

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6258549号
(P6258549)

(45) 発行日 平成30年1月10日(2018.1.10)

(24) 登録日 平成29年12月15日(2017.12.15)

(51) Int.Cl. F I
G06F 17/30 (2006.01) G06F 17/30 220Z
 G06F 17/30 110C

請求項の数 31 (全 26 頁)

<p>(21) 出願番号 特願2017-153599 (P2017-153599) (22) 出願日 平成29年8月8日(2017.8.8) (62) 分割の表示 特願2016-502660 (P2016-502660) の分割 原出願日 平成26年3月14日(2014.3.14) (65) 公開番号 特開2017-224331 (P2017-224331A) (43) 公開日 平成29年12月21日(2017.12.21) 審査請求日 平成29年8月30日(2017.8.30) (31) 優先権主張番号 13/835,590 (32) 優先日 平成25年3月15日(2013.3.15) (33) 優先権主張国 米国 (US) 早期審査対象出願</p>	<p>(73) 特許権者 515160389 インフォマティカ エルエルシー アメリカ合衆国 カリフォルニア州 94 063、レッドウッド シティ、シーポー ト ブルバード 2100 2100 Seaport Blvd, R edwood City, CA 9406 3 U. S. A. (74) 代理人 100120662 弁理士 川上 桂子 (72) 発明者 グロンディン、 リチャード カナダ ジェイ3イー 2ゼット1、ケベ ック、サンテージュリ、リュ ヴォ克蘭 3 最終頁に続く</p>
---	--

(54) 【発明の名称】 個別のデータ値に対する効率よい演算を行うための方法、装置、並びにコンピュータ可読媒体

(57) 【特許請求の範囲】

【請求項1】

一または複数の演算装置によって個別のデータ値に対して演算を効率良く行うための方法であって、前記方法は、

前記一または複数の演算装置の少なくとも一つが、テーブル内のデータカラムを宛先とした照会を受領する工程であって、前記データカラムはデータベースのドメインに対応し、前記照会は、前記照会に回答して検索されたデータをグループ化するための一または複数のグループセットを定義し、前記一または複数のグループセットのそれぞれのグループは、前記データベースの一または複数の他のドメインに関連付けられた一または複数の値の固有のグループに対応している、工程と、

前記一または複数の演算装置の少なくとも一つが、前記テーブルにおいて前記グループセットに対応する一または複数の値の固有のグループを含む列に存在する前記データカラム内のいずれかの値を特定することにより、前記一または複数のグループセット中の各グループセットについてエンティティマップベクトルを生成する工程であって、各エンティティマップベクトルの長さは前記データカラムに対応する前記ドメイン内の固有データ値の数と等しく、各エンティティマップベクトル内の各ビット位置は、固有データ値が語彙順に配列された一覧中の対応する固有データ値の語彙順位置に対応し、各エンティティマップベクトル中の各ビットの値は、前記グループセットにおける対応する固有データ値の有無を示す工程と、
 を含む、方法。

【請求項 2】

エンティティマップベクトル中のビットが「1」という値であれば前記グループセットに固有データ値が有ることを示し、エンティティマップベクトル中のビットが「0」という値であれば前記グループセットに固有データ値が無いことを示す、請求項 1 に記載の方法。

【請求項 3】

グループセット中の個別のデータ値の数は、当該グループセットに対応するエンティティマップベクトルの前記ビットの各々の値の合計に等しい、請求項 1 に記載の方法。

【請求項 4】

前記一または複数のグループセットは少なくとも二つのグループセットを含み、前記方法は、

前記グループセットの二つ以上からのデータを集計して、合成グループセットにするためのコマンドを、前記一または複数の演算装置の少なくとも一つが受領する工程と、

前記二つ以上のグループセットに対応するエンティティマップベクトルについて、前記一または複数の演算装置の少なくとも一つが、ブールOR演算を行い、合成エンティティマップベクトルを生成する工程と、

をさらに含む、請求項 1 に記載の方法。

【請求項 5】

前記合成グループセット中の個別のデータ値の数は、前記合成エンティティマップベクトルの前記ビットの各々の値の合計に等しい、請求項 4 に記載の方法。

【請求項 6】

前記テーブルはパーティション化され複数の演算ノードにわたって分散されており、前記データカラムはカラムパーティションのセットの中の第一カラムパーティションであり、前記ドメインはドメインパーティションのセットの中の第一ドメインパーティションであり、前記一または複数のエンティティマップベクトルは一または複数の第一パーティションエンティティマップベクトルであり、

前記方法は、

前記一または複数の演算装置の少なくとも一つが、前記第一ドメインパーティションを前記ドメインパーティションのセットの中の一または複数のその他のドメインパーティションと合成し、グローバルドメインを生成する工程と、

前記一または複数の演算装置の少なくとも一つが、前記一または複数の第一パーティションエンティティマップベクトルに対応する一または複数のグローバルエンティティマップベクトルを生成する工程であって、各グローバルエンティティマップベクトルの長さは前記グローバルドメインにおける固有データ値の数に等しく、前記グローバルエンティティマップベクトルにおける各ビット位置は、前記グローバルドメイン内の固有データ値が語彙順に配列された一覧中の対応する固有データ値の語彙順位置に対応し、前記エンティティマップベクトル内の各ビットの値は、前記グループセットにおける対応する固有データ値の有無を示す、工程と、

をさらに含む、請求項 1 に記載の方法。

【請求項 7】

前記一または複数の演算装置の少なくとも一つが、前記第一パーティションに対応する前記一または複数のグローバルエンティティマップベクトルと、第二パーティションに対応する一または複数の第二グローバルエンティティマップベクトルとに、ブールOR演算を行い、一または複数の合成グローバルエンティティマップベクトルを生成する工程をさらに含む、請求項 6 に記載の方法。

【請求項 8】

前記合成グローバルエンティティマップベクトルの各々の各ビットの値の合計は、前記第一パーティションと前記第二パーティションとにわたる当該グローバルエンティティマップベクトルに対応するあるグループセットの中の個別の要素の総数である、請求項 7 に記載の方法。

10

20

30

40

50

【請求項 9】

グループセットの数は、少なくとも部分的に、前記テーブル中の第二データカラム中の個別のデータ値の数によって決定される。請求項 1 に記載の方法。

【請求項 10】

前記テーブル中の一または複数の第二データカラムを参照して前記グループセットが定義される、請求項 1 に記載の方法。

【請求項 11】

個別のデータ値に対して演算を効率良く行うための装置であって、前記装置は、
一または複数のプロセッサと、

前記一または複数のプロセッサの少なくとも一つと動作的に結合され、命令を保存する、一または複数のメモリとを備え、前記命令は、前記一または複数のプロセッサの少なくとも一つによって実行される時に、前記一または複数のプロセッサの少なくとも一つに、テーブル内のデータカラムを宛先とした照会を受領し、前記データカラムはデータベースのドメインに対応し、前記照会は、前記照会に回答して検索されたデータをグループ化するための一または複数のグループセットを定義し、前記一または複数のグループセットのそれぞれのグループは、前記データベースの一または複数の他のドメインに関連付けられた一または複数の値の固有のグループに対応し、

前記テーブルにおいて前記グループセットに対応する一または複数の値の固有のグループを含む列に存在するデータカラム内のいずれかの値を特定することにより、前記一または複数のグループセット中の各グループセットについてエンティティマップベクトルを生成し、各エンティティマップベクトルの長さは前記データカラムに対応する前記ドメイン内の固有データ値の数と等しく、各エンティティマップベクトル中の各ビットの位置は、前記固有データ値が語彙順に配列された一覧中の対応する固有データ値の語彙順位置に対応し、各エンティティマップベクトル中の各ビットの値は、前記グループセットにおける対応する固有データ値の有無を示す、
ようにさせる装置。

【請求項 12】

エンティティマップベクトル中のビットが「1」という値であれば前記グループセットに固有データ値が有ることを示し、エンティティマップベクトル中のビットが「0」という値であれば前記グループセットに固有データ値が無いことを示す、請求項 11 に記載の装置。

【請求項 13】

グループセット中の個別のデータ値の数は、当該グループセットに対応する前記エンティティマップベクトルの前記ビットの各々の値の合計に等しい、請求項 11 に記載の装置。

【請求項 14】

前記一または複数のグループセットは複数のグループセットを含み、前記一または複数のメモリはさらに命令を保存し、前記命令は、前記一または複数のプロセッサの少なくとも一つによって実行される時に、前記一または複数のプロセッサの少なくとも一つに、
前記グループセットの二つ以上からのデータを集計して、合成グループセットにするためのコマンドを受領し、

前記二つ以上のグループセットに対応するエンティティマップベクトルについて、ブールOR演算を行い、合成エンティティマップベクトルを生成する、
ようにさせる、請求項 11 に記載の装置。

【請求項 15】

前記合成グループセット中の個別のデータ値の数は、前記合成エンティティマップベクトルの前記ビットの各々の値の合計に等しい、請求項 14 に記載の装置。

【請求項 16】

前記テーブルはパーティション化され複数の演算ノードにわたって分散されており、前記データカラムはカラムパーティションのセットの中の第一カラムパーティションであり

、前記ドメインはドメインパーティションのセットの中の第一ドメインパーティションであり、前記一または複数のエンティティマップベクトルは一または複数の第一パーティションエンティティマップベクトルであって、前記一または複数のメモリはさらに命令を保存し、前記命令は、前記一または複数のプロセッサの少なくとも一つによって実行される時に、前記一または複数のプロセッサの少なくとも一つに、

前記第一ドメインパーティションを前記ドメインパーティションのセットの中の一または複数のその他のドメインパーティションと合成し、グローバルドメインを生成し、

前記一または複数の第一パーティションエンティティマップベクトルに対応する一または複数のグローバルエンティティマップベクトルを生成し、

各グローバルエンティティマップベクトルの長さは前記グローバルドメイン内の固有データ値の数に等しく、前記グローバルエンティティマップベクトル中の各ビットの位置は、前記グローバルドメイン中の前記固有データ値が語彙順に配列された一覧中の対応する固有データ値の語彙順位置に対応し、前記エンティティマップベクトルの各ビットの値は、前記グループセット中における対応する固有データ値の有無を示す、

ようにさせる、請求項 1 1 に記載の装置。

【請求項 1 7】

前記一または複数のメモリはさらに命令を保存し、前記命令は、前記一または複数のプロセッサの少なくとも一つによって実行される時に、前記一または複数のプロセッサの少なくとも一つに、

