

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5525541号
(P5525541)

(45) 発行日 平成26年6月18日 (2014. 6. 18)

(24) 登録日 平成26年4月18日 (2014. 4. 18)

(51) Int. Cl. F I
G06T 1/00 (2006.01) G O 6 T 1/00 B
G06F 12/00 (2006.01) G O 6 F 12/00 5 1 7

請求項の数 98 (全 25 頁)

(21) 出願番号	特願2011-539631 (P2011-539631)	(73) 特許権者	509123208
(86) (22) 出願日	平成21年12月1日 (2009. 12. 1)		アビニシオ テクノロジー エルエルシー
(65) 公表番号	特表2012-510687 (P2012-510687A)		アメリカ合衆国 02421 マサチュー
(43) 公表日	平成24年5月10日 (2012. 5. 10)		セッツ州 レキシントン スプリング ス
(86) 国際出願番号	PCT/US2009/066210		トリート 201
(87) 国際公開番号	W02010/065511	(74) 代理人	100079108
(87) 国際公開日	平成22年6月10日 (2010. 6. 10)		弁理士 稲葉 良幸
審査請求日	平成24年11月30日 (2012. 11. 30)	(74) 代理人	100109346
(31) 優先権主張番号	61/119, 164		弁理士 大貫 敏史
(32) 優先日	平成20年12月2日 (2008. 12. 2)	(72) 発明者	ウェイクリング ティム
(33) 優先権主張国	米国 (US)		アメリカ合衆国 01810 マサチュー
			セッツ州 アンドーヴァー アボット ス
			トリート 11

最終頁に続く

(54) 【発明の名称】 データ管理システム内のデータセットのインスタンスのマッピング

(57) 【特許請求の範囲】

【請求項1】

少なくとも一つの入力データセットからデータを受取り、且つ、少なくとも一つの出力データセットにデータを出力するデータ処理プログラムの仕様を処理するステップと、

第1のデータセット候補を特定するステップであって、前記第1のデータセット候補は前記データ処理プログラムのうちの第1のデータ処理プログラムに表現され、前記第1のデータセット候補は第1のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、ステップと、

第2のデータセット候補を特定するステップであって、前記第2のデータセット候補は前記データ処理プログラムのうちの第2のデータ処理プログラムに表現され、前記第2のデータセット候補は前記第1のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、ステップと、

前記第1のデータセット候補と前記第2のデータセット候補との間のマッピングを受け取るユーザインターフェースを提供するステップであって、前記マッピングは、前記第1のデータセット候補と前記第2のデータセット候補の双方が同じデータセットを表現することを示す、ステップと、

前記マッピングの前記第1のデータセット候補もしくは前記第2のデータセット候補に対してデータを提供する、または、前記マッピングの前記第1のデータセット候補もしくは前記第2のデータセット候補からデータを受け取るデータ処理プログラムに関連して前記ユーザインターフェースを通して受け取られた前記マッピングを格納するステップと、

10

20

を含む方法。

【請求項 2】

前記セットを前記ユーザインターフェースを通して提示することを含む請求項 1 に記載の方法。

【請求項 3】

前記一つ以上の基準への一致の数量化に従って順位付けられた、可能性のあるマッピングのリストを前記ユーザインターフェースを通して提示することを含む請求項 1 に記載の方法。

【請求項 4】

前記可能性のあるマッピングのリストは、前記リスト内にてより高位に順序付けられた所与のデータセットのバージョンである可能性がより高い候補を含む請求項 3 に記載の方法。

10

【請求項 5】

前記基準の一つが前記第 1 のデータセット候補及び前記第 2 のデータセット候補を特定するマップに組み込まれている請求項 3 に記載の方法。

【請求項 6】

前記基準の一つは前記ユーザインターフェースから受け取られる請求項 3 に記載の方法

【請求項 7】

前記可能性のあるマッピングの少なくとも一つは前記データセット候補の一つを表現するデータフローグラフの構成要素を示し、且つ、前記可能性のあるマッピングの少なくとも一つは前記データセット候補の一つを表現しないデータフローグラフの構成要素を示す、請求項 3 に記載の方法。

20

【請求項 8】

複数の構成要素を含むデータフローグラフのサブグラフは前記データセット候補の一つを表現する請求項 7 に記載の方法。

【請求項 9】

前記サブグラフはデータ構成要素を含む請求項 8 に記載の方法。

【請求項 10】

前記サブグラフは実行可能な構成要素を含む請求項 8 に記載の方法。

30

【請求項 11】

前記第 1 のデータセット候補及び前記第 2 のデータセット候補を特定することは、前記第 1 のデータセット候補が前記第 2 のデータセット候補と共通の一つ以上の特徴を有するか否かを判別するためのヒューリスティクスを使用することを含む請求項 1 に記載の方法。

