

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2012-123790

(P2012-123790A)

(43) 公開日 平成24年6月28日(2012.6.28)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G 0 6 F 12/00 (2006.01)</b>	G 0 6 F 12/00 5 0 1 B	5 B 0 6 5
<b>G 0 6 F 3/06 (2006.01)</b>	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 1 Z	
	G 0 6 F 3/06 3 0 1 X	

審査請求 未請求 請求項の数 15 O L (全 20 頁)

(21) 出願番号 特願2011-245173 (P2011-245173)  
 (22) 出願日 平成23年11月9日 (2011.11.9)  
 (31) 優先権主張番号 12/961544  
 (32) 優先日 平成22年12月7日 (2010.12.7)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 390009531  
 インターナショナル・ビジネス・マシーンズ・コーポレーション  
 INTERNATIONAL BUSINESS MACHINES CORPORATION  
 アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード  
 (74) 代理人 100108501  
 弁理士 上野 剛史  
 (74) 代理人 100112690  
 弁理士 太佐 種一  
 (74) 代理人 100091568  
 弁理士 市位 嘉宏

最終頁に続く

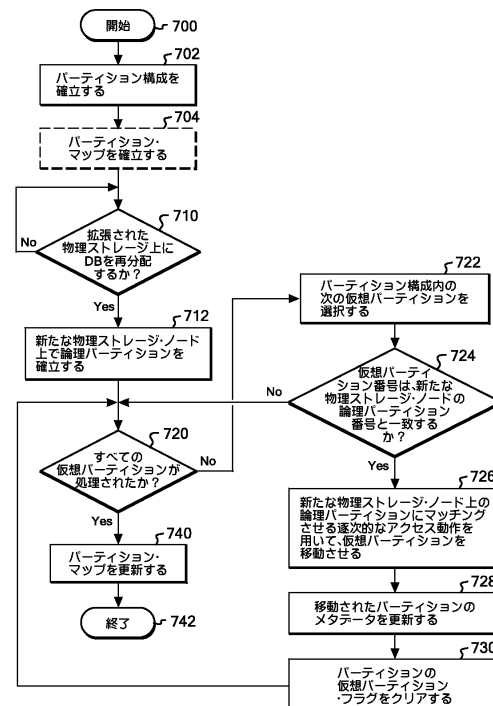
(54) 【発明の名称】 仮想パーティションを用いたデータベース再分配

## (57) 【要約】

【課題】 区分データベースを再分配する方法、プログラム及びデータ処理システムを提供すること。

【解決手段】 いくつかの実施形態において、区分データベースは、少なくとも論理的又は物理的な第1データ・ストレージ・ノード上の複数の論理又は物理パーティション内に格納され、複数の論理パーティションの中の第1のパーティションのサブセットは、仮想パーティションとして構成される。区分データベースを格納するように第2の物理データ・ストレージ・ノードの割り付けを指示する入力を受け取る。第2のパーティションが、第2のデータ・ストレージ・ノード上で構成される。入力にตอบสนองして、区分データベースは、第1のパーティション上の仮想パーティション内のデータを第2のパーティションに移動させることにより、第1及び第2のデータ・ストレージ・ノードにわたって再分配される。

【選択図】 図9



**【特許請求の範囲】****【請求項 1】**

区分データベースを、データ処理システムの少なくとも第 1 のデータ・ストレージ・ノード上の複数のパーティション内に格納するステップと、

前記複数のパーティションの中の第 1 のパーティションのサブセットを、仮想パーティションとして構成するステップと、

前記区分データベースを格納するように第 2 のデータ・ストレージ・ノードの割り付けを指示する入力を受け取るステップと、

前記第 2 のデータ・ストレージ・ノード上で、前記区分データベースの第 2 のパーティションを構成するステップと、

前記入力に応答して、前記第 1 のパーティション上の前記仮想パーティション内のデータを前記第 2 のパーティションに移動させることによって、前記区分データベースを前記第 1 及び第 2 のデータ・ストレージ・ノードにわたって再分配するステップと

を含む、データ処理方法。

10

**【請求項 2】**

前記仮想パーティションが、前記第 1 のパーティションの 1 つ又は複数のデータ・ブロックを含む、請求項 1 に記載の方法。

**【請求項 3】**

前記仮想パーティションが、前記第 1 のパーティション内のデータのみを含む、請求項 2 に記載の方法。

20

**【請求項 4】**

前記方法が、前記仮想パーティションと前記第 1 のパーティションとを関連付けるパーティション構成データ構造体を確立するステップを含み、

前記再分配するステップが、前記第 2 のパーティションに移動された前記データが仮想パーティション内に存在しないことを示すように、前記パーティション構成データ構造体を更新するステップを含む、

請求項 1 に記載の方法。

**【請求項 5】**

前記方法が、前記区分データベース内のデータを、複数のパーティション及び前記仮想パーティションのうちの特定の 1 つにマッピングする、パーティション・マップを確立するステップを含み、

30

前記再分配するステップが、前記第 2 のパーティションに移動されたデータが仮想パーティションにマッピングされないことを示すように、前記パーティション・マップを更新するステップを含む、

請求項 1 に記載の方法。

**【請求項 6】**

前記再分配するステップが、

前記仮想パーティション内のデータのバックアップを作成するステップと、

前記データを前記バックアップから前記第 2 のパーティションに復元するステップとを含む、請求項 1 に記載の方法。

40

**【請求項 7】**

前記仮想パーティション内の前記データが、第 1 のデータであり、

前記再分配の前に、前記第 1 のパーティションが、前記第 1 のデータと、前記仮想パーティション内に存在しない第 2 のデータとを格納し、

前記バックアップを作成するステップが、前記第 1 のデータを含み、かつ前記第 2 のデータを除外するバックアップを作成するステップを含む、

請求項 6 に記載の方法。

**【請求項 8】**

請求項 1 ～ 7 の何れか 1 項に記載の方法の各ステップをコンピュータに実行させる、コンピュータ・プログラム。

50

**【請求項 9】**

データ処理システムであって、  
プロセッサと、  
前記プロセッサに結合されたデータ・ストレージと、  
前記データ・ストレージ内に格納されたプログラム・コードであって、前記プロセッサによって実行されたときに、前記データ処理システムに、  
区分データベースを、前記データ・ストレージの少なくとも第 1 のデータ・ストレージ・ノード上の複数のパーティション内に格納することと、  
前記複数のパーティションの中の第 1 のパーティションのサブセットを、仮想パーティションとして構成することと、  
前記区分データベースを格納するように前記データ・ストレージの第 2 のデータ・ストレージ・ノードの割り付けを指示する入力を受け取ることと、  
前記第 2 のデータ・ストレージ・ノード上で、前記区分データベースの第 2 のパーティションを構成することと、  
前記入力に応答して、前記第 1 のパーティション上の前記仮想パーティション内のデータを前記第 2 のパーティションに移動させることによって、前記区分データベースを前記第 1 及び第 2 のデータ・ストレージ・ノードにわたって再分配することと  
を行わせる、プログラム・コードと  
を含む、データ処理システム。

10

**【請求項 10】**

20

前記仮想パーティションが、前記第 1 のパーティションの 1 つ又は複数のデータ・ブロックを含む、請求項 9 に記載のデータ処理システム。

**【請求項 11】**

前記仮想パーティションが、前記第 1 のパーティション内のデータのみを含む、請求項 10 に記載のデータ処理システム。

**【請求項 12】**

前記プログラムがさらに、前記コンピュータに、前記仮想パーティションと前記パーティションとを関連付けるパーティション構成データ構造体を確立することを行わせ、  
前記再分配することが、前記第 2 のパーティションに移動された前記データが仮想パーティション内に存在しないことを示すように、前記パーティション構成データ構造体を更新することを含む、  
請求項 9 に記載のデータ処理システム。