前記第一パーティションに対応する前記一または複数のグローバルエンティティマップベクトルと、第二パーティションに対応する一または複数の第二グローバルエンティティマップベクトルとに、ブールOR演算を行い、一または複数の合成グローバルエンティティマップベクトルを生成するようにさせる、

請求項 1 6 に記載の装置。

【請求項 1 8】

前記合成グローバルエンティティマップベクトルの各々の各ビットの値の合計は、前記第一パーティションと前記第二パーティションとにわたる当該グローバルエンティティマップベクトルに対応するあるグループセットの中の個別の要素の総数である、請求項 1 7 に記載の装置。

【請求項 1 9】

グループセットの数は、少なくとも部分的に、前記テーブル中の第二データカラム中の個別のデータ値の数によって決定される、請求項 1 1 に記載の装置。

【請求項 2 0】

前記テーブル中の一または複数の第二データカラムを参照して前記グループセットが定義される、請求項 1 1 に記載の装置。

【請求項 2 1】

コンピュータ可読命令を含むコンピュータプログラムであって、前記命令は、一または複数の演算装置によって実行される時に前記一または複数の演算装置の少なくとも一つに

、
テーブル内のデータカラムを宛先とした照会を受領し、前記データカラムはデータベースのドメインに対応し、前記照会は、前記照会に応答して検索されたデータをグループ化するための一または複数のグループセットを定義し、前記一または複数のグループセットのそれぞれのグループは、前記データベースの一または複数の他のドメインに関連付けられた一または複数の値の固有のグループに対応し、

前記テーブルにおいて前記グループセットに対応する一または複数の値の固有のグループを含む列に存在する前記データカラム内のいずれかの値を特定することにより、前記一または複数のグループセット中の各グループセットについてエンティティマップベクトルを生成し、各エンティティマップベクトルの長さは前記データカラムに対応する前記ドメイン内の固有データ値の数と等しく、各エンティティマップベクトル内の各ビット位置は、固有データ値が語彙順に配列された一覧中の対応する固有データ値の語彙順位置に対応

10

20

30

40

50

し、各エンティティマップベクトル中の各ビットの値は、前記グループセットにおける対応する固有データ値の有無を示す、
 ようにさせる、コンピュータプログラム。

【請求項 2 2】

エンティティマップベクトル中のビットが「1」という値であれば前記グループセットに固有データ値が有ることを示し、エンティティマップベクトル中のビットが「0」という値であれば前記グループセットに固有データ値が無いことを示す、請求項 2 1 に記載のコンピュータプログラム。

【請求項 2 3】

グループセット中の個別のデータ値の数は、当該グループセットに対応するビットマップベクトルの前記ビットの各々の値の合計に等しい、請求項 2 1 に記載のコンピュータプログラム。

10

【請求項 2 4】

前記一または複数のグループセットは複数のグループセットを含み、前記コンピュータプログラムはさらに追加の命令を含み、前記追加の命令は、前記一または複数の演算装置によって実行される時に、前記一または複数の演算装置の少なくとも一つに、

前記グループセットの二つ以上からのデータを集計して、合成グループセットにするためのコマンドを受領し、

前記二つ以上のグループセットに対応するエンティティマップベクトルについて、ブールOR演算を行い、合成エンティティマップベクトルを生成する、
 ようにさせる、請求項 2 1 に記載のコンピュータプログラム。

20

【請求項 2 5】

前記合成グループセット中の個別のデータ値の数は、前記合成エンティティマップベクトルの前記ビットの各々の値の合計に等しい、請求項 2 4 に記載のコンピュータプログラム。

【請求項 2 6】

前記テーブルはパーティション化され複数の演算ノードにわたって分散されており、前記データカラムはカラムパーティションのセットの中の第一カラムパーティションであり、前記ドメインはドメインパーティションのセットの中の第一ドメインパーティションであり、前記一または複数のエンティティマップベクトルは一または複数の第一パーティションエンティティマップベクトルであって、前記コンピュータプログラムはさらに追加の命令を含み、前記追加の命令は、一または複数の演算装置によって実行される時に前記一または複数の演算装置の少なくとも一つに、

30

前記第一ドメインパーティションを前記ドメインパーティションのセットの中の一または複数のその他のドメインパーティションと合成し、グローバルドメインを生成し、

前記一または複数の第一パーティションエンティティマップベクトルに対応する一または複数のグローバルエンティティマップベクトルを生成し、

各グローバルエンティティマップベクトルの長さは前記グローバルドメイン内の固有データ値の数に等しく、前記グローバルエンティティマップベクトル中の各ビットの位置は、前記グローバルドメイン中の前記固有データ値が語彙順に配列された一覧中の対応する固有データ値の語彙順位置に対応し、前記エンティティマップベクトルの各ビットの値は、前記グループセット中における対応する固有データ値の有無を示す、
 ようにさせる、請求項 2 1 に記載のコンピュータプログラム。

40

【請求項 2 7】

前記コンピュータプログラムはさらに追加の命令を含み、前記追加の命令は、一または複数の演算装置によって実行される時に前記一または複数の演算装置の少なくとも一つに、

前記第一パーティションに対応する前記一または複数のグローバルエンティティマップベクトルと、第二パーティションに対応する一または複数の第二グローバルエンティティマップベクトルとに、ブールOR演算を行い、一または複数の合成グローバルエンティテ

50

イマップベクトルを生成させる、請求項 2 6 に記載のコンピュータプログラム。

【請求項 2 8】

前記合成グローバルエンティティマップベクトルの各々の各ビットの値の合計は、前記第一パーティションと前記第二パーティションとにわたる当該グローバルエンティティマップベクトルに対応するあるグループセットの中の個別の要素の総数である、請求項 2 7 に記載のコンピュータプログラム。

【請求項 2 9】

グループセットの数は、少なくとも部分的に、前記テーブル中の第二データカラム中の個別のデータ値の数によって決定される、請求項 2 1 に記載のコンピュータプログラム。

【請求項 3 0】

前記テーブル中の一または複数の第二データカラムを参照して前記グループセットが定義される、請求項 2 1 に記載のコンピュータプログラム。

【請求項 3 1】

請求項 2 1 ~ 3 0 のいずれか一項に記載のコンピュータプログラムを記録した、コンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【関連出願データ】

【0 0 0 1】

本願は、2 0 1 3 年 3 月 1 5 日出願の米国特許出願第 1 3 / 8 3 5 , 5 9 0 号に対して優先権を主張するものであり、この参照によりその全内容が本明細書に組み込まれるものである。

【背景技術】

【0 0 0 2】

企業は頻繁にデータをデータウェアハウスに保存する。かかるデータウェアハウスは大量のデータを多数のデータベースおよびデータベーステーブルに分散させて有することが多い。したがって、このような多数のテーブルやデータベースからデータをまとめるためには、データ集計技術を利用しなければならない。

【0 0 0 3】

データウェアハウスの集計プロセスの性能を改善するために、下位のデータのサマリを保存する下位集計 (Low Level Aggregate、LLA) テーブルがよく使用される。多数のパーティションに保存される分散データベースの場合、各パーティションに L L A テーブルが利用でき、集計プロセスは各 L L A テーブルからのデータをくみあわせてデータベース全体についての集計されたテーブルを生成することができる。このような集計プロセスは、ロールアッププロセスとして知られ、ソースデータからではなく下位のサマリ (L L A テーブル) から上位のサマリが生成できるので有用である。

【0 0 0 4】

ロールアッププロセスは、超並列処理 (M P P) データベースアーキテクチャにおいて多数のデータパーティションに対して集計クエリを実行するために、つまりパーティション化されたテーブルについて照会を行うために利用できる。

【0 0 0 5】

ロールアッププロセス中に多数のテーブルからデータを正確に集計するためには、集計が行われている根底にある関数が加法でなければならない。例えば、第一のテーブルおよび第二のテーブル両方ともに出費額に対応するカラムがあり、ユーザが両方のテーブルの出費の総額の合計を計算しようとする、第一テーブルの出費額の合計と、第二テーブルの出費額の合計とを加算で集計し、正確な総額を出すことができる。しかしながら、ロールアッププロセスでユーザが集計しようとする関数すべてが、加法関数というわけではない。

【図面の簡単な説明】

【0 0 0 6】

10

20

30

40

50

【図1】図1A～1Gは、客の来店記録が入っている見本データベーステーブルと、見本データベーステーブルから生成される下位集計テーブルの例を示す。

【図2】図2A～2Bは、開示された実施形態にかかる、個別の(distinct)データ値について効率良く演算を行う方法のフローチャートである。

【図3】図3A～3Cは、開示された実施形態にかかる、テーブル中のデータの一または複数のカラムをトークン化する工程を示すテーブルである。

【図4】図4A～4Cは、トークン化されたテーブルから生成されたLLAテーブルと、LLAテーブルから生成された二つの見本上位集計テーブルとを示す。

【図5A】図5Aは、開示された実施形態にかかる、エンティティマップベクトルデータ構造を利用するトークン化されたテーブルから生成されたLLAテーブルを示す。

10

【図5B】図5Bは、開示された実施形態にかかる、エンティティマップベクトル中のビットの各々に値を割り当てる方法を示す。

【図5C】図5Cは、開示された実施形態にかかる、ドメインの語彙順配列の固有データ値のテーブルを示す。

【図5D】図5Dは、開示された実施形態にかかる、エンティティマップベクトル中のビットの各々に値を割り当てる方法を示す。

【図6A】図6Aは、開示された実施形態にかかる、図5AのLLAテーブルに保存されたエンティティマップベクトル上でロールアッププロセスを行う時の工程を示す。

【図6B】図6Bは、図6Aに示されるロールアッププロセスの結果である上位テーブルを示す。

20

【図7】図7A～7Bは、開示された実施形態にかかる、多数のパーティションにまたがってドメインを結合することにより分散データベース中の個別のデータ値について演算を効率良く行うためのフローチャートを示す。

【図8】図8A～8Dは、開示された実施形態にかかる、パーティションに特定されたエンティティIDの見本セットをグローバルエンティティIDへ変換するために必要な工程のそれぞれを示す。

【図9】図9は、開示された実施形態にかかる、多数のパーティションにまたがってグローバルエンティティマップベクトルを保存するためのデータキューブを示す。

【図10】図10は、開示された実施形態にかかる、経時的に割り当てられたエンティティIDを使用したデータベースウェアハウス更新演算を示す。

30

【図11A】図11Aは、開示された実施形態にかかる、ドメイン同期データ構造と、プロセスフローとを示す。

【図11B】図11Bは、マッピングテーブルを使用してローカルエンティティマップベクトルをグローバルエンティティマップベクトルにマッピングする時のマッピングを示す。