【請求項 12】

前記特徴は前記データセット候補の一つの表現におけるバイト及びレコードの量を含む請求項 11 に記載の方法。

【請求項 13】

前記特徴は前記データセット候補の一つの表現の名称を含む請求項 11 に記載の方法。

40

【請求項 14】

前記特徴は前記データセット候補の一つの表現の生成日を含む請求項 11 に記載の方法

【請求項 15】

前記特徴は前記データセット候補の一つの表現のデータフォーマットを含む請求項 11 に記載の方法。

【請求項 16】

前記マッピングの前記第 1 のデータセット候補及び前記第 2 のデータセット候補の少なくとも一つはデータ管理システムに知られているデータセットのグループに属する請求項 1 に記載の方法。

50

【請求項 17】

さらに、前記第 1 のデータセット候補及び前記第 2 のデータセット候補の間のフォーマットマッピングを提供することを含む請求項 1 に記載の方法。

【請求項 18】

前記マッピングは、前記第 1 のデータセット候補及び前記第 2 のデータセット候補を追跡する前記データ管理システムにおけるレコードを指し示す識別子を含む請求項 1 に記載の方法。

【請求項 19】

さらに、前記第 1 のデータセット候補及び前記第 2 のデータセット候補の一方又は双方の変化に基づいて前記マッピングを更新することを含む請求項 1 に記載の方法。

10

【請求項 20】

データストレージシステム内に格納されたデータをマッピングするシステムであって、命令を格納するメモリと、前記格納された命令を実行するプロセッサと、少なくとも一つの入力データセットからデータを受取り、且つ、少なくとも一つの出力データセットにデータを出力するデータ処理プログラムの仕様を格納するデータストレージシステムと、

第 1 のデータセット候補を特定し、前記第 1 のデータセット候補は前記データ処理プログラムのうちの第 1 のデータ処理プログラムに表現され、前記第 1 のデータセット候補は第 1 のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、マ

20

ッパであって、第 2 のデータセット候補を特定し、前記第 2 のデータセット候補は前記データ処理プログラムのうちの第 2 のデータ処理プログラムに表現され、前記第 2 のデータセット候補は前記第 1 のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、マッパと、

前記第 1 のデータセット候補と前記第 2 のデータセット候補との間のマッピングを受け取るユーザインターフェースであって、前記マッピングは、前記第 1 のデータセット候補と前記第 2 のデータセット候補の双方が同じデータセットを表現することを示し、前記ユーザインターフェースは、前記マッピングの前記第 1 のデータセット候補もしくは前記第 2 のデータセット候補に対してデータを提供する、または、前記マッピングの前記第 1 の

30

データセット候補もしくは前記第 2 のデータセット候補からデータを受け取るデータ処理プログラムに関連して前記データストレージシステム内の前記マッピングを格納する、ユーザインターフェースと

を含むシステム。

【請求項 21】

前記ユーザインターフェースが前記セットを提示する請求項 20 に記載のシステム。

【請求項 22】

前記ユーザインターフェースが、前記一つ以上の基準への一致の数量化に従って順位付けられた、可能性のあるマッピングのリストを提示する請求項 20 に記載のシステム。

【請求項 23】

前記可能性のあるマッピングのリストは、前記リスト内にてより高位に順序付けられた所与のデータセットのバージョンである可能性がより高い候補を含む請求項 22 に記載のシステム。

40

【請求項 24】

前記基準の一つが前記マッパに組み込まれている請求項 22 に記載のシステム。

【請求項 25】

前記基準の一つは前記ユーザインターフェースによって受け取られる請求項 22 に記載のシステム。

【請求項 26】

前記可能性のあるマッピングの少なくとも一つは前記データセット候補の一つを表現す

50

るデータフローグラフの構成要素を示し、且つ、前記可能性のあるマッピングの少なくとも一つは前記データセット候補の一つを表現しないデータフローグラフの構成要素を示す、請求項 2 2 に記載のシステム。

【請求項 2 7】

複数の構成要素を含むデータフローグラフのサブグラフは前記データセット候補の一つを表現する請求項 2 6 に記載のシステム。

【請求項 2 8】

前記サブグラフはデータ構成要素を含む請求項 2 7 に記載のシステム。

【請求項 2 9】

前記サブグラフは実行可能な構成要素を含む請求項 2 7 に記載のシステム。

10

【請求項 3 0】

前記マップは、前記第 1 のデータセット候補が前記第 2 のデータセット候補と共通の一つ以上の特徴を有するか否かを判別するためのヒューリスティクスを使用する請求項 2 0 に記載のシステム。