30

**【請求項 13】**

前記プログラムがさらに、前記コンピュータに、前記区分データベース内のデータを、複数のパーティション及び前記仮想パーティションのうちの特定の 1 つにマッピングする、パーティション・マップを確立することを行わせ、  
前記再分配することが、前記第 2 のパーティションに移動されたデータが仮想パーティションにマッピングされないことを示すように、前記パーティション・マップを更新することを含む、  
請求項 9 に記載のデータ処理システム。

40

**【請求項 14】**

前記再分配することが、  
前記仮想パーティション内のデータのバックアップを作成することと、  
前記データを前記バックアップから前記第 2 のパーティションに復元することと  
を含む、  
請求項 9 に記載のデータ処理システム。

**【請求項 15】**

前記仮想パーティション内の前記データが、第 1 のデータであり、  
前記再分配の前に、前記第 1 のパーティションが、前記第 1 のデータと、前記仮想パーティション内に存在しない第 2 のデータとを格納し、

50

前記バックアップを作成することが、前記第 1 のデータを含み、かつ前記第 2 のデータを除外するバックアップを作成することを含む、請求項 14 に記載のデータ処理システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般にデータ処理に関し、特に区分データベースの再分配に関する。

【背景技術】

【0002】

大量のデータが格納されるコンピュータ処理環境において、データは通常、リレーショナル・データベース管理システム (RDBMS) によって管理され、これを利用して、データの格納、アクセス及び管理のために 1 つ又は複数のデータベースをインスタンス化することができる。各データベースは、1 つ又は複数の表スペースを含み、これがリレーショナル・データ・モデルに従って表データを格納する。表形式の編成が含意するように、表データは、行及び列に論理的に配置され、各々の表の行は、関連付けられた行キーを有する。

10

【0003】

強化された管理性、性能及び/又は可用性を提供するために、リレーショナル・データベースは通常、各々が各自のデータ、インデックス、構成ファイル及びトランザクション・ログを有する、複数の論理又は物理パーティション (以降は、より明確な意味付けが必要とされる場合以外は単に「パーティション」と称する) に分けられている。所与のいずれの表の表データも、パーティションのうちの 1 つまたは複数の中に配置されることができ、表データが存在するパーティションは、典型的にはハッシュ関数によって定められる。データは、データベースのパーティションを横断して分配されているので、おそらくは多数のコンピュータ上にある、多数のプロセッサの力を連携して利用して、データベース内のデータを格納し、検索し、処理し及び管理することができる。

20

【0004】

オンライン・トランザクション処理 (OLTP) システム、データ・ウェアハウス企業、保険及び金融会社などのような、大量のデータ量を管理する企業では、格納データの量が増えるにつれて、そのデータ・ストレージ及び処理能力を拡張することが頻繁に必要となる。例えば、企業は、照会応答時間が長くなることを回避する一方で、増大するデータ量を扱うために、企業の既存の情報技術 (IT) インフラストラクチャに 1 つ又は複数の追加サーバ及びそれに関連付けられたストレージ・ノードを追加することがある。

30

【0005】

追加サーバを使えるようにするためには、RDBMS は、1 つ又は複数のデータベース・インスタンスを再分配及び再編成して、データベース・インスタンスが既存のサーバのストレージ・ノード上のみならず、新たに導入されたサーバのストレージ・ノード上にも存在するようにしなければならない。従来技術による、RDBMS がデータベースを再分配及び再編成するための従来のプロセスが、図 1 に示される。

【0006】

40

データベースの再分配及び再編成の従来のプロセスは、ブロック 100 で開始され、その後、ブロック 102 に進み、これは、RDBMS が、再分配されるデータベース全体のバックアップを作成することを示す。データベースのサイズによっては、データベースのバックアップの作成に、かなりの処理時間 (例えば、数日から数週間) を費やすことがある。プロセスは次にブロック 104 - 108 を含む反復ループに入り、ここで、データベースが、既存のストレージ・ノードと新たなストレージ・ノードとにわたって行ごとに再分配される。再分配は、ブロック 104 で開始し、ここで、RDBMS は、処理されるべき次のデータベース行のキー値を読み取る。RDBMS は次に、データベース行のキー値を再ハッシュして、再分配の後でそのデータベース行が存在することになるターゲット・パーティション番号を定める (ブロック 106)。ブロック 110 において、RDBMS

50

は、ターゲット・パーティション番号が既存のパーティション番号と同一かどうか、つまりデータベース行が移動しないかどうか、判断する。ターゲット・パーティション番号が既存のパーティション番号と一致する場合には、プロセスは、後述するブロック 118 に移る。しかしながら、ターゲット・パーティション番号が既存のパーティション番号と一致しない場合には、プロセスは、ブロック 112 - 116 へと進む。

【0007】

ブロック 112 - 116 において、RDBMS は、既存のストレージ・ノードから完全なデータベース行を読み出し、新たに追加されたストレージ・ノード上の新たなパーティションにデータベース行を挿入し、既存のストレージ・ノードからデータベース行を削除する。その後、ブロック 118 において、RDBMS は、データベースのすべての行が処理されたかどうか判断する。まだ処理されていなければ、プロセスは、前述のブロック 104 に戻る。しかしながら、ブロック 118 において、RDBMS が、データベースのすべての行が処理されたと判断した場合、プロセスは、ブロック 120 に進む。

【0008】

理解されるように、選択されたデータベース行を、ブロック 104 - 118 に示される再分配によって既存のストレージ・ノードから新たにインストールされたストレージ・ノードに移動させると、疎らにポピュレートされた、それゆえ十分に利用されていない既存のストレージ・ノードが残ることになる。したがって、ブロック 120 において、RDBMS は、既存のストレージ・ノード内のデータベース行を再編成して、データベースをコンパクトなストレージ編成に戻す。再編成が成功裡に完了したら、RDBMS は次に、ブロック 122 においてデータベース全体の第 2 のバックアップを作成する。さらに、ブロック 124 に示されるように、RDBMS は、ユーティリティを実行して、データベースに関連した統計量を収集し、表スペース、インデックス及びパーティションを再度特徴付けし、これらの統計量をカタログに記録する。最後に、ブロック 126 において、RDBMS は、パーティションを認識する任意のアプリケーション（例えば、Microsoft（登録商標）Internet Information Services（IIS））に、データベースが新たに追加されたストレージ・ノードにわたって再編成されたことを通知する。その後、データベースを再分配及び再編成するための従来プロセスは、ブロック 130 で終了する。

【0009】

図 2 - 図 4 は、従来技術による、新たに追加されたデータ・ストレージ・ノードにわたるデータベースを再分配及び再編成を示す。特に、図 2 は、データベースがポピュレートされた 4 つのデータベース・パーティション 202a - 202d を含むデータ・ストレージ・システム 200 を示す。データベースのサイズが現在インストールされているデータ・ストレージ・ノードの容量に近づいているので、データ・ウェアハウス企業は、追加のデータベース・パーティションをサポートするために、1 つまたは複数の追加ストレージ・ノードをデータ・ストレージ・システム 200 に追加することができる。

【0010】

図 3 に示される例において、データ・ウェアハウス企業は、4 つの追加データベース・パーティション 202e - 202h をサポートするために、1 つまたは複数の追加ストレージ・ノードをデータ・ストレージ・システム 200 に追加する。図 3 はさらに、図 1 のブロック 104 - 118 に示される従来の行ごとのデータベースの再分配に従うと、新たなデータベース・パーティション 202e - 202h に移動したデータベースの部分は密にコンパクト化されるが、元のデータベース・パーティション 202a - 202d 上に残ったデータベースの部分は疎らにポピュレートされているので、したがってデータ・ストレージ・システム 200 のストレージ容量の利用率が低下することを示している。したがって、図 1 のブロック 120 を参照して上述したように、RDBMS は、図 4 に示されるコンパクトで良好に分配されたデータベースを達成するためには、データベース部分 202a - 202d 上に存在する部分も再編成しなければならない。