【図12】図12は、開示された実施形態にかかる、個別のデータ値について演算を効率良く行うための方法を実施するために使用できるコンピューティング環境の一例を示す。

【発明を実施するための形態】

【0007】

方法、装置、およびコンピュータ可読媒体が実施例および実施形態として本明細書に記載されているが、個別のデータ値に対する効率よい演算を行うための方法、装置、およびコンピュータ可読媒体が、実施形態または図面に記載されたものに限定されないことを、当業者は認識する。図面および明細書が開示された特定の形態に限定されることを意図したものではないことを理解されたい。むしろ、添付の請求項の精神及び範囲から逸脱しないすべての変形、均等物、および代替物を包含することを意図している。本明細書中に使用された見出しはいずれも、整理するという目的のみのために使用されており、本明細書または請求項の範囲を制限することを企図していない。本明細書で使用されている通り、「～することができる(can)」という表現は、必須の意味(すなわち、「～でなければならない(must)」という意味)ではなく、許容の意味(すなわち、可能性があるという意味)で用いられている。同様に、「含む」(「include」「including」、および「incl

40

50

udes」) という用語は、「含む」ことを意味しており、それに限定されるという意味ではない。

【 0 0 0 8 】

図 1 A のテーブルを参照し、加法関数用のロールアッププロセスのいくつかの例を以下に説明する。図 1 A は、客の来店記録が入っている見本データベーステーブルを示す。このテーブルは、客の名前と、店舗 ID と、来店中の出費額とのカラムを有する。

【 0 0 0 9 】

本明細書全体で用いられる照会言語は、構造化照会言語 (SQL) であるが、いかなるデータベース照会言語も、本発明の範囲から逸脱することなく使用でき、その例としては、コンテクスチュアル照会言語 (CQL)、XQuery、YQL、Datalog、OQL、RDQL、多次元式 (MDX) など多数の言語が挙げられることが理解される。

10

【 0 0 1 0 】

図 1 B には、各店舗についての来店数の L L A テーブルが示されている。「来店数」カラムは、店舗のそれぞれについての記録数をカウントすることにより生成され、店舗 1 には来店数 3、店舗 2 には来店数 4 という結果になっている。全店舗の総来店数を特定するためにこのテーブルを使用してロールアッププロセスを実行しようとする、関数 Sum(来店数) でデータがロールアップでき、総数 7 という結果になる。この場合、来店の総数を特定するために元のテーブルへ戻る必要が無い。

【 0 0 1 1 】

同様に、図 1 C のテーブルは、店舗それぞれの総出費額の L L A テーブルである。この出費額は、各店舗用の元のテーブルの「出費額」カラムの値を合計することにより特定される。ユーザが全店舗の総出費額を特定するためにこのテーブルをロールアップしようすれば、図 1 C の L L A テーブル上の関数 Sum(出費額) を使用して、1 0 7 という正確な総額に至ることができる。

20

【 0 0 1 2 】

ロールアッププロセス中の L L A テーブルの集計は、単純な合計以外の関数についても行うことができる。図 1 D は、各店舗における来店ごとの平均出費額の L L A テーブルを示す。この額は、各店舗についての元のテーブルの「出費額」カラムの値を平均することによって特定される。さらに、L L A テーブルには、前述のように各店舗への来店数のカウントが入っている。ユーザが、全店舗での平均出費額を特定するためにこのテーブルを

30

【 0 0 1 3 】

【数 1】

$$\text{平均出費額(合計)} = \frac{\left((\text{平均}_{\text{店舗1}} \times \text{来店数}_{\text{店舗1}}) + (\text{平均}_{\text{店舗2}} \times \text{来店数}_{\text{店舗2}}) \right)}{\text{Sum (来店数)}}$$

【 0 0 1 4 】

繰り返しになるが、ロールアッププロセスは、元のテーブルを必要とせず、かつ、L L A テーブルのみを使用することによって行える。

40

【 0 0 1 5 】

図 1 E および 1 F は、それぞれ、各店舗における最小出費額が入っている L L A テーブルと、各店舗における最大出費額とが入っている L L A テーブルとを示す。いずれかの店舗での最小の出費額を特定するために図 1 E のテーブルをロールアップするためには、ユーザは L L A テーブルについて関数 Min(最小出費額) を使用でき、正確な答えである 5 に至ることができる。同様に、いずれかの店舗での最大の出費額を特定するために図 1 F のテーブルをロールアップするためには、ユーザは L L A テーブルについて関数 Max(最大出費額) を使用でき、正確な答えである 3 1 に至ることができる。

50

【 0 0 1 6 】

図 1 G を参照すると、各店舗を訪れたことのある個別の (distinct) 客の数のカウントが入っている L L A テーブルが示されている。これは、店舗 I D でグループ分けされた Count(distinct<客の名前>)関数を使用して行うことができる。したがって、この場合、3 人の個別の客 (Bill、George、および Larry) が店舗 1 に来店し、4 人の個別の客 (Bill、George、Max、および Larry) が店舗 2 に来店していた。

【 0 0 1 7 】

ユーザが図 1 G の L L A テーブルをロールアップしたいなら、根底にある主なテーブルである図 1 A のテーブルにアクセスしないで店舗 1 または店舗 2 によく来店する個別の客の総数を正しく特定する方法は無い。L L A テーブルの個別の客のカラム上で機能する count-distinct 関数は、2 つの個別の値の総数を返信するだけであり、合計として 7 を返すだろうが、両方の店舗を訪れた個別の客の正確な数は、4 人である。同様に、もし L L A テーブルに各々の店舗で出費された個別の出費ドル額の合計が入っており、ユーザが、全ての店舗で出費された個別のドル額のすべての合計を計算したければ、元のテーブルにアクセスしなければならないだろう。さもなければ、店舗 1 で出費された 1 1 ドルと店舗 2 で出費された 1 1 ドルが二重に計上されてしまうだろうからである。さらに、もし L L A テーブルに各店舗で出費された個別の出費額すべてについての平均が入っていれば、両店舗での個別の出費額すべてについての平均をロールアップで特定しようとしてもうまくいかず、根底にある主なテーブルにアクセスしなければならないだろう。

【 0 0 1 8 】

count<distinct>のような非加法コマンドを含んでいる照会は、L L A テーブルのロールアップ演算を使用することが出来ず、その結果、計算コストが非常に高くなる。なぜなら、これらコマンドを、根底にあるデータについて実行する必要があるからである。パーティション化データベースか、または分散データベースにおいては、正確な値を検索するためには根底にあるデータが集計されているかまたはパーティション間で交換されなければならないので、このコストは膨れ上がる。例えば、もしカラム X Y Z が二つのパーティションに散らばっており、ユーザがカラム X Y Z の個別の値の数を特定したいなら、第一パーティション上のカラム X Y Z の値のすべてを、第二パーティションのカラム X Y Z の値のすべてと共に、何らかの方法で一緒に管理して、個別の値の数を特定しなければならない。

【 0 0 1 9 】

出願人は、エンティティマップベクトル (EMV) と称される方法と、関連したデータ構造とを発見した。これは、ユーザが非加法集計関数を加法として使用できるようにし、それによって、以前は L L A テーブルよりも下位であるテーブルにある、根底にあるデータへのアクセスや処理を必要とした関数について、L L A テーブル上でロールアッププロセスが実行可能になった。出願人の技術によってさらに、パーティション化されたデータベースのノード間で交換されるデータの量が大幅に削減できるようになった。これは、各ノードについての L L A テーブルが、非加法集計関数についても、加法集計できるからである。

【 0 0 2 0 】

図 2 A は、開示された実施形態にかかる、EMV を生成する方法のフローチャートを示す。ステップ 2 0 1 で、データベースのテーブル中のデータのトークン化されたカラムが、EMV を生成する前に保存される。このトークン化プロセスを、図 3 A - 3 C を参照して説明する。さらに、図 2 B に示され、以下に記載されるように、EMV は、データのトークン化されたカラムを使用することなく、あるドメインでの固有データ値の語彙順に基づいて生成することができる。

【 0 0 2 1 】

図 3 A は、トークン化前の顧客取引テーブル 3 0 0 を示す。このテーブルは、タプル I D (TupleID) 3 0 1 と、週 (Week) 3 0 2、店舗 I D (SstoreID) 3 0 3、客 I D (CustID) 3 0 4、製品 I D (productID) 3 0 5、量 (Qty) 3 0 6、および出費額 (Spent)

10

20

30

40

50

307のデータカラムとを備える。データカラムの各々は、関連付けられたドメインに対応させてもよい。例えば、客IDを備えたドメインは、図示されたテーブルの客IDカラムに入力するために使用できるが、同じドメインを、客のアドレスと共に客IDを一覧にした別のテーブルに入力するために使用できる。もちろん、ドメインは、同じテーブルの複数のカラムに入力するために使用できる。例えば、都市名のドメインは、起点都市と目的地都市とに対応するテーブルのカラムに入力するために使用できる。したがって、週302のカラムはドメインD_Weekに対応でき、店舗ID303はドメインD_StoreIdに対応でき、客ID304はドメインD_CustIdに対応でき、製品ID305はドメインD_ProductIdに対応でき、量306はドメインD_Qtyに対応でき、出費額307はドメインD_Spentに対応できる。

10

【0022】

各データカラムのデータ値をトークン化するために、各ドメインに生じる固有データ値（異なるカラムに対応する）は、識別されて固有のトークン値にマッピングされなければならない。固有データ値をテーブル中の各ドメインの固有のトークン値にマッピングするが、図3Bのマッピングテーブル310はこのマッピングを示している。各ドメインの各固有データ値は、エンティティID、すなわちEIDと称される、そのドメイン用の固有のトークンにマッピングされる。したがって、例えば、D_Week312は3つの固有データ値を備え、したがってこれら3つの固有データ値が3つの対応する固有のEIDにマッピングされる。同様に、D_StoreId313もまた3つの固有データ値と3つのEIDとを備え、D_CustId314は4つの固有データ値と4つのEIDとを備え、D_ProductIdは4つ

20

【0023】

ドメイン中の固有データ値に対応するEIDにマッピングするために用いられるマッピングテーブルは、ルックアップ辞書として保存できるので、関数および関連付けられた処理は、必要であればEIDを使用してデータ値を検索できる。さらに、EIDに対する固有データ値のマッピングは、固有データ値の語彙順で生じるものとして示されており、これは割り当てプロセスや、その後のデータの保存や検索を簡便にするからであるが、EIDの固有データ値への割り当ては固有データ値のいかなる順序で生じてもよい。例えば、第一データ値が第二のデータ値より以前に遭遇すれば下位のEIDを有するよう、経時的な順序でデータ値に割り当てられるようにEIDを割り当ててもよい。