【請求項 3 1】

前記特徴は前記データセット候補の一つの表現におけるバイト及びレコードの量を含む請求項 3 0 に記載のシステム。

【請求項 3 2】

前記特徴は前記データセット候補の一つの表現の名称を含む請求項 3 0 に記載のシステム。

20

【請求項 3 3】

前記特徴は前記データセット候補の一つの表現の生成日を含む請求項 3 0 に記載のシステム。

【請求項 3 4】

前記特徴は前記データセット候補の一つの表現のデータフォーマットを含む請求項 3 0 に記載のシステム。

【請求項 3 5】

前記マッピングの前記第 1 のデータセット候補及び前記第 2 のデータセット候補の少なくとも一つはデータ管理システムに知られているデータセットのグループに属する請求項 2 0 に記載のシステム。

30

【請求項 3 6】

前記マップは、前記第 1 のデータセット候補と前記第 2 のデータセット候補との間のフォーマットマッピングを発生する請求項 2 0 に記載のシステム。

【請求項 3 7】

前記マッピングは、前記第 1 のデータセット候補及び前記第 2 のデータセット候補を追跡する前記データ管理システムにおけるレコードを指し示す識別子を含む請求項 2 0 に記載のシステム。

【請求項 3 8】

前記マップは前記第 1 のデータセット候補及び前記第 2 のデータセット候補の一方又は双方の変化に基づいて前記マッピングを更新する請求項 2 0 に記載のシステム。

40

【請求項 3 9】

データストレージシステムに格納されたデータをマッピングするシステムであって、少なくとも一つの入力データセットからデータを受取り、且つ、少なくとも一つの出力データセットにデータを出力するデータ処理プログラムの仕様を格納するデータストレージシステムと、

第 1 のデータセット候補を特定し、前記第 1 のデータセット候補は前記データ処理プログラムのうちの第 1 のデータ処理プログラムに表現され、前記第 1 のデータセット候補は第 1 のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、手段であって、

第 2 のデータセット候補を特定し、前記第 2 のデータセット候補は前記データ処理プロ

50

グラムの中の第2のデータ処理プログラムに表現され、前記第2のデータセット候補は前記第1のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、手段と、

前記第1のデータセット候補と前記第2のデータセット候補との間のマッピングを受け取るユーザインターフェースを提供する手段であって、前記マッピングは、前記第1のデータセット候補と前記第2のデータセット候補の双方が同じデータセットを表現することを示す、手段と、

前記マッピングの前記第1のデータセット候補もしくは前記第2のデータセット候補に対してデータを提供する、または前記マッピングの前記第1のデータセット候補もしくは前記第2のデータセット候補からデータを受け取るデータ処理プログラムに関連して前記ユーザインターフェースを通して受け取られた前記マッピングを格納する手段と、
を含むシステム。

【請求項40】

データストレージシステムに格納されたデータをマッピングするためのコンピュータプログラムを格納するコンピュータ読み取り可能ストレージデバイスであって、前記コンピュータプログラムは、コンピュータに、

少なくとも一つの入力データセットからデータを受取り、且つ、少なくとも一つの出力データセットにデータを出力するデータ処理プログラムの仕様を処理させる命令と、

第1のデータセット候補を特定させる命令であって、前記第1のデータセット候補は前記データ処理プログラムのうちの第1のデータ処理プログラムに表現され、前記第1のデータセット候補は第1のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、命令と、

第2のデータセット候補を特定させる命令であって、前記第2のデータセット候補は前記データ処理プログラムのうちの第2のデータ処理プログラムに表現され、前記第2のデータセット候補は前記第1のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、命令と、

前記第1のデータセット候補と前記第2のデータセット候補との間のマッピングを受け取るユーザインターフェースを提供させる命令であって、前記マッピングは、前記第1のデータセット候補と前記第2のデータセット候補の双方が同じデータセットを表現することを示す、命令と、

前記マッピングの前記第1のデータセット候補もしくは前記第2のデータセット候補に対してデータを提供する、または前記マッピングの前記第1のデータセット候補もしくは前記第2のデータセット候補からデータを受け取るデータ処理プログラムに関連して前記ユーザインターフェースを通して受け取られた前記マッピングを格納させる命令と
を含む、コンピュータ読み取り可能ストレージデバイス。