【発明の概要】

**【発明が解決しようとする課題】****【0011】**

本発明の目的は、区分データベースを再分配する方法、プログラム及びデータ処理システムを提供することである。

**【課題を解決するための手段】****【0012】**

いくつかの実施形態において、区分データベースは、少なくとも論理的又は物理的な第1のデータ・ストレージ・ノード上の複数の論理又は物理パーティション内に格納され、複数の論理的パーティションの中の第1のパーティションのサブセットは、仮想パーティションとして構成される。区分データベースを格納するように第2の物理データ・ストレージ・ノードの割り付けを指示する入力を受け取る。第2のパーティションが、第2のデータ・ストレージ・ノード上で構成される。入力に応答して、区分データベースは、第1のパーティション上の仮想パーティション内のデータを第2のパーティションに移動させることにより、第1及び第2のデータ・ストレージ・ノードにわたって再分配される。

**【図面の簡単な説明】****【0013】**

【図1】従来技術による、データベースの再分配及び再編成のための従来プロセスの高次フローチャートを示す。

【図2】従来技術による、新たなデータ・ストレージ・ノードにわたる従来のデータベースの再分配及び再編成を示す。

【図3】従来技術による、新たなデータ・ストレージ・ノードにわたる従来のデータベースの再分配及び再編成を示す。

【図4】従来技術による、新たなデータ・ストレージ・ノードにわたる従来のデータベースの再分配及び再編成を示す。

【図5】1つの実施形態による、例示的なデータ処理環境を示す。

【図6】図5のデータ処理企業の例示的なデータ・ストレージ・ノードを示す。

【図7】1つの実施形態による、パーティション構成データ構造体の例示的な実施形態を示す。

【図8】1つの実施形態による、パーティション・マップの例示的な実施形態を示す。

【図9】データベースの再分配方法の第1の例示的な実施形態の高次論理フローチャートを示す。

【図10】図9に示された第1の例示的な方法による、データベースの例示的な再分配を示す。

【図11】データベースの再分配方法の第2の例示的な実施形態の高次論理フローチャートを示す。

【図12】図11に示された第2の例示的な方法による、データベースの例示的な再分配を示す。

**【発明を実施するための形態】****【0014】**

ここで、図面、特に図5を参照すると、1つの実施形態による、例示的なデータ処理環境300の高次ブロック図が示される。図示されるように、例示的なデータ処理環境300は、大量のデータを管理する企業、政府機関、非営利団体又は教育機関などのような組織により又はそれに代わって運営されることがある、データ処理企業310を含む。データ処理企業310は、通信のために、有線又は無線の構内又は広域通信網、携帯電話網及び/又は公衆交換電話網(PSTN)のような、1つまたは複数の回線交換又はパケット交換通信ネットワーク304に結合されている。したがって、データ処理企業310は、通信ネットワーク304を介して、機器302a-302d(例えば、サーバ・コンピュータ・システム、パーソナル・コンピュータ・システム、携帯型コンピュータ・システム、移動体電話、スマートフォン、地上通信電話)と通信することができる。

**【0015】**

機器 302a - 302d とデータ処理システム 310 との間の通信は、例えば P S T N 又はインターネットプロトコル上の音声通信 ( V o I P ) 接続を介した音声通信、及び / 又は、例えばインスタント・メッセージ、簡易メール転送プロトコル ( S M T P ) 又はハイパーテキスト転送プロトコル ( H T T P ) を介したデータ通信を含むことができる。例えば、データ処理企業 310 と機器 302a - 302d との間の通信は、機器 302a - 302d からデータ処理企業 310 へのデータ要求の送信と、データ処理企業 310 から機器 302a - 302d への応答データ ( 例えば、プログラム・コード、画像、グラフィックス、テキスト、音声、映像及び / 又はそのようなデータを含むファイル ) の送信とを含むことができる。

#### 【 0016 】

10

さらに図 5 を参照すると、データ処理企業 310 は、サーバ 312a - 312n のような 1 つ又は複数の物理コンピュータ・システムを含むことができ、これらは、通信のために、例えば、イントラネット、仮想私設ネットワーク ( V P N ) 又はソケット接続のようなケーブル及び / 又はネットワーク接続を含む、通信ファブリック 314 によって結合することができる。図示された例示的な実施形態において、サーバ 312a は、サーバ 312a が通信ネットワーク 304 及び通信ファブリック 314 を介して通信することを可能にする、1 つ又は複数のネットワーク・インタフェース 316 を含む。サーバ 312a は、例えば、1 つ又は複数のデータベースに編成されたデータを管理し、アクセスし及び操作するために、データ及びプログラム・コードを処理する 1 つ又は複数のプロセッサ 320 をさらに含む。サーバ 312a は、ポート、ディスプレイ及び外付け機器などのような 20 入力 / 出力 ( I / O ) デバイス 322 も含み、これが、入力を受け取り、かつサーバ 312a によって実行された処理の出力を提供する。最後に、サーバ 312a は、メモリ、固体ドライブ、光学又は磁気ディスク・ドライブ、テープ・ドライブなどを含む 1 つ又は複数の揮発性又は不揮発性ストレージ・デバイスを含むことができる、データ・ストレージ 330 を含む。

#### 【 0017 】

図示された実施形態において、データ・ストレージ 330 は、サーバ 312a のハードウェア・リソースを管理し、かつサーバ 312a 上で実行されるその他のソフトウェアに対して共通サービスを提供する、オペレーティング・システム ( O S ) 332 を格納する。例えば、O S 332 は、A I X ( 登録商標 )、L i n u x ( 登録商標 )、A n d r o i d ( 登録商標 ) 又は W i n d o w s ( 登録商標 ) オペレーティング・システムのうちの 1 つで実装することができる。データ・ストレージ 330 はまた、ニューヨーク州、アーモ 30 ンクの I B M コーポレーションから入手可能な D B 2 ( 登録商標 ) リレーショナル・データベース管理システム ( R D B M S ) のようなデータベース・マネージャ 340 も含み、これが、例示的なデータベース 350 のような 1 つ又は複数のデータベース内のデータを管理し、これにアクセスし、これを操作する。いくつかの実施形態において、データベース・マネージャ 340 は、O S 332 又はその他のソフトウェア・プログラムと一体化されたものとして行うことができる。データベース・マネージャ 340 は、データベース 350 30 に加えて、データベース 350 の種々の論理パーティションを定め、かつパーティションをデータ処理企業 310 の物理ストレージ・リソースに対してマッピングする、1 つ又は 40 複数のパーティション構成データ構造体 352 を維持する。データベース・マネージャ 340 は、さらに後述するように、データベース 350 の仮想パーティションをデータベース 350 の論理パーティションに対してマッピングする、パーティション・マップ 354 を随意に維持することもできる。

#### 【 0018 】

種々の実施形態において、データベース・マネージャ 340 及び / 又は O S 332 は、通信ファブリック 314 及び通信ネットワーク 304 を介したサーバ 312a と機器 302a - 302d との通信をサポートする、コードを含むことができる。いくつかの実施形態においては、O S 332 及び / 又はデータベース・マネージャ 340 内に必ずしも適切な通信能力が組み込まれているわけではないので、データ・ストレージ 330 は、サーバ 50