30

【0024】

図3Cは、トークン化後の顧客取引テーブル320を示す。マッピングテーブル310からのEIDは、データのカラムの各々においてデータ値に置き換えられている。しかしながら、この高度な暗号化構造においてすら、テーブル320を介して生成されたLLAテーブルには、非加法の（distinctタイプの）集計関数を含むロールアッププロセスを実行するのに適切なフォーマットでデータが保存されることはない。

【0025】

例えば、LLAテーブルは、週と、店舗IDと、各週に各店舗に来店した個別の客とを一覧にした取引テーブル（「Trx」）320に基づいて、下記の定義を用いて、生成できる。

40

【0026】

```
Create Table LLA_TRX as
Select      week,
           StoreId,
           Count(distinct CustId) as CD_CustId,
From Trx
Group By week, StoreId;
```

この定義は、各週に各店舗に来店した個別の客に対応する値を、カラムCD_CustIdに割り当てる。さらに、「group-by」コマンドは、週-店舗IDのペアに対応するグループセ

50

ットを作る。例えば、グループセット1（週1、店舗1）、グループセット2（週1、店舗2）などである。グループセットの構成や数はどんなものでよい。例えば、グループセットは、一つのカラム、または二つ以上のカラムに基づいていてもよい。

【0027】

図4Aには、結果として得られたLLA_TRXテーブル400が示されている。このテーブルは、グループセットごとの個別の客の数を正確に表している。例えば、テーブルは、週1に店舗1を訪れた個別の客の数を2であると、または、週2に店舗1を訪れた個別の客の数を3であると、正確に一覧にしている。

【0028】

もしユーザが各週についてテーブルLLA_TRXをロールアップして各週に店舗のいずれかを訪れた個別の客の数を一覧にしようとしたなら、各グループセットについて個別の客の数の正確なカウントを作ることはできないであろう。この場合のグループセットとは、グループ1（週1）、グループ2（週2）、などである。例えば、上位LLAテーブルを作り出すために下記の定義が用いられたとする。

【0029】

```
Create Table LLA_TRX_WEEK as
Select      Week,
            Count(distinct CD_CustId) as CD_CD_CustId,
From LLA_Trx
Group By Week;
```

図4Bのテーブル410に図示される、結果として得られたテーブルの一覧では、各週に店舗のいずれかに来店した個別の客の数が不正確となるだろう。同様に、テーブルLLA_TRXを集計するためにSUM関数が用いられるとする。

【0030】

```
Create Table LLA_TRX_WEEK as
Select      Week,
            Sum(distinct CD_CustId) as Sum_CD_CustId,
From LLA_Trx
Group By Week;
```

すると、図4Cに図示される、結果として得られたテーブルの一覧でもまた、各週に店舗のいずれかに来店した個別の客の数が不正確になるだろう。例えば、週2の個別の客の数は、図3Cのテーブル320に基づくと4であり（客ID1, 2, 3および4）、図4Bのテーブル410の一覧では2であり、図4Cのテーブル420の一覧では5である。

【0031】

前述したように、これは、個別の値の集計が非加法であるからである。週ごとの個別の客の数について正確な値を特定するためには、非加法集計の照会に対してはLLA_TRXテーブルを無効にし、元のデータテーブルへ戻る必要がある。

【0032】

開示された実施形態にしたがって生成されたEMVデータ構造は、各グループセットの個別の値を、集計できるようなやり方で効率良く保存することによって、この問題を解決する。図2Aにもどり、ステップ202では、一または複数のグループセットを特定し、EMVデータ構造を生成するために使用される新たな関数やその他の同様のコマンドを備えた、特別な照会が受領できる。ステップ203では、演算システムが、一または複数のエンティティマップベクトルを生成できる。もちろん、ステップ201、202、および203は、異なる順序で実行できる。例えば、照会をまず受領し（ステップ202）、その後、トークン化されたデータカラムが保存され（ステップ201）、そして、一または複数のエンティティマップベクトルが生成できる（ステップ203）。

【0033】

ステップ202～203を容易にするために、値ではなくEMVを各グループセットについて作成するよう演算システムに指示する、新たなSQL関数count(Additive_distinc

10

20

30

40

50

t<カラムの名前>) を定義することができる。EMVは、異なるエンティティに対して各ビットがマッピングされたビットマップベクトルでありうる。EMVのグループセット中に個別の値に関連する情報を保存することにより、エンティティと称される各個別の値をビットマップベクトルの一つのビットにマッピングすることができる。したがって、図3Cのトークン化された取引テーブル(Trx)320の上記の例を使用して、LLAテーブルLLA__TRXを生成するために下記の定義を用いることができる。

【0034】

```
Create Table LLA_TRX as
Select      Week,
            StoreId,
            Count(Additive_distinct CustId) as CD_CustId,
From Trx
Group By Week, StoreId;
```

10

図5Aは、前記の結果得られるLLA__TRXテーブル500を示す。LLA__TRXテーブルのように、各グループセットに対応しているタプルID501、週502、および店舗ID503が保存される。ただし、個別の客のカウンタを数字で保存するのではなく、各ビットが異なるEIDを参照し、ビットの値がグループセット中のEIDの有無を示すEMVに、各グループセットについての個別の客が保存される。

【0035】

もちろん、Additive_distinctのようなSQL拡張子が、各グループセットについてEMVを作成するよう演算システムに指示する必要が無い。例えば、データベースエンジンが、Count(distinct<カラムの名前>)などの個別のデータ値の集計関数を含むカラムを備えたテーブルの作成を検出し、個別のデータ値のカウンタを保存するためにEMVを自動的に造ることができる。これによって、既存のフロントエンドツールでEMV構造が利用できるようにする。

20

【0036】

グループセット(週1、店舗ID2)に対応する第二のEMV513について、以下にEMVの構造を記載する。図5Bを参照すると、空のEMV505が示されている。先頭文字506の役割は、EMVにビットマップベクトルであると標識をつけることのみであり、それに適していればいかなる文字でもよく、また、省略することもできる。EMV505の長さ寸法は、ビットスロットの数であり、先頭文字を含まない。EMVの長さはEMVを生成するために用いられるカラムのドメイン中のEIDの総数に等しい。換言すれば、EMVの長さは、EMVを生成するために用いられているドメイン中の固有の値の数に等しい。したがって、例えば、図5Aのテーブル500用にEMVを生成するために用いられたドメインは、D_CustIdドメインであり、図3Bより、D_CustIdドメインは4つのEIDを有していることがわかる。したがってそのドメインについてのEMVの長さは、4という長さであり、この4はビットスロットに対応する。

30

【0037】

再度図5Bを参照すると、ビットスロットの各々は、異なるEIDに、EIDが割り当てられる順序で対応している。したがって、第一スロットはEID=1に対応し、第二スロットはEID=2に対応するなどである。EMV中のビットスロットの各々のビットの値を特定するために、EMVに対応するグループセットを調べ、EID値を含むかどうかを確認する。したがって、EMV513の第一ビットについては、グループセット(週1、店舗2)がCustId=1を含むかどうか、ステップ507で判断される。図3Cのテーブル320から、週1の間に店舗2にCustId=1が現れていないことがわかる。したがって、第一ビットの値は0である。

40

【0038】

同じやり方で、EMV513の第二ビットについては、グループセット(週1、店舗2)がCustId=2を含むかどうかステップ508で判断される。CustId=2も週1の間に店舗2に現れていないので、第二ビットの値もゼロである。EMV513の第三ビットについ

50

ては、グループセット（週1、店舗2）がCustId=3を含むかどうかステップ509で判断される。テーブル320から、週1の間に店舗2にCustId=3が2度現れていることがわかる。したがって、第三ビットの値は1に等しい。EMV513中の第四ビットについては、グループセット（週1、店舗2）がCustId=4を含むかどうかステップ510で判断される。CustId=4もまた週1の間に店舗2に現れているので、第四ビットも同様に1に設定される。その結果得られたEMV513は、店舗2に週1の間に来店した個別の客を、加法集計を可能にするフォーマットで、把握している。

【0039】

もしユーザがグループセットごとの個別の客の数を特定したければ、そのグループセットについてEMV中の各ビットを合計しさえすればよい。例えば、グループセット（週1、店舗1）についてのEMVはOB1100であり、個別の客の総数は $1 + 1 + 0 + 0 = 2$ である。グループセット（週2、店舗1）についてのEMVはOB1101であるので、そのグループセットの個別の客の総数は $1 + 1 + 0 + 1 = 3$ である。

【0040】

図2Bを参照し、以下にトークン化されていないデータカラムからEMVを生成する方法を記載する。ステップ210では、図2Aのステップ202と同様に、一または複数のグループセットを特定し、新たな拡張子またはその他の同様のコマンドを含み、データベースエンジンにEMVデータ構造を生成するよう指示する、照会が受領される。ステップ211で、データのカラムに対応するドメイン中の固有データ値の語彙順に基づいて、グループセットの各々についてエンティティマップベクトルが生成される。図5Cおよび5Dは、図3Aのトークン化されていないテーブル300中のグループセット（週201001、店舗1002）について生成されたEMVに関して、このプロセスをより詳細に示す。テーブル520は、客IDドメインの、語彙順に配置された固有データ値を示す。図5Dに示されたEMV生成プロセスは、EID生成のステップをバイパスするためにこの語彙順序を使用する。空のEMV525に示されるように、生成されたEMVの各々の長さは、ドメイン中の固有データ値の数に等しい。この場合、客IDドメインには4つの固有データ値があるので、あるグループセットについて生成されたEMVの長さは4となるであろう。各EMV中の各ビットの位置は、固有データ値が語彙順に配列された一覧520中の、対応する固有データ値の語彙順位置に対応し、EMV中の各ビットの値はグループセット中の対応する固有データ値の有無を示す。したがって、EMV533中の第一ビットについては、グループセット（週201001、店舗1002）が語彙順で最も若い位置の客IDを含むかどうか判断される。この客IDは165であり、図3Aのテーブル300から、この客IDがグループセット（週201001、店舗1002）には出現しないことがわかるので、ビットはゼロとなる。同様に、EMV533の第二ビットは語彙順で二番目に若い位置の固有データ値、167を参照し、これもグループセットには出現しないので、このビットはゼロとなる。EMV533の第三ビットは語彙順で酸番目に若い位置の固有データ値、268を参照し、図3Aのテーブル300から、この値はグループセット（週201001、店舗1002）に出現しているとわかるので、第三ビットは1に設定される。同様に、第四ビットに対応する値、301、もグループセットに出現するので、第四ビットは1に設定される。この結果、図5Bに生成されたものと同じEMVとなるが、これはトークン化されていないデータカラムから生成できる。