【請求項41】

前記第1のデータセットの各バージョンは、前記データストレージシステムに関連した異なる場所に格納される請求項1に記載の方法。

【請求項42】

前記第1のデータセットの各バージョンは、異なるデータストレージフォーマットを用いて解釈される請求項1に記載の方法。

【請求項43】

前記第1のデータセットの各バージョンは、前記データ処理プログラムの実行間で変化するパラメータを用いてアクセスされる請求項1に記載の方法。

【請求項44】

前記第1のデータセットの各バージョンは、前記データストレージシステムに関連した異なる場所に格納される請求項20に記載のシステム。

【請求項45】

前記第1のデータセットの各バージョンは、異なるデータストレージフォーマットを用いて解釈される請求項20に記載のシステム。

10

20

30

40

50

【請求項 4 6】

前記第 1 のデータセットの各バージョンは、前記データフローグラフの実行間で変化するパラメータを用いてアクセスされる請求項 2 0 に記載のシステム。

【請求項 4 7】

前記ユーザインターフェースが前記セットを提示する請求項 3 9 に記載のシステム。

【請求項 4 8】

前記ユーザインターフェースが、前記一つ以上の基準への一致の数量化に従って順位付けられた、可能性のあるマッピングのリストを提示する請求項 3 9 に記載のシステム。

【請求項 4 9】

前記可能性のあるマッピングのリストは、前記リスト内にてより高位に順序付けられた所与のデータセットのバージョンである可能性がより高い候補を含む請求項 4 8 に記載のシステム。

10

【請求項 5 0】

前記第 1 のデータセット候補及び前記第 2 のデータセット候補を特定する前記手段は、前記基準の一つを含む請求項 4 8 に記載のシステム。

【請求項 5 1】

前記基準の一つは、前記ユーザインターフェースによって受信される請求項 4 8 に記載のシステム。

【請求項 5 2】

前記可能性のあるマッピングの少なくとも一つは前記データセット候補の 1 つを表現するデータフローグラフの構成要素を示し、且つ、前記可能性のあるマッピングの少なくとも一つは前記データセット候補の 1 つを表現しないデータフローグラフの構成要素を示す、請求項 4 8 に記載のシステム。

20

【請求項 5 3】

複数の構成要素を含むデータフローグラフのサブグラフは前記データセット候補の 1 つを表現する請求項 5 2 に記載のシステム。

【請求項 5 4】

前記サブグラフはデータ構成要素を含む請求項 5 3 に記載のシステム。

【請求項 5 5】

前記サブグラフは実行可能な構成要素を含む請求項 5 3 に記載のシステム。

30

【請求項 5 6】

前記第 1 のデータセット候補及び前記第 2 のデータセット候補を特定する前記手段は、前記第 1 のデータセット候補が前記第 2 のデータセット候補と共通の一つ以上の特徴を有するか否かを判別するためのヒューリスティクスを使用する請求項 3 9 に記載のシステム。

【請求項 5 7】

前記特徴は前記データセット候補の一つの表現におけるバイト及びレコードの量を含む請求項 5 6 に記載のシステム。

【請求項 5 8】

前記特徴は前記データセット候補の一つの表現の名称を含む請求項 5 6 に記載のシステム。

40

【請求項 5 9】

前記特徴は前記データセット候補の一つの表現の生成日を含む請求項 5 6 に記載のシステム。

【請求項 6 0】

前記特徴は前記データセット候補の一つの表現のデータフォーマットを含む請求項 5 6 に記載のシステム。

【請求項 6 1】

前記マッピングの前記第 1 のデータセット候補及び前記第 2 のデータセット候補の少なくとも一つはデータ管理システムに知られているデータセットのグループに属する請求項

50

39に記載のシステム。

【請求項62】

前記第1のデータセット候補及び前記第2のデータセット候補を特定する前記手段は、前記第1のデータセット候補と前記第2のデータセット候補との間のフォーマットマッピングを発生する請求項39に記載のシステム。

【請求項63】

前記第1のデータセットの各バージョンは、前記データストレージシステムに関連した異なる場所に格納される請求項39に記載のシステム。

【請求項64】

前記第1のデータセットの各バージョンは、異なるデータストレージフォーマットを用いて解釈される請求項39に記載のシステム。

10

【請求項65】

前記第1のデータセットの各バージョンは、前記データフローグラフの実行間で変化するパラメータを用いてアクセスされる請求項39に記載のシステム。

【請求項66】

前記コンピュータプログラムは、コンピュータに、前記第1のデータセット候補及び前記第2のデータセット候補を前記ユーザインターフェースを通して提示させる命令をさらに含む、請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項67】

前記コンピュータプログラムは、コンピュータに、前記一つ以上の基準への一致の数量化に従って順位付けられた、可能性のあるマッピングのリストを前記ユーザインターフェースを通して提示させる命令をさらに含む、請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

20

【請求項68】

前記可能性のあるマッピングのリストは、前記リスト内にてより高位に順序付けられた所与のデータセットのバージョンである可能性がより高い候補を含む請求項67に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項69】

