年3月8日出願

10

20

30

40

50

）の優先権主張に基づくものであり、同出願の全記載内容は引用をもって本書に組み込み記載されているものとする。

【0002】

本発明は、分散ストレージに関し、特に、データ構造の制御が可能な分散ストレージシステム、および方法と装置に関する。

【背景技術】

【0003】

<分散ストレージシステム>

複数の計算機（データノード、あるいは単に「ノード」ともいう）をネットワーク結合し、各計算機のデータ格納部（HDD（Hard Disk Drive）やメモリ等）にデータを格納して利用するシステムを実現する分散ストレージシステム（Distributed Storage System）が利用されている。

10

【0004】

一般的な分散ストレージ技術では、

- ・データをどの計算機（ノード）に配置するか、
- ・処理をどの計算機（ノード）で行うか、

といった判断をソフトウェアや特別な専用ハードウェアにより実現し、システムの状態に対してその動作を動的に変更することでシステム内のリソース使用量を調整し、システム利用者（クライアント計算機）に対する性能を向上している。

20

【0005】

分散ストレージシステムにおいては、データが複数のノードに分散しているため、データにアクセスしようとするクライアントは、まず、データを持っているノードがどれであるかを知る必要がある。また当該データを持つノードが複数ある場合、どのノード（一つ以上）にアクセスするかを知る必要がある。

【0006】

分散ストレージシステムでは、一般に、ファイル管理として、ファイル本体と、当該ファイルのメタデータ（ファイルの格納場所、ファイルサイズ、オーナー等）を別々に保存する方式が用いられている。

【0007】

<メタサーバ方式>

30

分散ストレージシステムにおいて、クライアントがデータを保持しているノードを知るための技術の一つとしてメタサーバ方式が知られている。メタサーバ方式では、データの位置情報を管理する、一つ又は複数（ただし、少ない数）の計算機により構成されたメタサーバを設ける。

【0008】

しかしながら、メタサーバ方式では、分散ストレージシステムの構成の大規模化に伴って、データを格納しているノードの位置を検出する処理を行うメタサーバの処理性能が足りず（メタサーバ1台当りで管理するノード数が膨大となり、該メタサーバの処理性能が追いつかない）、導入したメタサーバが、かえってアクセス性能上のボトルネックになる、という問題がある。

40

【0009】

<分散KVS>

データを保持しているノードの位置を知るための別の手法（技術）として、分散関数（例えばハッシュ関数）を用いてデータの位置を求めるものがある。この種の手法は、例えば分散KVS（Key Value Store：キー・バリュー・ストア）と呼ばれている。

【0010】

分散KVSでは、全てのクライアントで、分散関数と、システムに参加しているノードのリスト（ノードリスト）とを共有する。

【0011】

50

また、格納データは、固定長あるいは任意長のデータ断片 ( V a l u e ) に分かれている。各データ断片 ( V a l u e ) 毎に一意に特定可能な識別子 ( K e y ) が付与され、 ( K e y 、 V a l u e ) のペアで保存される。例えばキーの値に応じて保存先のノード ( サーバ ) を変えることで、複数のノードにデータを分散保存することが可能となる。

【 0 0 1 2 】

各クライアントは、データにアクセスする際、キーを分散関数の入力値とし、分散関数の出力値とノードリストとを基に、データを格納しているノードの位置を算術的に求める。

【 0 0 1 3 】

クライアント間で共有する情報のうち、分散関数は、基本的に、時間が経過しても変化しない ( 時不変 ) 。一方、ノードリストの内容は、ノードの故障や追加に伴い、随時、変更される。このため、クライアントは、それらの情報に対して任意の方法でアクセス出来る必要がある。

【 0 0 1 4 】

< レプリケーション >

分散ストレージシステムにおいては、可用性 ( A v a i l a b i l i t y : システムが連続して動作できる能力 ) 確保のために、データの複製を複数ノードで保持し、データの複製を、負荷分散に活用することが一般的に行われている。

【 0 0 1 5 】

なお、作成するデータの複製を用いて負荷分散を実現する技術が特許文献 1 に記載されている。

【 0 0 1 6 】

本件に関して行われた先行文献サーチの結果サーチされた特許文献 2 には、サーバが情報構造定義部で情報構造定義体を定義し、登録用クライアントは情報構造定義体によりデータベースを構築し、データベースアクセスツールを生成し、このツールを用いてデータベースに情報を登録する構成が開示されている。また特許文献 3 には、分散型ストレージシステムにおいて、各複製がそれぞれ固有のロケータ値を介してアクセス可能なオブジェクトの複製を保存するストレージノードと、各オブジェクトに対するそれぞれのキーマップエントリを保存するキーマップインスタンスを含み、所定のオブジェクトについてはそれぞれのキーマップエントリは、オブジェクトの複製と、対応するキー値、各ロケータを含む構成が開示されている。

【 先行技術文献 】

【 特許文献 】

【 0 0 1 7 】

【 特許文献 1 】 特開 2 0 0 6 - 1 2 0 0 5 号公報 ( 特許第 4 5 2 8 0 3 9 号 )

【 特許文献 2 】 特開平 1 1 - 1 9 5 0 4 4 号公報 ( 特許第 3 9 1 1 8 1 0 号 )

【 特許文献 3 】 特表 2 0 0 9 - 5 2 2 6 5 9 号公報

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 1 8 】

以下に関連技術の分析を与える。

【 0 0 1 9 】

関連技術の分散ストレージシステムでは、可用性保持のため、複製データを複数のノードで同一の物理構造で保持している。これにより、アクセス応答性能と可用性保証を実現している。しかしながら、複製データを同一の物理構造で保持しているため、データの利用形態の特性が異なるアプリケーション等に対しては、別のデータ構造への変換、及び別のデータ構造を保持するためのストレージを用意しなければならない。

【 0 0 2 0 】

したがって、本発明の目的は、分散ストレージにおけるデータ複製において、可用性を確保するとともに、ストレージの利用効率の低下の回避、応答性能の低下の回避の少なく

10

20

30

40

50

とも１つを可能とする、分散ストレージシステムと方法を提供することにある。

【課題を解決するための手段】

【００２１】

本発明によれば、上記課題の少なくとも１つの解決を図るため、特に制限されるものではないが、概略以下の構成とされる。

【００２２】

本発明によれば、それぞれがデータ格納部を備え、ネットワーク結合される複数のデータノードを備え、データ複製先のデータノードが、前記データノード間で、論理的には同一であるが、物理的には異なるデータ構造をそれぞれの前記データ格納部に保持する、少なくとも二つのデータノードを含む分散ストレージシステムが提供される。本発明によれば、分散ストレージシステムを構成するデータノード装置として、他のデータノードとネットワーク結合され、更新対象のデータを複数のデータノードに複製する場合、前記データに関して、少なくとも一つの他のデータノードとの間で、論理的には同一であるが、物理的には異なるデータ構造を前記データ格納部に保持するデータノード装置が提供される。

10

【００２３】

本発明によれば、それぞれがデータ格納部を備え、ネットワーク結合される複数のデータノードを備えたシステムでの分散ストレージ方法であって、前記複数のデータノードの少なくとも二つのデータノードは、前記データノード間で論理的には同一であるが、物理的には異なる複数種のデータ構造の複製をそれぞれの前記データ格納部に保持する、分散ストレージ方法が提供される。

20

【００２４】

いくつかの実施形態によれば、前記複数のデータノードにおいて、ターゲットのデータ構造への変換をデータ更新要求とは非同期で行うようにしてもよい。あるいは、いくつかの実施形態によれば、前記データノードにおいて中間データ保持構造に受信データを保持し、前記更新要求の応答を返し、前記中間データ保持構造に保持されるデータ構造をターゲットのデータ構造に非同期に変換する。あるいは、いくつかの実施形態によれば、予め定められたテーブル単位でデータ配置先、配置先のデータ構造、データ分割を可変に制御する。

【発明の効果】

30

【００２５】

本発明によれば、分散ストレージにおけるデータ複製において、可用性を確保するとともに、ストレージの利用効率の低下の回避、応答性能の低下を回避の少なくとも１つ可能としている。

【図面の簡単な説明】

【００２６】

【図１】本発明の第１の例示的な実施の形態のシステム構成を示す図である。

【図２】本発明の第１の例示的な実施の形態を説明する図である。

【図３】本発明の第１の例示的な実施形態を説明する図である。

【図４】本発明の第１の例示的な実施形態のデータノードの構成例を示す図である。

40

【図５】本発明の第１の例示的な実施形態におけるデータ構造管理情報 ９ ２ １ を模式的に示す図である。

【図６】本発明の第１の例示的な実施形態におけるテーブルのデータ保持構造を説明する図である。

【図７】本発明の第１の例示的な実施形態におけるテーブルのデータ保持、非同期更新を模式的に説明する図である。

【図８】本発明の第１の例示的な実施形態におけるデータ配置特定情報 ９ ２ ２ の例を示す図である。

【図９】本発明の第１の例示的な実施形態における Write 処理の動作シーケンスを説明する図（１）である。

50



【図 10】本発明の第 1 の例示的な実施形態における W r i t e 処理の動作シーケンスを説明する図 ( 2 ) である。

【図 11】本発明の第 1 の例示的な実施形態における R E A D 系処理の動作シーケンスを説明する図である。

【図 12】本発明の第 1 の例示的な実施形態におけるクライアント実現手段 6 1 におけるアクセス処理の動作を説明するフローチャートである。

【図 13】本発明の第 1 の例示的な実施形態におけるデータノードにおけるアクセス処理の動作を説明するフローチャートである。

【図 14】本発明の第 1 の例示的な実施形態におけるデータ変換処理を説明するフローチャートである。

10

【図 15】本発明の第 2 の例示的な実施形態のデータ構造管理情報 9 2 1 を模式的に示す図である。

【図 16】本発明の第 2 の例示的な実施形態におけるデータ配置特定情報 9 2 2 の例を示す図である。

【図 17】本発明の第 2 の例示的な実施形態のデータノードの構成例を示す図である。

【図 18】本発明の第 3 の例示的な実施形態のデータノードの構成例を示す図である。

【図 19】本発明の第 3 の例示的な実施形態の全体の制御フローを説明するフローチャートである。

【図 20】本発明の第 3 の例示的な実施形態のデータ構造変換処理を説明するフローチャートである。

20

【図 21】本発明の第 3 の例示的な実施形態の変換処理を説明する図である。

【図 22】本発明の第 3 の例示的な実施形態のパーティショニング数の変更処理を説明するフローチャートである。

【図 23】本発明の第 3 の例示的な実施形態のパーティショニング数変更時の動作を説明するフローチャートである。

【図 24】本発明の第 3 の例示的な実施形態における分散テーブルのデータ配置を説明する図である。

【図 25】本発明の第 3 の例示的な実施形態における構造情報保持部 9 2 を説明する図である。

【図 26】本発明の第 4 の例示的な実施形態のコンシステント・ハッシング分割配置を説明する図である。

30

【図 27】本発明の第 4 の例示的な実施形態の情報記録形態を説明する図である。

【図 28】本発明の第 4 の例示的な実施形態におけるカラムベースのコンシステント・ハッシング分割配置を説明する図である。

【図 29】本発明の第 4 の例示的な実施形態において 1 カラムをパーティショニングした場合のコンシステント・ハッシング分割配置を説明する図である。

【発明を実施するための形態】

【0027】

本発明の好ましい態様 ( P r e f e r r e d M o d e s ) の一つによれば、複数種類のデータ構造を持ち、データ配置ノード ( 「データノード」という ) 間で論理的には同一であるが、物理的には異なる構造の複製 ( レプリカ ) を保持する。本発明においては、書き込み ( 更新 ) 要求とは非同期に行われるデータ構造変換の適用契機を制御可能としている。本発明においては、W r i t e の応答特性を優先した構造を中間構造 ( 中間データ保持構造 ) を備え、該中間構造に保持されるデータ構造をターゲットとなるデータ構造に非同期に変換する。

40

【0028】

本発明の好ましい態様においては、制御パラメータを変更するインタフェースを持つ。アクセス負荷に応じて制御パラメータを変更する。あるいは、処理負荷が増えたら、パーティショニング粒度を小さくする等の制御が行われる。

【0029】

50

本発明の好ましい態様によれば、複数種類のデータ構造を持つことが可能となるキー・バリュー・ストア (Key Value Store) を実現可能としている。本発明の好ましい態様によれば、論理的には同一内容であるが、物理的には異なるデータ構造の複製 (レプリカ) を持つ。この結果、

- ・異なる種類のアクセス負荷に対して高速に対応可能とし、
- ・可用性保持のための複製 (レプリカ) を他の用途に利用可能とし、データ容量の効率利用を可能としている。

#### 【0030】

本発明の好ましい態様において、データ送信元から該データを受け取る側のデータノードでは、受信データを複製に同期して直ちにターゲット構造に変換する代りに、中間構造形式で保持し、ターゲット構造への変換を非同期で行うようにしてもよい。例えば Write 要求に対してデータをバッファに保持して直ちに応答を返す等、アクセス要求に対する応答特性を優先した中間構造を備え、中間構造に保持されたデータ構造を、非同期でターゲット構造を変換することにより、データ構造の変換処理によって生じる、アクセス性能上のボトルネックを回避しながら、要求される高可用性の維持を可能としている。分散ストレージシステムの複数のデータノード上で複数種類のデータ構造への更新、変換を同時に行うことは、性能上、ボトルネックとなりやすい。本発明の好ましい態様においては、Write に特化した構造 (Write の応答性能を優先した中間データ保持構造) を用意し、可用性保証のための複製実行時には、同期式 (Sync) で中間構造で複製し、該中間構造で保持されるデータを非同期 (Async) で正式のターゲット構造に変換する。

#### 【0031】

さらに、本発明の好ましい態様によれば、データノードやデータ構造、非同期に構造変換を実行するための契機 (トリガー) を制御可能にすることで、様々なアプリケーションや負荷変動に対応可能としている。

#### 【0032】

本発明の好ましい態様によれば、特に制限されるものではないが、例えば、テーブル単位で、データ配置、データ構造、パーティショニング (分割) をコントロール可能としている。

#### 【0033】

データ構造として、例えば、

ロウストア (Row-store) :

- ・追記型 (データの格納領域に記録を追加)、
- ・更新型、

カラムストア (Column-store) :

- ・圧縮の有無、

ライトログ (例えばライト性能を優先するために更新情報を追記するための構造) :

インデックス (検索用の索引データ) の有無 :

データの格納順をソート (Sorting) しているか :

分割 (Partitioning) の有 / 無、分割数 :

分割 (Partitioning) 単位、アルゴリズム :

等の項目について組み合わせが選択される。

#### 【0034】

本発明の好ましい態様によれば、例えばデータをどのデータノードに置くかが制御の対象となるほか、どのデータ構造とするかも制御対象となる。

#### 【0035】

・例えば Write 要求のみが行われる場合、ライトログ (Write Log) や、追記型ロウテーブル (Row-table) とする。

#### 【0036】

- ・あるいは、Read と Write の組み合わせに対して、ロウテーブル (Row-table)

b 1 e ) が選択される。

【 0 0 3 7 】

・さらに、分析アプリケーションに対して、例えばカラムストア（あるいはカラム指向データベース）を選択する。カラムストア方式は、クエリ（Query）に対してストレージのリードアクセスを効率化する。

【 0 0 3 8 】

・あるいは、分散処理に対して、パーティショニング（データ分割）の粒度を相対的に小さくし、集中処理に対してパーティショニングを大きくするか、パーティショニングを止める等の制御を行うようにしてもよい。

【 0 0 3 9 】

・さらに、中間構造で保持されるデータを非同期（Async）でターゲット構造へ変換するためのトリガー（契機）を制御対象としてもよい。

【 0 0 4 0 】

・あるいは、分析アプリケーションが必要とするデータ鮮度（データの新しさの尺度）によって、データ変換の契機を調整するようにしてもよい。

【 0 0 4 1 】

本発明の好ましい態様によれば、それぞれがデータ格納部（図1の12）を備え、ネットワーク結合される複数のデータ配置ノード（データノード）を備えた分散ストレージシステムにおいて、クライアントからのデータ更新要求時等に行われる複製において、複製先の1つ又は複数のデータノードでは、更新要求を受けたデータベースにおけるデータ構造とは、異なる1つ又は複数種のデータ構造で複製データを前記データ格納部（図1の12）に格納する。その際、前記データノードは、複製データを一旦中間構造を保持して更新要求に対する応答をクライアントに返し、前記更新要求とは、非同期で目的のデータ構造に変換して格納する。