【0041】

エンティティマップベクトルを利用しないLLAテーブルの前記の例とは異なり、EMVを備えたLLAテーブルは、ロールアッププロセス中に正確に集計できる。例えば、いずれかの店舗に来店した個別の客の数についての上位テーブルを、下記の定義を用いて、LLAテーブルから週でグループ化して生成することができる。

【0042】

```
Create Table LLA_TRX_WEEK as
Select      Week,
            Count(Distinct CD_CustId) as EMV_CD_CustId,
```

10

20

30

40

50

From LLA_Trx

Group By Week;

図 6 A を参照すると、店舗にまたがるビットベクトルを集計し、週でグループ化するプロセスが示されている。ステップ 6 0 1 A、6 0 1 B、および 6 0 1 C において、集計されなければならないグループセットの各々が、各週について判断される。例えば、週 1 について合成すべきグループセットは、週 1、店舗 1 と、週 1、店舗 2 のグループセットに対応する客 ID である。週 2 について合成すべきグループセットは、週 2、店舗 1 と、週 2、店舗 2 と、週 2、店舗 3 のグループセットに対応する客 ID である。週 3 について合成すべきグループセットは、週 3、店舗 1 と、週 3、店舗 3 のグループセットに対応する客 ID である。

10

【 0 0 4 3 】

ステップ 6 0 2 A、6 0 2 B、および 6 0 2 C を参照すると、各週について合成すべきグループセットの各々の E M V が、ブール O R 演算 (Boolean OR) を使用して集計される。したがって、例えば、もし二つの E M V をブール O R 演算で合成し、第一 E M V のみが第一位置に「 1 」を有していれば、その結果得られる E M V は「 1 」を第一位置に有する。

【 0 0 4 4 】

その結果得られた合成 E M V がステップ 6 0 3 A、6 0 3 B、および 6 0 3 C に示されている。各 E M V は、異なる週のグループセットに対応する。前述のように、各週の個別の客のカウントは、E M V の各々のビットをすべて合計することにより特定できる。

20

【 0 0 4 5 】

図 6 B は、その結果得られたテーブル 6 0 0 を示し、テーブル 6 0 0 は図 6 A の合成 E M V に基づいて、各グループセットについてのタプル ID 6 0 1、週番号 6 0 2、および個別の客のカウント 6 0 3 を備えている。図 3 C の元のテーブル 3 2 0 を参照すると、テーブル 6 0 0 の個別の客のカウント 6 0 3 は、正確に示す t h a t 4 人の個別の客 (1 , 2 , 3 , および 4) が週 1 にいずれかの店舗に来店し、4 人の個別の客 (1 , 2 , 3 , および 4) が週 2 にいずれかの店舗に来店し、2 人の個別の客 (1 および 2) が週 3 にいずれかの店舗に来店したことが、わかる。

【 0 0 4 6 】

もちろん、非加法集計関数の集計は、グループセット中の個別の要素の数をカウントすることに限定されない。ユーザが、あるグループセット中の個別の値のセットを合計したい、またはグループセット中の個別の値のセットの平均値を求めたいなら、ルックアップ辞書中の E M V の各ビットに対応する固有データ値を探索し、加算して合計したり、またはもしビットが「 1 」であれば平均値に加えたりすることによって、達成できる。さらに、カスタム化計算やメトリクスなど、他の目的のためにエンティティマップベクトル中のビットに対応する固有データ値の各々を検索するために、ルックアップ辞書を用いることができる。エンティティマップベクトルが、固有データ値の語彙順の配列に基づいて、トークン化されていないデータ値から生成される状況では、同じ語彙順配列関係を、各ビットに対応する固有データ値を検索するために用いることができる。例えば、E M V の第一ビットが「 1 」なら、最も若い語彙順位置の固有データ値を、セット中のデータ値を含む、合計、平均、または他の計算に組み入れることができる。同様に、E M V の第二ビットが「 0 」なら、二番目に若い語彙順位置に対応する固有データ値は、合計、平均、または他の計算から除外することができる。

30

40

【 0 0 4 7 】

パーティション化されたデータベース中の照会処理およびパーティション化されたデータベース中のデータの表現に基づいて、E M V のプロセスの大要を以下に記載する。取引テーブルを含む 3 つのパーティションデータベースを仮定すると、そのテーブルの論理ビューは、View として下記のようにあらわされる。

Create View V1 as

Select * from Partition1.TRX

50

```

Union All
Select * from Partition2.TRX
Union All
Select * from Partition3.TRX

```

テーブルのパーティション化は、要求の効率的な分散並列処理を可能にするが、非加法集計関数は例外である。例えば、図3Aのテーブル300と同様のTRXという名の取引テーブルの元の照会が、「Select Week, Sum(Spent) from TRX group by Week (週を選択、週によるTRXグループからの(出費額)を合計)」であるが、テーブルが3つのパーティションに分散されていれば、照会は、分散された要求に変形することができる。

```

Select Week, Sum(P_Spent) from (
    Select Week, Sum(Spent) as P_Spent from Partition1.TRX group by Week
    Union All
    Select Week, Sum(Spent) as P_Spent from Partition2.TRX group by Week
    Union All
    Select Week, Sum(Spent) as P_Spent from Partition3.TRX group by Week
) T

```

10

Group By Week

その分散された要求については、「Select Week, Sum(Spent) as P_Spent from Partition1.TRX group by Week (週を選択、週によるパーティション1 . TRXグループからのP_Spentとして(出費額)を合計)」のような、各照会セグメントが独立して実行でき、その中間結果セットをまとめて付加して最後に再処理し、最終的な照会結果を作成できる。

20

【0048】

「Select Week, Count(Distinct CustId) from TRX group by Week (週を選択、週によるTRXグループからの(個別の客ID)をカウント)」のような非加法集計関数要求の場合、この要求は分散された要求に変形できる。

```

Select Week, Count(Distinct CustId) from (
    Select Week, CustId from Partition1.TRX group by Week, CustId
    Union All
    Select Week, CustId from Partition2.TRX group by Week, CustId
    Union All
    Select Week, CustId from Partition3.Trx group by Week, CustId
) T

```

30

Group By Week

分散された要求はこの場合、各パーティションから固有のCustIdの一覧を保存するが、膨大なデータをパーティションの各々の間で交換し、その後処理するよう命ずることがありうる。例えば、第一週についての第一パーティション中のCustId(客ID)の一覧を、第一週についての第二パーティション中のCustId(客ID)の一覧と比較し、重複が無いかが判断しなければならないだろう。その結果、かかる分散された要求は、要求実行プロセスに支障をきたす。

40

【0049】

パーティション間で交換されるデータ量を削減し性能を向上させるために、非加法集計関数を加法関数に変形する方法を一または複数のパーティション化されたデータベースに適用することができる。例えば、元の照会「Select Week, Count(Distinct CustId) from TRX group by Week (週を選択、週によるTRXグループからの(個別の客ID)をカウント)」が下記のように変形できる。

```

Select Week, Count(Distinct CustId) from
(
    Select Week, Count(Additive_distinct CustId) CustId from Partition1.TRX group by Week

```

50

Union All

```
Select Week, Count(Additive_distinct CustId) CustId from Partition2.TRX group by Week
```

Union All

```
Select Week, Count(Additive_distinct CustId) CustId from Partition3.Trx group by Week
```

) T

Group By Week

図7 Aは、開示された実施形態にかかる、この分散された要求を処理するステップを示す。3つの第一ステップ701は、パーティションの各々に対して平行して行うことができ、各パーティションに適用される図2 Aに示されるEMV生成ステップと同様である。もちろん、これらステップは平行して行われなくてもよく、なんらかの適した順序で行える。前記3つの第一ステップ701がパーティションの各々について実行され、その結果、各パーティションに一つずつ、計3つのEMVセットが得られ、EMVの各桁はローカルパーティションEIDを参照している。例えば、図8 Aは、3つのパーティションに対応して、3つの仮想EIDマッピングテーブル801 A、801 B、および801 Cと共に、3つの仮想EMVセット802 A、802 B、および802 Cを、パーティションの各々にある固有データ値の各々について示している。マッピングテーブル801 A~801 C中のEIDはすべて異なる固有データ値に対応しているので、また、各パーティション中のEIDの数が様々であるので、EMVの長さも、EMVセット802 A~802 C中の各EMVのビットの意味も、様々である。その結果、EMVセット802 A~802 Cは、このままではunion演算子によって合成することができない。

【0050】

図7 Aを参照すると、ステップ702で、ドメインパーティションの各々からグローバルドメインを構築するためにドメイン結合プロセスをおこない、EMVセットを標準化できる。図7 Bへ続き、ステップ703で、グローバルドメイン中の固有データ値の各々について、グローバルEIDが生成される。グローバルEIDがパーティションEIDとしてと同じ固有データ値に対応するので、図7 Bのステップ704でパーティションEIDの各々がグローバルEIDに変換できる。ステップ705では、パーティションに特定のEMVがグローバルEMVに変換され、そのすべてが各ビットにおいて、同じ長さを有し、同じ固有データ値を参照する。これはパーティション-グローバルEIDマッピングテーブルを使用して、パーティションに作り出されたEMVの各々を再構築し、グローバルEIDに対してマッピングする新たなEMVを作り上げることにより達成される。

【0051】

当該プロセスを示すために、図8 Bは、図8 Aのドメインパーティションのドメイン結合の結果得られたグローバルドメイン803と、対応するグローバルEID804とを示す。図8 Cは、パーティション1、811、パーティション2、812およびパーティション3、813用のパーティション-グローバルEIDマッピングテーブルを示す。前述のように、パーティション-グローバルEIDマッピングは、パーティションEIDとグローバルEIDとが固有データ値を相互参照することによって生成される。マッピングテーブルはその後、パーティション特定EMVをグローバルEMVに変換するために用いられる。変換の結果生成されたグローバルEMVは、パーティション1、821、パーティション2、822、およびパーティション3、823について図8 Dに示されている。

【0052】

パーティション特定EMVがグローバルEMVに変換されると、前述のようにブルOR演算で集計できる。その結果得られた合成グローバルEMVは、合成パーティションすべてにわたるあるグループセットにおいて、固有データ値が有るかどうかを示すことになろう。さらに、合成グローバルEMV中のビットの合計は、パーティションにわたってそのEMVに対応するグループセット中の個別のデータ値の総数と等しくなる。

