

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5710782号
(P5710782)

(45) 発行日 平成27年4月30日 (2015. 4. 30)

(24) 登録日 平成27年3月13日 (2015. 3. 13)

(51) Int. Cl.

F I

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

G O 6 F 12/00 5 1 3 J

G O 6 F 17/30 (2006. 01)

G O 6 F 17/30 1 1 0 C

G O 6 F 17/30 1 8 0 D

G O 6 F 17/30 3 4 0 C

請求項の数 19 (全 24 頁)

(21) 出願番号 特願2013-546104 (P2013-546104)
 (86) (22) 出願日 平成23年4月15日 (2011. 4. 15)
 (65) 公表番号 特表2014-505925 (P2014-505925A)
 (43) 公表日 平成26年3月6日 (2014. 3. 6)
 (86) 国際出願番号 PCT/US2011/032631
 (87) 国際公開番号 W02012/087366
 (87) 国際公開日 平成24年6月28日 (2012. 6. 28)
 審査請求日 平成26年2月3日 (2014. 2. 3)
 (31) 優先権主張番号 12/973, 668
 (32) 優先日 平成22年12月20日 (2010. 12. 20)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 506332063
 セールスフォース ドット コム インコ
 ーポレイティッド
 アメリカ合衆国 カリフォルニア州 94
 105, サンフランシスコ, ザ ランドマ
 ーク アット ワン マーケット, スイー
 ト 300
 (74) 代理人 100083806
 弁理士 三好 秀和
 (74) 代理人 100095500
 弁理士 伊藤 正和
 (74) 代理人 100111235
 弁理士 原 裕子

最終頁に続く

(54) 【発明の名称】 マルチテナントストアで横断的ストア結合を行う方法及びシステム

(57) 【特許請求の範囲】

【請求項 1】

リレーショナルデータストア及び非リレーショナルデータストアを有するマルチテナン
 トデータベースシステムからデータを取得する方法であって、

ホストシステムで前記マルチテナントデータベースシステムに対する要求を受信するス
 テップであって、前記要求は前記マルチテナントデータベースシステムから取得されるデ
 ータを特定するステップと、

前記ホストシステムを介する前記要求に基づいて、前記取得されるデータの1つ以上の
 位置を取得するステップと、

前記ホストシステムで前記要求に基づいて取得される複数のデータ要素を特定するデー
 タベースクエリを生成するステップであって、前記複数のデータ要素は前記非リレーショ
 ナルデータストア内に存在する1つ以上のデータ要素及び前記リレーショナルデータストア
 内に存在する1つ以上の他のデータ要素を含むステップと、

前記データを取得するために前記マルチテナントデータベースシステムに対して前記デー
 タベースクエリを実行するステップと

を含み、

前記データベースクエリは結合動作を特定し、

前記結合動作は、

前記非リレーショナルデータストアに対して実行される第1のサブクエリであって、前
 記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素を識別する第1

10

20

のサブクエリと、

前記リレーショナルデータストアに対して実行される第2のサブクエリであって、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素を識別する前記第1のサブクエリと前記リレーショナルデータストア内に存在する前記1つ以上の他のデータ要素との間のデータデルタを決定する第2のサブクエリと、

前記リレーショナルデータストア及び前記非リレーショナルデータストアに対して実行される第3のサブクエリであって、前記リレーショナルデータストアから前記非リレーショナルデータストアへと決定されたデータデルタに対応するデータを複製する第3のサブクエリと、

前記非リレーショナルデータストアに対して実行される第4のサブクエリであって、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素と前記リレーショナルデータストアから前記非リレーショナルデータストアに複製されることで前記非リレーショナルデータストア内から利用可能な前記1つ以上の他のデータ要素との双方をフェッチすることにより、前記非リレーショナルデータストアから取得される前記データをフェッチする第4のサブクエリと

を含む、方法。

【請求項2】

顧客スキーマが前記取得されるデータの1つ以上の位置を記述し、前記顧客スキーマは前記非リレーショナルデータストア内若しくは前記リレーショナルデータストア内に存在する又は前記非リレーショナルデータストア及び前記リレーショナルデータストアの双方から利用可能である前記データの複数のデータ要素の各々を特定し、

前記方法は、前記ホストシステムを介して、前記要求の受信に応答して前記顧客スキーマを取得するステップを更に含む、請求項1に記載の方法。

【請求項3】

前記データベースクエリは、複数のサブクエリを含み、

前記複数のサブクエリの中の少なくとも1つのサブクエリは、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素を前記非リレーショナルデータストアから取得するために行われ、

前記複数のサブクエリの中の少なくとも第2のサブクエリは、前記リレーショナルデータストア内に存在する前記1つ以上の他のデータ要素を前記リレーショナルデータストアから取得するために行われる、請求項1に記載の方法。

【請求項4】

前記マルチテナントデータベースシステムに対して前記データベースクエリを実行するステップは、前記リレーショナルデータストア及び前記非リレーショナルデータストアの双方に記憶されたデータ要素を参照することを含む、請求項1に記載の方法。

【請求項5】

前記結合動作は、前記ホストシステムの最適化エージェントを介して利用可能なクエリ最適化に基づいて生成される複数のサブクエリを含み、前記クエリ最適化は、

前記複数のサブクエリに対する特定の順序付け、

対応するサブクエリの実行に対する目的データストア、

前記取得されるデータに基づく1つ以上の事前クエリ評価、

前記リレーショナルデータストアから前記非リレーショナルデータストアへの複製命令、

前記リレーショナルデータストア及び前記非リレーショナルデータストアの各々から取得される前記複数のデータ要素の少なくとも1つ以上を特定し、且つ前記最適化エージェントにアクセス可能なメモリに配置されるメモリ内結合動作、及び

前記要求を充足するに際し前記最適化エージェントにアクセス可能な前記メモリから前記複数のデータ要素の前記少なくとも1つ以上を取得するための対応するサブクエリを含むグループから選択される、請求項1に記載の方法。

【請求項6】

前記結合動作は、

前記リレーショナルデータストアから2つ以上の関連テーブルを特定する結合動作と、
前記リレーショナルデータストアからの少なくとも1つの関連テーブル及び前記非リレーショナルデータストア内に存在する少なくとも1つ以上のデータ構造を特定する結合動作と、

前記非リレーショナルデータストア内に存在する2つ以上の離れた及び別個のデータ構造を特定する結合動作であって、前記2つ以上の離れた及び別個のデータ構造の各々が重複する共有キーを欠いている結合動作と

を含むグループから選択される、請求項1に記載の方法。

【請求項7】

10

前記結合動作は、

前記リレーショナルデータストアに対して実行される第5のサブクエリであって、決定されたデータデルタに基づいて前記リレーショナルデータストア内に存在する前記1つ以上の他のデータ要素を取得する第5のサブクエリと

を含む、請求項1に記載の方法。

【請求項8】

前記第5のサブクエリは、

前記リレーショナルデータストアに対して実行される時間ベースのサブクエリであって、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素内のデータ要素に対応する任意のタイムスタンプよりも遅いタイムスタンプを有する前記リレーショナルデータストア内のデータ要素に基づいて前記リレーショナルデータストアから取得される前記1つ以上の他のデータ要素を特定する時間ベースのサブクエリと、

20

前記リレーショナルデータストアに対して実行される記録識別子ベースのサブクエリであって、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素内のデータ要素に対応する任意の記録識別子よりも数値が大きい記録識別子を有する前記リレーショナルデータストア内のデータ要素に基づいて前記リレーショナルデータストアから取得される前記1つ以上の他のデータ要素を特定する記録識別子ベースのサブクエリと

の1つを含む、請求項7に記載の方法。

【請求項9】

前記マルチテナントデータベースシステムにおいて複数の新しいトランザクションを受信するステップであって、前記新しいトランザクションの各々は前記非リレーショナルデータストアに書き込まれる新しいデータを特定するステップと、

30

前記リレーショナルデータストアの追加ログに前記新しいデータを書き込むステップと
を更に含み、決定されたデータデルタに基づいて前記リレーショナルデータストア内に存在する前記1つ以上の他のデータ要素を取得する第5のサブクエリは、前記リレーショナルデータストアの前記追加ログから前記リレーショナルデータストア内に存在する前記1つ以上の他のデータ要素を取得する第5のサブクエリを含む、請求項7に記載の方法。

【請求項10】

前記追加ログがフラッシュ閾値に到達すると、前記リレーショナルデータストアの前記追加ログに書き込まれた前記新しいデータを前記非リレーショナルデータストアにフラッシュすることを更に含む、請求項9に記載の方法。

40

【請求項11】

前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素は、複数の圧縮フラットファイル若しくは複数のバイナリファイル又は前記圧縮フラットファイル及び前記バイナリファイルの組み合わせを含む、請求項1に記載の方法。

【請求項12】

前記リレーショナルデータストアは、リレーショナルデータベース管理システム(RDBMS)に従って実装されるリレーショナルデータベースを含み、前記リレーショナルデータベースの複数の関連テーブルは、前記リレーショナルデータベース内の2つ以上の関連テーブルの各々に対して1つ以上の重複する共通の特徴を介して相互に関連付けされる

50

、請求項 1 に記載の方法。

【請求項 1 3】

前記非リレーショナルデータストアは、各々が前記非リレーショナルデータストアに総記憶容量の少なくとも一部を提供する複数の基本ハードウェア記憶装置を有する分散構造型データベースを含み、前記非リレーショナルデータストア内のデータ要素は、主キーに基づいて参照可能であるが、2 つ以上の関連テーブルの間の 1 つ以上の重複する共通の特徴に基づいて参照可能ではない、請求項 1 に記載の方法。

【請求項 1 4】

前記リレーショナルデータストアは、Oracle 互換データベースの実装、IBM DB2 Enterprise Server 互換リレーショナルデータベースの実装、MySQL 互換リレーショナルデータベースの実装、及び Microsoft SQL Server 互換リレーショナルデータベースの実装を含むグループから選択されるリレーショナルデータベースの実装を含み、

前記非リレーショナルデータストアは、Vampire 互換非リレーショナルデータベースの実装、Apache Cassandra 互換非リレーショナルデータベースの実装、BigTable 互換非リレーショナルデータベースの実装、及び HBase 互換非リレーショナルデータベースの実装を含むグループから選択される NoSQL 非リレーショナルデータベースの実装を含む、請求項 1 に記載の方法。

【請求項 1 5】

前記マルチテナントデータベースシステムのインターフェースを介して前記要求を受信するステップは、前記マルチテナントデータベースシステムのウェブサーバを介して前記要求を受信することを含み、前記ウェブサーバは、前記要求の起点となる遠隔に設置されたエンドユーザクライアントマシンとのウェブベースインターフェースを提供し、