【 0 0 4 2 】

本発明の態様の1つによれば、データ構造情報（例えばデータ構造の管理情報やデータ配置特定情報）を保持管理する装置（図1の9）を備え、データアクセス手段（図1の611）およびデータノードをアクセスする手段（図4の112）は、データ構造情報を用いて、複製対象のデータに対するデータ構造（物理構造）を決定する。このため、分散ストレージノード毎に、複製データを異なるデータ構造で保持することができる。

【 0 0 4 3 】

本発明の態様の1つによれば、分散ストレージシステムにおいて、複製先のデータノードは、クライアントからの更新要求に対して、更新処理性能を優先する中間構造（中間データ保持構造、中間バッファ構造ともいう）に、データを一旦保持して、該更新要求に対して応答し、データ構造管理情報で指定されるデータ構造への変換を非同期で実行する。このため、複数種のデータ構造をそれぞれ中間データ保持構造に保持しつつ、更新処理の応答性能を維持できる。

【 0 0 4 4 】

本発明の態様の1つによれば、複数種のデータ構造を持ち、クライアント側が、アクセス内容に応じて適切なデータ構造に処理の振分け（適切なデータ構造を保持するデータノードをアクセスするように振り分ける）を行うようにしてもよい。このため、アクセス処理性能を向上することができる。

【 0 0 4 5 】

前記した関連技術を、上記本発明の態様の観点から分析する。

【 0 0 4 6 】

前述したとおり、関連技術の分散ストレージシステムにおいては、可用性保持のため、複製データを複数のノードで同一の物理構造で保持している。このため、可用性保持のための複製データの保持格納形式を制御することができない。

【 0 0 4 7 】

例えば、

- ・データの配置場所、
- ・データ配置（内部）構造、
- ・データを分散して格納するか、集中的に格納するかという格納方式、

等の複製データの保持格納形式について、可変に制御することができない。

**【 0 0 4 8 】**

データ移行等において、データ移行元のストレージ/データベースと、移行先のストレージ/データベースとは、同一データを異なるデータ構造で表現したものである、ということができる。例えば複製データを同一のデータ構造（物理構造）にて複数ノードで保持する構成において、互いに異なるデータ構造の各々について、各ノードで複製を保持する場合、ストレージ容量が過剰に必要とされる（この場合の複製に必要なストレージ容量は、データ容量×複製数×データ構造の種類の数）。そのため、計算機やディスク等のハードウェアを多く用意して利用することによって、購入コストや消費電力等の運用コストが増大してする（大量のデータコピー、大量のデータ構造の変換処理が必要とされる）。

10

**【 0 0 4 9 】**

また、関連技術において、分散ストレージシステムを利用するユーザ（例えばアプリケーション開発者）が、実現したいアプリケーション・ロジックを踏まえた上で、

- ・適切なデータ構造の選択、
- ・適切なスキーマの設計、
- ・適切なデータベースソフトウェア、設定の使い分け

を行う必要がある、ということである。いずれも、データベースシステムおよびストレージシステムに対して高い知見がユーザに要求されることから、これらをユーザ側で行うことは、實際上、困難である。

20

**【 0 0 5 0 】**

また、複製にあたり、適切なデータ構造を選択した場合であっても、複数のデータベースシステムを用意し、データの移行を行う必要がある、ということである。これらの処理は、計算機（サーバ）等において、データの入出力等の負荷が大きい。このため、移行先のデータベースのデータは、移行元のデータベースより古いデータとならざるを得ない。また、前述したように、同一内容のデータを互いに異なる複数のデータ構造として保持する場合、ストレージ利用効率が悪化してしまう。

**【 0 0 5 1 】**

本発明の態様の1つによれば、複製データを複数種のデータ構造（物理構造）で保持することで、要求される高可用性と、高速応答等性能を確保しつつ、データ構造変換のボトルネックを解消し、ストレージの利用効率を高めることが出来る。

30

**【 0 0 5 2 】**

以下、添付図面を参照して、いくつかの例示的な実施形態について説明する。

**【 0 0 5 3 】****< 実施形態 1 >**

本発明の第1の例示的な実施形態について図面を参照して説明する。図1は、本発明の第1の実施形態のシステム構成の一例を示す図である。データノード1～4、ネットワーク5、クライアントノード6、構造情報管理手段（構造情報管理装置）9を備える。

40

**【 0 0 5 4 】**

データノード1～4は、分散ストレージを構成するデータ格納ノードであり、1つ以上の任意の数によって構成される。ネットワーク5は、データノード1～4を含むネットワークノード間の通信を実現する。クライアントノード6は、分散ストレージにアクセスする計算機ノードである。クライアントノード6は必ずしも独立して存在しなくてもよい。なお、データノード1～4がクライアント計算機を兼ねる例は、図2を参照して後述される。

**【 0 0 5 5 】**

データノード1～4は、それぞれ、データ管理・処理手段（データ管理・処理部）11、21、31、41、データ格納部12、22、32、42を備える。

50

## 【 0 0 5 6 】

データ管理・処理手段 X 1 ( X = 1、2、3、4 ) は、分散ストレージに対するアクセス要求を受け付け、処理を実行する。

## 【 0 0 5 7 】

データ格納部 X 2 ( X = 1、2、3、4 ) はデータノードの担当するデータの保持、記録を行う。

## 【 0 0 5 8 】

クライアントノード 6 は、クライアント機能実現手段 ( クライアント機能実現部 ) 6 1 を備える。

## 【 0 0 5 9 】

クライアント機能実現手段 6 1 は、データノード 1 ~ 4 によって構成される分散ストレージにアクセスする。

## 【 0 0 6 0 】

クライアント機能実現手段 6 1 はデータアクセス手段 ( データアクセス部 ) 6 1 1 を備える。

## 【 0 0 6 1 】

データアクセス手段 ( データアクセス部 ) 6 1 1 は、構造情報管理手段 9 から構造情報 ( データ構造管理情報とデータ配置特定情報 ) を取得し、その構造情報を用いて、アクセス先のデータノードを特定する。

## 【 0 0 6 2 】

なお、各データノード 1 ~ 4 やネットワーク 5 内の任意の装置 ( スイッチ、中間ノード ) において、構造情報管理手段 9 の構造情報保持部 9 2 に格納される構造情報の一部又は全てを自装置内又は他の装置内のキャッシュ ( 不図示 ) に保持するようにしてもよい。

## 【 0 0 6 3 】

すなわち、以下の実施形態の動作の説明において、構造情報保持部 9 2 に格納される構造情報に対するアクセスは、自装置内又は予め定められた所定の場所に配設されたキャッシュに対してアクセスするようにしてもよい。キャッシュに格納された構造情報の同期については、公知の分散システムの技術が適用できるため、ここでは詳細は省略する。よく知られているように、キャッシュを利用することでストレージ性能を高速化することが出来る。

## 【 0 0 6 4 】

構造情報管理手段 ( 構造情報管理装置 ) 9 は、構造情報を変更する構造情報変更手段 9 1 と、構造情報を保持する構造情報保持部 9 2 を備える。構造情報保持部 9 2 は、データ構造管理情報 9 2 1 ( 図 4 参照 ) とデータ配置特定情報 9 2 2 を含む ( 図 4 参照 ) 。データ構造管理情報 9 2 1 は、後に図 5 を参照して説明されるが、テーブル識別子に対して、複製を特定するレプリカ識別子と、前記レプリカ識別子に対応したデータ構造の種類を特定するデータ構造情報と、指定されたデータ構造として格納されるまでの時間情報である更新契機からなるエントリをデータの複製数分有する。データ配置特定情報 9 2 2 は、後に図 8 を参照して説明されるが、テーブル識別子に対応して、前記レプリカ識別子と、前記レプリカ識別子に対応した 1 つ又は複数のデータ配置先のデータノード情報を有する。

## 【 0 0 6 5 】

本実施形態において、クライアントノード 6 は、データノード 1 ~ 4 とは独立に ( 別々に ) 設けることは必ずしも必要とされない。つまり、以下、変形例として説明するように、データノード 1 ~ 4 のうち、任意の 1 つ以上のノードに、クライアント機能実現手段 6 1 を備えた構成としてもよい。

## 【 0 0 6 6 】

< 実施形態 1 の変形例 >

図 2 は、本発明の第 1 の実施形態の変形例の構成を示す図である。図 2 に示す通り、データノード 1、2、3、4 の各々に、クライアント機能実現手段 6 1 が配設されている。

## 【 0 0 6 7 】

図2を参照すると、データノード1、2、3、4に配設されるクライアント機能実現手段61は、図1のデータアクセス手段611の他に、構造情報キャッシュ保持部612を備える。

【0068】

構造情報キャッシュ保持部612は、構造情報保持部92に格納される構造情報の一部又は全てを格納するキャッシュメモリである。

【0069】

構造情報同期手段(構造情報同期装置)93は、構造情報のキャッシュの同期を制御する。構造情報保持部92のデータを取得し、データノードのクライアント機能実現手段61の構造情報キャッシュ保持部612の情報を更新する。

10

【0070】

構造情報同期手段93は、システムを構成する任意の機器に、任意の数、具備するようにしてもよい。例えば、各データノード1~4の少なくとも1つを実現する計算機上でソフトウェアとして動作させるようにしてもよい。

【0071】

図2において、データノード1~4をそれぞれ個別の計算機として実現した場合の例を図3に示す。図3の例では、1つ以上の任意の数のデータノード計算機101~104と、ネットワーク105から構成される。

【0072】

データノード計算機101~104は、それぞれCPU101a、データ記憶装置101b、データ転送装置101cを備える。CPU101aにより、データ管理・処理手段11、クライアント機能実現手段61の機能の全て又は一部を実現する。

20

【0073】

データ記憶装置101bは、例えば、ハードディスクドライブ、フラッシュメモリ、DRAM(Dynamic Random Access Memory)、MRAM(Magnetoresistive Random Access Memory)、FeRAM(Ferroelectric Random Access Memory)、PRAM(Phase change RAM)、RAIDコントローラに結合された記憶装置、磁気テープのようにデータを記録可能な物理媒体、又は、ストレージノードの外部に設置された媒体にデータを記録する制御装置である。ネットワーク105及びデータ転送装置101cは、例えばEthernet(登録商標)、Fibre ChannelやFCoE(Fibre Channel over Ethernet(登録商標))、InfiniBand(Intel社その他による団体が推進する高速IOバスアーキテクチャ)、QsNet(Quadrics社製品)、Myrinet(Myricom社製品)、Ethernet(登録商標)、あるいはこれらを利用するTCP/IP(Transmission/Control Protocol/Internet Protocol)やRDMA(Remote Direct Memory Access)のような上位プロトコルによって実現しうる。ただし、ネットワーク105の実現方法は、これらに限られない。Ethernet(登録商標)で実現する場合の例としては、データ転送装置101cは計算機に接続されるネットワークカード、ネットワーク105はEthernet(登録商標)ケーブルおよびスイッチ等から構成される。

30

40

【0074】

データノード1~4の実現は、仮想化された計算機(Virtual Machine)であってもよい。代表的な例としてVMWare(VMWare社製品)、Xen(Citrix社商標)等がある。

【0075】

<データノードの詳細の一例>

図4は、本発明の第1の実施形態の構成例をより詳細に説明する図である。図4には、図1のデータノード1~4を中心に示した構成が示されている。なお、図4等の図面において、簡単化のため、構造情報保持部92に格納される構造情報は参照符号92で参照さ

50

れる場合がある。

【 0 0 7 6 】

データノードのデータ管理・処理手段 1 1 は、アクセス受付手段（アクセス受付部） 1 1 1、アクセス処理手段（アクセス処理部） 1 1 2、データ構造変換手段（データ構造変換部） 1 1 3 を備えている。

【 0 0 7 7 】

アクセス受付手段 1 1 1 は、データアクセス手段 6 1 1 からアクセス要求を受け付け、処理完了後にデータアクセス手段 6 1 1 に応答を返す。

【 0 0 7 8 】

アクセス処理手段 1 1 2 は、構造情報保持部 9 2 の構造情報（あるいはその任意の場所に保持されるキャッシュ情報）を用い、アクセス処理を、該当するデータ格納部 1 2 X（X = 1、2、3）に対して処理を行う。

【 0 0 7 9 】

データ構造変換手段 1 1 3 は、一定契機毎に構造別データ格納部 1 2 1 のデータを用いて、構造別データ格納部 1 2 X（X = 1、2、3）に変換する。

【 0 0 8 0 】

データ格納部 1 2 は、複数種の構造別データ格納部を備えている。図 4 では、構造別データ格納部 1 2 1（データ構造 A）、構造別データ格納部 1 2 2（データ構造 B）、構造別データ格納部 1 2 3（データ構造 C）を備える。

【 0 0 8 1 】

どのようなデータ構造を選択するかは、構造別データ格納部 1 2 X（X = 1、2、3）単位で任意である。

【 0 0 8 2 】

本実施形態では、構造別データ格納部 1 2 1（例えばデータ構造 A）は、データの書き込みを伴う処理（データの追加や更新）に対する応答性能に特化した構造をとる。具体的には、データ変更内容をキュー（例えば F I F O（F i r s t I n F i r s t O u t））として高速なメモリ（デュアルポート R A M 等）上に保持するソフトウェア、アクセス要求処理内容を任意の記憶媒体にログとして追記するソフトウェア等が実装される。データ構造 B、データ構造 C は、データ構造 A とは異なるデータ構造であり、互いに異なるデータアクセス特性を持つ。

【 0 0 8 3 】

データ格納部 1 2 は、必ずしも単一の記憶媒体でなくてもよい。図 4 のデータ格納部 1 2 を複数のデータ配置ノードからなる分散ストレージシステムとして実現し、各構造別データ格納部 1 2 X を分散して格納する方式であってもよい。

【 0 0 8 4 】

データ配置特定情報 9 2 2 は、分散ストレージに格納するデータ、あるいはデータ断片の格納先を特定するための情報（および情報を格納、取得する手段）である。データの分散配置方式は、前述した通り、例えばメタサーバ方式や分散 K V S 方式が一般的に利用される。

【 0 0 8 5 】

メタサーバ方式の場合、データの位置情報を管理する情報（例えばブロックアドレスとその対応するデータノードアドレス）がデータ配置特定情報 9 2 2 である。メタサーバは、この情報（メタデータ）を参照することで、必要なデータの配置先を知ることが出来る。

【 0 0 8 6 】

分散 K V S 方式の場合、システムに参加するノードのリストが、このデータ配置特定情報に該当する。データを格納する識別子と、ノードリスト情報を用いることによって、データ格納先のデータノードを決定することが出来る。

【 0 0 8 7 】

データアクセス手段 6 1 1 は、構造情報管理手段 9 におけるデータ配置特定情報 9 2 2

10

20

30

40

50

、あるいは、予め定められた所定の場所に記憶されるデータ配置特定情報 9 2 2 のキャッシュ情報を用いて、アクセスすべきデータノード 1 ~ 4 を特定し、データノードのアクセス受付手段 1 1 1 に対してアクセス要求を発行する。

【 0 0 8 8 】

< データ構造管理情報 >

データ構造管理情報 9 2 1 は、データの集合毎にデータの格納方式を特定するためのパラメータ情報である。図 5 は、図 4 のデータ構造管理情報 9 2 1 の一例を示す図である。特に制限されるものではないが、本実施形態では、データの格納方式を制御する単位を、テーブルとする。そして、テーブル毎（テーブル識別子毎）に、レプリカ識別子、データ構造の種別、更新契機の各情報を、データ複製の複製数分、用意する。

10

【 0 0 8 9 】

図 5 ( A ) では、各テーブルは、可用性確保（保持）のために、3 つの複製を保持する。レプリカ識別子は、それぞれの複製を特定する情報であり、図 5 ( A ) では、0、1、2 として付与されている。

【 0 0 9 0 】

データ構造は、データの格納方式を示す情報である。図 5 ( A ) では、3 種類のデータ構造（A、B、C）をレプリカ識別子毎に異なる方式を指定している。

【 0 0 9 1 】

図 5 ( B ) に、データ構造 A、B、C の例を示す。データの格納方式の種類として、  
A：キュー、  
B：ロウストア、  
C：カラムストア  
が指定されている。