前記基準の一つが前記第1のデータセット候補及び前記第2のデータセット候補を特定するマップに組み込まれている請求項67に記載のコンピュータ読み取り可能ストレージデバイス。

30

【請求項70】

前記基準の一つは前記ユーザインターフェースから受け取られる請求項67に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項71】

前記可能性のあるマッピングの少なくとも一つは前記データセット候補の一つを表現するデータフローグラフの構成要素を示し、且つ、前記可能性のあるマッピングの少なくとも一つは前記データセット候補の一つを表現しないデータフローグラフの構成要素を示す、請求項67に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項72】

複数の構成要素を含むデータフローグラフのサブグラフは前記データセット候補の一つを表現する請求項71に記載のコンピュータ読み取り可能ストレージデバイス。

40

【請求項73】

前記サブグラフはデータ構成要素を含む請求項72に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項74】

前記サブグラフは実行可能な構成要素を含む請求項72に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項75】

前記第1のデータセット候補及び前記第2のデータセット候補を特定することは、前記

50

第1のデータセット候補が前記第2のデータセット候補と共通の一つ以上の特徴を有するか否かを判別するためのヒューリスティクスを使用することを含む請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項76】

前記特徴は前記データセット候補の一つの表現におけるバイト及びレコードの量を含む請求項75に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項77】

前記特徴は前記データセット候補の一つの表現の名称を含む請求項75に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項78】

前記特徴は前記データセット候補の一つの表現の生成日を含む請求項75に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項79】

前記特徴は前記データセット候補の一つの表現のデータフォーマットを含む請求項75に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項80】

前記マッピングの前記第1のデータセット候補及び前記第2のデータセット候補の少なくとも一つはデータ管理システムに知られているデータセットのグループに属する請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項81】

前記コンピュータプログラムは、コンピュータに、前記第1のデータセット候補と前記第2のデータセット候補との間のフォーマットマッピングを提供させる命令をさらに含む、請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項82】

前記マッピングは、前記第1のデータセット候補及び前記第2のデータセット候補を追跡するデータ管理システムにおけるレコードを指し示す識別子を含む請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項83】

前記コンピュータプログラムは、コンピュータに、前記第1のデータセット候補及び前記第2のデータセット候補の一方又は双方の変化に基づいて前記マッピングを更新させる命令をさらに含む、請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項84】

前記第1のデータセットの各バージョンは、前記データストレージシステムに関連した異なる場所に格納される請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項85】

前記第1のデータセットの各バージョンは、異なるデータストレージフォーマットを用いて解釈される請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項86】

前記第1のデータセットの各バージョンは、前記データフローグラフの実行間で変化するパラメータを用いてアクセスされる請求項40に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項87】

前記データ処理プログラムの少なくとも一つは、データのフローを表現するリンクによって相互接続された計算を表現するノードを含むデータフローグラフを含み、前記データフローグラフは、前記少なくとも一つの入力データセットからデータのフローを受け取り、且つ、前記少なくとも一つ出力データセットにデータのフローを提供する請求項1に記載の方法。

【請求項88】

前記第1のデータ処理プログラムはデータフローグラフのサブグラフであり、前記第2

10

20

30

40

50

のデータ処理プログラムは前記データフローグラフのサブグラフである、請求項 1 に記載の方法。

【請求項 8 9】

前記データ処理プログラムの少なくとも一つは、データのフローを表現するリンクによって相互接続された計算を表現するノードを含むデータフローグラフを含み、前記データフローグラフは、前記少なくとも一つの入力データセットからデータのフローを受け取り、且つ、前記少なくとも一つ出力データセットにデータのフローを提供する請求項 2 0 に記載のシステム。

【請求項 9 0】

前記第 1 のデータ処理プログラムはデータフローグラフのサブグラフであり、前記第 2 のデータ処理プログラムは前記データフローグラフのサブグラフである、請求項 2 0 に記載のシステム。

10

【請求項 9 1】

前記データ処理プログラムの少なくとも一つは、データのフローを表現するリンクによって相互接続された計算を表現するノードを含むデータフローグラフを含み、前記データフローグラフは、前記少なくとも一つの入力データセットからデータのフローを受け取り、且つ、前記少なくとも一つ出力データセットにデータのフローを提供する請求項 3 9 に記載のシステム。

【請求項 9 2】

前記第 1 のデータ処理プログラムはデータフローグラフのサブグラフであり、前記第 2 のデータ処理プログラムは前記データフローグラフのサブグラフである、請求項 3 9 に記載のシステム。

20

【請求項 9 3】

前記データ処理プログラムの少なくとも一つは、データのフローを表現するリンクによって相互接続された計算を表現するノードを含むデータフローグラフを含み、前記データフローグラフは、前記少なくとも一つの入力データセットからデータのフローを受け取り、且つ、前記少なくとも一つ出力データセットにデータのフローを提供する請求項 4 0 に記載のコンピュータ読み取り可能ストレージデバイス。