前記要求は、前記マルチテナントデータベースシステムに関するホスト組織内で動作する前記マルチテナントデータベースシステムからのサービスに対する要求を含む、請求項 1 に記載の方法。

【請求項 1 6】

前記マルチテナントデータベースシステムは、複数の離れた別個の顧客組織によって共有されるハードウェア及びソフトウェアの要素を更に含み、前記離れた別個の顧客組織の各々は前記マルチテナントデータベースシステムが実行されるホスト組織から遠隔に設置される、請求項 1 に記載の方法。

【請求項 1 7】

前記非リレーショナルデータストアは複数の分散型コンピュータノードを含み、各コンピュータノードは少なくとも 1 つのメモリ、1 つ以上のプロセッサ及び 1 つ以上の通信可能にインターフェース接続されたハードディスクドライブを備え、前記分散型コンピュータノードの各々は、中央トランザクション局からの承認又は制御無しで非リレーショナルデータベーストランザクションを読み取り、書き込み、及び更新する機能を有する分離した非リレーショナルデータベースインスタンスを含み、

前記リレーショナルデータストアは、モノリシックリレーショナルデータベースインスタンスへの更新又は変更が前記モノリシックリレーショナルデータベースインスタンスに通信可能にインターフェース接続され且つ制御される持続的記憶装置に対する持続的記憶にコミットされるかどうかを制御する中央トランザクション局と計算リソースを調整するプロセッサ及びメモリを含むモノリシックリレーショナルデータベースインスタンスを含む、請求項 1 に記載の方法。

【請求項 1 8】

ホストシステム内のプロセッサによって実行されると方法を実行する命令が記憶された一時的でないコンピュータ可読記憶媒体であって、前記方法は、

マルチテナントデータベースシステムから取得されるデータを特定する要求を受信するステップと、

前記要求に基づいて、前記取得されるデータの 1 つ以上の位置を取得するステップと、

前記要求に基づいて取得される複数のデータ要素を特定するデータベースクエリを生成するステップであって、前記複数のデータ要素は前記マルチテナントデータベースシステムの非リレーショナルデータストア内に存在する1つ以上のデータ要素及び前記マルチテナントデータベースシステムのリレーショナルデータストア内に存在する1つ以上の他のデータ要素を含むステップと、

前記データを取得するために前記マルチテナントデータベースシステムに対して前記データベースクエリを実行するステップと

を含み、

前記データベースクエリは結合動作を特定し、前記結合動作は、

前記非リレーショナルデータストアに対して実行される第1のサブクエリであって、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素を識別する第1のサブクエリと、

前記リレーショナルデータストアに対して実行される第2のサブクエリであって、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素を識別する前記第1のサブクエリと前記リレーショナルデータストア内に存在する前記1つ以上の他のデータ要素との間のデータデルタを決定する第2のサブクエリと、

前記リレーショナルデータストア及び前記非リレーショナルデータストアに対して実行される第3のサブクエリであって、前記リレーショナルデータストアから前記非リレーショナルデータストアへと決定されたデータデルタに対応するデータを複製する第3のサブクエリと、

前記非リレーショナルデータストアに対して実行される第4のサブクエリであって、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素と前記リレーショナルデータストアから前記非リレーショナルデータストアに複製されることで前記非リレーショナルデータストア内から利用可能な前記1つ以上の他のデータ要素との双方をフェッチすることにより、前記非リレーショナルデータストアから取得される前記データをフェッチする第4のサブクエリと

を含む、一時的でないコンピュータ可読記憶媒体。

【請求項19】

プロセッサ及びメモリと、

リレーショナルデータストア及び非リレーショナルデータストアが実装されるマルチテナントデータベースシステムとの通信インターフェースと、

前記マルチテナントデータベースシステムから取得されるデータを特定する要求を受信する要求プロセッサと、

前記要求に基づいて、前記取得されるデータの1つ以上の位置を取得する顧客スキーマプロセッサと、

前記要求に基づいて取得される複数のデータ要素を特定するデータベースクエリを生成するサブクエリ生成器であって、前記複数のデータ要素は前記マルチテナントデータベースシステムの前記非リレーショナルデータストア内に存在する1つ以上のデータ要素及び前記マルチテナントデータベースシステムの前記リレーショナルデータストア内に存在する1つ以上の他のデータ要素を含むサブクエリ生成器と、

前記データを取得するために前記マルチテナントデータベースシステムに対して前記データベースクエリを実行するクエリ実行器と

を備え、

前記要求に基づいて前記データベースクエリを生成する前記サブクエリ生成器は、前記要求を充足するために前記要求に基づいて複数のサブクエリを生成するサブクエリ生成器を含み、前記サブクエリは、

前記非リレーショナルデータストアに対して実行される第1のサブクエリであって、前記非リレーショナルデータストア内に存在する前記1つ以上のデータ要素を識別する第1のサブクエリと、

前記リレーショナルデータストアに対して実行される第2のサブクエリであって、前記

10

20

30

40

50

非リレーショナルデータストア内に存在する前記１つ以上のデータ要素を識別する前記第１のサブクエリと前記リレーショナルデータストア内に存在する前記１つ以上の他のデータ要素との間のデータデルタを決定する第２のサブクエリと、

前記リレーショナルデータストア及び前記非リレーショナルデータストアに対して実行される第３のサブクエリであって、前記リレーショナルデータストアから前記非リレーショナルデータストアへと決定されたデータデルタに対応するデータを複製する第３のサブクエリと、

前記非リレーショナルデータストアに対して実行される第４のサブクエリであって、前記非リレーショナルデータストア内に存在する前記１つ以上のデータ要素と前記リレーショナルデータストアから前記非リレーショナルデータストアに複製されることで前記非リレーショナルデータストア内から利用可能な前記１つ以上の他のデータ要素との双方をフェッチすることにより、前記非リレーショナルデータストアから取得される前記データをフェッチする第４のサブクエリと

を含む、システム。

【発明の詳細な説明】

【技術分野】

【０００１】

本明細書に記載の主題は、概してコンピュータの分野に関し、より詳細にはマルチテナントストアで横断的ストア結合を行う方法及びシステムに関する。

【０００２】

（優先権の主張）

本願は、出願番号１２／９７３，６６８及び代理人整理番号８９５６Ｐ００６／３２０ＵＳを有し、２０１０年１２月２０日に提出され、“METHODS AND SYSTEMS FOR PERFORMING CROSS STORE JOINS IN A MULTI-TENANT STORE”と題された米国特許出願、並びに出願番号６１／３２５，７０９及び代理人整理番号８９５６Ｐ００６Ｚ／３２０PROVを有し、２０１０年４月１９日に提出され、“METHODS AND SYSTEMS FOR PERFORMING CROSS STORE JOINS IN A MULTI-TENANT STORE”と題された仮特許出願に関し、これらへの優先権を主張すると共に、全ての内容が参照により本明細書に組み込まれる。本願は、出願番号６１／３２５，９５１及び代理人整理番号８９５６Ｐ００７Ｚ／３２１PROVを有し、２０１０年４月２０日に提出され、“METHODS AND SYSTEMS FOR OPTIMIZING QUERIES IN A MULTI-TENANT STORE”と題された仮特許出願に更に関し、これへの優先権を主張し、全ての内容が参照により本明細書に組み込まれる。

【０００３】

（著作権表示）

本文書の開示の一部は著作権によって保護される資料を含んでいる。著作権者は、特許文献又は特許開示が特許商標庁の特許出願又は記録に現れるものについては、何人によるこれらの複製に対して異議を持たないが、それ以外はどのようなものであっても全ての著作権を保有する。

【背景技術】

【０００４】

背景欄で検討される主題は、単に背景欄で言及された結果として先行技術であるとみなされるべきではない。同様に、背景欄で言及される又は背景欄の主題と関連する課題は、先行技術で過去に認識されていたとみなされるべきではない。背景欄における主題は、請求項に記載された主題の実施形態に対応し得る異なるアプローチを表しているだけである。

【０００５】

コンピュータ環境内では、データを持続的に記憶するために様々なデータ記憶環境が選

10

20

30

40

50

扱われる場合がある。例えば、データは、ハードドライブにファイルシステムデータを持続的に記憶するオペレーティングシステムによって管理されるファイルシステム内に記憶されてもよい。又はデータは、データベース内に持続的に記憶されてもよい。各々が固有の利点及び欠点を有する様々な種類のデータベースが利用可能である。例えば、所謂リレーショナルデータベースは、各テーブルによって共有される共通の特徴を使用して、様々なデータテーブルをデータベース内で互いに「関連付け」する能力を提供する。例えば、リレーショナルデータベースでは、従業員識別子は1つより多くのテーブルを関連付ける共通の特徴として使用されてもよい。しかしながら、このようなデータベース構造はいくつかの欠点を有している。そのうちの1つは、関連が高水準の計算オーバーヘッドコスト及びリレーショナルデータベースのスケラビリティの範囲を制限する計算の複雑性を伴うという点である。

10

【0006】

非リレーショナルデータベースのモデル及び実装も存在し、一般的にはより良好なスケラビリティを示すが、リレーショナルデータベースのモデル及び実装と関連しない異なる欠点も示す。例えば、非リレーショナルデータベースの実装は、大きなファイル又はオブジェクトを記憶するための改善されたスケラビリティを示すことが多いが、選択的なデータセットを記憶し又は急速に変化するデータセットに対してデータ保障を実装する等の他の点ではあまり適切でない場合がある。

【0007】

残念ながら、同時に複数のデータから情報を参照するデータベースクエリは著しく非効率的であり、複数のデータストアの実装から得られる可能性がある利益を損なう。更に、分散データベースモデルの個別の実装を同時に参照するデータベースクエリは、以前のデータベースクエリ機構の使用では完全に実行不可能な場合がある。

20

【0008】

実施形態は例示目的で示されており、限定を目的とするものではなく、以下の詳細な説明を参照して図面と併せて検討すればより完全に理解され得る。

【図面の簡単な説明】

【0009】

【図1】実施形態が動作可能な例示的なアーキテクチャを示す。

【図2】実施形態が動作可能な代替の例示的なアーキテクチャを示す。

30

【図3】実施形態が動作可能な代替の例示的なアーキテクチャを示す。

【図4】実施形態が動作可能な代替の例示的なアーキテクチャを示す。

【図5】実施形態が動作し、設置され、統合され、又は構成され得るシステムの図表現を示す。

【図6】一実施形態に従ってマルチテナントストアで横断的ストア結合を実行する方法の流れ図を示す。

【図7】一実施形態に従うコンピュータシステムの例示的態様で機械の図表現を示す。

【発明を実施するための形態】

【0010】

本明細書にはマルチテナントストアで横断的ストア結合を実行するシステム、装置、及び方法が記載される。一実施形態では、このような方法は、リレーショナルデータストア及び非リレーショナルデータストアを有するマルチテナントデータベースシステムからデータを取得することを含む。例えば、このような方法では、マルチテナントデータベースシステムに関するホストシステムは、マルチテナントデータベースシステムから取得されるデータを特定する要求を受信し、要求に基づいて、ホストシステムを介して取得されるデータの1つ以上の位置を取得し、ホストシステムにおいて、要求に基づいて取得される複数のデータ要素を特定するデータベースクエリを生成し、複数のデータ要素は非リレーショナルデータストア内に存在する1つ以上のデータ要素及びリレーショナルデータストア内に存在する1つ以上の他のデータ要素を含み、データを取得するためにマルチテナントデータベースシステムに対してデータベースクエリを実行する。