20

【 0 0 9 2 】

この場合、テーブル識別子「S t o c k s」のレプリカ識別子 0 は、データ構造 B（ロウストア）として格納される。

【 0 0 9 3 】

データ構造は、それぞれデータを格納するための方式であり、

A：キュー（Q U E U E）は、リンクトリスト（L i n k e d L i s t）である。

【 0 0 9 4 】

B：ロウストア（R O W S T O R E）は、テーブルのレコードを行（R O W）順に格納する。

30

【 0 0 9 5 】

C：カラムストア（C O L U M N S T O R E）は、列（C O L U M N）順に格納する。

【 0 0 9 6 】

図 6 に、テーブルのデータ保持構造の一例を示す。図 6 の（A）のテーブルは、Key カラムと、3 つの V a l u e カラムを備え、各ローは、Key と 3 つの V a l u e のセットからなる。

【 0 0 9 7 】

カラムストア、ロウストアは、それぞれ図 6 に示すように、記憶媒体上の格納順序を行（ロー）ベース、列（カラム）ベースに格納されている形式のことを指す。

40

【 0 0 9 8 】

図 6 では、テーブル（図 6 の（A）参照）の格納方式として、

レプリカ識別子 0 と 1 のデータとして、データ構造 B（ロウストア）で保持し（図 6 の（B）、（C）参照）、

レプリカ識別子 2 のデータとして、データ構造 C（カラムストア）として保持する（図 6 の（D）参照）。

【 0 0 9 9 】

再び図 5 ( A ) を参照すると、データ構造管理情報 9 2 1（図 4 参照）における更新契

50



機は、データを指定されたデータ構造として格納されるまでの時間契機である。S t o c k s のレプリカ識別子 0 の例では 3 0 s e c と指定されている。したがって、S t o c k s のレプリカ識別子 0 のデータ構造 B (ロウストア) を格納するデータノードにおいて、ロウストア方式の構造別データ格納部 1 2 2 に対して、データの更新が反映されるのが 3 0 s e c 契機であることを示す。データ更新が反映されるまでの間は、キュー等の中間構造としてデータが保持される。また、データノードでは、クライアントからの要求に対しても、中間構造に格納して応答が行われる。本実施形態では、指定されたデータ構造への変換は、更新要求とは、非同期 (A s y n c h r o n o u s ) で行われる。

#### 【 0 1 0 0 】

図 7 は、テーブルのデータ保持、非同期更新の例を模式的に説明する図である。更新契機が「0」より大きい場合には、各データノードは、W r i t e (更新要求) の応答速度に優れた構造を中間構造として持ち、更新内容を受け付ける。中間構造に書き込みを行った時点で、更新要求元のクライアントに対して処理完了の応答を返す。

#### 【 0 1 0 1 】

各データノードの中間構造 (W r i t e 向け中間構造、W r i t e 優先中間構造、あるいは「中間データ保持構造」ともいう) に書き込まれた更新データは、各データノードにおいて、それぞれ、データ構造 B、C に、それぞれ非同期 (A s y n c) に更新される。図 7 に示す例では、W r i t e により、レプリカ識別子が 0 のデータノードにおいて、W r i t e 向け中間構造には、データ構造 A が格納保持され、レプリカ識別子 1、2 のデータノードに対して同期方式 (S y n c h r o n o u s) で、W r i t e 向け中間構造に保持されたデータ構造 A のデータがレプリケート (複製) され、レプリカ識別子 1、2 のデータノードの各々において、W r i t e 向け中間構造にはデータ構造 A のデータが一旦格納保持される。レプリカ識別子 0、1、2 に対応するデータ構造にそれぞれ対応するデータノードにおいて、ターゲットのデータ構造 B、B、C への変換は、図 5 (A) に示すようなデータ構造管理情報 9 2 1 の更新契機情報により指定される。

#### 【 0 1 0 2 】

図 7 に示すように、一つのデータノードの W r i t e 向け中間構造に書き込まれた更新データ (データ構造 A) のデータノード間での複製は、書き込み (更新) と同期 (S y n c) して行われる。このような構成をとることによって、W r i t e (書き込み) データに対して、すぐに R E A D (読み出し) 系のアクセスがないデータに対しては、W r i t e の応答速度を高めることが出来る。

#### 【 0 1 0 3 】

また、(後の) R E A D 系アクセス時には、当該 R E A D アクセスに必要なデータ構造に既に変換されているため、変換されたデータ構造を用いて、R E A D 系アクセスを処理することで、処理の高速化を実現することができる。さらに、R E A D 系アクセスの種類によって、適切なデータ構造を選んでアクセス先ノードを使い分けることも出来る。

#### 【 0 1 0 4 】

本実施形態では、単に説明の簡易化のため、データ構造の種類数を A、B、C の 3 つとしたが、データ構造の種類数は 3 つに制限されるものでないことは勿論であり、特性の異なる任意の複数種類であってもよい。また、データ構造の例として、キュー、カラムストア、ロウストアの 3 種を例示したが、かかる例に制限されるものでないことは勿論である。例えば、

- ・ロウストア構造におけるインデックスの有無、
  - ・インデックスを作成したカラムの種類の違い、
  - ・更新を追記構造で格納するロウストア形式、
- 等であってもよい。

#### 【 0 1 0 5 】

図 5 に示した例とは異なる方式として、データ構造管理情報 9 2 1 において、データ構造の種類数の代わりに、データ格納プログラムを指定するようにしても良い。例えば、図 5 (A) のデータ構造 A としてデータをキューに格納するプログラム A、データ構造 B、C

10

20

30

40

50

として異なるデータベース・ソフトウェアを指定する。この場合、データ構造 A が指定されているテーブルのレプリカ識別子を格納するデータノードでは、受け付けたデータをプログラム A を実行することで処理する。

【0106】

<データ配置特定情報>

図8は、図4のデータ配置特定情報922の例を示す。各テーブル識別子のレプリカ識別子0、1、2毎に、配置ノードが指定されている。これは、前述したメタサーバ方式に対応している。

【0107】

<分散KVS>

分散KVS方式の場合、データ配置特定情報922は、分散ストレージに参加しているノードリスト情報（不図示）が該当する。このノードリスト情報をデータノード間で共有することによって、「テーブル識別子」+「レプリカ識別子」をキー情報として、コンシステント・ハッシング方式により、配置ノードを特定することが出来る。また、レプリカの配置先として、コンシステント・ハッシング方式における隣接ノードに格納することができる。コンシステント・ハッシング方式は第4の実施形態で説明する。

【0108】

再び図8を参照すると、データ配置特定情報922において、配置ノードは、可用性を保証するためには、同一のテーブルが同一ノードに保持されることがないように指定されなければならない。

【0109】

例えば、図5(A)のStocksテーブルのレプリカ識別子0と1と2の配置ノードは互いに重複してはならない。なお、可用性の考慮を無視するのであれば、この制限はこの限りではない。つまり、複数種類のレプリカを同一ノードに保持してもよい。

【0110】

<Write処理のシーケンス>

本発明の第1の実施形態の動作について説明する。図9は、図1乃至図8を参照して説明した本発明の第1の実施形態におけるWrite処理（更新を伴う処理）のシーケンスを示す図である。

【0111】

クライアント機能実現手段61は、構造情報管理手段9の構造情報保持部92に保持されているデータ配置特定情報922（図4、図8参照）の情報を取得する（あるいは任意場所のキャッシュメモリから情報を取得する）。

【0112】

クライアント機能実現手段61は、取得した情報を用いて、Write処理を行うデータの配置先のデータノード（図9では、レプリカ識別子0のデータノード1）に対して、Writeアクセス命令を発行する。

【0113】

データノード1のアクセス受付手段111は、Writeアクセス要求（Write処理要求）を受け付け、レプリカ識別子1、2に指定されているデータノード2、3に対してWriteアクセスを転送する。レプリカ識別子1、2のデータノードを特定する方法としては、データノード1が構造情報保持部92（あるいは適切なキャッシュ）にアクセスしても良いし、クライアント実現手段61が発行するWriteアクセス命令にデータ構造管理情報921の全部あるいは一部の情報をともに渡すようにしてもよい。

【0114】

各データノードのアクセス処理手段112は、受け取ったWriteアクセス要求の処理を行う。

【0115】

アクセス処理手段112は、データ構造管理情報921の情報を参照して、Write処理を実行する。

10

20

30

40

50

## 【0116】

更新契機が「0」より大きい場合には、Write 処理内容をデータ構造 A の構造別データ格納部 121 に格納する。

## 【0117】

更新契機が「0」の場合には、データ構造管理情報 921 に指定されているデータ構造の構造別データ格納部 12X に対して格納する。

## 【0118】

アクセス処理手段 112 は、Write 処理完了後、アクセス受付手段 111 に、完了通知を発行する。

## 【0119】

レプリカ先のデータノード (2、3) は、レプリカ元のデータノード 1 のアクセス受付手段 111 に Write 完了応答を返答する。

## 【0120】

アクセス受付手段 111 は、データノード 1 のアクセス処理手段 112 からの完了通知と、各レプリカ先のデータノード 2、3 の完了通知を待ち合わせ、全て受け取った後に、クライアント機能実現手段 61 に対して応答する。

## 【0121】

データ構造変換手段 113 (図 4 参照) は、定期的に構造別データ格納部 121 (データ構造 A) のデータを、構造別データ格納部 12X (データ構造管理情報 921 に指定されている、最終格納先データ構造) に変換して格納する。

## 【0122】

なお、図 9 の例では、データノード 1 が、レプリカ先のデータノード 2、3 に対して、Write アクセスを転送しているが、図 10 に示すように、クライアント機能実現手段 61 が、格納先のデータノードの全てに対して、Write アクセスを発行するようにしても良い。

## 【0123】

図 10 の例では、図 9 と比較して、Write アクセス要求の待ち合わせをクライアント機能実現手段 61 が行うことが異なる。

## 【0124】

<参照系処理のシーケンス>

図 11 は、本発明の第 1 の実施形態における参照系処理 (READ 処理) のシーケンスを示す図である。

## 【0125】

クライアント計算機 (クライアントノード) 6 は、データ構造管理情報 921 の情報を取得して、命令の実行先ノードを特定する。レプリカデータを配置するノードは、レプリカ識別子のいずれを用いてもよいが、行う処理によって適切なノードを選択することが望ましい。

## 【0126】

参照系処理とは、データの読み込みを伴う処理をいい、例えば SQL (Structured Query Language) 文における Select 文による命令等に対応する。

## 【0127】

また、

あるテーブル A からデータを読み出し、

当該データを用いた演算結果をテーブル B に更新する場合、

テーブル A からのデータ読み出しは参照系処理に該当する。

## 【0128】

あるいは、テーブル A を参照した後、テーブル A を更新するような処理の場合、一括して Write 処理 (図 9、図 10 記載) として扱っても良い。あるいは、テーブル A の参照処理は参照系処理として扱い、テーブル A の更新を、更新処理として扱ってもよい。

10

20

30

40

50

## 【 0 1 2 9 】

< クライアント機能実現手段の動作 >

図 1 2 は、クライアント機能実現手段 6 1 の視点によるアクセス処理の動作を説明するフローチャートである。図 1 2 を参照して、クライアントのアクセスフローについて説明する。

## 【 0 1 3 0 】

まず、クライアント機能実現手段 6 1 が、構造情報保持部 9 2 の情報をマスタ、あるいは任意の箇所のキャッシュにアクセスすることで取得する ( 図 1 2 のステップ S 1 0 1 ) 。

## 【 0 1 3 1 】

次に、クライアントが発行する命令内容が W r i t e 処理であるか参照処理 ( R e a d ) であるかを識別する ( ステップ S 1 0 2 ) 。

## 【 0 1 3 2 】

これは、発行命令のコマンドにより指定したり、命令の実行コードを解析したりすることで特定することが出来る。例えば、S Q L を処理するストレージシステムの場合、

- ・ I N S E R T 命令 ( テーブルへレコードを追加する S Q L 命令 ) であれば、W r i t e 処理、

- ・ S E L E C T 命令 ( テーブルからレコードを削除する S Q L 命令 ) であれば、参照系処理、

である。

## 【 0 1 3 3 】

あるいは、クライアント機能実現手段 6 1 を用いて、命令を呼び出す際に、明示的に指定するようにしても良い ( そのような A P I ( A p p l i c a t i o n P r o g r a m I n t e r f a c e ) を準備する ) 。

## 【 0 1 3 4 】

ステップ S 1 0 2 の結果、W r i t e 処理であれば、ステップ S 1 0 3 以降に進む。

## 【 0 1 3 5 】

W r i t e 処理の場合、クライアント機能実現手段 6 1 は、更新が必要なノードをデータ配置特定情報 9 2 2 の情報を用いて特定する。この処理は、図 9 を参照して説明した通りである。

## 【 0 1 3 6 】

クライアント機能実現手段 6 1 は、特定したデータノードに対して、命令実行要求 ( 更新要求 ) を発行する ( ステップ S 1 0 3 ) 。

## 【 0 1 3 7 】

クライアント機能実現手段 6 1 は、更新要求発行先のデータノードからの応答通知を待ち合わせ、更新要求が、各データノードに保持されたことを確認する ( ステップ S 1 0 4 ) 。

## 【 0 1 3 8 】

図 1 2 は、クライアント機能実現手段 6 1 が、更新先のデータノードに対して命令を発行し、応答を待ち合わせるといふ図 1 0 のシーケンスに対応するクライアント機能実現手段 6 1 の動作を説明するためのフローチャートである。

## 【 0 1 3 9 】

ステップ S 1 0 2 の結果、参照処理である場合には、ステップ S 1 0 5 へ進む。

## 【 0 1 4 0 】

ステップ S 1 0 5 では、まず、クライアント機能実現手段 6 1 は、処理内容の特性を特定 ( 認識 ) する ( ステップ S 1 0 5 ) 。

## 【 0 1 4 1 】

クライアント機能実現手段 6 1 は、特定した処理特性と、その他のシステム状況を踏まえて、アクセス対象のデータノードを選択し、命令要求を発行する処理を行う ( ステップ S 1 0 6 ) 。

10

20

30

40

50

## 【 0 1 4 2 】

クライアント機能実現手段 6 1 は、その後、データノードからアクセス処理結果を受け取る（ステップ S 1 0 7 ）。

## 【 0 1 4 3 】

以下、ステップ S 1 0 5、ステップ S 1 0 6 の処理について説明を補充する。

## 【 0 1 4 4 】

まず、クライアント機能実現手段 6 1 は、データ構造管理情報 9 2 1 に格納されている情報から、アクセス対象のデータが保持されているデータ構造の種類を知ることが出来る。例えば、図 5（A）の例の場合、WORKERS テーブルにアクセスする場合、レプリカ識別子 0、1 は、データ構造 B、レプリカ識別子 2 は、データ構造 C である。

10

## 【 0 1 4 5 】

そして、クライアント機能実現手段 6 1 では、データノードに対して行われるデータアクセスが、どちらのデータ構造に適しているかを判断し、適している方を選択する。

## 【 0 1 4 6 】

より詳しくは、例えば、クライアント機能実現手段 6 1 では、アクセス要求である SQL 文を解析し、テーブル識別子が「WORKERS」のテーブル内のあるカラムの総和をとる命令の場合には、データ構造 C（カラムストア）を選択し、ある特定のレコードを取り出す命令の場合には、データ構造 B（ロウストア）が向いていると判断する。

## 【 0 1 4 7 】

ある特定のレコードを取り出す命令であった場合、レプリカ識別子 0、1 では、どちらを選択しても良い。なお、必ずしも「最新のデータで処理を行う必要が無い場合」、レプリカ識別子 1（更新契機は 3 0 s e c）を用いることが望ましい。

20

## 【 0 1 4 8 】

この「最新のデータで処理を行う必要が無い場合」であることの特定は、アプリケーション・コンテキストに依存する。このため、クライアント機能実現手段 6 1 に受け渡される命令に、利用するデータ構造や、必要なデータの鮮度（データの新しさ）を特定する情報を、明示的に指定する形式としても良い。

## 【 0 1 4 9 】

アクセスすべきレプリカ識別子（データ構造）を特定した後、アクセスすべきデータノードを算出する。このとき、分散ストレージシステムの状況に応じて、アクセスノードの選択を変更できるようにしても良い。例えば、あるテーブルが同一のデータ構造 B として、データノード 1、2 に格納されている際に、データノード 1 のアクセス負荷が大きい場合に、データノード 2 を選択するような動作に変更してもよい。

30

## 【 0 1 5 0 】

また、別のデータ構造 C として、データノード 3 に格納されている場合に、データノード 3 のアクセス負荷が、データノード 1、2 と比較して小さければ、処理するアクセス内容がデータ構造 B の方が向いていたとしても、データノード 3（データ構造 C）に対して、アクセス要求を発行するようにしても良い。