3 1 2 a が通信ファブリック 3 1 4 及び通信ネットワーク 3 0 4 を介して他のサーバ 3 1 2 及び機器 3 0 2 a - 3 0 2 d と通信することを可能にする、ウェブ・サーバ（例えば、Apache、IIS など）、音声自動応答（IVR）及び／又はその他のプログラム・コードのような、通信コード 3 4 2 をさらに含むことができる。特に、通信コード 3 4 2 は、実装されている場合、データベース・マネージャ 3 4 0 に対するデータベース照会の通信と、データベース・マネージャ 3 4 0 から要求者に対する応答データの通信とをサポートする。

#### 【0019】

データ・ストレージ 3 3 0 の内容は、いくつかの実施形態においてはサーバ 3 1 2 a 上にローカライズされていてもよく、他の実施形態においては、多重のサーバ 3 1 2 a - 3 1 2 n のデータ・ストレージ 3 3 0 にわたって分配されることを理解されたい。さらに、サーバ 3 1 2 a のデータ・ストレージ 3 3 0 内に図示された内容は、随意に、データ処理企業 3 1 0 のストレージ・エリア・ネットワーク（SAN）3 6 0 上に部分的に又は全体が常駐することができることを理解されたい。図示されるように、SAN 3 6 0 は、ストレージ要求を受け取り、かつ供給するスイッチ／コントローラ（SW/C）と、その各々が 1 つ又は複数の物理的不揮発性メモリドライブ、ハードディスクドライブ、光記憶ドライブ、テープドライブなどを含むことができる、多重のデータ・ストレージ・ノード 3 7 0 a - 3 7 0 k とを含む。いくつかの実施形態において、データ・ストレージ・ノード 3 7 0 a - 3 7 0 k は、このような物理ストレージ・リソースの仮想化された抽象を示す論理エンティティとすることができる。

#### 【0020】

上記の説明を検討すると、データ処理企業 3 1 0 が実体化される形態は、例えば、編成の種類、データベース 3 5 0 のサイズ、データベース 3 5 0 に照会することができる機器 3 0 2 a - 3 0 2 d の数などの 1 つ又は複数の要因に基づいて、実施形態ごとに様々であり得ることが理解されよう。例えば、1 つ又は複数の手持ち型、ノートブック型、デスクトップ型又はサーバのコンピュータ・システムを含むことができる、このような実装のすべてが、添付の特許請求の範囲において述べられる本発明の実施形態として企図される。

#### 【0021】

図 6 は、図 5 のデータ処理企業 3 1 0 内のデータ・ストレージ・ノード 4 0 0（例えば、SAN 3 6 0 のデータ・ストレージ・ノード 3 7 0 又はサーバ 3 1 2 のデータ・ストレージ 3 3 0 内のデータ・ストレージ・ノード）のより詳細な図を示す。図示された例において、データ・ストレージ・ノード 4 0 0 は、8 つの論理又は物理パーティションをホストし、以後、これらを、それぞれ LP 0 - LP 7 と番号付けされた論理パーティションと仮定する。RDBMS において、データ・ブロック B 0 - B 1 5 の各々は、共通の行キー・ハッシュを有する 1 つ又は複数のデータベース行に対応することができる。

#### 【0022】

本開示によれば、データベース・マネージャ 3 4 0 は、データ・ブロック B 0 - B 1 5 のサブセットを仮想パーティションに割り当てる。例えば、データベース・マネージャ 3 4 0 は、データ・ブロック B 0 - B 1 5 の各々を、VP 8 - VP 1 5 と番号付けされた 8 つの仮想パーティションのうちのそれぞれ 1 つに割り当てることができる。種々のシナリオにおいて、各仮想パーティションは、好ましくはすべてが共通の論理パーティション上に存在する、1 つ又は複数のデータ・ブロックを含むことができる。図 9 - 図 1 2 を参照して以下でさらに論じるように、データベース・マネージャ 3 4 0 は、仮想パーティションを参照することによって、データベース 3 5 0 を効率的に再分配することができる。

#### 【0023】

ここで図 7 を参照すると、1 つの実施形態によるパーティション構成データ構造体 3 5 2 の例示的な実施形態が示される。図示された実施形態において、例えば 1 つ又は複数のデータベース構成ファイルに実装することができる、パーティション構成データ構造体 3 5 2 は、データベース 3 5 0 の複数の論理パーティションを定め、かつ論理パーティションをデータ処理企業 3 1 0 の物理ストレージ・リソースに対してマッピングする、複数の



構成エントリ 500 を含む。

【0024】

例示的な実施形態において、パーティション構成データ構造体 352 の各構成エントリ 500 は、ノード番号フィールド 502、ホスト名フィールド 504、論理パーティション番号フィールド 506 及び仮想パーティション・フラグ 508 を含む、多数のフィールドを含む。ノード番号フィールド 502 は、データベース 350 のパーティションを一意に識別する整数を指定する。ノード番号を論理パーティションに限定する従来の区分データベースとは対照的に、ノード番号フィールド 502 は、好ましくは、データベース 350 の論理パーティション及び仮想パーティションの各々に対する一意のノード番号を含む。ホスト名フィールド 504 は、ノード番号フィールド 502 内で識別されたデータベース・パーティションの TCP/IP ホスト名（例えば、「サーバ A」）を識別する。さらに、論理ポート・フィールド 506 は、ノード番号フィールド 502 内で識別されたデータベース・パーティションに割り当てられた論理ポート（例えば、論理パーティション）を指定し、仮想パーティション・フラグ 508 は、ノード番号フィールド 502 内で指定されたパーティションが仮想パーティションであるかどうかを指定する。構成エントリ 500 は、論理パーティションへの通信パス及び / 又はオペレーティング・システム固有の情報といった追加の構成情報を提供する、1 つ又は複数の追加のフィールドを含むことができることを理解されたい。

10

【0025】

図 7 で示されるパーティション構成データ構造体 352 の例示的な実施形態を想定すると、図 6 のデータ・ストレージ・ノード 400 を記述するパーティション構成データ構造体 352 の部分は、以下の表 I に示されるように表すことができる。

20

【表 1】

表 I

ノード番号	ホスト名	論理ポート番号	VP
0	サーバ A	0	—
1	サーバ A	1	—
2	サーバ A	2	—
3	サーバ A	3	—
4	サーバ A	4	—
5	サーバ A	5	—
6	サーバ A	6	—
7	サーバ A	7	—
8	サーバ A	0	V
9	サーバ A	1	V
10	サーバ A	2	V
11	サーバ A	3	V
12	サーバ A	4	V
13	サーバ A	5	V
14	サーバ A	6	V
15	サーバ A	7	V

30

40

【0026】

ここで図 8 を参照すると、1 つの実施形態による例示的なパーティション・マップ 354 が示される。図示された実施形態において、データベース・マネージャ 340 は、パーティション・マップ 354 を、その各々がハッシュ値フィールド 602、仮想パーティション番号フィールド 604 及び論理パーティション番号フィールド 606 を含む複数の行

50

600を含む、ルックアップ表として実装する。したがって、各行600は、それぞれのハッシュ値（例えば、データベース350の行の行キーに由来するハッシュ値から導出される）を、論理パーティション番号と、適用できる場合には仮想パーティション番号とに関連付ける。例えば、0から4095までの範囲のハッシュ値と、図6に示されるような8つの論理パーティションLP0-LP7を実装するデータ・ストレージ・ノード400とを仮定すると、パーティション・マップ354は、以下の表IIにまとめられる値を格納する4096個の行600を含むことができる。

【表2】

表II

ハッシュ値	仮想パーティション番号	論理パーティション番号
0	—	0
1	—	1
2	—	2
3	—	3
...	...	...
7	—	7
8	8	0
9	9	1
10	10	2
...	...	...
15	15	7
16	—	0
...	...	...
23	—	7
24	8	0
...	...	...
31	15	7
...	...	...
4095	—	7