40

50

【0011】

1つより多くのデータベースを検索又は参照する連合クエリは、異なるデータストアに記憶されたテーブル間の結合動作を要求すること等によって、特にデータベースの最下行レベルでデータを参照する場合に、著しく非効率的である。その理由は、この動作は非常に多くのネットワーク帯域幅を消費するので、このような結合動作はうまくスケールせず、より大きなデータベースの実装で実行できるように実装できないからである。このような結合動作の課題は、例えば、リレーショナルデータベースの実装と非リレーショナルデータベースの実装との結合等、分散モデルで動作するデータベースの複数の実装間でデータ結合を要求する場合に更に悪化する。本明細書に記載の手法は、より大きなデータベースシステムで実行できるように実装可能なように、特に、リレーショナル及び非リレーショナルモデル等の分散動作モデルで動作する複数のデータストアの実装を利用するシステムで実行できるように実装可能なように、このような結合動作を実行する能力を促進する。

10

【0012】

例えば、本明細書に記載の手法を用いて、結合動作は、非リレーショナルデータベースに記憶されたオブジェクトに対する非リレーショナルデータベースクエリを開始することによって実行されてもよく、ここでは1つ以上の外部キーペアレントがOracleTM等のリレーショナル型データベースの実装で記憶されるオブジェクトである。例えば非リレーショナルデータベースに記憶されたオブジェクトが非リレーショナルデータベースの実装に存在すると共にOracleTMに記憶されたオブジェクトがリレーショナルデータベースの実装に存在するにもかかわらず、非リレーショナルデータベースに記憶されたチャイルドテーブルは、マスターテーブルとしてOracleTMに記憶された“Account”テーブルを有してもよい。

20

【0013】

以下の記載では、様々な実施形態の十分な理解を提供するために、多くの特定の詳細が特定のシステム、言語、コンポーネント等の例として説明される。しかしながら、特定の詳細は開示の実施形態を実施するために用いられる必要がないことが当業者には明らかであろう。他の例では、開示の実施形態を不必要に分かりにくくするのを避けるために、周知の材料又は方法は詳細には記載されていない。

【0014】

図面に描かれた及び本明細書に記載された様々なハードウェア要素に加えて、実施形態は以下に記載の様々な動作を更に含む。このような実施形態に従って記載された動作は、ハードウェア要素によって実行され、又は命令と共にプログラミングされた汎用若しくは専用プロセッサに動作を実行させるために使用可能な機械実行可能命令で具現化されてもよい。代替的に、動作はハードウェア及びソフトウェアの組み合わせによって実行されてもよい。

30

【0015】

実施形態は、本明細書に記載の動作を実行するためのシステム又は装置にも関する。開示のシステム又は装置は、要求された目的で特別に構築されてもよく、又はコンピュータに記憶されたコンピュータプログラムによって選択的に起動又は再構成される汎用コンピュータを含んでもよい。このようなコンピュータプログラムは、限定されないが、フロッピー（登録商標）ディスク、光学ディスク、CD-ROM及び磁気光学ディスクを含む任意の種類のディスク、ROM (read-only memory)、RAM (random access memory)、各々がコンピュータシステムバスに結合されるEPROM、EEPROM、磁気若しくは光学カード又は一時的でない電子命令を記憶するのに適した任意の種類の媒体等の一時的でないコンピュータ可読記憶媒体に記憶されてもよい。一実施形態では、命令が記憶されたコンピュータ可読記憶媒体は、マルチテナントデータベース環境内の1つ以上のプロセッサに本明細書に記載の方法及び動作を実行させる。別の実施形態では、このような方法及び動作を実行するための命令は、後で実行するために一時的でないコンピュータ可読媒体に記憶される。

40

50

【 0 0 1 6 】

本明細書に提示されるアルゴリズム及びディスプレイは、本質的に特定のコンピュータ又は他の装置に関連せず、特定のプログラミング言語を参照して記載された実施形態でもない。様々なプログラミング言語が本明細書に開示された実施形態の教示を実装するために利用可能なことが理解されるであろう。

【 0 0 1 7 】

図 1 は、実施形態が動作可能な例示的なアーキテクチャ 1 0 0 を示す。アーキテクチャ 1 0 0 は、ネットワーク 1 2 5 を介して複数の顧客組織 (1 0 5 A、1 0 5 B 及び 1 0 5 C) と通信可能にインターフェース接続されたホストシステム 1 1 0 を表す。ホストシステム 1 1 0 内でデータベース機能及びコード実行環境を実装する複数の基本ハードウェア、ソフトウェア及び論理要素 1 2 0 を有するマルチテナントデータベースシステム 1 3 0 がホストシステム 1 1 0 内にあり、マルチテナントデータベースシステム 1 3 0 のハードウェア、ソフトウェア及び論理要素 1 2 0 は、ネットワーク 1 2 5 を介してホストシステム 1 1 0 に通信可能にインターフェース接続することによりホストシステム 1 1 0 によって提供されるサービスを利用する複数の顧客組織 (1 0 5 A、1 0 5 B 及び 1 0 5 C) とは離れており別個である。このような実施形態では、離れた別個の顧客組織 (1 0 5 A - 1 0 5 C) の各々は、マルチテナントデータベースシステム 1 3 0 が実行されるホストシステム 1 1 0 を介して顧客組織 (1 0 5 A - 1 0 5 C) にサービスを提供するホスト組織から遠く離れて設置されてもよい。代替的に、1 つ以上の顧客組織 1 0 5 A - 1 0 5 C が、基本データが持続的に記憶されるマルチテナントデータベースシステム 1 3 0 を提供する同一のホスト組織内等、ホストシステム 1 1 0 と共同設置されてもよい。

【 0 0 1 8 】

一実施形態では、マルチテナントデータベースシステム 1 3 0 のハードウェア、ソフトウェア及び論理要素 1 2 0 は少なくとも非リレーショナルデータストア 1 5 0 及びリレーショナルデータストア 1 5 5 を含み、これらはホストシステム 1 1 0 内でデータベース機能及びコード実行環境を実装するハードウェア、ソフトウェア及び論理要素 1 2 0 に従って動作する。ホストシステム 1 1 0 は、ネットワークを介して複数の顧客組織 1 0 5 A - 1 0 5 C の 1 つ以上から要求 1 1 5 を更に受信してもよい。例えば、着信要求 1 1 5 は、マルチテナントデータベースシステム 1 3 0 内の顧客組織 1 0 5 A - C の 1 つのためにサービスに対する要求又はデータを取得若しくは記憶する要求に対応してもよい。

【 0 0 1 9 】

図 2 は、実施形態が動作可能な代替の例示的なアーキテクチャ 2 0 0 を示す。一実施形態では、ホストシステム 1 1 0 は、リレーショナルデータストア 1 5 5 及び非リレーショナルデータストア 1 5 0 を有するマルチテナントデータベースシステム 1 3 0 からデータを取得する方法を実装する。

【 0 0 2 0 】

例えば、このような実施形態では、ホストシステム 1 1 0 でマルチテナントデータベースシステム 1 3 0 に対する要求 1 1 5 が受信される。この要求 1 1 5 はマルチテナントデータベースシステム 1 3 0 から取得されるデータ 2 1 8 を特定する。一部の実施形態では、ホストシステム 1 1 0 内で動作する別個のウェブサーバ 2 1 0 がネットワーク 1 2 5 を介して着信要求 1 1 5 を受信する。例えば、ウェブサーバ 2 1 0 は、ネットワーク 1 2 5 を介して様々な顧客組織 1 0 5 A - C から要求 1 1 5 を受信することを担当する。ウェブサーバ 2 1 0 は、要求 1 1 5 の起点となるエンドユーザクライアントマシン (例えば、顧客組織 1 0 5 A - C 内に設置されるエンドユーザ装置等) とのウェブベースインターフェースを提供してもよい。この要求 1 1 5 は、例えば、遠隔に実装されたクラウドコンピュータサービスを提供するホストシステム 1 1 0 等のホスト組織内で動作するマルチテナントデータベースシステム 1 3 0 からのサービスに対する要求を構成する。また、最適化エージェント 2 4 5 が所定の実施形態に従って事前クエリを展開すること及びデータクエリを最適化すること等の追加の機能を提供してもよい。

【 0 0 2 1 】

一実施形態では、ホストシステム 110 は、要求 115 に基づいて、取得されるデータ 218 の 1 つ以上の位置 216 を取得する。一実施形態では、顧客スキーマ 240 が取得されるデータ 218 の 1 つ以上の位置 216 を記述する。ここで、顧客スキーマ 240 は、非リレーショナルデータストア 150 内若しくはリレーショナルデータストア 155 内に存在する又は非リレーショナルデータストア 150 及びリレーショナルデータストア 155 の双方から利用可能である取得されるデータ 218 の複数のデータ要素の各々を特定する。一実施形態では、ホストシステム 110 は要求 115 の受信に応答して顧客スキーマ 240 を取得する。代替的に、ホストシステム 110 は、顧客スキーマ 240 から取得されるデータ 218 の 1 つ以上の位置 216 を取得する。

【0022】

例えば、特定の実施形態において、取得されるデータ 218 の 1 つ以上の位置 216 は、取得されるデータ 218 を構成する複数のデータ要素の各々がマルチテナントデータベースシステム内のどこに設置されているかを特定する顧客スキーマ 240 に記憶され及びこれから取得される。このような顧客スキーマ 240 は、例えば、上記マルチテナント記憶能力を実装又は提供するホストシステム 110 の様々な要素に高速で効率的なアクセスを提供するグローバルキャッシュレイヤを介してアクセス可能であってもよい。代替的な実施形態では、取得されるデータ 218 の 1 つ以上の位置 216 は、ホストシステム 110 によって、最適化エージェント 245 によって、ホストシステム 110 のクエリレイヤ 260 によって、又は描かれている非リレーショナルデータストア 150 及びリレーショナルデータストア 155 に散在する複数のデータ要素を有するデータ 218 等、分散データベースの実装に散在するマルチテナントデータベースシステム 130 から取得されるデータ 218 の位置 216 を決定することに関与するホストシステム 110 の他の要素によって顧客スキーマ 240 から取得されてもよい。