## 【 0 1 5 1 】

クライアント機能実現手段 6 1 では、このようにして算出・選択されたデータノードに対して、アクセス要求を発行し（ステップ S 1 0 6）、該データノードから、アクセス処理結果を受け取る（ステップ S 1 0 7）。

40

## 【 0 1 5 2 】

< データノードの動作 >

図 1 3 は、図 4 のデータノードにおけるアクセス処理を説明するフローチャートである。図 1 3、図 4 を参照して、データノードの動作について詳細に説明する。

## 【 0 1 5 3 】

まず、データノードのデータ管理・処理手段 1 1 のアクセス受付手段 1 1 1 がアクセス処理要求を受け付ける（図 1 3 のステップ S 2 0 1）。

## 【 0 1 5 4 】

50

次に、データノードのデータ管理・処理手段 1 1 のアクセス受付手段 1 1 1 は、受け付けた処理要求の内容が W r i t e 処理であるか、参照処理であるか判定する（ステップ S 2 0 2 ）。

【 0 1 5 5 】

ステップ S 2 0 2 の結果、W r i t e 処理であった場合、データノードのデータ管理・処理手段 1 1 のアクセス処理手段 1 1 2 は、構造情報保持部 9 2 におけるデータ構造管理情報 9 2 1 の情報を取得する（ステップ S 2 0 3 ）。データ構造管理情報 9 2 1 の情報取得は、マスタデータにアクセスしてもよいし、任意の箇所にあるキャッシュデータにアクセスするようにしてもよいし、あるいは、図 1 又は図 2 のクライアント機能実現手段 6 1 が、データノードに対して発行する要求に情報（マスタデータ又はキャッシュデータへのアクセス）を付与し、アクセス処理手段 1 1 2 では、その情報を用いてアクセスするようにしてもよい。

10

【 0 1 5 6 】

次に、アクセス処理手段 1 1 2 は、データ構造管理情報 9 2 1 の情報から、該データノードに対する処理の更新契機が「 0 」( 零 ) であるかどうかを判定する（ステップ S 2 0 4 ）。

【 0 1 5 7 】

ステップ S 2 0 4 の結果、更新契機が「 0 」の場合、アクセス処理手段 1 1 2 は、構造情報保持部 9 2 の構造情報に指定されたデータ構造を、直接、更新する（ステップ S 2 0 5 ）。すなわち、更新データを指定されたデータ構造に変換し対応する構造別データ格納部 1 2 X ( X = 1、2、3 ) に格納する。

20

【 0 1 5 8 】

更新契機が「 0 」でない場合、アクセス処理手段 1 1 2 は、W r i t e 向け中間構造（構造別データ格納部 1 2 1 ）に、更新データを格納する（ステップ S 2 0 6 ）。

【 0 1 5 9 】

ステップ S 2 0 5、2 0 6 の場合、いずれも、処理完了後、アクセス受付手段 1 1 1 は、要求元のクライアント実現手段 6 1 に対して、処理完了通知を応答する（ステップ S 2 0 7 ）。

【 0 1 6 0 】

ステップ S 2 0 2 の結果、データの参照処理であった場合、参照処理の実行を行う（ステップ S 2 0 8 ）。

30

【 0 1 6 1 】

参照処理の実行方式として、特に制限されるものでないが、代表的には、以下の 3 種類の方法を挙げることができる。

【 0 1 6 2 】

( 1 ) 第 1 の方法は、データ構造管理情報 9 2 1 に指定されているデータ構造のデータ格納部のデータを利用して処理する。これは最も性能が優れるが、更新契機が大きい場合には、W r i t e 向け中間構造のデータが参照処理に反映されていない可能性がある。このため、データの不整合が生じる可能性がある。ただし、アプリケーション開発者が事前に認識して利用する場合や、W r i t e 後に、データの読み出しが更新契機内に起きないことがわかっているか、もし新しいデータアクセスが必要な場合には、更新契機が「 0 」のレプリカ識別子データにアクセスすると決めている場合には、特に、問題はない。

40

【 0 1 6 3 】

( 2 ) 第 2 の方法は、別途行われる変換処理の適用を待ってから処理する方法である。これは、実装が容易であるが、応答性能が劣化する。応答性能を求めないアプリケーションの場合、問題はない。

【 0 1 6 4 】

( 3 ) 第 3 の方法は、データ構造管理情報 9 2 1 に指定されているデータ構造と、W r i t e 向け中間構造に保持されているデータの両方を読んで処理する。この場合、常に、最新のデータを応答できるが、第 1 の方法より性能が劣化する。

50

## 【 0 1 6 5 】

上記第 1 乃至第 3 のいずれの方法をとってもよい。また、複数の種類を実現し、システムの設定ファイルとして記述する、クライアント機能実現手段 6 1 から発行される処理命令の中に、どの方法で実行するかを指定するようにしてもよい。

## 【 0 1 6 6 】

< データ構造変換手段の変換動作 >

図 1 4 は、図 4 のデータ構造変換手段 1 1 3 におけるデータ変換処理の動作を示すフローチャートである。図 1 4、図 4 を参照して、データ変換処理を説明する。

## 【 0 1 6 7 】

データ構造変換手段 1 1 3 は、定期的に変換処理の必要の有無を判定するため、データノード内のタイマ（図 4 では不図示）でのタイムアウト発生による呼び出しを待つ（図 1 4 のステップ S 3 0 1）。なお、このタイマは、専用タイマとしてデータ構造変換手段 1 1 3 内に備えるようにしてもよい。タイマのタイムアウト時間は、図 5（A）の更新契機（sec）に対応する。

10

## 【 0 1 6 8 】

次に、構造情報保持部 9 2 の構造情報（データ情報）を取得し（ステップ S 3 0 2）、変換が必要なデータ構造があるか否かを判定する（ステップ S 3 0 3）。例えば、タイマで判定が 1 0 秒毎に行われるときに、更新契機が 2 0 秒のデータ構造は、2 0 秒毎に変換処理を実行するため、1 0 秒時点では、変換処理を行わなくても良い。

## 【 0 1 6 9 】

変換処理が必要でない場合には、タイマ呼び出し待ち（タイマでのタイムアウト発生により呼び出されるまでウェイト）に戻る（ステップ S 3 0 1）。

20

## 【 0 1 7 0 】

一方、変換処理が必要な際には、更新向け中間データ構造から、変換対象のデータに対する更新処理内容を読み出し（ステップ S 3 0 4）、変換先の構造別データ格納部 1 2 X（X = 1 ~ 3）へ更新情報を反映する処理を行う（ステップ S 3 0 5）。

## 【 0 1 7 1 】

< 実施形態 2 >

本発明の第 2 の実施の形態について説明する。本発明の第 2 の実施の形態では、データを、所定単位で複数に分割して、複数のデータノードに格納できるようにしている。本実施形態のシステムの基本構成は、図 1、図 2、図 4 等に示した構成とされるが、図 1 5、図 1 6 を参照して説明されるように、本実施形態においては、データ構造管理情報 9 2 1、データ配置特定情報 9 2 2 の内容が拡張されている。また、図 1 7 を参照して説明されるように、本実施形態においては、データノードのアクセス受付手段が、アクセス処理手段にアクセス要求を発行するときに、他のデータノードのアクセス処理手段に対しても、アクセス要求を発行し、さらに、データ構造変換手段が、他のデータノードのデータ構造変換手段に対して、変更要求を発行する構成とされていることが、前記第 1 の実施形態と相違している。なお、本実施形態におけるデータノードの構成も、基本的には、図 4 に従うが、その詳細は図 1 7 を参照して後述される。

30

## 【 0 1 7 2 】

本実施形態では、格納対象とするデータ（テーブル識別子）を、複製の格納単位（レプリカ識別子）毎に、パーティショニング（分割）して、分割した格納単位を、各データノードでそれぞれ格納することができる。

40

## 【 0 1 7 3 】

図 1 5 は、データ構造管理情報 9 2 1（図 4 参照）の例を示す図である。データ構造管理情報 9 2 1 は、テーブル識別子に対して、複製数分、レプリカ識別子と、該レプリカ識別子に対応したパーティション数を備える。

## 【 0 1 7 4 】

パーティション数が「1」であるレプリカ識別子は、複製（レプリカ）を 1 つのデータノードに格納する。その場合の動作は、前記第 1 の実施形態と同一である。

50

## 【 0 1 7 5 】

パーティション数が「1」よりも大きい場合、そのレプリカ識別子のデータを、複数のデータノードに分割して格納する。図16は、その場合のデータ配置特定情報922の例を示す図である。

## 【 0 1 7 6 】

データ構造管理情報921において、あるレプリカ識別子のパーティション数が「1」よりも大きい場合、データ配置特定情報922（図4参照）において、当該レプリカ識別子に対して、図16に示すように、配置ノードのリスト（分割して格納する複数のデータノードのリスト）を記録する。

## 【 0 1 7 7 】

図15のデータ構造管理情報921の例では、テーブル識別子「WORKERS」のレプリカ識別子2のパーティション数が「4」である。図16のデータ配置特定情報922では、テーブル識別子「WORKERS」のレプリカ識別子2の「配置ノード」として、ノード番号2、3、5、6が指定されている。

## 【 0 1 7 8 】

配置ノードの決定は、テーブル識別子毎に、システム全体として想定される要求可用性レベルを保つように決める。マニュアル（人手）で行ってもよいし、図15のデータ構造管理情報921、図16のデータ配置特定情報922の内容をプログラムで自動生成するようにしてもよい。

## 【 0 1 7 9 】

例えば、一般的に、可用性レベルは、複製数（レプリカ数）に応じて決定される。求める可用性レベルが3レプリカであれば、レプリカ識別子を3つ用意し、それぞれの配置ノードが互いに重複しないように決定する。

## 【 0 1 8 0 】

図16の例では、テーブル識別子「WORKERS」のレプリカ識別子の各配置ノードは、互いに重複しないよう指定されている。なお、レプリカ識別子を4つ以上用意してもよいことは勿論である。例えばレプリカ識別子が4つの場合、求める可用性レベルが「3」のままであれば、同一のテーブル識別子のレプリカ識別子の配置ノードとして、1つまで重複して選ぶことが出来る（例えば、4つのレプリカ識別子のうち、配置ノードが重複するレプリカ識別子が2つあってもよい）。

## 【 0 1 8 1 】

各レプリカ識別子のデータ格納構造と、分割配置戦略（パーティショニング・ストラテジ）により、パーティショニング時の配置ノードの重複を許すか否かが異なる。

## 【 0 1 8 2 】

例えば、次のような場合には、パーティショニング時の配置ノードを重複して格納することが出来る。ノード番号1 - 18のデータノードに、ロウストア形式（データ構造B）で、12分割のレプリカを、2つ格納する場合、互いに重複を許さない場合には格納が不可能である。しかし、この場合、次のようにすれば、2レプリカ・レベルの可用性を満たしつつ、配置ノードを重複させて割り当てることが出来る。

## 【 0 1 8 3 】

レプリカ識別子0は、ノード番号1 - 12、  
レプリカ識別子1は、ノード番号7 - 18、  
に分割して格納するものとする。

## 【 0 1 8 4 】

このとき、レプリカ識別子0と1の同一レコードのデータが、同一ノードに格納されないように、分割配置戦略が決定されていれば、可用性レベルを満たすことが出来る。具体的には、下記のように、テーブルをパーティショニングする際に、ある任意のカラムの値によって、分散配置する場合（カラムの値の前半、後半で分割）、

・レプリカ識別子0のノード番号1 - 6には、カラムの値の前半、ノード番号7 - 12にはカラムの値の後半、

10

20

30

40

50



・レプリカ識別子 1 のノード番号 7 - 1 2 には、カラムの値の前半、ノード番号 1 3 - 1 8 には、カラムの値の後半  
というように格納することで、同一のレコードが、同一のノードに格納されることは回避される。このようにすることで、配置ノードの割り当てを重複させながら、可用性を満たすことが出来る。

【 0 1 8 5 】

配置ノード先の決定は、システムあるいはテーブル識別子毎に指定される可用性レベルを満たすように行う。

【 0 1 8 6 】

パーティション数が「 1 」よりも大きいレプリカ識別子に対する更新時のアクセス先は、配置ノード群のいずれを選んでも良い。あるいは、常に、リストの最初のノードを選ぶようにしてもよい（例えばテーブル識別子「WORKERS」のレプリカ識別子「2」の場合、ノード番号 2 のデータノード）。後者の方が、データ構造変換手段 1 1 3 における、構造別データ格納部 1 2 1 から構造別データ格納部 1 2 2、1 2 3 への変換処理がやや簡略化される。

【 0 1 8 7 】

パーティション時には、コンシステント・ハッシング法等を用いて分散配置してもよいし、前述したようなテーブルのあるカラムの値や、ユニークな Key の範囲などで格納先を決定してもよい。

【 0 1 8 8 】

分割配置戦略を複数用意する場合には、データ配置特定情報 9 2 2（図 4 参照）に、レプリカ識別子毎に選択された分割配置戦略の情報を記録する必要がある。

【 0 1 8 9 】

本実施形態において、パーティショニングを行う際には、前記第 1 の実施の形態と比較して、データ構造変換手段 1 1 3（図 1 7 参照）における変換処理（図 1 4 ステップ S 3 0 5）や、更新契機が「0」の場合のデータ構造の更新処理（図 1 3 ステップ S 2 0 5）が異なり、指定された配置ノード先のデータ格納部を更新する点が相違している。

【 0 1 9 0 】

また、データノードのアクセス処理時において、アクセス先が、パーティショニングにより、複数ノードにまたがる場合には、アクセス受付手段 1 1 1（図 1 7 参照）は、配置先の他のデータノードのアクセス処理手段 1 1 2（図 1 7 参照）に対してアクセス要求を発行する必要がある。

【 0 1 9 1 】

更新処理時に、更新契機（図 5（A）参照）が「0」の場合、更新処理対象のレコードが格納されるデータノード全てのアクセス処理手段 1 1 2 に対してアクセス要求を発行する必要がある。

【 0 1 9 2 】

参照処理についても、処理対象のレコードが格納されるデータノードの全てのアクセス処理手段 1 1 2 に要求を発行する。必要なデータノードの選択については、分散配置戦略に依存する。

【 0 1 9 3 】

図 1 7 は、本発明の第 2 の実施形態の構成を示す図であり、データノード 1 ~ X の構成が示されている。本実施形態においては、前記第 1 の実施形態のアクセス受付手段 1 1 1 と相違して、アクセス受付手段 1 1 1 は、自ノード内のアクセス処理手段 1 1 2 に対して、アクセス要求を発行する際に、他ノードのアクセス処理手段 1 1 2 にも発行する場合がある。同様に、データ構造変換手段 1 1 3 は、定期的に変換処理の必要の有無を判定し、データ構造の変換を行う場合、パーティショニングされたデータを格納する他のデータノードのデータ構造変換手段 1 1 3 に対してデータ変換要求を発行する。本発明の第 2 の実施形態によれば、データを分割して複数のデータノードに格納することができる。

【 0 1 9 4 】

## &lt; 実施形態 3 &gt;

次に、本発明の第3の実施形態について説明する。本実施形態では、データ構造管理情報921をアクセス負荷に応じて変更するようにしている。変更された値をシステムのデータ構造に反映することで、データ構造の設定内容(図5に示したようなレプリカ識別子毎のデータ構造の割り当て)の不適切さの修正や、システム運用後のアクセスパターンの変化などに対応可能とする。これを実現する制御パラメータの自律変更の動作について説明する。

## 【0195】

図18は、本発明の第3の実施形態のデータノードの構成を示す図である。図1、図2、図4を参照して説明した前記第1の実施形態と比較して、本実施形態においては、履歴記録部71と変更判定手段(変更判定部)72が追加されている。本実施形態の各データノードのアクセス受付手段111(あるいは他の任意の手段において)は受け付けたアクセス要求を履歴記録部71に記録するよう動作する。履歴記録部71は、各テーブル識別子のレプリカ識別子毎のアクセス要求(あるいはアクセス処理内容)を記録する。

## 【0196】

履歴記録部71は、システム全体で1つ備えた構成としてもよい。あるいは、各データノードに履歴記録部71を備え、各データノードで個別に各テーブル識別子のレプリカ識別子毎のアクセス要求を記録していき、各データノードで個別に集められたアクセス履歴を、任意の方法で、集約する仕組みを設けてもよい。

## 【0197】

変更判定手段(変更判定部)72は、履歴記録部71に格納された履歴情報を用いて、データ構造を変換するか否かについて判定する。変更判定手段72は、システム全体で1つ備えた構成としてもよいし、あるいは、各データノードで変更判定手段72を分散して動作させ、変更判定を行うような構成としてもよい。