10

20

30

40

50

## 【0027】

ここで図9を参照すると、第1の実施形態による、データベースを再分配する例示的な方法の高次論理フローチャートが示される。図示された方法は、例えば、サーバ312の1つ又は複数のプロセッサ320によるデータベース・マネージャ340の実行を通じて、実施することができる。本明細書で提示される他の論理フローチャートと同様に、ステップは、厳密に経時的な順序ではなく論理的に記述されており、少なくともいくつかの実施形態において、1つ又は複数のステップは、同時に行われることもあり、又は図示されているのとは異なる順序で行われることもあることを理解されたい。

## 【0028】

図9に示されたプロセスは、ブロック700で開始し、その後、ブロック702に進み、これは、データベース・マネージャ340が、例えば、管理者の入力に応答して又は所定のデフォルトに基づいて自動的に、データベース350内で所望の数の仮想パーティションを構成することを示す。表Iで示され、かつ図6に描かれた例示的な区分データベース350において、データベース・マネージャ340は、仮想パーティションVP8-VP15をそれぞれ論理パーティションLP0-LP7内に確立するために、ブロック702において、パーティション構成データ構造体352の最後から8つのエントリ500を入力することができる。上記のように、仮想パーティションは、データベース350を格

納するために割り付けられた物理ストレージ容量が拡張したときに再分配されることになる、データベース 350 のデータを含む。仮想パーティションの数及び位置を構成すると共に、データベース・マネージャ 340 は、データのハッシュ値（例えば、行キー）とパーティション構成データ構造体 352 によって構成された論理及び仮想パーティションとの間ですばやくマッピングするために、随意にパーティション・マップ 354 を確立する（ブロック 704）。データベース・マネージャ 340 は、代替的に、各ハッシュ値に関連付けられる論理及び仮想パーティションを必要なときに計算することもできるので、ブロック 704 は随意である。

#### 【0029】

プロセスは、ブロック 704 からブロック 710 に進み、これは、受け取った入力、拡張された物理ストレージ容量にわたってデータベース 350 を再分配すべきことを指示しているかどうかを、データベース・マネージャ 340 が判断することを示す。理解されるように、データベース 350 を格納するために利用可能な拡張された物理ストレージ容量は、サーバ 312 をデータ処理企業 310 に追加すること、追加のデータ・ストレージ・ノード 370 を SAN 360 に追加すること、及び / 又はデータ処理企業 310 の既存のデータ・ストレージ・ノードを、データベース 350 を格納するように再割り付けすることによって、利用可能となり得る。拡張された物理ストレージ容量にわたってデータベース 350 を再分配すべきであることを指示する入力をデータベース・マネージャ 340 が検出しない場合には、プロセスはブロック 710 に留まる。プロセスがブロック 710 に留まっている間は、データベース・マネージャ 340 は、当該分野で公知のように、データベース 350 の構造化照会言語（SQL）クエリに応答するデータを提供すること、及び任意の要求された管理又は構成関数などを実行することを含む、従来のデータベース処理を行う。ブロック 710 におけるデータベース・マネージャ 340 による、拡張された物理ストレージ容量にわたってデータベース 350 を再分配すべきであることを指示する入力（例えば、ユーザ・コマンド）を受け取ったとの判断に応答して、プロセスは、ブロック 712 へと移る。

#### 【0030】

ブロック 712 は、データベース・マネージャ 340 が、データベース 350 を格納するために割り付けられた新たな物理ストレージ・ノード上に論理パーティションを確立することを示す。プロセスは次に、ブロック 720 - 730 を含むループへと入り、ここで、仮想パーティションが、既存の論理パーティションからブロック 712 で確立された新たな論理パーティションに再分配される。最初にブロック 720 を参照すると、データベース・マネージャ 340 は、例えばパーティション構成データ構造体 352 を参照することにより、データベース 350 のすべての仮想パーティションが処理されたか否かを判断する。データベース 350 のすべての仮想パーティションがすでに処理されたとブロック 720 においてデータベース・マネージャ 340 が判断したことに応答して、プロセスは、ブロック 720 から、後述される 740 に進む。しかしながら、まだデータベース 350 のすべての仮想パーティションが処理されてはいないとブロック 720 においてデータベース・マネージャ 340 が判断した場合には、データベース・マネージャ 340 は、処理のための仮想パーティション、例えば、パーティション構成データ構造体 352 内にリストされた次の仮想パーティションを選択する（ブロック 722）。

#### 【0031】

ブロック 724 において、データベース・マネージャ 340 は、例えば、仮想パーティション番号が、新たに割り付けられたストレージ・ノード上に確立された論理パーティションのうちの 1 つに割り当てられた論理パーティション番号と一致するか否かを判断することによって、処理のために選択された仮想パーティションを移動させるか否かを判断する。現在選択された仮想パーティションを移動させないとの判断に応答して、プロセスは、すでに説明したブロック 720 に戻る。しかしながら、ブロック 724 においてデータベース・マネージャ 340 が、選択された仮想パーティションを移動すべきであると判断した場合、プロセスはブロック 726 へと移る。ブロック 726 は、データベース・マネ

ージャ 3 4 0 が、既存の論理パーティションから、一致する論理パーティション番号を有する論理パーティションへの逐次的なアクセス動作を用いて、仮想パーティションのデータを移動させることを示す。データベース・マネージャ 3 4 0 は次に、移動されたパーティションに関連してデータ・ストレージ・ノード上に格納されたメタデータを更新し（ブロック 7 2 8）、パーティション構成データ構造体 3 5 2 内の関連付けられた構成エントリ 5 0 0 の仮想パーティション・フラグ 5 0 8 をクリアする（ブロック 7 3 0）。結果として、移動されたパーティションは、もはや仮想パーティションではなくなり、新たに割り付けられたデータ・ストレージ・ノード上の論理パーティションのうちの 1 つのデータ・ブロックへと変換される。プロセスは、ブロック 7 3 0 からブロック 7 2 0 へと戻り、これは、データベース・マネージャ 3 4 0 が、もしあれば次の仮想パーティションを処理することを示す。

10

#### 【 0 0 3 2 】

ブロック 7 2 0 においてデータベース・マネージャ 3 4 0 が、すべての仮想パーティションがすでに処理されたと判断したことに応答して、データベース・マネージャ 3 4 0 は、ハッシュ値と論理及び仮想パーティション番号との間の修正された関係を反映するようにパーティション・マップ 3 5 4 を更新する（ブロック 7 4 0）。ブロック 7 4 0 に続いて、図 9 に示されたプロセスは、ブロック 7 4 2 で終了する。

#### 【 0 0 3 3 】

図 1 0 は、図 9 に示された例示的な方法による、データベース 3 5 0 の例示的な再分配を示す。図 1 0 に示される例において、データベース 3 5 0 は、図 6 を参照して先に論じたように、元は、データ・ストレージ・ノード 4 0 0 の 8 つの論理パーティション（すなわち L P 0 - L P 7）上のみに格納されている。論理パーティション L P 0 - L P 7 上で、データベース・マネージャ 3 4 0 は、データ・ブロック B 8 - B 1 5 をそれぞれ仮想パーティション V P 8 - V P 1 5 として構成する。

20

#### 【 0 0 3 4 】

データベース 3 5 0 を収容するために利用可能なデータ処理環境 3 1 0 の物理データ・ストレージ容量は、その後、追加のデータ・ストレージ・ノード 8 0 0 を含むように拡張される。ブロック 7 1 2 に関して述べたように、データベース・マネージャ 3 4 0 は、データ・ストレージ・ノード 8 0 0 を、8 つの論理パーティション番号 L P 8 - L P 1 5 で構成する。さらに、図 9 のブロック 7 2 0 - 7 3 0 に従って、データベース・マネージャ 3 4 0 は、仮想パーティション V P 8 - V P 1 5（それぞれデータ・ブロック B 8 - B 1 5 に対応する）の各々を、データ・ストレージ・ノード 8 0 0 上の論理パーティション L P 8 - L P 1 5 のそれぞれ 1 つの上に再分配し、データ・ブロック B 0 - B 7 をデータ・ストレージ・ノード 4 0 0 の論理パーティション L P 0 - L P 7 上に残す。