【0023】

一実施形態では、ホストシステム 110 は、要求 115 に基づいてデータベースクエリ 217 を生成する。ここで、データベースクエリ 217 は、取得される複数のデータ要素であって、非リレーショナルデータストア 150 内に存在する 1 つ以上のデータ要素及びリレーショナルデータストア 155 内に存在する 1 つ以上の他のデータ要素を含む複数のデータ要素を特定する。特定の実施形態では、データベースクエリ 217 は、取得されるデータ 218 の取得される 1 つ以上の位置 216 に更に基づく。このようなデータベースクエリ 217 は、クエリレイヤ 260 又は最適化エージェント 245 等のホストシステム 110 のサブシステムによって生成するためにホストシステム 110 によって更に生成又は代理されてもよい。

【0024】

一実施形態では、ホストシステム 110 は、マルチテナントデータベースシステム 130 に対して生成されたデータベースクエリ 217 を実行して、図 2 によって描かれているようなデータ 218 を取得する。ここで、マルチテナントデータベースシステム 130 の複数の基本ハードウェア、ソフトウェア及び論理要素 120 に向かう下向き矢印は、マルチテナントデータベースシステム 130 の実装機能に渡されるデータベースクエリ 217 を描く。また、上向きに曲がった矢印によって描かれるマルチテナントデータベースシステムによって応答して戻されるデータ 218 は、分散データストア、非リレーショナルデータストア 150 及びリレーショナルデータストア 155 の各々を起点とする複数のデータ要素をホストシステム 110 に返信する。

【0025】

一実施形態では、データベースクエリ 217 は複数のサブクエリを含む。このような実施形態では、複数のサブクエリの中の少なくとも 1 つは、非リレーショナルデータストア 150 内に存在する 1 つ以上のデータ要素を非リレーショナルデータストア 150 から取得するために行われる。また、複数のサブクエリの中の少なくとも第 2 のサブクエリは、リレーショナルデータストア 155 内に存在する 1 つ以上の他のデータ要素をリレーショナルデータストア 155 から取得するために行われる。例えば、「非リレーショナルデー

タストアからデータ要素「a」を取得」(例えば、150)及び「リレーショナルデータストアからデータ要素「b」を取得」(例えば、155)等の複数のサブクエリ列並びに一般的なSQL(Structured Query Language)型のクエリを反映して「select 'x' from 'y' where 'z'」を記述する別のサブクエリ列が、データベースクエリ217の拡大図内に図2によって描かれる。このようなクエリは、選択された実装クエリ言語又は文法に依存して基本データストア(例えば、150及び155)にクエリを行うのに適切であってもよく、又はそうでなくてもよい。

【0026】

従って、このような実施形態に従って、マルチテナントデータベースシステム130に対してデータベースクエリ217を実行することは、必要なデータ218を取得するようにリレーショナルデータストア155及び非リレーショナルデータストア150の双方に記憶されたデータ要素を参照することを含む。

10

【0027】

図3は、実施形態が動作可能な代替の例示的なアーキテクチャ300を示す。特に、所定の実施形態に従ってデータベースクエリ217により特定される結合動作が更に描かれる。

【0028】

例えば、一実施形態によれば、結合動作305はデータベースクエリ217によって特定される。

20

【0029】

特定の実施形態では、結合動作305は複数のサブクエリを含む。例えば、このような実施形態では、第1のサブクエリ306が非リレーショナルデータストア150に対して実行されて、非リレーショナルデータストア150内に存在する1つ以上のデータ要素を識別する。これは非リレーショナルデータストア150への湾曲した破線の矢印によって描かれている。

【0030】

このような実施形態では、第2のサブクエリ307がリレーショナルデータストア155に対して実行されて、非リレーショナルデータストア150内に存在する1つ以上のデータ要素を識別する第1のサブクエリとリレーショナルデータストア155内に存在する1つ以上の他のデータ要素との間のデータデルタ310を決定する。

30

【0031】

この実施形態では、第3のサブクエリ308がリレーショナルデータストア155及び非リレーショナルデータストア150に対して実行される。第3のサブクエリはリレーショナルデータストア155から非リレーショナルデータストア150へと決定されたデータデルタ310に対応するデータを複製する。例えば、このような第3のサブクエリ308は、リレーショナルデータストア155内に存在する1つ以上の他のデータ要素を取得して、例えば、一時テーブル、ファイルに入れ、データを一時的にキャッシュ等してもよい。その後、このような第3のサブクエリ308は、非リレーショナルデータストア150に対してデータデルタ310に対応する取得データの挿入又は書き込みコマンドを発行して、非リレーショナルデータストア150にデータデルタ310のデータを書き込み、記憶し、又は挿入する。このようにして、複製を終了して、更にリレーショナルデータストア155に存在していたが以前に利用できなかったデータ要素を今度は非リレーショナルデータストア150から利用可能にする。図3の破線を参照すると、リレーショナルデータストア155から非リレーショナルデータストア150に識別されたデータデルタ310を複製するために、両データストア(リレーショナルデータストア155及び非リレーショナルデータストア150)に対して実行される第3のサブクエリが描かれている。

40

【0032】

一方のデータストアから他方のデータストアにデータを複製又は同期するための決定は、様々な考慮事項に基づいてもよい。例えば、リレーショナルデータストア155から非

50

リレーショナルデータストア 150 に複製するための決定は、初期の場所からの小さなデータセットをより大きなデータセットを有する場所に複製するための決定又はポリシーに基づいてもよい。例えば、要求されたデータの一部である 1 つ以上のデータ要素は、リレーショナルデータストア 155 から非リレーショナルデータストア 150 の複製を行うためにネットワーク帯域幅の観点からその逆よりも効果的であってもよい。

【0033】

一部の実施形態では、その反対も同様に当てはまり、データの複製は、非リレーショナルデータストア 150 からリレーショナルデータストア 155 へと反対方向に行ってもよい。このような決定は、例えば、最適化エージェント 245 によって実行され又は行われてもよい。非リレーショナルデータベースの実装（例えば、150）を用いる特定の実施形態では、リレーショナルデータベース型のオブジェクトがリレーショナルデータストア 155（例えば、OracleTM）に記憶され、結合動作 305 を特定する 1 つのサブクエリを介して非リレーショナルデータストア 150 に複製される。その後、複製によって非リレーショナルデータストア 150 から全ての必要なデータが利用可能になるので、データ取得動作を特定する別のサブクエリが非リレーショナルデータストア 150 から全ての必要なデータを引き出してもよい。それにもかかわらず、このような例では、データの少なくとも一部は持続的に記憶されており、最初はリレーショナルデータストア 155（例えば、OracleTM）からのみ利用可能である。

【0034】

他の複製の決定及び考慮事項は、最適化エージェント 245 によって同様に考慮及び実装されてもよい。例えば、1 つの複製ポリシーは、複製されたデータが常に同期されているか同期されることを保証されているかどうか、又は一部の逸脱が許容可能なリスクであるかどうか等、複製動作が提供する一貫性保証に基づいてもよい。

【0035】

非リレーショナルデータストア 150 から 10 億個のチャイルド行がクエリを行われ得るように、リレーショナルデータストア 155 から非リレーショナルデータストア 150 への 1 千万 “Account” テーブルの複製を要求する複製動作を例にする。このような例では、非リレーショナルデータベースのクエリエンジンが、OracleTM（例えば、リレーショナル）データを取得するために JDBC（Java（登録商標）Database Connectivity）を介して大量のコールアウトを行うために利用されてもよい。このような例では、OracleTM RAC（Oracle Real Application Cluster）が利用されてもよい。又は、一貫性保証が重要な考慮事項であり、非常に大きなデータの複製が開始されている OracleTM 11g データガードに基づく保証機構が必要な一貫性を提供するために利用されてもよい。データの着信要求に対して応答性のオンザフライでクエリ時にこのようにすることは着信要求への応答又は充足に際し許容できない長い遅延を必要とする場合があるので、このような大きな複製は最も適切には前もって行われてもよい。

【0036】

逆に、小テーブルと対応する少ないデータ転送を例にする。このような例では、現在のデータを有するリレーショナルデータストア 155（例えば、OracleTM）内に存在する 1 つ又は複数の小テーブルの全内容が、例えば、そのデータの少なくとも一部が小テーブル内に存在し且つリレーショナルデータストア 155（例えば、OracleTM）によって持続的に記憶されるデータに対する着信要求に応答して、クエリ時に非リレーショナルデータストア 150 に複製されてもよい。このようなポリシーは、例えば、最適化エージェント 245 によって決定されるこの種のクエリ及び複製の固定コストが重大ではない大規模解析に適切であってもよい。

【0037】

別の考慮事項は、例えば、特定の組織ごとに関連付けられるテーブル又はオブジェクトのサイズに依存する OrgID ベースによる OrgID に関してもよい。例えば、中規模から大規模の組織に対応するデータ（例えば、既定のサイズ閾値に基づく）が前もって複

10

20

30

40

50

製されるポリシーが採用されてもよい。一実施形態では、事前複製は、非リレーショナルデータストア 150 内の事前複製テーブル（例えば、事前複製が選択されている）等、リレーショナルデータストア 155 内の特定のオブジェクトへの変更を補足して別のテーブルにこうした変更をプッシュするスキニーテーブル複製コード（skinny table replication code）を利用してもよい。

【0038】

所定の実施形態では、特定の組織が大量の変更を引き起こす場合があるので、非リレーショナルデータストア 150 で要求された解析が実行されることを許可しつつも、リアルタイム同期は必ずしも変更ごとに要求されず又は適切ではない。従って、このような実施形態では、特定の間隔で更新を同期するオプションが（例えば、以下で検討される最適化エージェント 245 及びハードウェアベースのクエリレイヤエージェント 501 及び 734 を介して）提供される。このような実施形態では、特定の間隔での更新を可能にするポリシーは、結果未決の参照及び許容可能な逸脱であり得る又は採用された複製ポリシー若しくはデータベースクエリの基本目的に依存して後続のデータチェック及び検証を要求し得る他の「乱雑なデータ（sloppy data）」にもかかわらず、より効果的な書き込み及び更新を提供する。

【0039】

複製に関する別の考慮事項は、テーブルの濃度等のデータストアの統計であってもよい。統計は、最適化エージェント 245 により生成され、利用可能であり、又は集められてもよい。例えば、比較的小規模な行の組が必要とされる所定の実施形態では、処理されるデータベースクエリ全体が大規模であり且つクエリが行われる行の組全体を送信することが合理的な送信コストを必要とすると決定される場合（例えば、既定のサイズ比又は固定閾値等に基づいて）、サブクエリ（例えば、306 - 309）がリレーショナルデータストアから 155 要求される行の組全体にクエリを行って、非リレーショナルデータストア 150 にクエリが行われた行の組全体を送信してもよい。このようなポリシーは、データの実際の要求を受信することなく事前複製又は前もって複製を行う必要性を回避してもよく、上記のデータの非一貫性の問題を更に回避してもよい。