【0053】

10

20

30

40

50

このドメイン結合プロセスは、データベース更新演算に利用できる。データウェアハウス環境では、更新プロセスが新たなデータを付加する工程を含むことがよくある。関連付けられたパーティションLLAテーブルを生成するために新たなパーティションとして付加されたデータを処理し、上述した技術を使用してパーティションLLAテーブルを既存のLLAテーブルと合成することにより、前記新たなデータを既存のLLAテーブルと結合することができる。

【0054】

例えば、パーティションの各々についてのグローバルEMVが、データキューブに保存される。図9は、開示された実施形態にかかる、グローバルEIDを保存するデータキューブ900の例を示す。図示されているように、次元には、EMV次元901、グループセット次元902、およびパーティション次元903がある。新たなデータ更新が受領されると、パーティションとして扱い、前述のドメイン結合マージプロセスを使用してデータキューブに加えることができる。

10

【0055】

前述のように、EIDも、データ値の語彙順ではなく、データ値が出現した経時的な順序で割り当てられてもよい。このようにEIDを割り当てることにより、データウェアハウスの更新のためのドメイン結合ステップを省略できる。図10を参照すると、既存のデータウェアハウスが1001に示されている。この例では、データ値が処理された経時的な順序を用いてEIDを各データ値に割り当てている。テーブル1002は、以下の順序で処理された4つのデータ値用のEIDを示している：368、167、165、268。さらに、テーブル1003は、データ値を含む3つのグループセットについてのEMVを示している。したがって、例えば、第二のEMV「OB1110」は、データ値368、167、および165を含むグループセットに対応する。

20

【0056】

図10はまた、データウェアハウスのデータ更新1004を示す。データ更新は、上述したパーティションと同様である。ローカルEIDマッピングテーブル1005は、同様に経時的順序で割り当てられた、更新におけるデータ値についてのEIDを示している。さらに、二つのグループセットについてのEMVが1006に示されている。この場合、更新におけるデータ値のうち二つは新規であり、データ値の一つはすでにデータウェアハウスに存在していたものである。

30

【0057】

更新されたデータウェアハウスが1007に示されている。EIDが経時的順序で割り当てられているので、ドメイン結合プロセスを行う必要がない。データ更新における二つの新たなデータ値、392および163が、次に使用可能な二つのEID、この場合5および6、に割り当てられる。既存のEIDの再マッピングが無いので、各EMVの最後に新たなEIDの数と等しい数のゼロを加えることによって各EMVの長さを伸ばすだけで、EMVの既存のテーブル1003が更新される。この結果、更新されたEMVを備えたテーブル1009となる。もちろん、ある位置に文字が無ければゼロと解釈することができるので、ゼロが必ずしもEMVに付加されなくてもよい。

【0058】

データ値の経時的な順序に基づいてEIDが割り当てられている場合にデータ更新に必要な再マッピングのみが、実際のデータ更新においてEMVに対して行われる。この場合、更新1006のEMVは、テーブル1010の更新されたEMVを生成するためにはローカルEIDテーブル1005ではなくグローバルEIDテーブル1008を使用するために、再マッピングされなければならない。したがって、例えば、更新1006の第一EMVが、値392および163（EMVの一番目および三番目の位置の「1」）を含むグループセットに対応し、これらデータ値に対応するグローバルEIDが5および6なので、データ更新内の第一EMVに対応する更新されたデータテーブル1010中の第一EMVは、五番目および六番目の位置に「1」を有することになる。

40

【0059】

50

パーティションに特定された一または複数のEMVが第一データパーティション中のトークン化されていないデータカラムから生成され、一または複数の他のパーティションと結合されるというシナリオでは、第一データパーティション中のトークン化されていないデータカラムに対応する第一ドメインパーティションを、グローバルドメインを生成するためのドメインパーティションの群の中の、一または複数の他のドメインパーティションと合成できる。EMV生成のプロセスと同様に、グローバルドメイン中のすべての固有データ値の語彙順配列に基づいて、パーティションに特定された一または複数のEMVに対応する一または複数のグローバルEMVを生成するために、このグローバルドメインを用いることができる。各グローバルEMVの長さはグローバルドメイン中の固有データ値の数と等しくてもよく、グローバルEMV中の各ビットの位置はグローバルドメイン中の固有データ値の語彙順の一覧における固有データ値に対応する語彙順の位置に対応してもよく、グローバルEMV中の各ビットの値は、グループセット中における、対応する固有データ値の有無を示すことができる。

10

【0060】

前述のドメインマージプロセスをバイパスするために使用することができ、分散システムでのデータベースの演算をより効率的にできる、ドメイン同期プロセスを、下記に記載する。図11Aは、分散データベース中の二つのパーティション、1101および1102を示す。図8Aの分散データベースと同様に、パーティションの各々はローカルEIDテーブルを含み、パーティション1101にはテーブル1104、パーティション1102にはテーブル1106を備える。この例では、EIDはデータ値の語彙順に基づいて割り当てられているものとして示されているが、EIDは前述のように経時的に割り当てられられてもよい。さらに、分散データベースは、グローバルEIDテーブル1103を備える。グローバルEIDテーブル1103には、データ値すべてをグローバルEIDにマッピングしたものが入っている。繰り返しになるが、これらグローバルEIDはデータ値の語彙順に基づいて割り当てられているが、データ値の経時的な順序にもとづいて割り当てられることもできる。グローバルEIDテーブル1103はパーティション1101およびパーティション1102の外に保存されているものとして示されているが、グローバルEIDテーブルは、パーティションの一つ、または両方に保存されることも可能である。

20

【0061】

各パーティション1101および1102もまた、ローカルEID - グローバルEIDマッピングテーブル、パーティション1101についてのテーブル1105、およびパーティション1102についてのテーブル1107を含む。EMV1108および1109の例を使用して下記に記載するように、これらマッピングテーブルは、ローカルEMVをグローバルEMVに変えるために使用できる。

30

【0062】

パーティション1101またはパーティション1102のどちらかで更新が受領されると、新たな固有データ値がグローバルEIDテーブル1103へ送られる。これら新たな固有データ値は、グローバルEIDテーブルを更新するために用いることができる。例えば、EIDが固有データ値の語彙順で割り当てられる時、グローバルEIDテーブル1103中のデータ値が再ソートされ、受領された新たな固有データ値を考慮してグローバルEIDが固有データ値の各々に再割り当てされる。経時的に割り当てられたEIDが用いられるなら、新たな固有データ値が次に使用可能なEIDに割り当てられて、既存の固有データ値のすべてに対してEIDを再割り当てる必要がないので、このプロセスは大幅に簡略化される。

40

【0063】

グローバルEIDテーブル1103が更新された後、新たなグローバルEIDがパーティション1101および1102へ戻され得る。新たなグローバルEIDを用いて、ローカルEID - グローバルEIDマッピングテーブル1105および1107の各々を更新することができる。繰り返しになるが、このプロセスには経時的に割り当てられたEIDが効率的である。新たなグローバルEIDのみが新たな固有データ値について割り当てら

50

れるからである。語彙順で割り当てられた E I D では、既存の固有データ値についてのすべての E I D が再シャッフルされうる。

【 0 0 6 4 】

照会が複数のパーティションから E M V を求める場合、ローカル E I D - グローバル E I D マッピングテーブル、1 1 0 5 および 1 1 0 7、は、なんらかのローカル E M V をグローバル E M V へ変換するために用いることができる。図 1 1 B を参照すると、パーティション 1 1 0 1 中の見本 E M V 1 1 0 8 についての変換プロセスが示されている。図 1 1 B に示されるように、ローカル E M V 中の各ビットは、そのローカル位置から、等価のグローバル E M V 1 1 1 0 中のそのグローバル位置へ、マッピングされる。したがって、例えば、ローカル E I D - グローバル E I D マッピングテーブル 1 1 0 5 がローカル E I D 「 3 」をグローバル E I D 「 5 」へマッピングするので、ローカル E M V 1 1 0 8 中の三番目の位置ビットの値は、グローバル E M V 1 1 1 0 の五番目の位置へ割り当てられる。グローバル E M V 1 1 1 0 中の四番目の位置など、グローバル E M V 中の位置へのマッピングが無ければ、その値は、グループセット中のある固有データ値が無いことを反映してゼロに設定される。なぜなら、もしその値が存在すれば、ローカル E I D - グローバル E I D マッピングテーブル 1 1 0 5 中にそれがあつたはずだからである。同様に、パーティション 1 1 0 2 中のローカル E M V 1 1 0 9 からグローバル E M V 1 1 1 1 を生成するためには、ローカル E I D - グローバル E I D マッピングテーブル 1 1 0 7 を使用できる。

【 0 0 6 5 】

このドメイン同期プロセスを通じて、上述されたドメインマージ手順をバイパスすることができ、その結果処理時間が改善でき、一度にエクスポートする必要のある固有データ値の数を減らすことができる。もちろん、各パーティションに E M V が一つ示されているのは、明瞭化だけのためである。実際は、あるパーティションにローカル E I D および値が存在すれば、いくつかのグループセットおよび E M V にその値が入るはずである。さらに、二つのパーティションが示されているのも単に明瞭化のためだけであり、データベースが含むパーティションの数はどんな数であってもよい。

【 0 0 6 6 】

上記の技術の一または複数は、一または複数のコンピュータシステム中に実現することができ、または、それらで構成されることができる。図 1 2 は、コンピューティング環境 1 2 0 0 の一般例を示す。コンピューティング環境 1 2 0 0 は、記載された実施形態の用途または機能の範囲についていかなる限定も示唆するものではない。

【 0 0 6 7 】

図 1 2 を参照すると、コンピューティング環境 1 2 0 0 は、少なくとも一つの処理ユニット 1 2 1 0 およびメモリ 1 2 2 0 を含む。処理ユニット 1 2 1 0 はコンピュータ実行可能な命令を実行し、現実のまたは仮想のプロセッサであってもよい。多重処理システムにおいては、多数の処理ユニットがコンピュータ実行可能な命令を実行して処理パワーを増加させている。メモリ 1 2 2 0 は、一時的なメモリ（例えば、レジスタ、キャッシュ、RAM など）、非一時的なメモリ（例えば、ROM、EEPROM、フラッシュメモリなど）、またはその二つの組合せであってもよい。メモリ 1 2 2 0 は、上記技術を実現するソフトウェア 1 2 8 0 を記憶していてもよい。