30

#### 【 0 0 3 5 】

サーバ 3 1 2 上にあるデータ・ストレージ・ノード 8 0 0 がホスト名「サーバ B」を有すると仮定すると、データベース・マネージャ 3 4 0 は、パーティション構成データ構造体 3 5 2 を上記の表 I でまとめられた状態から下記の表 I I I に示される状態に更新する。

。

【表 3】

表 I I I

ノード番号	ホスト名	論理ポート番号	V P
0	サーバA	0	—
1	サーバA	1	—
2	サーバA	2	—
3	サーバA	3	—
4	サーバA	4	—
5	サーバA	5	—
6	サーバA	6	—
7	サーバA	7	—
8	サーバB	0	—
9	サーバB	1	—
10	サーバB	2	—
11	サーバB	3	—
12	サーバB	4	—
13	サーバB	5	—
14	サーバB	6	—
15	サーバB	7	—

10

20

さらに、データベース・マネージャ 3 4 0 は、パーティション・マップ 3 5 4 を上記の表 I I にまとめられた状態から下記の表 I V に示される状態に更新する。

【表 4】

表 I V

ハッシュ値	仮想パーティション番号	論理パーティション番号
0	—	0
1	—	1
2	—	2
3	—	3
...	...	...
7	—	7
8	—	0
9	—	1
1 0	—	2
...	...	...
1 5	—	7
1 6	—	0
...	...	...
2 3	—	7
2 4	—	0
...	...	...
3 1	—	7
...	...	...
4 0 9 5	—	7

10

20

30

40

50

## 【 0 0 3 6 】

図 1 と図 9 - 図 1 0 とを比較することによって、図 9 に示される例示的なプロセスが、処理が集中する図 1 のステップの多くを不要にすることに留意されたい。例えば、図 9 において、ブロック 1 0 6 に示されるようにデータベース 3 5 0 の行を再ハッシュする必要はない。さらに、図 1 のブロック 1 0 2 及び 1 2 2 に示されるように、データベース 3 5 0 の再分配の前又は後にデータベース 3 5 0 をバックアップする必要もない。その上、ブロック 1 2 0 に示されるようにデータベース 3 5 0 の再編成を行う必要もなく、ブロック 1 2 4 に示されるようにデータベース統計量を更新する必要もない。最後に、ブロック 1 2 6 に示されるようにパーティションを認識するアプリケーションを更新する必要もない。

## 【 0 0 3 7 】

図 1 1 を参照すると、第 2 の実施形態による、データベースを再分配する例維持的な方法の高次論理フローチャートが示される。特に、図 1 1 に示されるプロセスは、データベース 3 5 0 の仮想パーティションのバックアップ及び復元を通じて、データベース 3 5 0 を新たに割り付けられた 1 つ又は複数の物理ストレージ・ノード上に再分配する。

## 【 0 0 3 8 】

図 1 1 に示されるプロセスは、ブロック 9 0 0 で開始し、その後、以前に説明したように、例えば管理者の入力に応答して、データベース 3 5 0 内で所望の数の仮想パーティションを構成することを示す、ブロック 9 0 2 に進む。仮想パーティションの数及び位置を構成すると共に、データベース・マネージャ 3 4 0 は、データのハッシュ値（例えば、行キー）とパーティション構成データ構造体 3 5 2 によって構成された論理及び仮想パーティションとの間ですばやくマッピングするために、随意にパーティション・マップ 3 5 4 を確立する（ブロック 9 0 4）。プロセスは、ブロック 9 0 4 からブロック 9 1 0 に進み、これは、受け取った入力、拡張された物理ストレージ容量にわたってデータベース 3

50を再分配すべきことを指示しているかどうかを、データベース・マネージャ340が判断することを示す。拡張された物理ストレージ容量にわたってデータベース350を再分配すべきであることを指示する入力をデータベース・マネージャ340が検出しない場合には、プロセスはブロック910に留まる（この時間の間は、データベース・マネージャ340は、他の従来のデータベース管理動作を行うことができる）。

#### 【0039】

ブロック910におけるデータベース・マネージャ340による、拡張された物理ストレージ容量にわたってデータベース350を再分配すべきであることを指示する入力を受け取ったとの判断に応答して、プロセスは、ブロック912へと移る。ブロック912は、データベース・マネージャ340が、データベース350を格納するために割り付けられた新たな物理ストレージ・ノード上に論理パーティションを確立することを示す。プロセスは次に、ブロック920 - 930を含むループへと入り、ここで、仮想パーティションが、ブロック912で確立された既存の論理パーティションからバックアップされる。最初にブロック920を参照すると、データベース・マネージャ340は、例えばパーティション構成データ構造体352を参照することにより、データベース350のすべての仮想パーティションが処理されたか否かを判断する。データベース350のすべての仮想パーティションがすでに処理されたとブロック920においてデータベース・マネージャ340が判断したことに応答して、プロセスは、ブロック920から、後述される940に進む。しかしながら、まだデータベース350のすべての仮想パーティションが処理されてはいないとブロック920においてデータベース・マネージャ340が判断した場合には、データベース・マネージャ340は、処理のための仮想パーティション、例えば、パーティション構成データ構造体352内にリストされた次の仮想パーティションを選択する（ブロック922）。次に、データベース・マネージャ340は、選択された仮想パーティションのバックアップを作成するが、好ましくは、そのバックアップから、仮想パーティションをホストする残りの論理パーティションを除外する（ブロック926）。データベース・マネージャ340は次に、選択された仮想パーティションに関連付けられたパーティション構成データ構造体352内の仮想パーティション・フラグ508をクリアする（ブロック930）。プロセスは、ブロック930からブロック920へと戻り、これは、データベース・マネージャ340が、もしあれば次の仮想パーティションを処理することを示す。

#### 【0040】

ブロック920においてデータベース・マネージャ340が、既存の物理データ・ストレージ・ノード上のデータベース350のすべての仮想パーティションがすでに処理されたと判断したことに応答して、データベース・マネージャ340は、仮想パーティションの各々を、ブロック926で作成したバックアップから、データ処理企業310の新たに割り付けられた物理ストレージ・ノードのそれぞれの論理パーティション（例えば、バックアップされた仮想パーティションの仮想パーティション番号と一致する論理パーティション番号を有する論理パーティション）へと復元する。結果として、移動されたパーティションは、もはや仮想パーティションではなくなり、新たに割り付けられたデータ・ストレージ・ノードの論理パーティション上のデータ・ブロックへと変換される。データベース・マネージャ340は次に、復元されたパーティションに関連してデータ・ストレージ・ノード上に格納されたメタデータを更新し（ブロック942）、移動されたパーティションを既存の物理ストレージ・ノードから削除する（ブロック944）。データベース・マネージャ340は、もしあれば、ハッシュ値と論理及び仮想パーティション番号との間の修正された関係を反映するようにパーティション・マップ354を更新する（ブロック946）。ブロック946に続いて、図11に示されたプロセスは、ブロック950で終了する。

#### 【0041】

図12は、図11に示される第2の例示的な方法による、データベースの例示的な再分配を示す。図12に示される例において、データベース350は、図6及び図10を参照

10

20

30

40

50

して先に論じたように、元は、データ・ストレージ・ノード 400 の 8 つの論理パーティション（すなわち LP0 - LP7）上のみに格納されている。論理パーティション LP0 - LP7 上で、データベース・マネージャ 340 は、データ・ブロック B8 - B15 をそれぞれ仮想パーティション VP8 - VP15 として構成する。