【0040】

データをどこに持続的に記憶すべきか、そしてデータをどこから取得すべきかの考慮事項は、特定のデータストアの基本的な実装ハードウェアに更に基づいてもよい。例えば、非リレーショナルデータストア 150 は、例えば、ギガバイトあたりのコストを単位として費用の安い記憶装置に大規模なフラットファイル及びバイナリファイルを記憶するように最適化されてもよい。非リレーショナルデータストア（例えば、150）が圧縮フラットファイル及びバイナリファイルの大量の読み込みに最適化されており、OracleTM等のリレーショナルモデルのデータストア（例えば、155）とは対照的にギガバイトあたりの要求を処理するのに計算コストが余りかからないので、このような基本のハードウェア実装が可能である。OracleTM等のリレーショナルデータストア 155 は、例えば、全ての記憶されたデータに必須の解析、トランザクション処理及びロールバック機能等の企業レベルのデータ保護を実装するためのリレーショナルデータストア 155 の要件によりギガバイトあたりではより高価な記憶ハードウェアを必要としてもよい。データベーストランザクションが最終的に完了する前に失敗又は中断する非一貫性の状態において、実行可能にデータストアがトランザクション処理を欠いたままであり得る直接挿入モデルとは対照的に、このような企業レベルのデータ保護は、トランザクション処理を介して、故障時に特定のトランザクションの「やり直し（redo）」を行う能力を更に可能にする。

【0041】

従って、所定の実施形態では、記憶されたデータへの最近の編集が、解析、トランザクション処理及びロールバックを実装するリレーショナルデータストアによって処理及び記憶される。また、リレーショナルデータストア 155 に書き込まれた更新の少なくとも一部は、その後、対応するデータを持続的に記憶するギガバイトあたりのコストを減らすた

10

20

30

40

50

めに、非リレーショナルデータストア 150 に複製、移転、処理又は移動される。このようにして、ホストシステム 110 は、所定のリレーショナルデータストア 155 に関連付けられる企業レベルデータ保護を利用すると同時に一部の非リレーショナルデータストア 150 を介して利用可能なあまり費用のかからない持続的記憶装置の利益を享受してもよい。

【0042】

一実施形態では、最適化エージェント 245 が非リレーショナルデータストア 150 及びリレーショナルデータストア 155 の双方から利用可能なデータの「ビュー」を有するので、最適化エージェント 245 は、データ 218 に対する着信要求 115 に対応する所定のデータ要素が 1 つより多くのソースから 1 つより多くの手法で取得され得る場合に、

「選択的クエリ」をもたらすことができる。最適化エージェント 245 は、このような要求 115 を充足するために発行される複数のサブクエリの改善されたシーケンス又は順序付けを更にもたらしすることができる。

【0043】

従って、所定の実施形態に従って、様々な利用可能な考慮事項を考慮すると、第 4 のサブクエリ 309 が結合動作 305 内に更に含まれ、非リレーショナルデータストア 150 に対して実行される。ここで、第 4 のサブクエリ 309 は、非リレーショナルデータストア 150 内に存在する 1 つ以上のデータ要素とリレーショナルデータストア 155 から非リレーショナルデータストア 150 に複製されることで非リレーショナルデータストア 150 内から利用可能な 1 つ以上の他のデータ要素との双方をフェッチすることにより、非

リレーショナルデータストア 150 から取得されるデータをフェッチする。このようにして、複数のデータ要素 315 は、結合動作 305 によってトリガされたデータ複製の前に複数のデータ要素 315 の一部が非リレーショナルデータストア 150 から最初に利用可能でないにもかかわらず、非リレーショナルデータストア 150 等のデータストアの 1 つから完全に取得されてもよい。

【0044】

一の実施形態では、要求されたデータの 1 つ以上のデータ要素が持続的に記憶される複数のデータストア間でデータを複製する代わりに、データを持続的に記憶するクエリを受ける 2 つ以上のデータストア（例えば、150 及び 155）の各々から離れた場所に全てのデータ要素を取得するようにポリシーが用いられてもよい。例えば、データは、クエリレイヤ 260 へのメモリ内結合動作を利用して取得され、又はグローバルキャッシュレイヤ（例えば、図 5 の要素 550）へのメモリ内結合動作を介して取得されてもよい。このようなメモリ内結合動作は、最適化エージェント 245 から利用可能な既知の統計に基づいて、又は特定されたサイズ閾値（例えば、行の数、サイズ（例えば、データのメガバイト数）を単位とするデータ量、要求されたデータの濃度等）に基づいて選択されてもよい。

【0045】

例えば、ホストシステム 110 内の既知の統計及び解析に基づいて利用可能な他の考慮事項は、例えば、要素の最大数が知られており、最大又は推定クエリコストが決定可能若しくは知られていて最適化エージェント 245 から利用可能である特定のクエリに対する既知の選択リスト量から導かれる特定のデータベースクエリ 217 又はサブクエリ（例えば、306 - 309）に対するクエリコストを含んでもよい。それは更に、既に行われた解析に基づいて知られ、又は複数の利用可能なデータストア（例えば、150 及び 155）の中で、最少量の時間で又は最小の計算リソースを利用 / 消費して最少数を有する結果をもたらすことができる最適化エージェント 245 を介して（例えば、1 つ以上の事前クエリを介して）決定可能であってもよい。例えば、大きなデータベースクエリ 217 では、各データストア（例えば、150 及び 155）から必要なデータの小部分に対して事前クエリを行って、どの事前クエリがより効果的な結果をもたらすかを決定し、このような決定に基づいて、より効果的なデータストア（例えば、事前クエリの結果によって（150 又は 155）を目標にする一次データベースクエリ 217 を充足するために必要とされ

る様々なサブクエリ(306 - 309)を生成することが望ましくてもよい。最適化エージェント245によって適切な解析が前もって行われる場合、事前クエリを発行しなくても、クエリポリシーが簡単に要求されてもよい。例えば、このような解析決定が行われて、顧客スキーマ240を介して1つ以上のデータの場所に対して記憶及び特定されてもよい。

【0046】

代替的な実施形態では、異なる又は追加の結合動作305が行われてもよい。例えば、一実施形態では、データベースクエリ217を介してマルチテナントデータベースシステム130に対して実行される結合動作305は、リレーショナルデータストア155から2つ以上の関連テーブルを特定する結合動作305、リレーショナルデータストア155及び非リレーショナルデータストア150内に存在する少なくとも1つ以上のデータ構造から少なくとも1つの関連テーブルを特定する結合動作305、及び非リレーショナルデータストア150内に存在する2つ以上の離れた及び別個のデータ構造を特定する結合動作305からなる結合動作305のグループから選択される結合動作305を含んでもよい。ここで、2つ以上の離れた及び別個のデータ構造の各々は、非リレーショナルデータストア150内で2つの別個のデータ構造を関連付け又は関係付けるための共有された特徴、ストリング、バイナリ又は英数字キー等の重複した共有キーを欠いている。

【0047】

例えば、一実施形態では、非リレーショナルデータストア150は、多くのデータ構造、ファイル、オブジェクト及び他のこうした情報を記憶する能力を提供するが、このようなデータ構造、ファイル、オブジェクト及び他の情報を「関連付け」する機能を実装しない。しかしながら、2つの別個のデータ構造の各々を特定する結合動作305は、例えば、所望されるが以前に非関連付けされた全ての情報を有する単一データ構造を形成し、又は代替の場所でこのような結合動作によって特定される所望の情報を取得及び一時的にキャッシュするために、非リレーショナルデータストア150に対して適切に形成された結合動作305を実行することができるマルチテナントデータベース130内の複数の基本ハードウェア、ソフトウェア、及び論理要素120等、実装された非リレーショナルデータストア150の外部の機能及び論理に依存することにより、各々のこのようなデータ構造を識別及びリンク又は関係付けることができる。

【0048】

代替的な実施形態では、特定された結合動作305は、非リレーショナルデータストア150に対して実行され且つ非リレーショナルデータストア内に存在する1つ以上のデータ要素を取得するための第1のサブクエリ(例えば、306)、リレーショナルデータストア155に対して実行され且つ非リレーショナルデータストア150内に存在する1つ以上のデータ要素とリレーショナルデータストア155内に存在する1つ以上の他のデータ要素との間のデータデルタ310を決定する第2のサブクエリ(例えば、307)、及びリレーショナルデータストア155に対して実行され且つ決定されたデータデルタ310に基づいてリレーショナルデータストア155内に存在する1つ以上の他のデータ要素を取得するための第3のサブクエリ(例えば、308)を含む。

【0049】

このような実施形態では、決定されたデータデルタ310に基づいてリレーショナルデータストア155内に存在する1つ以上の他のデータ要素を取得する第3のサブクエリ(例えば、308)は、時間ベースのクエリフィルタ機構又は記録ベースのフィルタ機構の何れかを含んでもよい。

【0050】

例えば、一実施形態では、時間ベースのサブクエリは、リレーショナルデータストア155に対して実行される。ここで、時間ベースのサブクエリは、非リレーショナルデータストア150内に存在する1つ以上のデータ要素内のデータ要素に対応するどのタイムスタンプよりも遅いタイムスタンプを有するリレーショナルデータストア155内のデータ要素に基づいてリレーショナルデータストア155から取得される1つ以上の他のデータ

要素を特定する。

【 0 0 5 1 】

代替的な実施形態では、記録識別子ベースのサブクエリは、リレーショナルデータストア 1 5 5 に対して実行される。ここで、記録ベースのサブクエリは、非リレーショナルデータストア 1 5 0 内に存在する 1 つ以上のデータ要素内のデータ要素に対応するどの記録識別子よりも数値が大きい記録識別子を有するリレーショナルデータストア 1 5 5 内のデータ要素に基づいてリレーショナルデータストア 1 5 5 から取得される 1 つ以上の他のデータ要素を特定する。

【 0 0 5 2 】

図 4 は、実施形態が動作可能な代替の例示的なアーキテクチャ 4 0 0 を示す。特に、所定の実施形態に従うマルチテナントデータベースシステムによって受信される新しいトランザクションの処理が更に詳細に描かれている。