## 【0198】

変更判定手段72は、構造変換が必要な際に、構造情報変更手段91に対して、データ構造の変換処理要求を発行する。

## 【0199】

構造情報変更手段91は、変更判定手段72からの変換処理要求に応答して、構造情報保持部92の情報を変更し、さらに、対象データノードのデータ管理・処理手段11内のデータ構造変換手段113に対して変換処理を要求する。

## 【0200】

本発明の第3の実施形態における制御パラメータの自律変更およびデータ構造の自律変換動作の流れについて、図19、図20、図21を用いて説明する。

## 【0201】

## &lt; 制御動作 &gt;

図19は、図18に示した本実施形態における制御動作を説明するフローチャートである。図19の動作を、例えば定期的に行うことによって、システムのデータ構造を自律的に変更・反映することが出来る。実行周期は、任意であるが、例えば周期を長くした場合、実行中の変更処理と、整合を取る必要がある。また、周期的な実行以外にも、所定のイベント検出に応答して変更処理を行うようにしてもよい。イベントとしては、例えばシステムの任意のいずれかの構成要素により、負荷の変更を検出(例として、一部のデータノードのCPU、ディスクなどのハードウェア利用率の大きな変化など)した場合等である。

## 【0202】

図19の動作フローは、テーブル識別子毎の構造変換処理の必要の有無の判定と、変換処理を示すものである。システムが保持管理する全てのテーブル識別子について、図19のフローを行う必要がある。

## 【0203】

変更判定手段72は、履歴記録部71のアクセス履歴情報の取得を行う(ステップS4

10

20

30

40

50

01)。

【0204】

次に、変更判定手段72は、取得したアクセス履歴情報を利用して、最近の一定期間（例えば最近1日以内、あるいは最近1週間以内等）に受け付けた全てのアクセス内容が、該当テーブル識別子のいずれかのレプリカとして適したデータ構造を持っているか否かを判定する（ステップS402）

【0205】

ステップS402において受け付けたアクセス内容に対して、レプリカ識別子のいずれかに適したデータ構造を持っている場合には、ステップS403に進む。ここで、レプリカ識別子のいずれかに適したデータ構造を持っている場合とは、例えば、列（カラム）アクセスが必要なアクセス要求を受け付けている際に、任意のレプリカ識別子のデータ構造として、カラムストア構造を持っている場合等である。

10

【0206】

ステップS403では、変換判定手段72は、各レプリカ識別子が不要なデータ構造を持っているかどうか判定する。例えば、列アクセスが必要なアクセス要求が履歴として全く無いのに、カラムストア構造を多数持つ場合、不要なデータ構造といえる。

【0207】

不要なデータ構造が無い場合には、特に変換処理をする必要が無いので、変更判定手段72は、フローを終了する。一方、不要なデータ構造がある場合、ステップS404に進む。

20

【0208】

ステップS404において、変更判定手段72は、各レプリカ識別子のデータ構造と、アクセス要求量・内容から、データ構造の変更の可否の判断を行う。データ構造の変更の可否の判断は、例えば予め定義したルール等に基づいて行われる。

【0209】

ルールとしては、以下が挙げられる。特に制限されるものでないが、ルールは、if <条件> then <アクション>（条件成立時アクションを実行）のif-then構造とされる。

【0210】

（R1）列アクセスのアクセス要求数が一定以下、且つ、行アクセスの総アクセス要求数が一定数以上の場合、カラムストア構造をロウストア構造に変換する（またはその逆）。

30

【0211】

（R2）テーブル識別子に対するアクセス要求総数が一定以上の場合、レプリカ数を増やす。

【0212】

（R3）テーブル識別子に対し、あるカラムの値による検索クエリーが一定数以上ある場合、いずれかのレプリカ識別子にインデックスを付与する。逆にアクセスが無い場合に、インデックスを削除する。

【0213】

（R4）テーブル識別子に対し、リード処理要求が一定数以上ある場合に、パーティショニング数を増加する（あるいは、この逆）。

40

【0214】

（R5）テーブル識別子に対し、複数レコードにまたがる更新処理要求が一定数以上ある場合に、パーティショニング数を削減する。あるいは、パーティショニング数を「1」にする。

【0215】

なお、ルールは上記に制限されず、任意のものを動作させてよい。

【0216】

ステップS404によってデータ構造やレプリカ数を変更する必要がある場合、ステップS405へ進む。その必要が無い場合、変更判定手段72は、フローを終了する。

50

## 【 0 2 1 7 】

ステップ S 4 0 5 において、変更判定手段 7 2、構造情報変更手段 9 1、データ構造変換手段 1 1 3 等により、データ構造を実際に変換する。レプリカを増やす場合、構造情報管理手段 9 のデータ構造管理情報 9 2 1 に、レプリカを増やすテーブル識別子のレコードを 1 つ増やし、ユニークなレプリカ識別子を付与し、その配置ノード先を決定する。配置ノードの決定は、前記第 1 の実施形態と同様に行われるが、可用性レベル以上のレプリカ数を保持していれば、他の配置ノードと重複しても良い。

## 【 0 2 1 8 】

また、レプリカは、新しいレプリカ識別子と同一のレプリカから配置ノード先へデータを複製する。

10

## 【 0 2 1 9 】

ステップ S 4 0 5 のデータ構造を変換する動作について、図 2 0、図 2 1 を参照して、より詳細に説明する。簡単化のため、図 2 0、図 2 1 については、レプリカ識別子は、パーティショニングされていない。以下では、図 1 8 のデータ構造変換手段 1 1 3 の変換処理は、データ構造を B から C に変換する例に即して説明する。

## 【 0 2 2 0 】

< データ構造変換の動作 >

図 2 0 は、本実施形態における、データ構造変換の動作を説明するフローチャートである。

## 【 0 2 2 1 】

まず、構造情報保持部 9 2 ( 図 1 6 ) のデータ構造管理情報 9 2 1 ( 図 4 ) に対して、変更判定手段 7 2 ( 図 1 6 ) が変更要求を発行する ( ステップ S 5 0 1、すなわち図 1 9 のステップ S 4 0 5 )。これにより、構造情報変更手段 9 1 は、変更先のデータノード X のデータ構造変換手段 1 1 3 に対して、変換処理要求を行う。

20

## 【 0 2 2 2 】

ステップ S 5 0 2 において、変更先のレプリカ識別子のデータをもつデータノード X では、該当レプリカ識別子のローカル複製 ( 局所的複製 ) を作成する。このローカル複製は、物理コピーではなく、ストレージによるスナップショット技術を用いてもよい。また、複製を取らず、変換元のデータとして、他ノードのレプリカ識別子のデータを用いても良い。この複製処理は、変換処理の実装方式によっては、必ずしも必要が無い。

30

## 【 0 2 2 3 】

さらに、ステップ S 5 0 3 において、構造変換処理として、データ構造変換手段 1 1 3 は、変換元のデータをデータ格納部から読み出し、変換先のデータとして異なるデータ構造として書き込む処理を行う。

## 【 0 2 2 4 】

データ構造変換手段 1 1 3 による構造変換の完了後に、変換処理中 ( あるいは変換処理開始の時点で ) 蓄積されているデータ構造 A のデータ格納部にデータ構造 A のデータ構造で格納されている更新データを、変換先のデータ構造に適用する ( ステップ S 5 0 4 )。

## 【 0 2 2 5 】

最後に、データ構造管理情報 9 2 1 ( 図 4 参照 ) の内容を変更し、クライアントノード 6 のデータアクセス手段 6 1 1 ( 図 1 参照 ) がアクセス要求の応答後に変換先のデータを用いるようにする ( ステップ S 5 0 5 )。

40

## 【 0 2 2 6 】

データ構造管理情報 9 2 1 ( 図 4 ) の変更後、変換元のデータを削除する。なお、変換元のデータは必ずしも削除しなくてもよいが、削除することで、メモリ利用効率が向上する。

## 【 0 2 2 7 】

< データ構造変換処理時のデータノードの処理 >

図 2 1 は、図 1 8 に示した本実施形態における変換処理中のデータノード内の処理を説明する図である。図 1 8 のデータ構造変換手段 1 1 3 でデータ構造の変換処理中 ( ステッ

50

プ S 5 0 2 - 5 0 4 ) において、アクセス処理手段 1 1 2 は、アクセス要求を、データ構造 A とデータ構造 B を用いて、アクセス要求を応答する。このとき、更新処理は、データ構造 A ( W r i t e 向け中間構造 ) に保持しておき、データ構造変換手段 1 1 3 で変換処理中は、データ構造 B ( R o w - S t o r e ) への適用を行わない。

【 0 2 2 8 】

データ構造変換手段 1 1 3 でのデータ構造変換処理が完了後 ( ステップ S 5 0 5 ) に、アクセス処理手段 1 1 2 は、W r i t e 向け中間構造であるデータ構造 A と、変換先のデータ構造 C ( C o l u m n S t o r e ) を用いて、アクセス要求を処理する。

【 0 2 2 9 】

なお、クライアント機能実現手段 6 1 ( 図 1 参照 ) から、アクセス先のデータノードを決定する際に、データ構造変換処理中のデータノードにはアクセスせず、他のレプリカ識別子のデータを用いるようにした場合、図 2 1 に示すように、データ構造変換処理中におけるアクセス処理手段 1 1 2 の排他処理の一部は不要になり、システム構成が簡略化される。逆に、図 2 1 のような制御機構を具備することで、データ構造変換処理中のレプリカ識別子データでも処理を行うことが出来る。

【 0 2 3 0 】

< パーティション数の変更動作 >

図 2 2、図 2 3 は、本実施形態において、パーティション数を変更する動作を説明するフローチャートである。パーティション数の変更処理は、図 1 9 と同一のフローチャートとして表現できる。以下では、図 2 2 について、図 1 9 との相違点に着目して説明する。また、パーティション数だけでなく、分散戦略を変更してもよい。分散戦略の変更の一例として、例えばラウンドロビンによる分散から、任意のカラムの値範囲による分散への変更、あるいはその逆等があげられる。

【 0 2 3 1 】

ステップ S 6 0 2 ( 図 1 9 のステップ S 4 0 2 に相当 ) は、変更判定手段 7 2 は、アクセス要求処理数に対し、必要性能に十分な分散数を保持しているか否かを判定する ( 例えば、全データをスキャンするような処理のような、データ並列の処理に対しては分散されている方が性能として有利なことが多い )。必要十分な分散数であれば、ステップ S 6 0 3 に進む。必要十分な分散数でなければ、ステップ S 6 0 4 に進む。

【 0 2 3 2 】

ステップ S 6 0 3 において、変更判定手段 7 2 は、レプリカ識別子毎に不要な分割がされていないか判定する。例えば、データ並列のアクセス処理要求が少ないのに、過剰に分散配置されているレプリカ識別子が該当する。

【 0 2 3 3 】

不要な分割がされていれば、ステップ S 6 0 4 へ進み、無ければフローを終了する。

【 0 2 3 4 】

ステップ S 6 0 4 において、変更判定手段 7 2 は、パーティション数の変更の要否判断を行う。前述したように、任意に指定されたルールに基づき、パーティション数の変更内容を決定する。変更が不要の場合には、変更判定手段 7 2 は、フローを終了する。変更が必要な場合には、変更判定手段 7 2 は、パーティション数を変更する ( ステップ S 6 0 5 )。ステップ S 6 0 5 は、パーティション数を実際に変更する処理である。

【 0 2 3 5 】

< パーティション数の変更処理 >

図 2 3 に、図 2 2 のステップ S 6 0 5 ( 変更判定手段 7 2 によるパーティショニング数変更処理 ) のフローを示す。以下では、図 2 3 について、図 2 0 と異なる点に着目して説明する。

【 0 2 3 6 】

ステップ S 7 0 2 のローカル複製は、図 2 1 に示したような変換処理中のアクセス要求の応答に利用するために準備する。

【 0 2 3 7 】

10

20

30

40

50

ステップS703では、パーティション数の変更により、配置ノードが変更されるレコードについて、データを変更先のデータノードにコピーする処理である。

【0238】

ステップS704は、図20のS504とほぼ同等であるが、データ構造Aに格納されているデータ構造変換中の更新処理内容の適用先が、別のデータノードになることがある点異なる。

【0239】

ステップS705は、図20のS505とほぼ同等である。

【0240】

パーティショニングされたデータについて、配置先ノードを変更したり、一部のデータをディスクに書き出したり、別途用意したアーカイブストレージに格納することにより、システムの容量効率やストレージコストを低減することが出来る。

10

【0241】

例えば、図24に示すように、注文履歴のような追記的にレコードを記録するような履歴記録型テーブル(A)に対して、分割配置戦略を、時系列に決定し、古いデータ(B1、B2)を、ディスクに書き出すか(C1、C2)。あるいは別のアーカイブに書き出し、新しいデータ(B3：最も新しいパーティショニング・テーブル)のみをメモリ(C3)上に保持するようにしてもよい。

【0242】

本実施形態において、構造情報保持部92のデータ配置特定情報922は、例えば図25に示すようなものとなる。データ配置特定情報922は、テーブル識別子に関して各レプリカ識別子に対応して、配置ノード、分割配置戦略、配置物理媒体の各情報を有する。なお、図24の履歴記録型テーブル(A)は、テーブル識別子の順番で格納される。

20

【0243】

分割配置戦略として、配置戦略の情報(ラウンドロビン、カラム1の値分散、時系列等)が指定されている。

【0244】

データ配置特定情報922では、テーブル識別子“orders”のレプリカ識別子2が、時系列に配置ノード2-10に分散配置され、配置先の物理媒体(memory、disk等)が指定されている。

30

【0245】

<実施形態4>

本発明の第4の実施形態としてコンシステント・ハッシングへの適用例を説明する。以下では、テーブルAをカラムストア形式でコンシステント・ハッシング分割配置する場合の例について、図26を用いて説明する。なお、本実施形態において、コンシステント・ハッシングで、データが配置されるデータノード(データ配置ノード)を決める処理は、図18の変更判定手段72で行うようにしてもよい。ノード情報は、変更判定手段72により、構造情報保持部92に記録される。特に制限されないが、本実施形態においては、キー値(Key)と、前記キー値に対応してカラム毎に1又は複数のデータレコードを有するセットをロウ方向の単位とし、ロウの識別はキー値(Key)で行われ、各カラムにカラム識別子(Value1、Value2、...)が付与されたテーブルに関して、前記キー値と、カラム識別子と、テーブル識別子を組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定する。

40

【0246】

なお、ロウストア形式において、コンシステント・ハッシング分割する場合には、レコードのKey値でハッシングして、データ配置ノードを決めればよい。Key値あるいはユニーク(一義的)なレコードIDを用いてデータ配置ノードを決定する。

【0247】

図26に模式的に示すように、コンシステント・ハッシング法においてハッシュ関数へ

50

引数を、

テーブル識別子 + カラム名 + K e y 値

を組み合わせた文字列（テーブル識別子：t a b l e A + カラム識別子 V a l u e 2 + K e y 値：a c c）を渡し、ハッシュ値が算出される。

【 0 2 4 8 】

該引数に対するハッシュ関数の出力（ハッシュ値）と、格納先ノードリスト（例えばデータノード 1 ~ 4）の情報から、コンシステント・ハッシング法により、データノードを決定することが出来る。

【 0 2 4 9 】

また、レコード毎にユニークなレコード I D を付与しておき、

テーブル識別子 + カラム名 + レコード I D

をハッシュ関数に渡す引数としてもよい。

【 0 2 5 0 】

図 2 7 ( A )、( B ) は、本実施形態におけるデータ配置ノードの記録方式について説明するための図である。カラムストア形式であるため、カラム毎にデータを記録する。外側の四角形は、データ配置ノードの記録領域の管理単位であり、例えばメモリや H D D（ハードディスクドライブ）のページに対応する。ページのサイズは任意としてよい。ページ内の任意の場所（図では末尾）に、テーブル識別子（t a b l e A）とカラム名（v a l u e 1）を指定するための管理情報を記録する。1つのカラム列全てが1つのページに収まらない場合には、他のユニットに記録する必要があるが、その他のユニットへのポインタ情報等を、この場所（記憶領域）に記録してもよい。セルの値は、ページ内の任意のアドレスに格納する。図 2 7 ( A ) では、ページの先頭側から順にセルの値（カラム名 v a l u e 1 の各値）を記録している。