【 0 0 6 8 】

コンピューティング環境は、さらに追加の構成を有することができる。例えば、コンピューティング環境 1 2 0 0 は、記憶装置 1 2 4 0、一または複数の入力装置 1 2 5 0、一または複数の出力装置 1 2 6 0、および一または複数の通信接続部 1 2 9 0 を備える。バス、制御器、またはネットワークなどの相互接続機構 1 2 7 0 は、コンピューティング環境 1 2 0 0 の構成部分を相互接続する。通常、オペレーティングシステムソフトウェアまたはファームウェア（図示省略）は、コンピューティング環境 1 2 0 0 で作動するその他のソフトウェア用の操作環境を提供し、コンピューティング環境 1 2 0 0 の構成部分の動作を調整する。

【 0 0 6 9 】

10

20

30

40

50

記憶装置 1240 は着脱可能であっても、着脱不可であってもよく、磁気ディスク、磁気テープ、カセット、CD-ROM、CD-RW、DVD など、情報を記憶するために用いることができ、コンピューティング環境 1200 中でアクセスすることができる媒体を含む。記憶装置 1240 は、ソフトウェア 1280 への命令を記憶できる。

【0070】

入力装置（複数も可）1250 は、キーボード、マウス、ペン、トラックボール、タッチ画面、またはゲーム制御器などのタッチ入力装置、音声入力装置、走査装置、デジタルカメラ、遠隔制御、コンピューティング環境 1200 に入力を行う他の装置でありうる。出力装置（複数も可）1260 は、表示器、テレビ、モニタ、プリンタ、スピーカ、コンピューティング環境 1200 からの出力を行う他の装置でありうる。

10

【0071】

通信接続部（複数も可）1290 は、通信媒体を介して別のコンピューティングエンティティへの通信を可能にする。通信媒体は、コンピュータ実行可能な命令、音声または画像情報、変調データ信号のデータ、等の情報を伝達する。変調データ信号は、信号の特徴の一つまたは複数が信号中の情報を暗号化するように設定または変更された信号である。例であって限定ではないが、通信媒体としては、電気、光学、RF、赤外線、音響などのキャリアによって実現される有線または無線技術などがある。

【0072】

インプリメンテーション（実現）は、コンピュータ可読媒体の一般的なコンテキストで説明できる。コンピュータ可読媒体は、コンピューティング環境内でアクセスできる媒体であれば、いずれの媒体であってもよい。例であって限定ではないが、コンピューティング環境 1200 内においては、コンピュータ可読媒体としては、メモリ 1220、記憶装置 1240、通信媒体、および上記のいずれかの組み合わせなどが挙げられる。

20

【0073】

図 12 は、コンピューティング環境 1200、表示装置 1260、および入力装置 1250 を個別の装置として示しているが、これは識別のしやすさのためである。コンピューティング環境 1200、表示装置 1260、および入力装置 1250 は個別の装置（例えば、モニタやマウスに有線で接続されたパーソナルコンピュータなど）でもよく、単一の装置に統合されたもの（例えば、スマートフォンやタブレットなどタッチディスプレイを備えたモバイル装置など）でありえ、または装置の組み合わせ（例えば、タッチスクリーン表示装置と動作的に連結される演算装置、単一の表示装置および入力装置に取り付けられた複数の演算装置など）でありえる。コンピューティング環境 1200 はセットトップボックスや、パーソナルコンピュータ、または、一または複数のサーバであってもよく、例えばネットワーク接続されたサーバの形式、クラスタ化サーバ環境、または演算装置のクラウドネットワークでありうる。

30

【0074】

記載された実施形態を参照して本発明の原理を記載し説明したが、記載された実施形態は、構成や詳細においてかかる原理から逸脱することなく変更できることは、認識されるであろう。特段の記載の無い限り、本明細書に記載されたプログラム、プロセス、または方法は、特定の種別のコンピューティング環境に関連するわけではなく、限定もされないことを、理解されたい。各種の汎用のまたは特殊なコンピューティング環境は、本明細書に記載された教示に従って、使用されたり操作が実行されたりできる。ソフトウェアに示される本実施形態の要件は、ハードウェアで実現でき、その逆も可である。

40

【0075】

本発明の原理が適用できる、考えられうる多くの実施形態を鑑み、以下の請求項および均等物の範囲および精神に該当するような実施形態はすべて、本発明であると主張する。

【図 1 A】

顧客名	店舗 ID	出費額
Bill	1	10
Bill	2	5
George	1	11
George	2	11
Max	2	21
Larry	1	31
Larry	2	18

【図 1 B】

店舗 ID	来店数
1	3
2	4

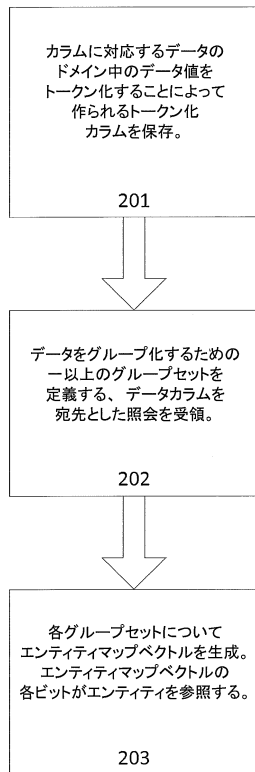
【図 1 C】

店舗 ID	出費額
1	52
2	55

【図 1 D】

店舗 ID	平均出費額	来店数
1	17 1/3	3
2	13 3/4	4

【図 2 A】



【図 1 E】

店舗 ID	最小出費額
1	10
2	5

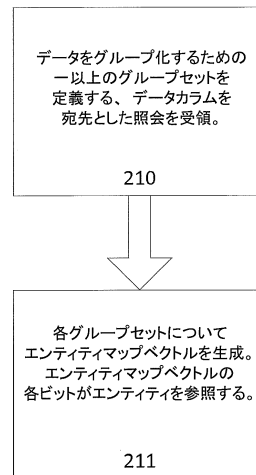
【図 1 F】

店舗 ID	最大出費額
1	31
2	21

【図 1 G】

店舗 ID	個別の客の数
1	3
2	4

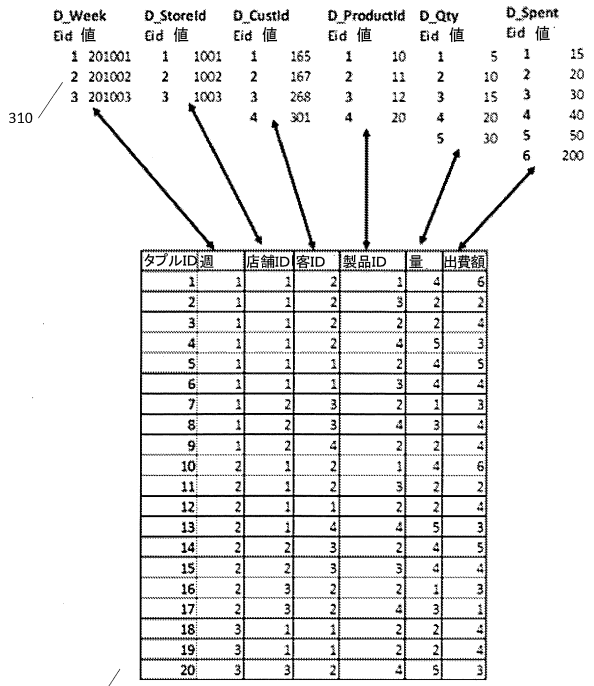
【図 2 B】



【図3A】

301	302	303	304	305	306	307
ダブルID	週	店舗ID	客ID	製品ID	量	出費額
1	201001	1001	167	10	20	200
2	201001	1001	167	12	10	20
3	201001	1001	167	11	10	40
4	201001	1001	167	20	30	30
5	201001	1001	165	11	20	50
6	201001	1001	165	12	20	40
7	201001	1002	268	11	5	30
8	201001	1002	268	20	15	15
9	201001	1002	301	11	10	44
10	201002	1001	167	10	20	200
11	201002	1001	167	12	10	20
12	201002	1001	165	11	10	40
13	201002	1001	301	20	30	30
14	201002	1002	268	11	20	50
15	201002	1002	268	12	20	40
16	201002	1003	167	11	5	30
17	201002	1003	167	20	15	15
18	201003	1001	165	11	10	44
19	201003	1001	165	11	10	40
20	201003	1003	167	20	30	30

【図3C】



【図3B】

312	313	314	315	316	317				
D_Week Eid 値	D_Storeid Eid 値	D_Custid Eid 値	D_Productid Eid 値	D_Qty Eid 値	D_Spent Eid 値				
1	201001	1	165	1	10	5	1	15	
2	201002	2	167	2	11	2	10	2	20
3	201003	3	268	3	12	3	15	3	30
		4	301	4	20	4	20	4	40
				5	30	5	30	5	50
				6	200	6	200	6	200

【図4A】

ダブルID	週	店舗ID	客ID
1	1	1	2
2	1	2	2
3	2	1	3
4	2	2	1
5	2	3	1
6	3	1	1
7	3	3	1

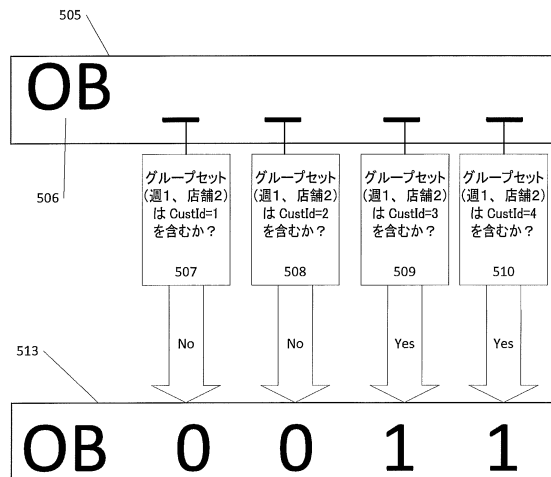
【図5A】

501	502	503	504	513
ダブルID	週	店舗ID	客ID	
1	1	1	1	001100
2	1	1	2	000011
3	2	1	1	001101
4	2	2	2	000010
5	2	3	3	000100
6	3	1	1	001000
7	3	3	3	000100