【請求項 9 4】

前記第 1 のデータ処理プログラムはデータフローグラフのサブグラフであり、前記第 2 のデータ処理プログラムは前記データフローグラフのサブグラフである、請求項 4 0 に記載のコンピュータ読み取り可能ストレージデバイス。

30

【請求項 9 5】

前記第 1 のデータセットの各バージョンは、異なるデータフローグラフ、データフローグラフサブセットまたは実行可能な構成要素に関連する請求項 8 7 に記載の方法。

【請求項 9 6】

前記第 1 のデータセットの各バージョンは、異なるデータフローグラフ、データフローグラフサブセットまたは実行可能な構成要素に関連する請求項 8 9 に記載のシステム。

【請求項 9 7】

前記第 1 のデータセットの各バージョンは、異なるデータフローグラフ、データフローグラフサブセットまたは実行可能な構成要素に関連する請求項 9 1 に記載のシステム。

40

【請求項 9 8】

前記第 1 のデータセットの各バージョンは、異なるデータフローグラフ、データフローグラフサブセットまたは実行可能な構成要素に関連する請求項 9 3 に記載のコンピュータ読み取り可能ストレージデバイス。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明はデータ管理システム内のデータセット（データセット）のインスタンス（instance）のマッピングに関する。

50

【背景技術】

【0002】

近年のデータ管理システムは、そのシステムの異なる特徴を表す多数の要素を含む。より複雑でないシステムは、しばしば、正確な視覚化の目的のために追加の処理なしでデータを直接見られることを可能にする。より複雑なシステムは、意味あるようにデータを見ることができるようにするために追加のメカニズムを必要とするであろう。多くの要素よりなる複雑なデータ管理システムは、データを多く異なる形式で格納し、且つ、データを多くの異なる方法で処理することができる。これらの格納形式及び処理形式の多くは、相互関係を解析する方法なしでは明白とならない態様において相互に関係する。

【発明の概要】

【0003】

< 関連出願 >

本願は、2008年12月2日出願の米国特許出願第61/119,164号の優先権を主張するものであり、ここに引用して組み込まれる。

概して、コンピュータシステムによって使用される、データストレージシステム（データ格納システム）内に格納されたデータをマッピングする方法は、

ノードを含むデータフローグラフの仕様を処理するステップであって、前記ノードがデータのフローを表現するリンクによって相互接続された計算を表現し、前記データフローグラフの少なくとも一つが少なくとも一つの入力データセットからデータのフローを受取り、且つ、前記データフローグラフの少なくとも一つが少なくとも一つの出力データセットにデータのフローを提供する、ステップ；

データセットの一つ以上のセットを特定するステップであって、所与のセット内の各データセットが単一のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、ステップ；

所与のセットにおける少なくとも二つのデータセットの間のマッピングを受け取るユーザインターフェースを提供するステップ；及び、

データを前記マッピングのデータセットに提供するかまたはデータを前記マッピングのデータセットから受け取る一つのデータフローグラフに関連して前記ユーザインターフェースを通して受け取られた前記マッピングを格納するステップ；

を含む。

【0004】

他の概要において、データストレージシステム内に格納されたデータをマッピングするシステムは、

ノードを含むデータフローグラフの仕様を格納するデータストレージシステム（データ格納システム）であって、前記ノードがデータのフローを表現するリンクによって相互接続された計算を表現し、前記データフローグラフの少なくとも一つが少なくとも一つの入力データセットからデータのフローを受取り、且つ、前記データフローグラフの少なくとも一つが少なくとも一つの出力データセットにデータのフローを提供する、データストレージシステム；

データフローグラフと関連したデータセットの一つ以上のセットを特定するマップ（mapper）であって、所与のセット内の各データセットが単一のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、マップ；及び、

所与のセットにおける少なくとも二つのデータセットの間のマッピングを受け取るユーザインターフェースであって、データを前記マッピングのデータセットに提供するかまたはデータを前記マッピングのデータセットから受け取る一つのデータフローグラフに関連して前記データストレージシステム内の前記マッピングを格納するユーザインターフェース；

を含む。

【0005】

他の概要において、データストレージシステムに格納されたデータをマッピングするシ

10

20

30

40

50

ステムは、

ノードを含むデータフローグラフの仕様を処理する手段であって、前記ノードがデータのフローを表現するリンクによって相互接続された計算を表現し、前記データフローグラフの少なくとも一つが少なくとも一つの入力データセットからデータのフローを受取り、且つ、前記データフローグラフの少なくとも一つが少なくとも一つの出力データセットにデータのフローを提供する、手段；