【 0 2 5 1 】

また、セルの値が、どの K e y に相当する情報であることを示す情報を、別途、任意の場所に記録しておく必要がある。図 2 7 ( A ) では、同一ユニット内の管理情報の直前に記録しておく。そこには、K e y の情報（あるいはユニークなレコード I D）とそれがどのアドレスに格納されているかの情報（ポインタ）を記録する。情報（K e y : c c # 8）は、K e y : c c のセルの値がアドレス # 8、（K e y : a b # 4）は、K e y : a b のセルの値がアドレス # 4、（K e y : a a # 0）は、K e y : a a のセルの値がアドレス # 0 に格納されていることを記録するものである。

【 0 2 5 2 】

また、図 2 7 ( B ) のように、同一テーブルの別のカラム（v a l u e 2）の情報を別の記録管理ユニット（メモリ又は H D D）に記録するようにしてもよい。あるいは、さらに簡単な方法で分割配置としてもよい。

【 0 2 5 3 】

本実施形態におけるパーティショニングの第 1 の例として、キー値と、前記キー値に対応してカラム毎に 1 又は複数のデータレコードを有するセットをロウ方向の単位とし、ロウの識別はキー値で行われ、各カラムにカラム識別子が付与されたテーブルのパーティショニング（カラムストア）を行う場合、テーブル識別子とカラム識別子とを組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定し、カラム単位で別々のデータノードに分散配置するようにしてもよい。別々のデータノード間でパーティショニング単位に異なるデータ構造で格納してもよい。

【 0 2 5 4 】

図 2 8 は、テーブルのパーティショニングとして、テーブルのカラム毎に、データ配置ノードを分散して配置する場合を模式的に示す図である。ハッシュ関数へ与える値として、テーブル識別子とカラム名称（例えば、（t a b l e A : v a l u e 2）あるいは、（t a b l e A : v a l u e 3））を渡すだけでよい。該引数に対するハッシュ関数の出力（ハッシュ値）から格納ノードが算出される。

10

20

30

40

50

## 【 0 2 5 5 】

あるいは、本実施形態におけるパーティショニングの第2の例として、キー値と、前記キー値に対応してカラム毎に1又は複数のデータレコードを有するセットをロウ方向の単位とし、ロウの識別はキー値で行われ、各カラムにカラム識別子が付与されたテーブルに関して1つのカラムを、パーティショニングする場合、テーブル識別子とカラム識別子と一義的な接尾子とを組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定し、1つのカラムを、複数のデータノードに分散配置するようにしてもよい。配置先の複数のデータノード間でパーティショニング単位に異なるデータ構造で格納してもよい。

10

## 【 0 2 5 6 】

図29は、図28において、テーブルの1つのカラムを二つにパーティショニングする場合を模式的に示す図である。この場合、カラムをパーティショニングするために、ハッシュ関数の引数として与える値として、テーブル識別子とカラム名称に加えて、数字等のユニークな接尾子を付与することで、複数種類のデータ配置ノード（格納ノード）を取得する。

## 【 0 2 5 7 】

この結果、Key値がa\_b、a\_c\_cの場合、データ配置ノード（格納ノード）1に配置し、Key値がd\_d、e\_eの場合、データ配置ノード（格納ノード）2に配置する。

## 【 0 2 5 8 】

このKey値と接尾子の組み合わせ（あるいはそれを計算できる値）を、図18の構造情報保持部92に格納する。また、Key値が数字の場合、数値範囲毎に接尾子を指定するようにしてもよい。例えば、1 - 100は識別子を0とする（結果として、格納ノード1に格納される）。このようにすることで、構造情報保持部92への保持管理するデータ容量を削減することが出来る。

20

## 【 0 2 5 9 】

なお、上記実施形態の第1、第2の例のテーブル・パーティショニングでは、カラムストア方式のパーティショニングを説明したが、ロウストア方式についても同様に適用可能である。この場合、カラム識別子の代わりにキー値等が用いられる。

## 【 0 2 6 0 】

コンシステント・ハッシング方式において、例えば、分散ストレージシステムへ参加する複数のデータ配置ノードを、システムの動作状態に対応したグループに分け、データの書き込み要求を受けたデータ配置ノードでは、分散ストレージシステムへ参加する複数のデータ配置ノードに対して、グループごとに規定されるデータ複製数分、データの複製を作成するようにしてよい。この場合、各グループに対応して、データの複製作成数を決定し、複数のデータ配置ノードを論理的に配置したハッシュリングを辿り、グループごとの規定されるデータ複製数に達成するまで、複製先を探索し、複製先データ配置ノードのリストを作成するようにしてもよい。あるいは、複製先データ配置ノードのリストを受け、前記リストの各データ配置ノードに対して、複製命令を発行するようにしてもよい。クライアントからのデータの書き込み要求に対して複製先データ配置ノードのリストを作成し、ハッシュリング上に配置される複数のデータ配置ノードが属する所属グループに対応して、各所属グループに対応するデータ複製数のデータを複製するようにしてもよい。

30

40

## 【 0 2 6 1 】

分散ストレージシステムやデータベースシステムを利用して企業の情報システムが実現されており、企業の業務内容の中心となるサービスを提供するシステムは「基幹系システム」あるいは「基幹系業務システム」と呼ばれ、販売や在庫管理システム、レジのPOSシステム（Point of sale system）等が含まれる。これら基幹系システムの情報を（時には集約して）、企業の意思決定に用いるためにデータ分析を行うシステムが、データウェアハウスとして知られている。これらのシステム（基幹系システム、データウェアハウス）では、一般的にデータに対するアクセス特性が異なるため、それ

50



それぞれのアクセス特性に向くように（高速処理を行うために）、データベースシステムを用意し、データ構造を特化させることが行われている。データウェアハウス・システムにおいては、例えば複数の基幹系システムからデータ（例えばトランザクション・データ等）を抽出し再構成し情報分析、意思決定のための大規模データベースを含む。基幹系システムのデータベースからデータウェアハウス・データベースへ、データの移行を行う必要があり、この工程は、E T L（E x t r a c t / T r a n s f o r m / L o a d）と呼ばれている。E T Lは、基幹系システムとデータウェアハウス・システム双方のデータ量の増大に伴い、高負荷になることが知られているが、本発明を適用することでデータ構造変換のボトルネックを解消し、ストレージの利用効率を高めることができる。

【 0 2 6 2 】

10

本発明に係るデータ記憶システムは、並列データベースや並列データ処理システム、分散ストレージ、並列ファイルシステム、分散データベース、データグリッド、クラスタコンピュータに適用することができる。

【 0 2 6 3 】

前記開示された実施形態の全部又は一部は、特に制限されないが、以下に記載される。

【 0 2 6 4 】

（付記 1）

それぞれがデータ格納部を備え、ネットワーク結合される複数のデータノードを備え、データ複製先のデータノードが、前記データノード間で、論理的には同一であるが、物理的には異なるデータ構造をそれぞれの前記データ格納部に保持する、少なくとも二つのデータノードを含む、分散ストレージシステム。

20

【 0 2 6 5 】

（付記 2）

複製先の前記データノードにおいて、目的のデータ構造への変換を複製データの受付とは非同期で行う、付記 1 記載の分散ストレージシステム。

【 0 2 6 6 】

（付記 3）

複製先の前記データノードにおいて、中間データ保持構造に前記複製データを保持して応答を返し、前記中間データ保持構造に保持されるデータ構造を、目的のデータ構造に非同期で変換する、付記 2 記載の分散ストレージシステム。

30

【 0 2 6 7 】

（付記 4）

予め定められたテーブル単位でデータの配置先のデータノード、配置先でのデータ構造、データ分割を可変に制御する手段を備えた付記 2 記載の分散ストレージシステム。

【 0 2 6 8 】

（付記 5）

データが配置されるデータノードを、コンシステント・ハッシングで求める手段を備えた、付記 1 乃至 4 のいずれか 1 に記載の分散ストレージシステム。

【 0 2 6 9 】

（付記 6）

データ更新時に行われるデータの複製において、前記複製先のデータノードでは、更新要求対象のデータを、それぞれ、指定されたデータベースでのデータ構造とは異なるデータ構造に変換してデータを前記データ格納部に格納し、その際、前記データノードは、更新対象のデータを、一旦、中間データ保持構造を保持して前記更新に対する応答を返し、前記更新要求とは非同期で目的のデータ構造に変換して格納する、付記 1 乃至 5 のいずれか 1 に記載の分散ストレージシステム。

40

【 0 2 7 0 】

（付記 7）

格納対象のデータを識別する識別子であるテーブル識別子に対応させて、複製を特定するレプリカ識別子と、前記レプリカ識別子に対応したデータ構造の種類を特定するデータ

50

構造情報と、指定されたデータ構造に変換して格納されるまでの時間情報である更新契機情報と、を、前記データ構造の種類の数に対応させて備えたデータ構造管理情報と、

前記テーブル識別子に対応して、前記レプリカ識別子と、前記レプリカ識別子に対応した1つ又は複数のデータ配置先のデータノード情報と、を備えたデータ配置特定情報と、を記憶管理する構造情報保持部を有する構造情報管理装置と、

前記データ構造管理情報と前記データ配置特定情報とを参照して、更新処理及び参照処理のアクセス先を特定するデータアクセス部を備えたクライアント機能実現部と、

それぞれが前記データ格納部を備え、前記構造情報管理装置と前記クライアント機能実現部とに接続される複数の前記データノードと、

を備え、

前記データノードは、

前記クライアント機能実現部からのアクセス要求に基づき、更新処理を行う場合に、中間データ保持構造にデータを保持して前記クライアント機能実現部に応答を返すデータ管理・処理部と、

前記データ構造管理情報を参照し、指定された更新契機に応答して、前記中間データ保持構造に保持されるデータを、前記データ構造管理情報で指定されたデータ構造に変換する処理を行うデータ構造変換部と、

を備えることを特徴とする、付記1乃至6のいずれか1に記載の分散ストレージシステム。

【0271】

(付記8)

前記中間データ保持構造は、指定された目的のデータ構造としてデータが前記データ格納部に格納されるまでの間、前記データを保持する、付記7記載の分散ストレージシステム。

【0272】

(付記9)

前記クライアント機能実現部が、前記更新処理又は前記参照処理の内容に応じてアクセス先のデータノードを、前記データ構造管理情報と前記データ配置特定情報より選択する、付記7記載の分散ストレージシステム。

【0273】

(付記10)

前記クライアント機能実現部は、前記構造情報管理装置の前記構造情報保持部に保持されている前記データ配置特定情報、又は、前記構造情報保持部に保持される情報をキャッシュする構造情報キャッシュ保持部に保持されているデータ配置特定情報を取得し、データ配置先のデータノードに対して、アクセス命令を発行する、付記7記載の分散ストレージシステム。

【0274】

(付記11)

前記データノードは、アクセス受付部、アクセス処理部、データ構造変換部を備え、

前記データノードの前記データ格納部は、構造別データ格納部を備え、

前記アクセス受付部は、前記クライアント機能実現部からの更新要求を受け付け、前記データ配置特定情報においてレプリカ識別子に対応して指定されているデータノードに対して更新要求を転送し、

前記データノードの前記アクセス処理部は、受け取った更新要求の処理を行い、前記データ構造管理情報の情報を参照して更新処理を実行し、その際、前記データ構造管理情報の情報から、前記データノードに対する前記更新契機が零の場合、更新データを、前記データ構造管理情報に指定されるデータ構造に変換して前記構造別データ格納部を更新し、

前記更新契機が零でない場合、前記中間データ保持構造に、一旦、更新データを書き込み、処理完了を応答し、

前記アクセス受付部は、前記アクセス処理部からの完了通知と、レプリカ先のデータノ

10

20

30

40

50

ードの完了通知を受けると、前記クライアント機能実現部に対して応答し、

前記データ構造変換部は、前記中間データ保持構造のデータを、前記データ構造管理情報に指定されているデータ構造に変換し変換先の前記構造別データ格納部に格納する、付記 7 又は 10 記載の分散ストレージシステム。

【 0 2 7 5 】

(付記 1 2 )

前記クライアント機能実現部は、参照系アクセスの場合、データノードに対して行われるデータアクセスに適しているデータ構造を選択し、レプリカ識別子を特定した後、アクセスすべきデータノードを算出し、選択されたデータノードに対してアクセス要求を発行し前記データノードからアクセス処理結果を受け取る、付記 7 記載の分散ストレージシステム。

10

【 0 2 7 6 】

(付記 1 3 )

前記クライアント機能実現部が、前記データノード内に配設されている、付記 7 記載の分散ストレージシステム。

【 0 2 7 7 】

(付記 1 4 )

前記クライアント機能実現部が、前記構造情報保持部に保持される情報をキャッシュする構造情報キャッシュ保持部を備えた付記 1 3 記載の分散ストレージシステム。

【 0 2 7 8 】

20

(付記 1 5 )

前記クライアント機能実現部の前記構造情報キャッシュ保持部の構造情報と、前記構造情報管理装置の前記構造情報保持部に保持される構造情報を同期させる構造情報同期部を備えた付記 1 4 記載の分散ストレージシステム。

【 0 2 7 9 】

(付記 1 6 )

前記データ構造管理情報が、データを複数のデータノードに分割して格納する分割数であるパーティション数をレプリカ識別子に対応して備え、

前記データ配置特定情報は、前記データ構造管理情報においてパーティション数が 2 以上に対応するレプリカ識別子に対応した配置ノードとして、複数のデータノードを含み、

30

アクセス要求を受けた前記データノードの前記アクセス受付部は、パーティショニングされたデータの配置先が複数のデータノードにまたがる場合に、前記複数のデータノードを構成する他のデータノードのアクセス処理部にアクセス要求を発行する、付記 7 記載の分散ストレージシステム。

【 0 2 8 0 】

(付記 1 7 )

アクセス要求を受けた前記データノードの前記データ構造変換部は、前記更新契機が零のとき、他のデータノードの前記データ構造変換部に対してアクセス要求を発行する、付記 7 又は 11 記載の分散ストレージシステム。

【 0 2 8 1 】

40

(付記 1 8 )

アクセス要求の履歴を記録する履歴記録部と、

前記履歴記録部の履歴情報を用いてデータ構造の変換を行うか否かを判定する変更判定部と、

を備えた付記 7 記載の分散ストレージシステム。

【 0 2 8 2 】

(付記 1 9 )

前記変更判定部は、データ構造の変換が必要と判定した場合、前記構造情報管理装置の前記構造情報変更部に変換要求を出力し、

前記構造情報管理装置の前記構造情報変更部は、前記構造情報保持部の情報を変更し、

50

前記データノードの前記データ構造変換部に変換要求を出力し、

前記データノードの前記データ構造変換部は前記データノードの前記データ格納部に保持されるデータ構造の変換を行う、付記 18 記載の分散ストレージシステム。

【0283】

(付記 20)

それぞれがデータ格納部を備え、ネットワーク結合される複数のデータノードを備えたシステムでの分散ストレージ方法であって、

データ複製先のデータノードの少なくとも二つのデータノードが、前記データノード間で、論理的には同一であるが、物理的には異なるデータ構造をそれぞれの前記データ格納部に保持する、分散ストレージ方法。

10

【0284】

(付記 21)

複製先の前記データノードにおいて、目的のデータ構造への変換を複製データの受付とは非同期で行う、付記 20 記載の分散ストレージ方法。

【0285】

(付記 22)

複製先の前記データノードにおいて、中間データ保持構造に複製データを保持して応答を返し、前記中間データ保持構造に保持されるデータ構造を、目的のデータ構造に非同期で変換する、付記 21 記載の分散ストレージ方法。

【0286】

20

(付記 23)

予め定められたテーブル単位でデータの配置先のデータノード、配置先でのデータ構造、データ分割を可変に制御する、付記 21 記載の分散ストレージ方法。

【0287】

(付記 24)

データが配置されるデータノードをコンシステント・ハッシングで求める、付記 20 乃至 23 のいずれか 1 に記載の分散ストレージ方法。

【0288】

(付記 25)

データ更新時に行われるデータの複製において、前記複製先のデータノードでは、更新要求対象のデータを、それぞれ、指定された目的のデータベースでのデータ構造とは異なるデータ構造に変換してデータを前記データ格納部に格納し、その際、前記データノードは、更新対象のデータを一旦、中間構造を保持して前記更新に対する応答を返し、前記更新要求とは非同期で、目的のデータ構造に変換して格納する、付記 20 乃至 24 のいずれか 1 に記載の分散ストレージ方法。

30

【0289】

(付記 26)