#### 【0042】

データベース 350 を収容するため割り付けられたデータ処理環境 310 の物理データ・ストレージ容量は、その後、追加のデータ・ストレージ・ノード 800 を含むように拡張され、データベース・マネージャ 340 は、その上で LP8 - LP15 で番号付けされる 8 つの論理パーティションを構成する。図 11 のブロック 920 - 930 に従って、データベース・マネージャ 340 は、仮想パーティション VP8 - VP15（それぞれ、データ・ブロック B8 - B15 に対応する）の各々のバックアップを含み、好ましくは論理パーティション LP0 - LP7 上に存在するその他のデータを除外する、仮想パーティション・バックアップ 1000 を作成する。ホストの論理パーティションに戻す従来の復元を行うのではなく、データベース・マネージャ 340 は、各仮想パーティションを、仮想パーティション・バックアップ 1000 から、データ・ストレージ・ノード 800 上の論理パーティション LP8 - LP15 のうちのそれぞれ 1 つへと復元する。さらに、データベース・マネージャ 340 は、対応する仮想パーティションをデータ・ストレージ・ノード 400 から削除し、データ・ブロック B0 - B7 をデータ・ストレージ・ノード 400 上の論理パーティション LP0 - LP7 上に残す。この方法において、データベース・マネージャ 340 は、図 10 に示すようにデータを直接移動させることによるのではなく、そのバックアップ能力を強化することによって、データベース 350 の仮想パーティションを既存の物理データ・ストレージ・ノード 400 から新たに割り付けられた物理データ・ストレージ・ノード 800 上に再分配する。しかしながら、得られるパーティション構成データ構造体 352 及びパーティション・マップ 354 は、上記の表 IIII 及び表 IV にまとめられたものと同一になる。

#### 【0043】

説明してきたように、少なくともいくつかの実施形態において、区分データベースは、少なくとも論理的又は物理的な第 1 データ・ストレージ・ノード上の複数の論理又は物理パーティション内に格納され、複数の論理パーティションの中の第 1 のパーティションのサブセットは、仮想パーティションとして構成される。区分データベースを格納するように第 2 の物理データ・ストレージ・ノードの割り付けを指示する入力を受け取る。第 2 のパーティションが、第 2 のデータ・ストレージ・ノード上で構成される。入力に応答して、区分データベースは、第 1 のパーティション上の仮想パーティション内のデータを第 2 のパーティションに移動させることにより、第 1 及び第 2 のデータ・ストレージ・ノードにわたって再分配される。

#### 【0044】

本発明を、1 つ又は複数の好ましい実施形態を参照して説明しながら具体的に示してきたが、当業者であれば、本発明の真意及び範囲から逸脱することなく、その形態及び詳細における種々の変更を行うことができることが理解されよう。例えば、本発明の機能を命令するプログラム・コードを実行するコンピュータ・システムに関して態様を説明してきたが、本発明は代替的に、データ処理システムで処理されて、本発明の機能を実行することができるプログラム・コードを格納した、有形の一過性ではないデータ・ストレージ媒体（例えば、光又は磁気ディスク又はメモリ）を含む、プログラム製品として実装することもできることを理解されたい。

#### 【符号の説明】

#### 【0045】

200：データ・ストレージ・システム  
202：データベース・パーティション  
300：データ処理環境  
400、800：データ・ストレージ・ノード

10

20

30

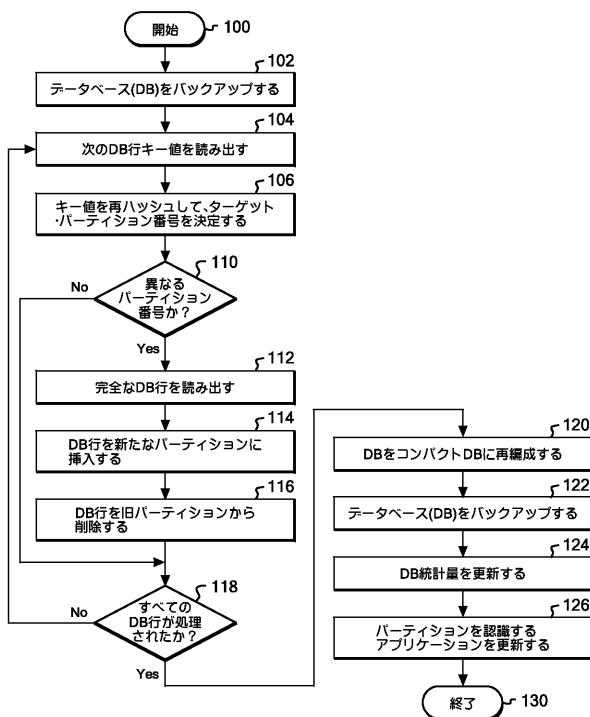
40

50

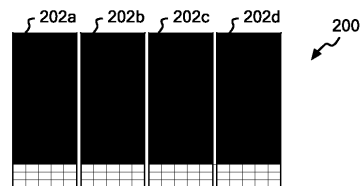


3 5 2 : パーティション構成データ構造体  
 3 5 4 : パーティション・マップ  
 5 0 0 : 構成エントリ

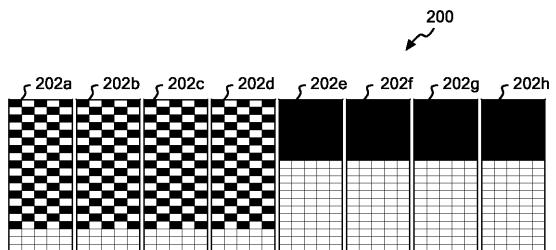
【 図 1 】



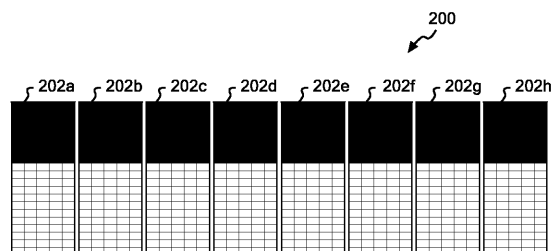
【 図 2 】



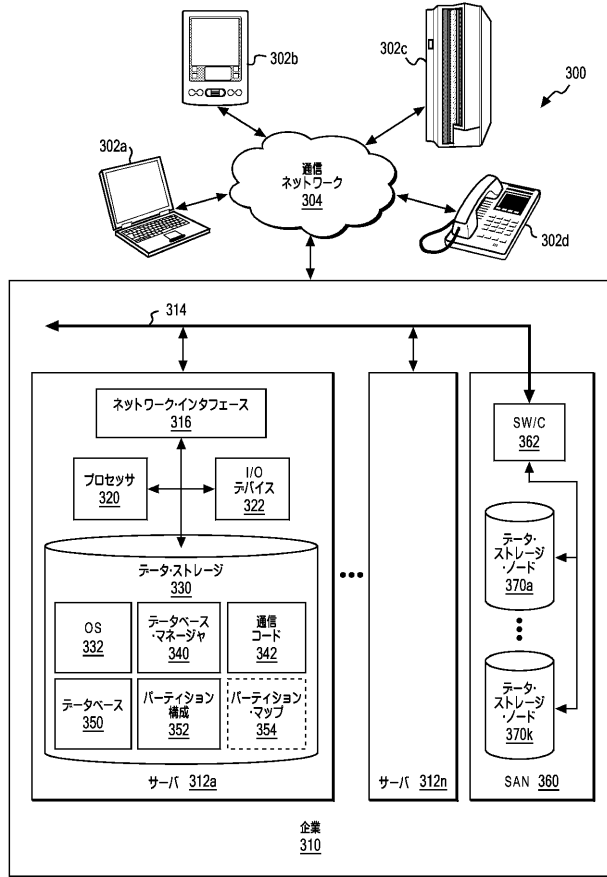
【 図 3 】



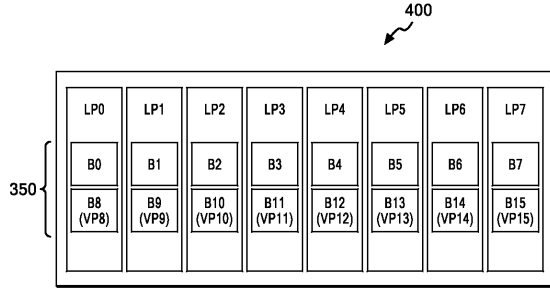
【 図 4 】



【図 5】



【図 6】



【図 7】

ノード番号 502	ホスト名 504	論理ポート番号 506	VP 508
⋮	⋮	⋮	⋮

500 (bracketed on the right)