データセットの一つ以上のセットを特定する手段であって、所与のセット内の各データセットが単一のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、手段；

所与のセットにおける少なくとも二つのデータセットの間のマッピングを受け取るユーザインターフェースを提供する手段；及び、

データを前記マッピングのデータセットに提供するかまたはデータを前記マッピングのデータセットから受け取る一つのデータフローグラフに関連して前記ユーザインターフェースを通して受け取られた前記マッピングを格納する手段；

を含む。

【0006】

他の概要において、データストレージシステムに格納されたデータをマッピングするためのコンピュータプログラムは、コンピュータ読み出し可能媒体に格納され、且つ、コンピュータに、

ノードを含むデータフローグラフの仕様を処理させる命令であって、前記ノードがデータのフローを表現するリンクによって相互接続された計算を表現し、前記データフローグラフの少なくとも一つが少なくとも一つの入力データセットからデータのフローを受取り、且つ、前記データフローグラフの少なくとも一つが少なくとも一つの出力データセットにデータのフローを提供する、命令；

データセットの一つ以上のセットを特定させる命令であって、所与のセット内の各データセットが単一のデータセットの異なるバージョンを特定するための一つ以上の基準に一致する、命令；

所与のセットにおける少なくとも二つのデータセットの間のマッピングを受け取るユーザインターフェースを提供させる命令；及び、

データを前記マッピングのデータセットに提供するかまたはデータを前記マッピングのデータセットから受け取る一つのデータフローグラフに関連して前記ユーザインターフェースを通して受け取られた前記マッピングを格納させる命令；

を含む。

【0007】

態様は一つ以上の以下の特徴を含むことができる。

【0008】

前記セットは前記ユーザインターフェース上に提示される。

【0009】

前記一つ以上の基準への一致の数量化に従って順位付けられた可能性のあるマッピングのリストは、前記ユーザインターフェースを通して提示される。

【0010】

前記可能性のあるマッピングのリストは、前記リスト内にてより高位に順序付けられた所与のデータセットのインスタンスである可能性がより高い候補を含む。

【0011】

前記基準の一つは、前記一つ以上のデータセットを特定するマップに組込まれている。

【0012】

前記基準の一つは前記ユーザインターフェースから受け取られる。

【0013】

前記可能性のあるマッピングの少なくとも一つは一つのデータセットを表現するデータフローグラフの構成要素を示し、且つ、前記可能性のあるマッピングの少なくとも一つは

10

20

30

40

50

データセットを表現しないデータフローグラフの構成要素を示す。

【0014】

複数の構成要素を含むデータフローグラフのサブグラフは、ひとつのデータセットを表現する。

【0015】

前記サブグラフはデータ構成要素を含む。

【0016】

前記サブグラフは実行可能な構成要素を含む。

【0017】

前記データセットの一つ以上のセットを特定することは、所与のセット内の一つのデータセットが他のデータセットと共通の一つ以上の特徴を有するか否かを判別するためのヒューリスティクスを使用することを含む。

10

【0018】

前記特徴は前記データセットの表現におけるバイト及びレコードの量を含む。

【0019】

前記特徴は前記データセットの表現の名称を含む。

【0020】

前記特徴は前記データセットの表現の生成日を含む。

【0021】

前記特徴は前記データセットの表現のデータフォーマットを含む。

20

【0022】

前記マッピングのデータセットの少なくとも一つはデータ管理システムに知られているデータセットのグループに属する。

【0023】

所与のセットのデータセット間のフォーマットマッピングが提供される。

【0024】

前記マッピングは、前記データセットを追跡する前記データ管理システムにおけるレコードを指し示す識別子を含む。

【0025】

データセットの変化に基づいて前記マッピングが更新される。

30

【0026】

本発明の態様は以下の一つ以上の利点を有することができる。

【0027】

バージョンを特定(識別)する基準に従ってデータセットの集合(セット)を特定することによって、一つのデータセットの二つのインスタンスの間的一致が純粋な手作業よりもより効率的になされ得る。さらに、少なくとも二つのデータセット間のマッピングを受取るユーザインターフェースを提供することにより、そのマッピングは、システムが純粋に自動化された場合よりもより正確になる。