【 0 0 5 3 】

例えば、所定の実施形態では、持続的な記憶のためにマルチテナントデータベースシステム 1 3 0 に書き込み又は挿入される新しい情報は、長期間にわたって非リレーショナルデータストア 1 5 0 に現在時に記憶されるように指定されてもよいが、それにもかかわらず一時的にリレーショナルデータストア 1 5 5 に書き込まれてもよい。例えば、一時的に一方のデータストア（リレーショナルデータストア 1 5 5 等）にデータを書き込んで、後に別のデータストア（非リレーショナルデータストア 1 5 0 等）にデータを移行するための考慮事項は、例えば、一方のデータストア対他方の書き込み応答時間の改善、更に代替のデータストアからの取得時間の改善を含んでもよい。一方のデータストアは、より低い計算又はより低い動作コストと関連してもよい。非リレーショナルデータストア 1 5 0 等の特定のデータストアが、滅多に更新されないが頻繁に取得されるデータでより効果的に動作してもよい。代替的に、リレーショナルデータストア 1 5 5 等の他のデータストアが、著しく断片化されたデータ又は滅多に更新されないデータを有する前の例示と比較して非常に頻繁に更新され又は追加されるデータでより良好な動作効率を示してもよい。

【 0 0 5 4 】

従って、所定の実施形態によれば、（例えば、先に描かれた要求 1 1 5 内の）マルチテナントデータベースシステム 1 3 0 において受信される新しいトランザクションは、非リレーショナルデータストア 1 5 0 に書き込まれる新しいデータ 4 1 6 を含み又は特定する。一部の実施形態では、新しいデータ 4 1 6 は、新しいデータ 4 1 6 が非リレーショナルデータストア 1 5 0 に書き込まれるという指示にかかわらず、リレーショナルデータストア 1 5 5 の追加ログ 4 1 0 に書き込まれる。新しいデータ 4 1 6 が書き込まれるというこのような指示は、新しいトランザクション 4 1 5 によって、例えば、新しいトランザクション 4 1 5 内のターゲット 4 1 9 の属性によって特定されてもよい。代替的に、例えば、最適化エージェント 2 4 5 によって決定される新しいデータ 4 1 6 の特徴に基づいて、又は新しいトランザクション 4 1 5 に対応する O r g I D と関連付けられるフラグ若しくは記憶された設定に基づいて、ホストシステム 1 1 0 によって決定が行われてもよい。

【 0 0 5 5 】

一部の実施形態では、決定されたデータデルタ（例えば、3 1 0）に基づいてリレーショナルデータストア 1 5 5 内に存在する 1 つ以上の他のデータ要素を取得するためのサブクエリを含む結合動作（例えば、3 0 5）は、リレーショナルデータストア 1 5 5 の追加ログ 4 1 0 からリレーショナルデータストア 1 5 5 内に存在する 1 つ以上の他のデータ要素を取得するためのサブクエリを含む。例えば、追加ログに書き込まれた新しいデータ 4 1 6 が取得されてもよく、又は追加ログに記憶された新しいデータ 4 1 6 の要素が取得されてもよい。

【 0 0 5 6 】

一実施形態では、ホストシステム 1 1 0 は、追加ログがフラッシュ閾値に到達すると、追加ログ 4 1 0 のフラッシュをトリガして、リレーショナルデータストア 1 5 5 の追加ログ 4 1 0 に書き込まれた新しいデータ 4 1 6 を非リレーショナルデータストア 1 5 0 にフ

10

20

30

40

50

ラッシュする。その結果、例えば、新しいデータはフラッシュされたデータ 4 1 7 として非リレーショナルデータストアに存在し、リレーショナルデータストア 1 5 5 の追加ログ 4 1 0 に以前に存在した新しいデータ 4 1 6 に対応することになる。

【 0 0 5 7 】

異なる種類のデータがマルチテナントデータベースシステム 1 3 0 によって記憶されてもよい。例えば、一実施形態では、非リレーショナルデータストア 1 5 0 内に存在する 1 つ以上のデータ要素は、複数の圧縮フラットファイル又は複数のバイナリファイル又は圧縮フラットファイル及びバイナリファイルの組み合わせに対応する。このようなファイルは、非リレーショナルデータベースアーキテクチャ（例えば、1 5 0）を介してより効果的に記憶されてもよい。

10

【 0 0 5 8 】

別の実施形態では、リレーショナルデータストア 1 5 5 は、リレーショナルデータベース管理システム（R D B M S）に従うリレーショナルデータベースを実装する。ここで、リレーショナルデータベースの複数の関連テーブルは、リレーショナルデータベース内の 2 つ以上の関連テーブルの各々に対して 1 つ以上の重複する共通の特徴を介して相互に関連付けられ、これによりリレーショナル型データストア 1 5 5 と共通して関連付けられる「リレーションシップ」を形成する。

【 0 0 5 9 】

一実施形態では、非リレーショナルデータストア 1 5 0 は、各々が非リレーショナルデータストア 1 5 0 に総記憶容量の少なくとも一部を提供する複数の基本ハードウェア記憶装置を有する分散構造型データベースを実装する。このような実施形態では、非リレーショナルデータストア 1 5 0 内のデータ要素は、主キーに基づいて参照可能であるが、リレーショナルデータストア 1 5 5 内のデータ要素の場合におけるように、2 つ以上の関連テーブルの間の 1 つ以上の重複する共通の特徴に基づいて参照可能ではない。

20

【 0 0 6 0 】

一実施形態では、リレーショナルデータストア 1 5 5 は、O r a c l e 互換データベースの実装、I B M D B 2 E n t e r p r i s e S e r v e r 互換リレーショナルデータベースの実装、M y S Q L 互換リレーショナルデータベースの実装、及び M i c r o s o f t S Q L S e r v e r 互換リレーショナルデータベースの実装の中から選択されるリレーショナルデータベースモデルを実装する。

30

【 0 0 6 1 】

一実施形態では、非リレーショナルデータストア 1 5 0 は、V a m p i r e 互換非リレーショナルデータベースの実装、A p a c h e C a s s a n d r a 互換非リレーショナルデータベースの実装、B i g T a b l e 互換非リレーショナルデータベースの実装、及び H B a s e 互換非リレーショナルデータベースの実装の中から選択される N o S Q L 非リレーショナルデータベースを実装する。

【 0 0 6 2 】

一実施形態では、非リレーショナルデータストア 1 5 0 は、各コンピュータノードが少なくとも 1 つのメモリ、1 つ以上のプロセッサ及び 1 つ以上の通信可能にインターフェース接続されたハードディスクドライブを含む複数の分散型コンピュータノードを含む。このような実施形態では、分散型コンピュータノードの各々は、中央トランザクション局からの承認又は制御無しで非リレーショナルデータベーストランザクションを読み込み、書き込み、及び更新する機能を有する分離した非リレーショナルデータベースインスタンスを更に含んでもよい。

40

【 0 0 6 3 】

特定の実施形態では、リレーショナルデータストア 1 5 5 は、モノリシックリレーショナルデータベースインスタンス（m o n o l i t h i c r e l a t i o n a l d a t a b a s e i n s t a n c e）への更新又は変更がモノリシックリレーショナルデータベースインスタンスに通信可能にインターフェース接続され且つ制御される持続的記憶装置に対する持続的記憶にコミットされるかどうかを制御する中央トランザクション局と計

50

算リソースを調整するプロセッサ及びメモリを含むモノリシックリレーショナルデータベースインスタンスを実装する。

【 0 0 6 4 】

図 5 は、実施形態が動作し、設置され、統合され、又は構成され得るシステム 5 0 0 の図表示を示す。

【 0 0 6 5 】

一実施形態では、システム 5 0 0 は、メモリ 5 9 5 及び 1 つ以上のプロセッサ 5 9 0 を含む。例えば、メモリ 5 9 5 は実行される命令を記憶してもよく、プロセッサ 5 9 0 はこのような命令を実行してもよい。システム 5 0 0 は、バス 5 1 5 と通信可能にインターフェース接続される複数の周辺機器の間でシステム 5 0 0 内のトランザクション及びデータを転送するためのバス 5 1 5 を含む。システム 5 0 0 は、例えば、要求を受信し、応答を返し、顧客組織 1 0 5 A C 内に設置されたクライアント装置等のリモートクライアントとインターフェース接続するウェブサーバ 5 2 5 を更に含む。

【 0 0 6 6 】

システム 5 0 0 は、データベースクエリ及びデータベースサブクエリを最適化して、基本データストアにクエリを行う最適な又は好ましい手法を決定するために選択的に事前クエリを調整するように設計される最適化エージェント 5 3 5 を有するように更に描かれる。システム 5 0 0 は、通信可能にインターフェース接続された装置及びシステムにキャッシュサービスを提供し、特に、顧客スキーマデータ（例えば、メタデータ等）をキャッシュをするグローバルキャッシュレイヤ 5 5 0 を更に含む。顧客スキーマデータは、グローバルキャッシュレイヤ 5 5 0 と連動して動作可能な顧客スキーマ 5 3 0 によって提供され、例えば、必須データ要素がマルチテナントデータベースシステム内のリレーショナルデータベース若しくは非リレーショナルデータベースの実装又はその両方によって記憶されるかどうかを特定し、対応する要求に対するデータセットを構成する 1 つ以上のデータ要素に関して基本データストア内の場所を特定する。顧客スキーマ 5 3 0 は、システム 5 0 0 内のハードドライブ、持続性データストア又は他の記憶場所に記憶されてもよい。

【 0 0 6 7 】

ハードウェアベースのクエリレイヤエージェント 5 0 1 はシステム 5 0 0 内で別個にあり、要求プロセッサ 5 7 0、顧客スキーマプロセッサ 5 7 5、サブクエリ生成器 5 8 0 及びクエリ実行器 5 8 5 を含む。一実施形態によれば、要求プロセッサ 5 7 0 は、（例えば、ウェブサーバ 5 2 5 から、上述のようにホストシステム 1 1 0 から、又は接続されたネットワークインターフェースから）取得するデータを特定する要求を受信する。要求プロセッサ 5 7 0 は、顧客スキーマプロセッサ 5 7 5 と連動して、基本データストアから取得される要求データの 1 つ以上の位置を取得する。要求プロセッサ 5 7 0 は、更にサブクエリプロセッサ 5 8 0 と連動して、決定された 1 つ以上のこのようなデータの位置に基づいて適切な基本データストアから要求された 1 つ以上のデータ要素を取得するために必要なサブクエリを構築及び生成し、又は一方のデータストアから他方にデータサブセットを同期、フラッシュ又は複製させる結合動作を開始するために必要なサブクエリを生成し、その結果、後続のサブクエリは 1 つのデータストアから要求されたデータセットを取得することができる。サブクエリ生成器 5 8 0 によって生成されるこのようなサブクエリは、最適化エージェント 5 3 5 から利用可能な統計及び事前クエリの結果に依存してもよい。クエリ実行器 5 8 5 は、通信可能にインターフェース接続されたデータベースの実装に対して生成されたクエリ及びサブクエリを実行する。