格納対象のデータを識別する識別子であるテーブル識別子に対応させて、複製を特定するレプリカ識別子と、前記レプリカ識別子に対応したデータ構造の種類を特定するデータ構造情報と、指定されたデータ構造に変換して格納されるまでの時間情報である更新契機情報と、を前記データ構造の種類の数に対応させて備えたデータ構造管理情報と、

40

前記テーブル識別子に対応して、前記レプリカ識別子と、前記レプリカ識別子に対応した 1 つ又は複数のデータ配置先のデータノード情報と、を備えたデータ配置特定情報と、

を含む構造情報を構造情報管理部で記憶管理し、

クライアント側では、前記データ構造管理情報と前記データ配置特定情報を参照して、更新処理及び参照処理のアクセス先を特定し、

前記データノードは、

前記クライアント側からのアクセス要求に基き、更新処理を行う場合に、中間データ保持構造にデータを保持して前記クライアントに応答を返し、

前記データ構造管理情報を参照し、指定された更新契機に応じて、前記中間データ保持

50

構造から指定されたデータ構造に変換する、ことを特徴とする、付記 2 5 記載の分散ストレージ方法。

【 0 2 9 0 】

(付記 2 7 )

前記データ構造管理情報が、データを複数のデータノードに分割して格納する分割数であるパーティション数を、レプリカ識別子に対応して備え、

前記データ配置特定情報は、前記データ構造管理情報においてパーティション数が 2 以上に対応するレプリカ識別子に対応した配置ノードとして、複数のデータノードを含み、

アクセス要求を受けた前記データノードでは、パーティショニングされたデータの配置先が複数のデータノードにまたがる場合に、前記複数のデータノードを構成する他のデータノードに対してアクセス要求を発行する、付記 2 6 記載の分散ストレージ方法。

10

【 0 2 9 1 】

(付記 2 8 )

アクセス要求に履歴を記録する履歴記録部での履歴情報を用いて、データ構造の変換を行うか否かを判定し、変換が必要な場合、前記構造情報を変換し、さらに前記データノードのデータ構造を変換する、付記 2 6 記載の分散ストレージ方法。

【 0 2 9 2 】

(付記 2 9 )

キー値と、前記キー値に対応して 1 又は複数のデータレコードを 1 又は複数のカラムに有するセットをロウ方向の単位とし、各カラムにカラム識別子が付与されたテーブルに関して、前記キー値と、前記カラム識別子と、前記テーブル識別子を組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定する、付記 5 記載の分散ストレージシステム。

20

【 0 2 9 3 】

(付記 3 0 )

キー値と、前記キー値に対応して 1 又は複数のデータレコードを 1 又は複数のカラムに有するセットをロウ方向の単位とし、各カラムにカラム識別子が付与されたテーブルに関して、前記テーブル識別子と前記カラム識別子とを組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定し、カラム単位で別々のデータノードに分散配置する、付記 5 記載の分散ストレージシステム。

30

【 0 2 9 4 】

(付記 3 1 )

キー値と、前記キー値に対応して 1 又は複数のデータレコードを 1 又は複数のカラムに有するセットをロウ方向の単位とし、各カラムにカラム識別子が付与されたテーブルに関して、前記テーブル識別子と前記カラム識別子と一義的な接尾子とを組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定し、1つのカラムを、複数のデータノードに分散配置する、付記 5 記載の分散ストレージシステム。

40

【 0 2 9 5 】

(付記 3 2 )

1 又は複数のデータレコードを 1 又は複数のカラムに有するセットをロウ方向の単位とし、各カラムにカラム識別子が付与され、レコード毎に一義的なレコード識別子が付与されたテーブルに関して、前記テーブル識別子と前記カラム識別子と前記レコード識別子を組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定する、付記 5 記載の分散ストレージシステム。

【 0 2 9 6 】

(付記 3 3 )

50

キー値と、前記キー値に対応して1又は複数のデータレコードを1又は複数のカラムに有するセットをロウ方向の単位とし、各カラムにカラム識別子が付与されたテーブルに関して、前記キー値と、前記カラム識別子と、前記テーブル識別子を組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定する、付記24記載の分散ストレージ方法。

【0297】

(付記34)

キー値と、前記キー値に対応して1又は複数のデータレコードを1又は複数のカラムに有するセットをロウ方向の単位とし、各カラムにカラム識別子が付与されたテーブルに関して、前記テーブル識別子と前記カラム識別子とを組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定し、カラム単位で別々のデータノードに分散配置する、付記24記載の分散ストレージ方法。

【0298】

(付記35)

キー値と、前記キー値に対応して1又は複数のデータレコードを1又は複数のカラムに有するセットをロウ方向の単位とし、各カラムにカラム識別子が付与されたテーブルに関して、前記テーブル識別子と前記カラム識別子と一義的な接尾子とを組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定し、1つのカラムを、複数のデータノードに分散配置する、付記24記載の分散ストレージ方法。

【0299】

(付記36)

1又は複数のデータレコードを1又は複数のカラムに有するセットをロウ方向の単位とし、各カラムにカラム識別子が付与され、レコード毎に一義的なレコード識別子が付与されたテーブルに関して、前記テーブル識別子と前記カラム識別子と前記レコード識別子を組み合わせた文字列を引数としてハッシュ関数でハッシュ値を求め、前記ハッシュ値と、格納先ノードリスト情報から、コンシステントハッシュにより、データ配置先のデータノードを決定する、付記24記載の分散ストレージ方法。

【0300】

なお、上記の特許文献の各開示を、本書に引用をもって繰り込むものとする。本発明の全開示（請求の範囲を含む）の枠内において、さらにその基本的技術思想に基づいて、実施形態の変更・調整が可能である。また、本発明の請求の範囲の枠内において種々の開示要素（各請求項の各要素、各実施形態の各要素、各図面の各要素等を含む）の多様な組み合わせないし選択が可能である。すなわち、本発明は、請求の範囲を含む全開示、技術的思想にしたがって当業者であればなし得るであろう各種変形、修正を含むことは勿論である。

【符号の説明】

【0301】

1～4 データノード

5 ネットワーク

6 クライアントノード

9 構造情報管理手段（構造情報管理装置）

11、21、31、41 データ管理・処理手段（データ管理・処理部）

12、22、32、42 データ格納部

61 クライアント機能実現手段（クライアント機能実現部）

71 履歴記録部

72 変更判定手段（変更判定部）

91 構造情報変更手段（構造情報変更部）

10

20

30

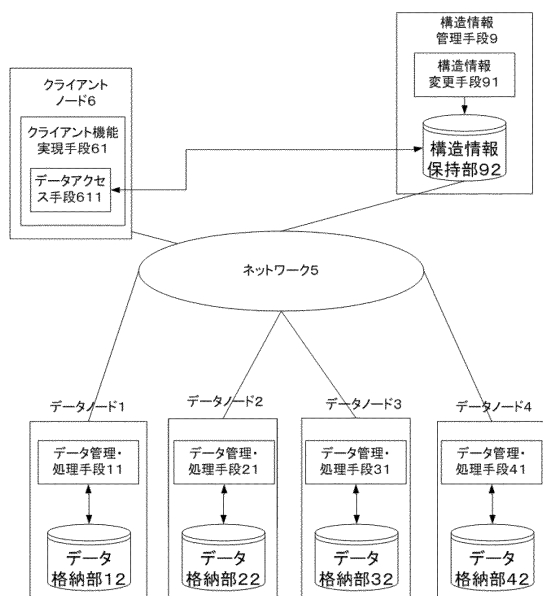
40

50

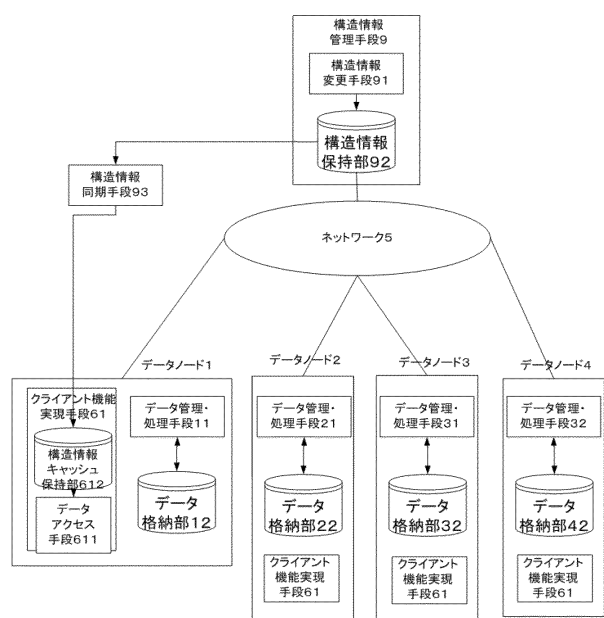
- 9 2 構造情報保持部
- 9 3 構造情報同期手段（構造情報同期部）
- 1 0 1 ~ 1 0 4 データノード計算機
- 1 0 1 a C P U
- 1 0 1 b データ記憶装置
- 1 0 1 c データ転送装置
- 1 0 5 ネットワーク
- 1 1 1 アクセス受付手段（アクセス受付部）
- 1 1 2 アクセス処理手段（アクセス処理部）
- 1 1 3 データ構造変換手段（データ構造変換部）
- 1 2 1、1 2 2、1 2 3、1 2 X 構造別データ格納部
- 6 1 1 データアクセス手段（データアクセス部）
- 6 1 2 構造情報キャッシュ保持部
- 9 2 1 データ構造管理情報
- 9 2 2 データ配置特定情報

10

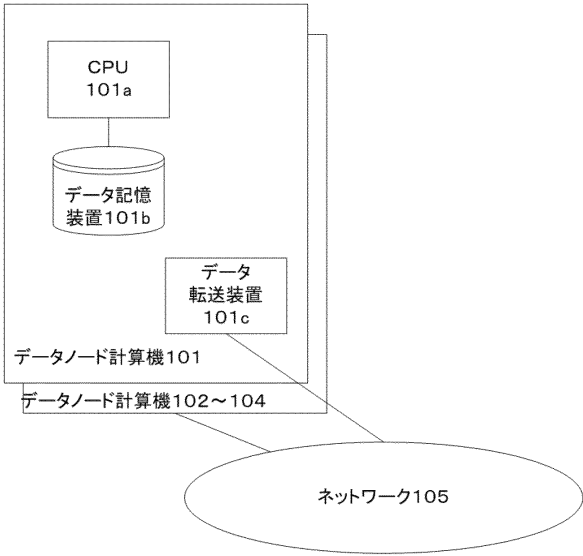
【図 1】



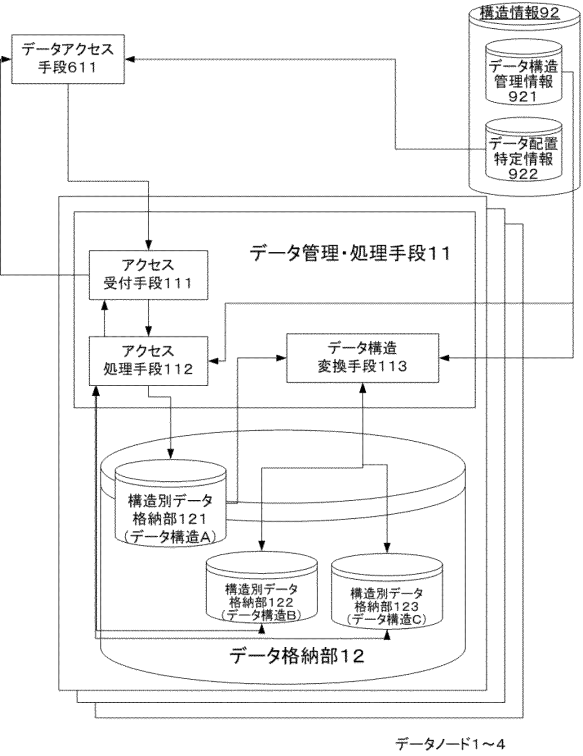
【図 2】



【図 3】



【図 4】



【図 5】

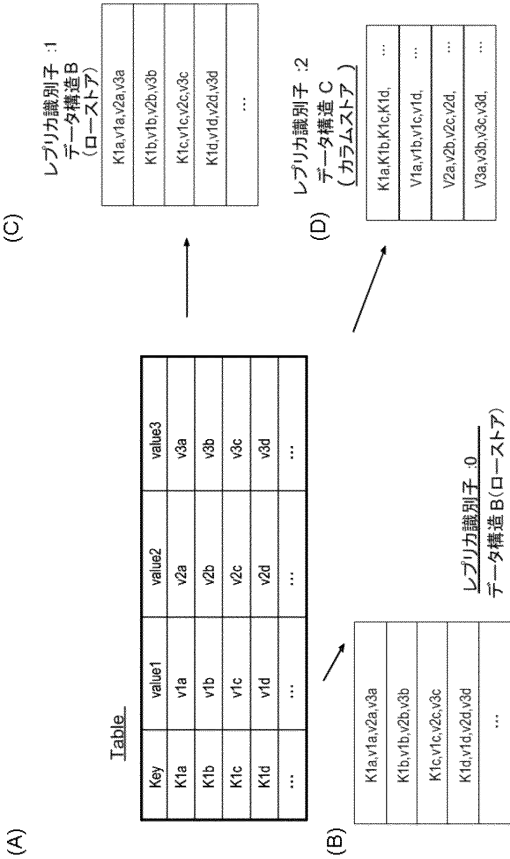
(A)

テーブル識別子	レプリカ識別子	データ構造	更新回数 (sec)
Stocks	0	B	30
	1	B	60
	2	C	60
WORKERS	0	B	0
	1	B	30
	2	C	20
Shops	0	A	0
	1	B	30
	2	C	240
Salary	0	B	30
	1	B	30
	2	C	240

(B)

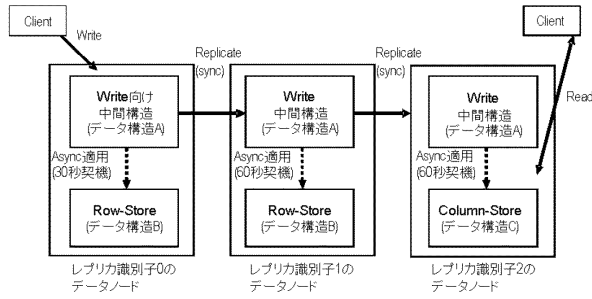
データ構造	データ格納方式の種類
A	キュー
B	ローストア
C	カラムストア
...	...