【図4B】

ダブルID	週	CD_CD_CustId
1	1	1
2	2	2
3	3	1

【図5B】



【図4C】

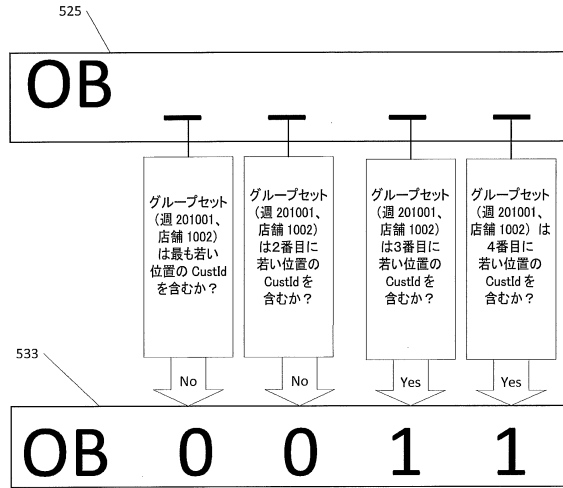
ダブルID	週	Sum_CD_CustId
1	1	4
2	2	5
3	3	2

【図5C】

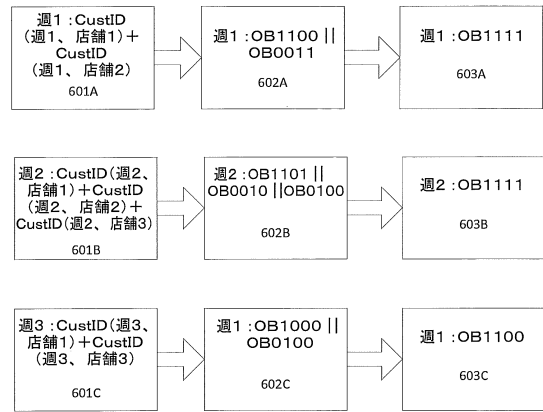
520

D_CustId 中の固有データ値の語彙順の一覧
165
167
268
301

【図5D】



【図6A】

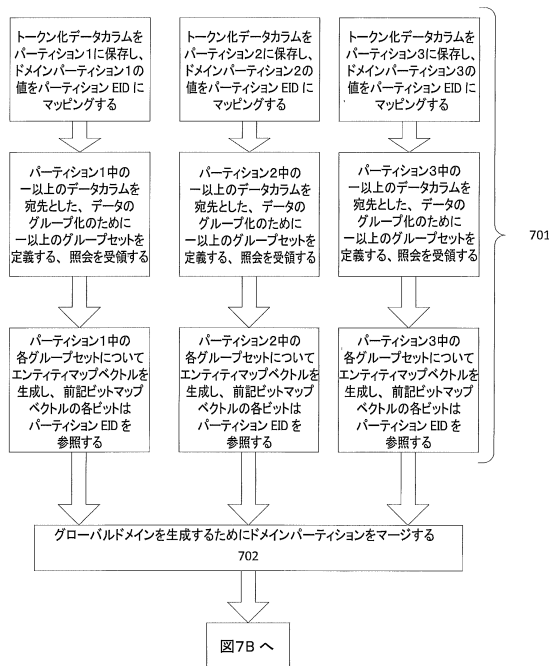


【図6B】

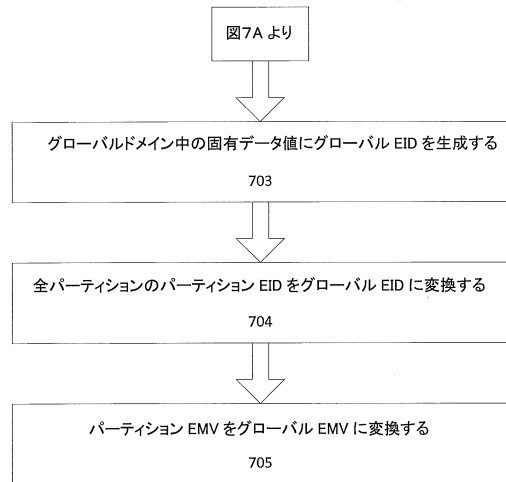
610

タブID	週	EMV_CD_CustId
1	1	4
2	2	4
3	3	2

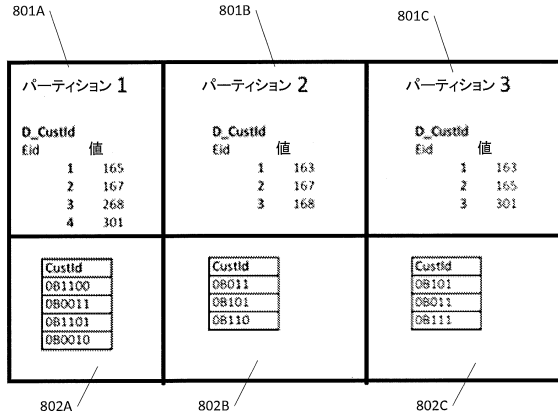
【図7A】



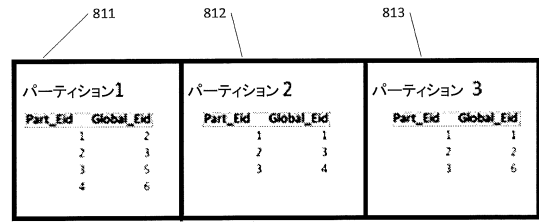
【図7B】



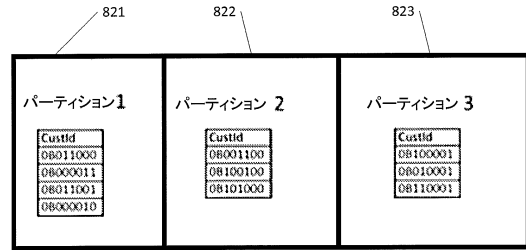
【図 8 A】



【図 8 C】



【図 8 D】



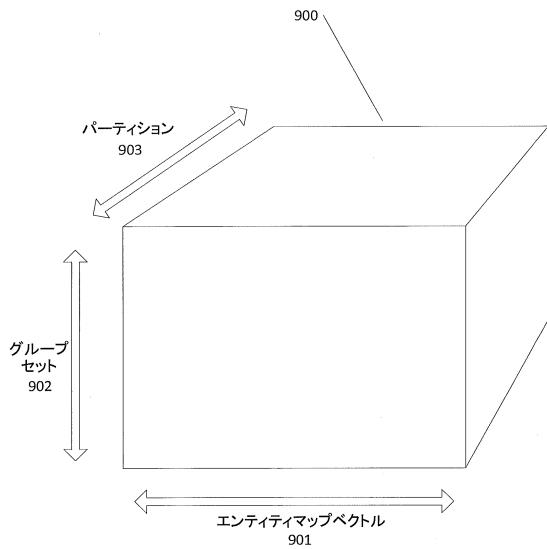
【図 8 B】

D_CustId	値
Eid	
1	163
2	165
3	167
4	168
5	268
6	301

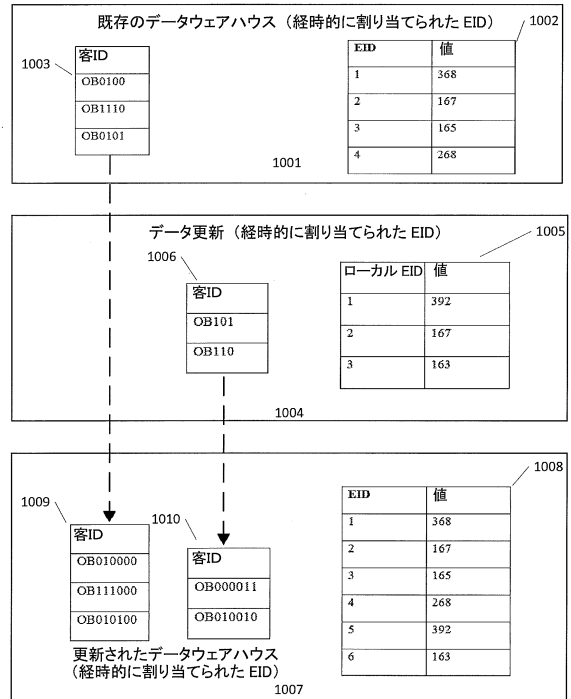
803

804

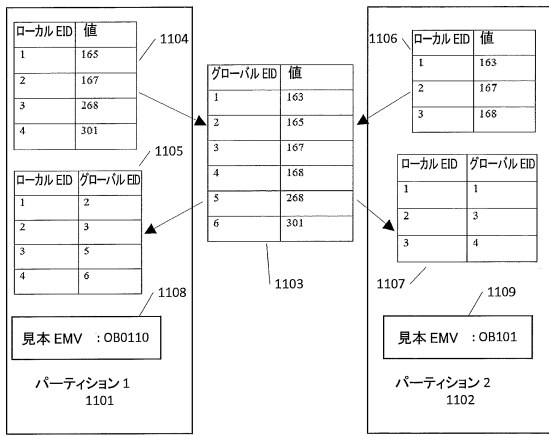
【図 9】



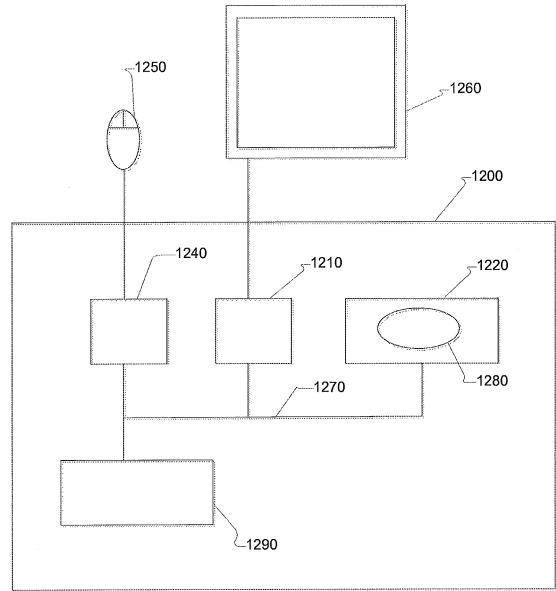
【図 10】



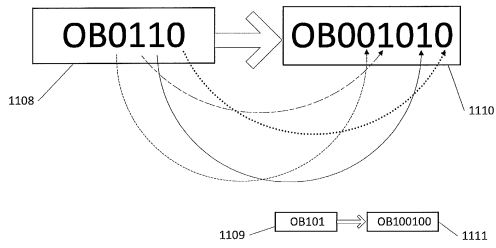
【図 11A】



【図 12】



【図 11B】



フロントページの続き

(72)発明者 ファデエイチエフ、 ユージェニー
カナダ エイチ４ブイ ３エヌ２、ケベック、モントリオール、プリンスオブウェールズ、５１３
９

審査官 吉田 誠

(56)参考文献 特開平 10 - 97544 (JP, A)
米国特許第 6647372 (US, B1)
米国特許第 7653605 (US, B1)
米国特許第 9218379 (US, B1)

(58)調査した分野(Int.Cl., DB名)
G06F 17/30