【 0 0 6 8 】

図 6 は、一実施形態に従ってマルチテナントストアで横断的ストア結合を実行する方法 6 0 0 を示す流れ図であり、所定の実施形態に従ってデータベースクエリ（例えば、2 1 7）によって特定される結合動作を特定及び実行することを含む。方法 6 0 0 は、ハードウェア（例えば、回路、専用論理、プログラム可能論理、マイクロコード等）、ソフトウェア（例えば、読み込み、書き込み、更新、最適化、事前クエリの開始、サブクエリの開始等、又はこれらの組み合わせ等の様々なクエリ動作を行うために処理装置上で実行され

る命令)を含み得る論理を処理することによって実行されてもよい。一実施形態では、方法600は、図5の要素501で描かれるハードウェアベースのクエリレイヤ等のハードウェア論理によって実行される。所定の実施形態によれば以下にリストされるブロック及び/又は動作の一部は選択的である。提示されたブロックの番号は明確のためであって、様々なブロックが行われなければならない動作の順番を規定することを意図していない。

【0069】

方法600は、ホストシステムでマルチテナントデータベースシステムに対する要求を受信する処理論理で開始する。この要求はマルチテナントデータベースシステムから取得されるデータを特定する(ブロック605)。ブロック610では、ホストシステムを介する要求に基づいて、処理論理は取得されるデータの1つ以上の位置を取得する。

10

【0070】

ブロック615では、ホストシステムを介して、処理論理は要求の受信に応答して顧客スキーマを取得する。例えば、顧客スキーマが取得されるデータの1つ以上の位置216を記述してもよい。ここで、顧客スキーマは、非リレーショナルデータストア内若しくはリレーショナルデータストア内に存在する又は非リレーショナルデータストア及びリレーショナルデータストアの双方から利用可能であるデータの複数のデータ要素の各々を特定する。

【0071】

ブロック620では、ホストシステムを介して、処理論理は要求に基づいてデータベースクエリを生成する。例えば、データベースクエリは、取得される複数のデータ要素であって、非リレーショナルデータストア内に存在する1つ以上のデータ要素及びリレーショナルデータストア内に存在する1つ以上の他のデータ要素を含む複数のデータ要素を特定してもよい。データベースクエリは複数のサブクエリを更に含んでもよい。一実施形態では、データベースクエリはサブクエリの1つを介して結合動作を特定する。同様に、パーティション、フラッシュ、同期又は複製動作がサブクエリを介して特定されてもよい。

20

【0072】

ブロック625において、処理論理はデータを取得するためにマルチテナントデータベースシステムに対してデータベースクエリを実行する。

【0073】

ブロック630では、処理論理はマルチテナントデータベースシステムで新しいトランザクションを取得する。新しいトランザクションの各々は非リレーショナルデータストアに書き込まれる新しいデータを特定する。ブロック635において、処理論理はリレーショナルデータストアの追加ログに新しいデータを書き込む。例えば、一実施形態では、データベースクエリのサブクエリは、取得されるデータがリレーショナルデータストアの追加ログから取得されることを特定する。

30

【0074】

ブロック640では、追加ログがフラッシュ閾値に到達すると、処理論理はリレーショナルデータストアの追加ログに書き込まれた新しいデータを非リレーショナルデータストアにフラッシュする。

【0075】

40

図7は、コンピュータシステムの例示的形態のマシン700の図形表現を示す。一実施形態によれば、このコンピュータシステムの中で、マシン700に本明細書で検討された1つ以上の方法を実行させるために、一組の命令が実行されてもよい。代替的な実施形態では、マシンは、LAN(Local Area Network)、イントラネット、エクストラネット、又はインターネットで他のマシンと接続(例えば、ネットワーク化)されてもよい。マシンは、クライアントサーバネットワーク環境におけるサーバ又はクライアントマシンの能力で、又はピアツーピア(又は分散ネットワーク環境)におけるピアマシンとして又はマルチテナントデータベース記憶サービスを提供するオンデマンド環境を含むオンデマンドサービス環境内のサーバ若しくは一連のサーバとして動作してもよい。マシンの所定の実施形態は、パーソナルコンピュータ(PC)、タブレットPC、セッ

50

トトップボックス (STB)、パーソナルデジタルアシスタント (PDA)、携帯電話、ウェブ装置、サーバ、ネットワークルータ、スイッチ若しくはブリッジ、コンピュータシステム、又はそのマシンによって取り行われる動作を特定する一組の命令 (連続的又はそれ以外) を実行することが可能な任意のマシンの形態であってもよい。更に、1つのマシンのみが示されているが、「マシン」という用語は、本明細書で検討された1つ以上の方法を実行する一組 (又は複数の組) の命令を個別の又は共同で実行する複数のマシン (例えば、コンピュータ) の集合を含むようにも解釈される。

【0076】

例示的なコンピュータシステム700は、プロセッサ702、主メモリ704 (例えば、ROM (read-only memory)、フラッシュメモリ、SDRAM (synchronous DRAM) 若しくはRDRAM (Rambus DRAM) 等のDRAM (dynamic random access memory)、フラッシュメモリ等のスタティックメモリ、SRAM (static random access memory)、揮発性であるが高データ速度のRAM等)、及び二次メモリ718 (例えば、ハードディスクドライブを含む持続性記憶装置及び持続性マルチテナントデータベースの実装) を含み、これらはバス730を介して互いに通信する。主メモリ704は、(例えば、リレーショナルデータストア及び非リレーショナルデータストアの双方に散在し、ハードウェアベースのクエリレイヤエージェント734を介して取得可能なデータ要素の場所等、2つ以上の異なったデータストアの中で特定のデータ又はデータセットを構成するデータ又はデータ要素の1つ以上の位置を特定する) 顧客スキーマ724を含む。主メモリ704は、顧客スキーマ724を介して提供される情報の種類等の大規模データセットの複数のデータ要素間のメタデータ及び他の関連又は対応情報を提供するために、システム全体でアクセス可能なグローバルキャッシュレイヤ等のグローバルキャッシュレイヤ723を更に含む。主メモリ704及びそのサブ要素 (例えば、723及び724) は、本明細書で検討された方法を実行するために処理論理726及び702と連動して動作可能である。

【0077】

プロセッサ702は、マイクロプロセッサ又は中央処理装置等の1つ以上の汎用処理装置を表す。より詳細には、プロセッサ702は、CISC (complex instruction set computing) マイクロプロセッサ、RISC (reduced instruction set computing) マイクロプロセッサ、VLIW (very long instruction word) マイクロプロセッサ、他の命令セットを実装するプロセッサ、又は命令セットの組み合わせを実装するプロセッサであってもよい。また、プロセッサ702は、ASIC (application specific integrated circuit)、FPGA (field programmable gate array)、DSP (digital signal processor)、又はネットワークプロセッサ等の1つ以上の専用処理装置であってもよい。プロセッサ702は、動作及び本明細書で検討された機能を実行するための処理論理726を実行するように構成される。

【0078】

コンピュータシステム700は、ネットワークインターフェースカード708を更に含んでもよい。また、コンピュータシステム700は、ユーザインターフェース710 (ビデオディスプレイ装置、LCD (liquid crystal display) 若しくはCRT (cathode ray tube) 等)、英数字入力装置712 (例えば、キーボード)、カーソル制御装置714 (例えば、マウス) 及び信号生成装置716 (例えば、統合スピーカ) を含んでもよい。コンピュータシステム700は、周辺装置736 (例えば、無線若しくは有線通信装置、メモリ装置、記憶装置、オーディオ処理装置、ビデオ処理装置等) を更に含んでもよい。コンピュータシステム700は、データベースクエリ及びサブクエリを管理し、マルチテナントデータベースシステム等の基本データストアとトランザクションを調整するハードウェアベースのクエリレイヤエージェント73

10

20

30

40

50

4 を更に含んでもよい。

【 0 0 7 9 】

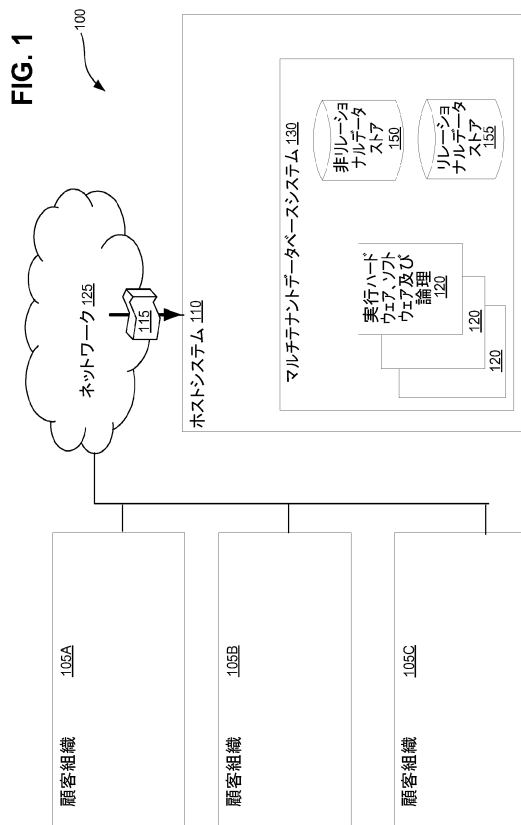
二次メモリ 7 1 8 は、本明細書で検討された 1 つ以上の方法又は機能を実現化する 1 つ以上の命令の組（例えば、ソフトウェア 7 2 2）が記憶される一時的でないマシン可読記憶媒体（又は、より詳細には、一時的でないマシンアクセス可能記憶媒体）7 3 1 を含んでもよい。また、ソフトウェア 7 2 2 が、コンピュータシステム 7 0 0 によってその実行中に主メモリ 7 0 4 内に及び / 又はプロセッサ 7 0 2 内に完全に又は少なくとも部分的に存在してもよい。ここで、主メモリ 7 0 4 及びプロセッサ 7 0 2 はマシン可読記憶媒体も構成する。更に、ソフトウェア 7 2 2 は、ネットワークインターフェースカード 7 0 8 を介してネットワーク 7 2 0 上で送信又は受信されてもよい。