【図 6】

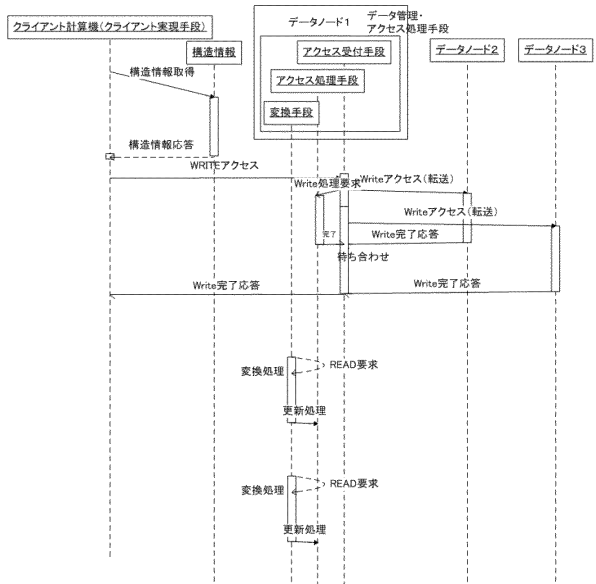




【図 7】



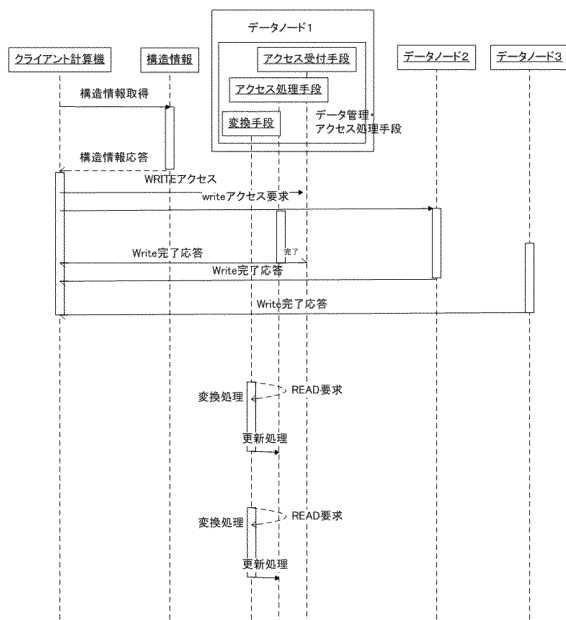
【図 9】



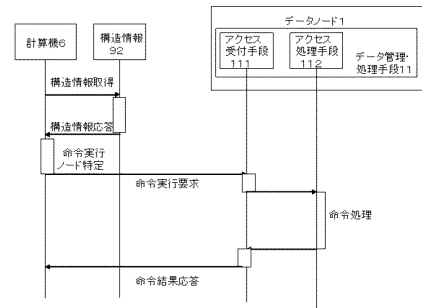
【図 8】

テーブル識別子	レプリカ識別子	配置ノード
Stocks	0	1
	1	2
	2	3
WORKERS	0	4
	1	1
	2	2
Shops	0	3
	1	4
	2	1
Salary	0	2
	1	3
	2	4

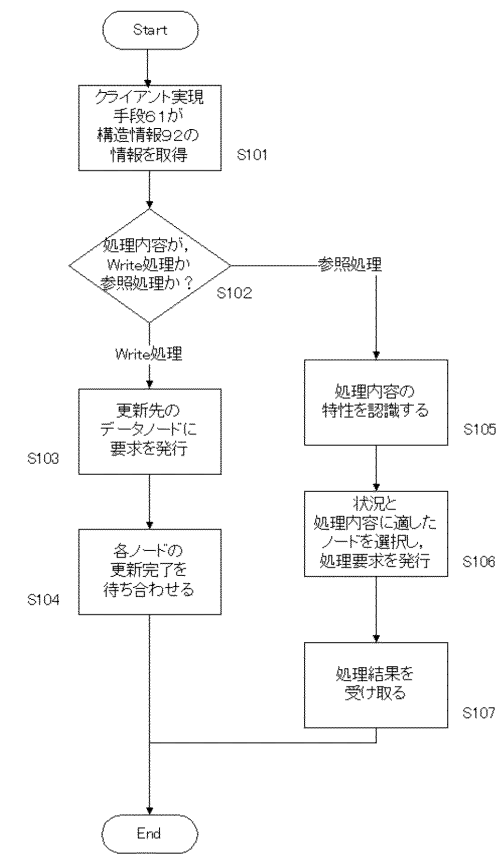
【図 10】



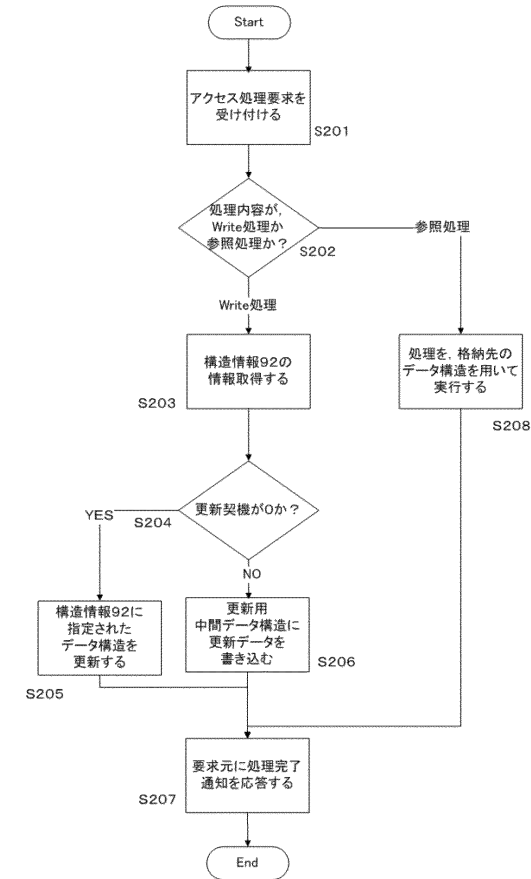
【図 11】



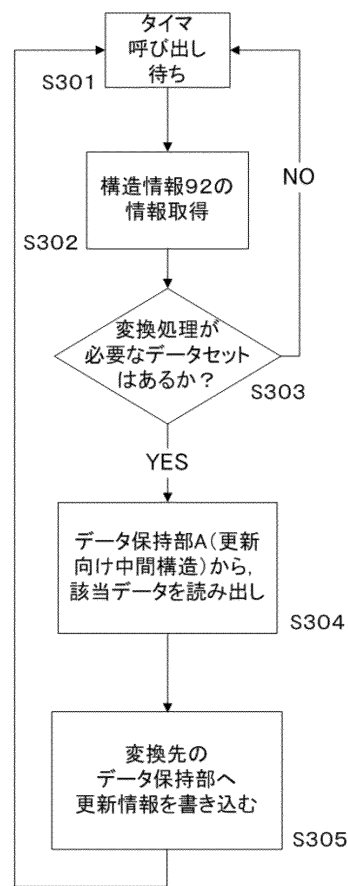
【図 1 2】



【図 1 3】



【図 1 4】



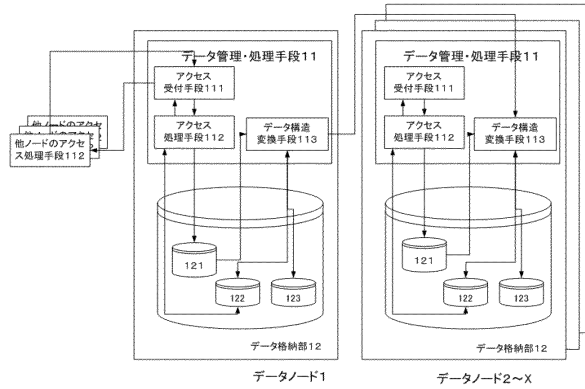
【図 1 5】

テーブル識別子	レプリカ識別子	パーティション数
Stocks	0	1
	1	1
	2	1
WORKERS	0	1
	1	1
	2	4
Shops	0	1
	1	1
	2	1
Salary	0	1
	1	16
	2	1

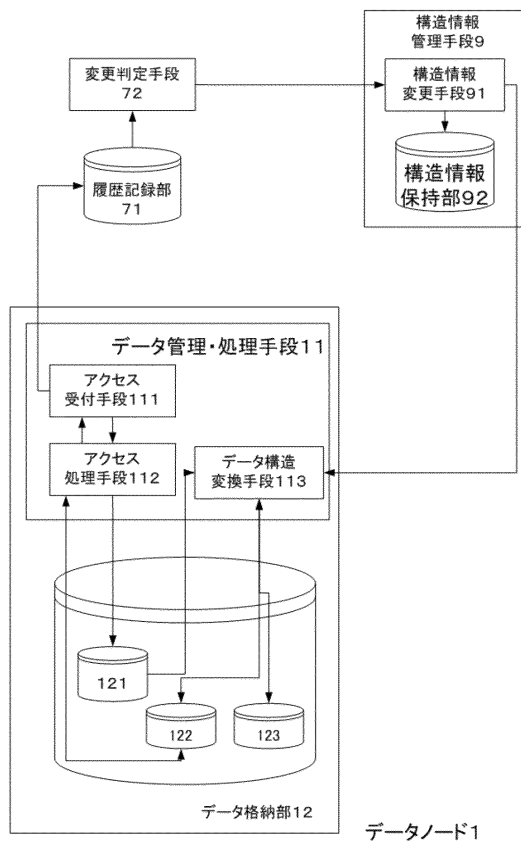
【図 16】

テーブル識別子	レプリカ識別子	配置ノード
Stocks	0	1
	1	2
	2	3
WORKERS	0	4
	1	1
	2	2,3,5,6
Shops	0	3
	1	4
	2	1
Salary	0	2
	1	3,5-15
	2	4

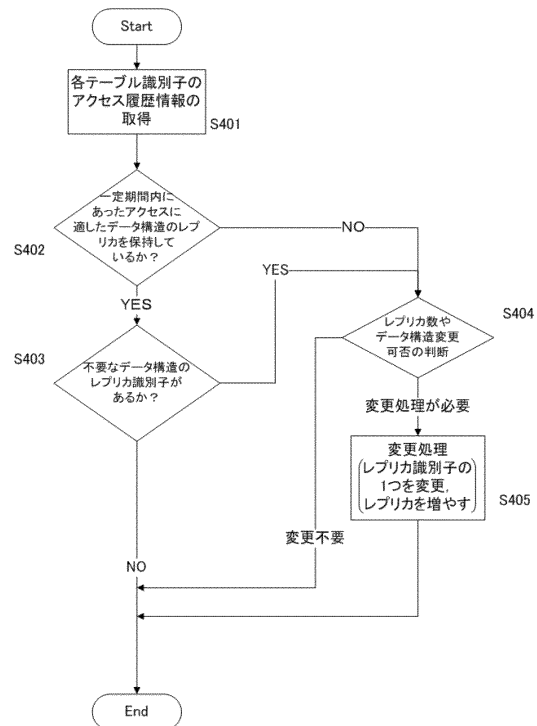
【図 17】



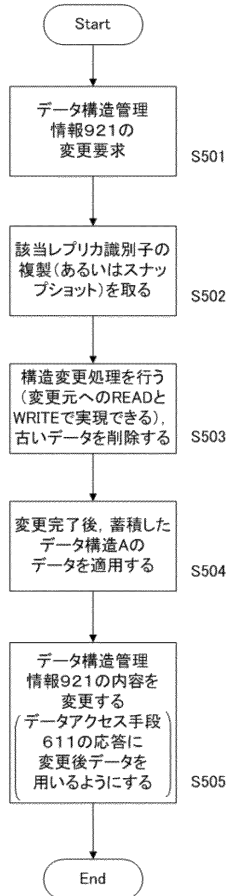
【図 18】



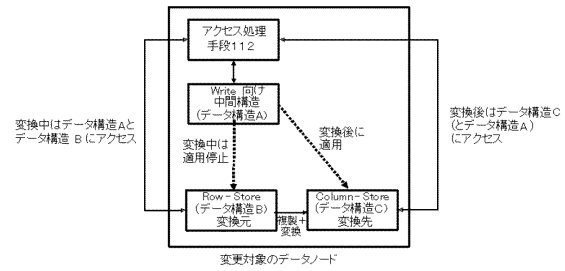
【図 19】



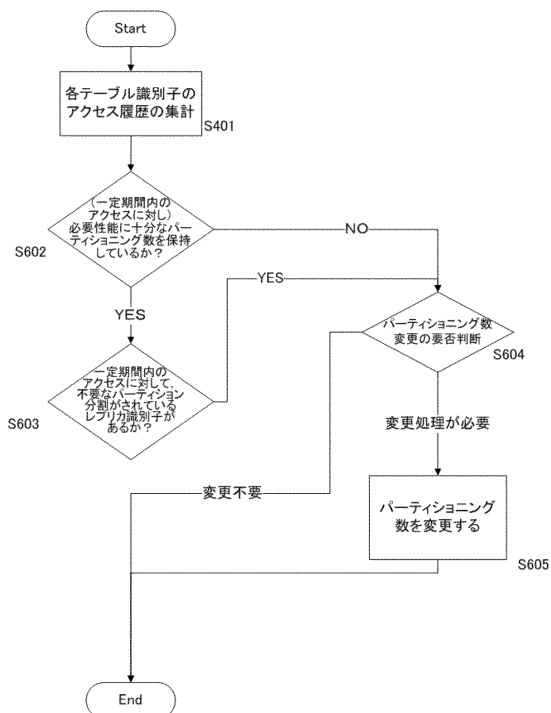
【図 20】



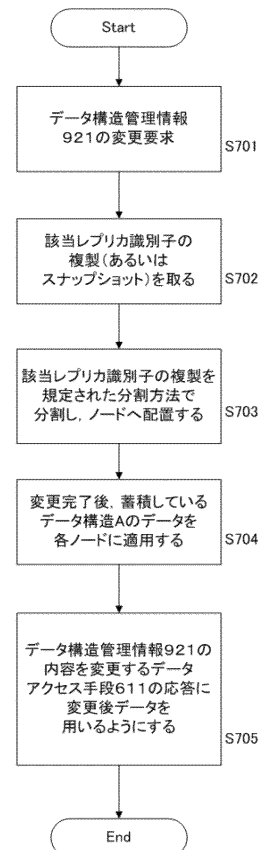
【図 21】



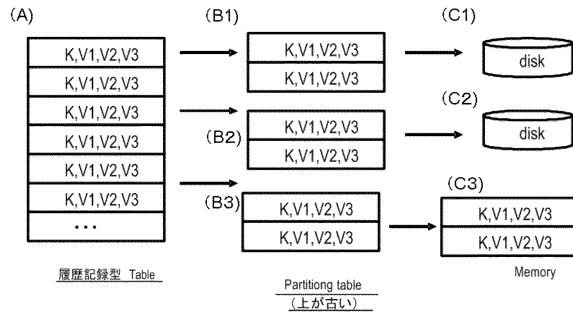
【図 22】



【図 23】



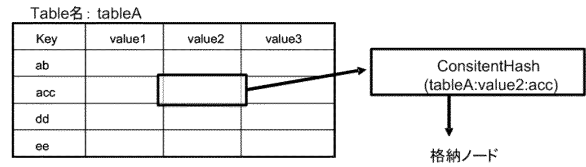
【図 24】



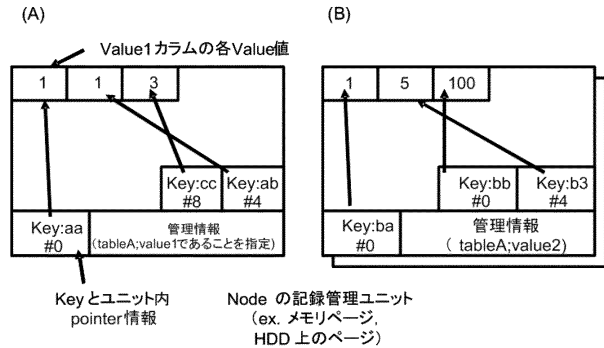
【図 25】

テーブル識別子	レプリカ識別子	配置ノード	分散配置戦略	配置物理媒体
Stocks	0	1	---	memory
	1	2	---	memory
	2	3	---	memory
WORKERS	0	4	---	memory
	1	1	---	memory
	2	2,3,5,6	roundRobin	memory
Shops	0	3	---	memory
	1	4	---	memory
	2	1	---	memory
Salary	0	2	---	memory
	1	3,5-15	コラム1の値分散	memory
	2	4	---	memory
orders	0	0	---	memory
	1	1	---	memory
	2	2-10	時系列	memory, disk,disk,...

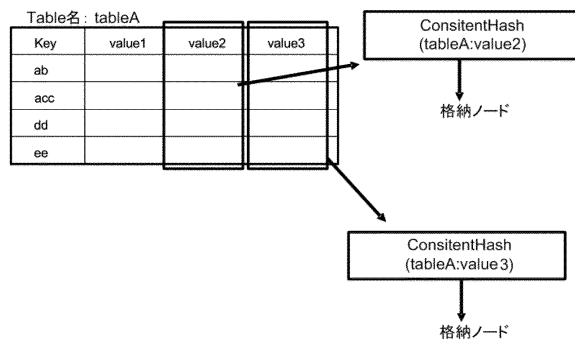
【図 26】



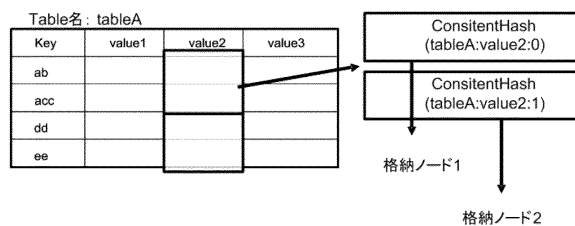
【図 27】



【図 28】



【図 29】



## フロントページの続き

(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)/エネルギー利用最適化データセンタ基盤技術の研究開発/サーバの最適構成とクラウド・コンピューティング環境における進化するアーキテクチャーの開発/クラウド・コンピューティング技術」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

(56)参考文献 特開2011-008711(JP,A)

特開2010-146067(JP,A)

中村俊介 ほか1名,読み出し性能と書き込み性能を選択可能なクラウドストレージ,情報処理学会研究報告[CD-ROM],日本,一般社団法人情報処理学会,2011年 2月15日,第2011-OS-116巻 第9号,pp.1~7

Avinash Lakshman ほか1名,Cassandra - A Decentralized Structured Storage System,ACM SIGOPS Operating Systems Review,米国,ACM,2010年 4月,Volume 44, Issue 2, pp.35~40

(58)調査した分野(Int.Cl.,DB名)

G06F 12/00

G06F 3/06

JSTPlus(JDreamIII)