352 (bracketed on the right)

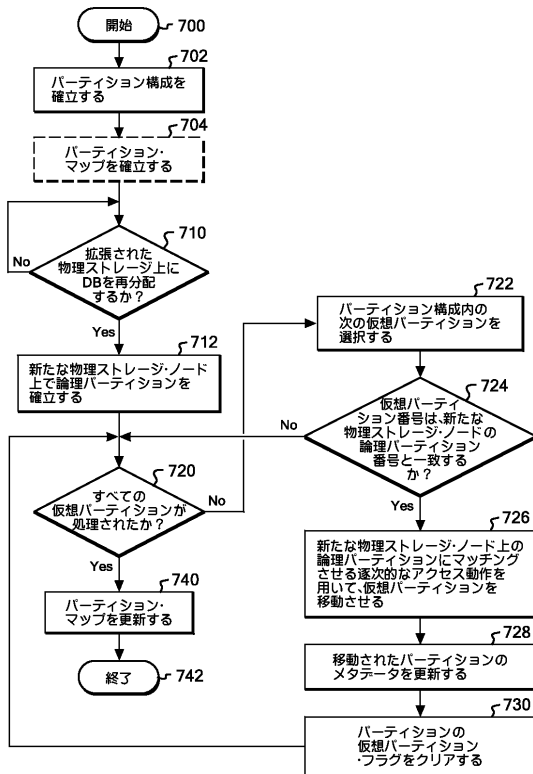
【図 8】

ハッシュ値 602	仮想パーティション番号 604	論理パーティション番号 606
⋮	⋮	⋮

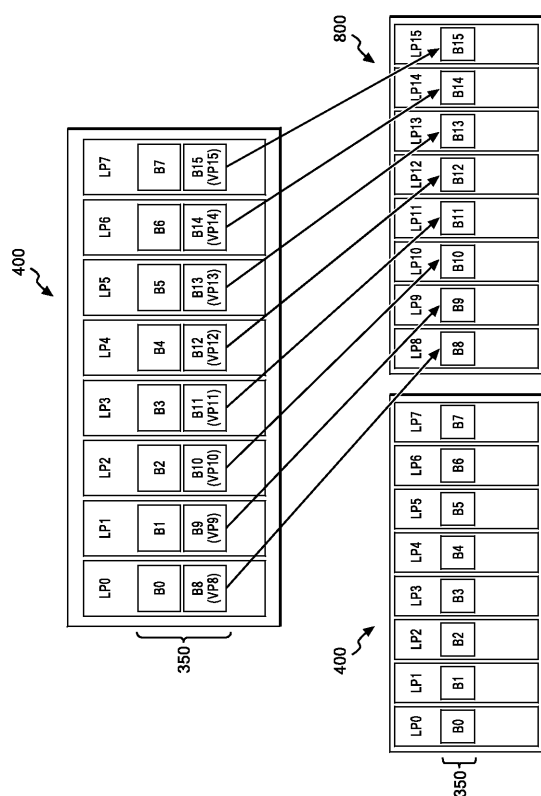
600 (bracketed on the right)

354 (bracketed on the right)

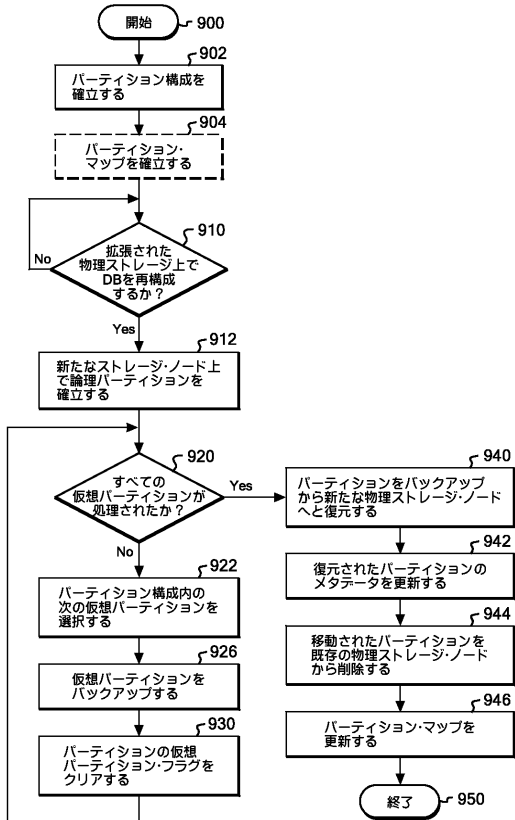
【図 9】



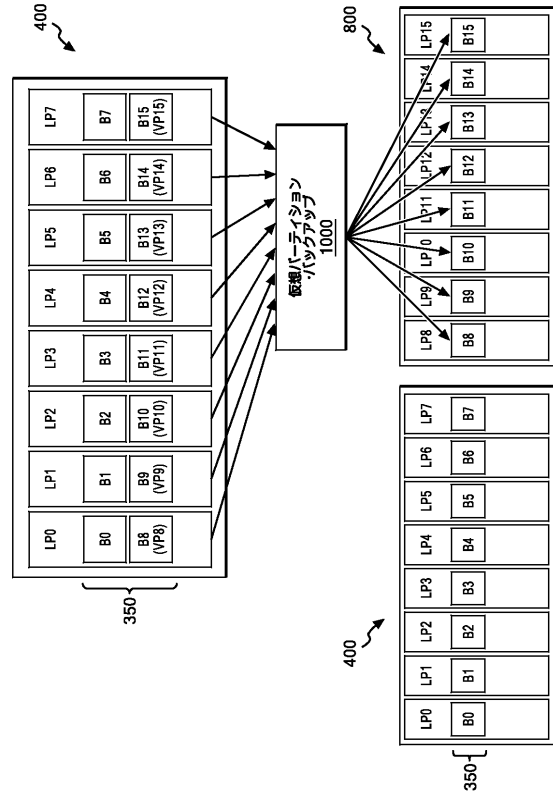
【図 10】



【図 1 1】



【図 1 2】



---

フロントページの続き

(72)発明者 ニーラジュ・シャルマ

インド共和国 5 6 0 0 4 5 バンガロール アウター・リング・ロード ラクエナハリ・ナガワ  
ラ マニヤタ・ビジネス・パーク ディー 2 郵便受け 2 - 1 エー - 1 9 3

(72)発明者 サウラブ・ジェイン

インド共和国 5 6 0 0 4 5 バンガロール アウター・リング・ロード ラクエナハリ・ナガワ  
ラ マニヤタ・ビジネス・パーク ディー 2 郵便受け 2 - 1 エー - 1 9 3

F ターム(参考) 5B065 ZA01