【 0 0 8 0 】

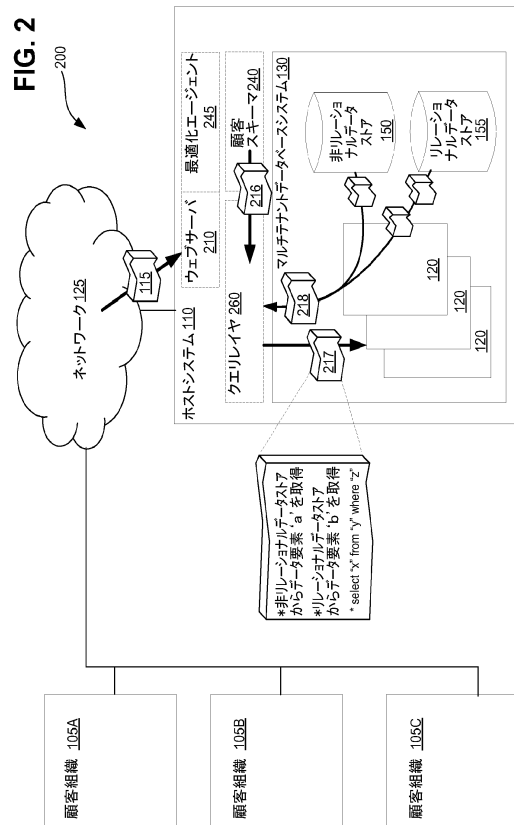
本明細書で検討された主題は例示として特定の実施形態に関して記載されているが、本発明の実施形態は明示的に列挙された開示の実施形態に限定されないことが理解されるべきである。それとは反対に、本開示は当業者に明白な様々な修正及び類似の構成を含むことが意図されている。従って、添付の特許請求の範囲は、全てのこのような修正及び類似の構成を包含するように最も広い解釈に一致するべきである。先の記載は例示的であり限定的ではないことが理解されるべきである。先の記載を読んで理解すると、多くの他の実施形態が当業者には明白であろう。従って、本開示の主題の範囲は、添付の請求項が享受する均等物の全範囲と共に、このような請求項を参照することで決定されるべきである。

10

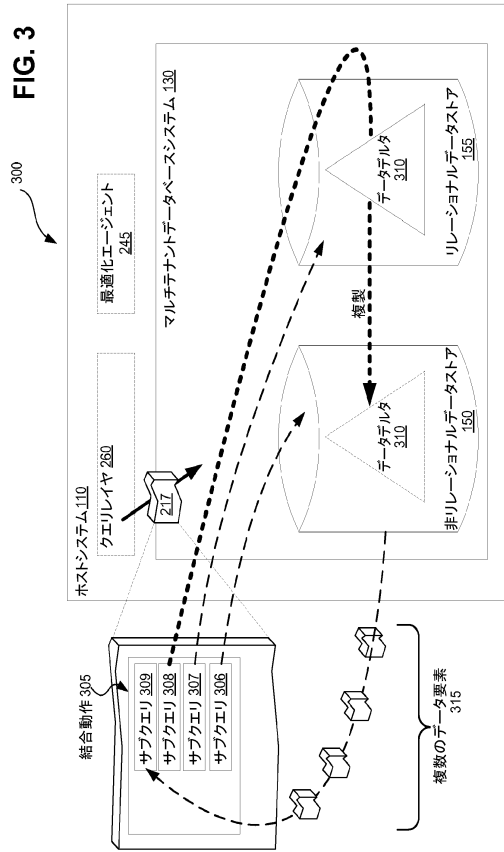
【 図 1 】



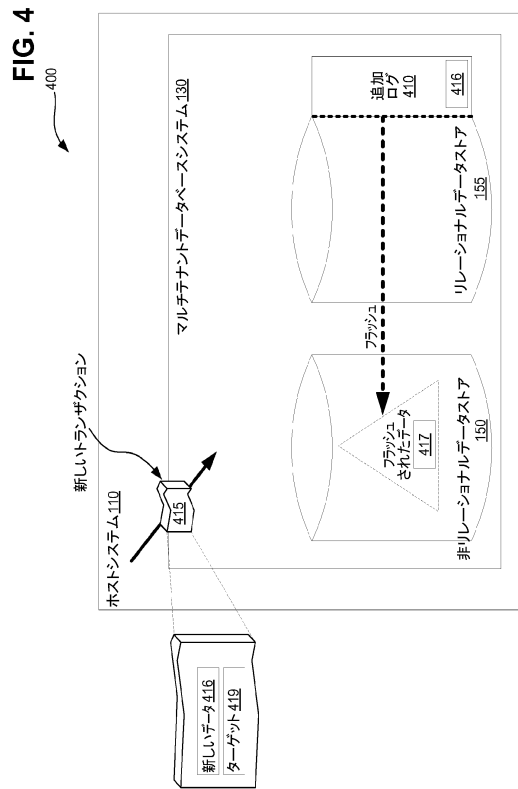
【 図 2 】



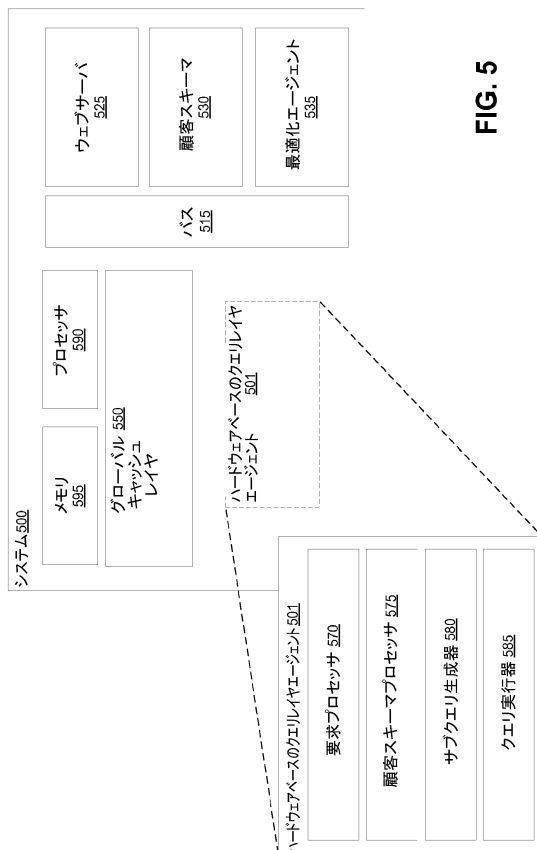
【 図 3 】



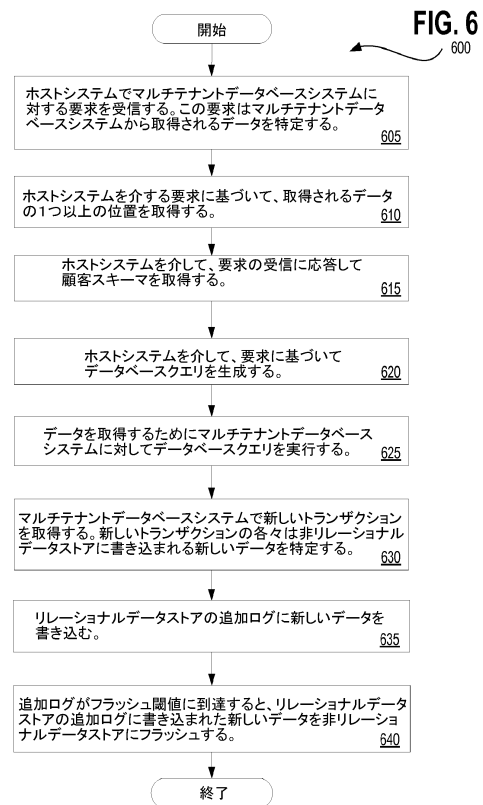
【 図 4 】



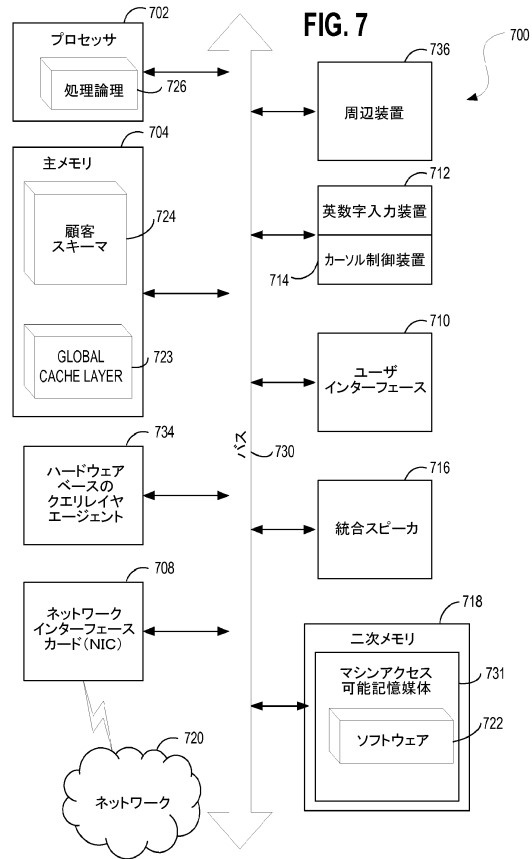
【 図 5 】



【 図 6 】



【図 7】



フロントページの続き

- (72)発明者 イードスン、 ビル シー .
アメリカ合衆国 9 4 3 0 3 カリフォルニア州 パロアルト コリーナ ウェイ 3 8 7 9
- (72)発明者 ワイスマン、 クレイグ
アメリカ合衆国 9 4 1 1 5 カリフォルニア州 サンフランシスコ サクラメント ストリート
2 8 3 8
- (72)発明者 オリバー、 ケヴィン
アメリカ合衆国 9 4 1 2 7 カリフォルニア州 サンフランシスコ 1 5 ス アベニュー 2 5
7 9
- (72)発明者 テイラー、 ジェームズ
アメリカ合衆国 9 4 1 1 4 カリフォルニア州 サンフランシスコ 2 5 ス ストリート 4 3
8 5
- (72)発明者 フェル、 サイモン ズィー .
アメリカ合衆国 9 4 9 2 5 カリフォルニア州 コーテ マデラ フライング クラウド コー
ス 1 4
- (72)発明者 シュナイダー、 ドノヴァン エー .
アメリカ合衆国 9 4 1 2 7 カリフォルニア州 サンフランシスコ アプトス アベニュー 2
5

審査官 原 秀人

- (56)参考文献 特開2001-051879(JP, A)
特開2010-224824(JP, A)
特開2000-222430(JP, A)
特開平09-146804(JP, A)
特開平11-096055(JP, A)
米国特許出願公開第2011/0258630(US, A1)
中田 敦, 最新のクラウド技術を解剖する設計者が明かす基盤の実像, 日経コンピュータ, 日本
, 日経BP社, 2010年 9月 1日, 第764号, p. 68 - 73
萩原 正義, クラウドの設計セオリー, 日経SYSTEMS, 日本, 日経BP社, 2010年1
1月26日, 第212号, p. 99 - 103

(58)調査した分野(Int.Cl., DB名)

G 0 6 F 1 2 / 0 0

G 0 6 F 1 7 / 3 0