【0028】

本発明の他の特徴及び利点は以下の説明及び特許請求の範囲の請求項から明らかになる。

40

【図面の簡単な説明】

【0029】

【図1】データフローグラフである。

【図2】データセットのマappa及び関連する構成要素の概観図である。

【図3A】データセットのマappaによって取扱われるシナリオの線図である。

【図3B】データセットマappaによって取扱われる異なるシナリオの線図である。

【図3C】データセットマappaによって取扱われる異なるシナリオの線図である。

【図3D】データセットマappaによって取扱われる異なるシナリオの線図である。

【図3E】データセットマappaによって取扱われる異なるシナリオの線図である。

50

【図4】データセットのマッパの動作フローチャートである。

【図5】データセットリンケージマッピングである。

【図6】データセットフォーマットマッピングである。

【発明を実施するための形態】

【0030】

概観

データ処理要素はグラフの形式とすることができる。グラフに基づく計算は、構成要素（格納されたデータに対応するデータストレージ構成要素、または、実行可能な処理に対応する計算の構成要素の何れか）を表現するグラフにおける頂点（vertices）を有する有向グラフにより表現される「データフローグラフ」であってそのグラフにおける有向リンクまたは「辺（edges）」が構成要素間のデータのフロー（流れ）を表現する「データフローグラフ」を用いて実行される。データフローグラフ（これは、また、単に「グラフ」と呼ばれる。）はモジュールの統一体（modular entity）である。各グラフは、一つ以上の他の複数のグラフから構成されることができ、また、特定のグラフはより大きなグラフにおける一つの構成要素であることができる。グラフ開発環境（GDE）は、実行可能なグラフを特定するとともにそのグラフの構成要素のパラメータを定義するユーザインターフェースを提供する。

【0031】

図1を参照すると、データフローグラフ101の一例は、データフローグラフ101の実行可能構成要素104a - 104jによって処理されるべきデータの集まりを提供する入力構成要素102を含む。たとえば、データセット102は、データベースシステムに関連するデータレコードあるいはトランザクション処理システムに関連するトランザクションを含むことができる。実行可能な構成要素の各々は、全体のデータフローグラフ101によって規定される計算の一部に関連している。作業要素（たとえばデータの集まりからの個別的なデータレコード）は一つの構成要素の一つ以上の入力ポートに入り、また、出力作業要素（幾つかの場合には、入力作業要素または入力作業要素の処理されたバージョン）は一般にその構成要素の一つ以上の出力ポートから出る。グラフ101において、構成要素104e、104g及び104jからの出力作業要素は、出力データ構成要素102a - 102c内に格納される。

【0032】

データセットは、特定のデータの集まりを表わす（たとえばオブジェクト指向データベース内に格納された）オブジェクトである。データフローグラフのシステムのコンテキストにおいて、一つの構成要素は一つのデータセットを表わすことができる。このような場合、グラフはデータセットを表現する構成要素（即ち、単に「データセット構成要素」と一つ以上の方法で相互に作用することができる。一つのデータセット構成要素は、所与のデータセットによって表現された物理的データをアクセスするための命令を含み、従って、一つのグラフはデータセットからの入力をデータセット構成要素を用いて受け取り、出力をデータセットにデータセット構成要素を用いて提供し、且つ、中間ステップにてデータセットのデータをデータセット構成要素を用いて処理することができる。データセット構成要素は、データセットオブジェクトの一つのインスタンスを含む所与のデータセットオブジェクトに関連する種々の情報を含むことができる。このようなシステムは、数十、数百、あるいは数千のグラフ及びこれに関連するデータセット構成要素を有することになる。このようなシステムの複雑度が増大するにつれて、種々のグラフとデータセット構成要素との間の関係を管理することは益々困難となる。システムにおける二つ以上のデータセット構成要素は同一のデータソースを表現することができ、また、そのようなデータセット構成要素は、各々、別のグラフ、グラフサブセットまたは実行可能な構成要素に関連付けることができる。

【0033】

たとえば、一つの可能なシナリオにおいては、単一のデータセットはデータ管理システムに関連した二つ以上の場所に格納され得る。このシナリオにおいては、二つ以上のデー

10

20

30

40

50

タソースが、同一データの類似バージョンまたは同一バージョンを含む。このシステムにおける二つのグラフはこの単一のデータセットを取り扱うことになるが、各グラフは、別のデータベーステーブルまたは他の形式のデータセット構成要素を、別のデータファイルから読み取り且つ別のデータファイルへと書き込む。

【 0 0 3 4 】

類似したシナリオにおいて、所与のデータセットによって表現されたデータ（たとえばデータファイル）は二つ以上の場所に格納されるだけでなく、異なるデータストレージフォーマットを用いて解釈され得る。上述の例と同様、二つのグラフは、（フォーマットのみ異なる）同一のデータを含む二つの別のデータファイルに作用することができる。各データファイルは、同一のデータのインスタンスを含んでいるにもかかわらず、異なるデータ形式の構成（配列、arrangement）を有することがある。

10

【 0 0 3 5 】

別のシナリオにおいては、一つのグラフはデータセットのインスタンスを含む一つのデータファイルに作用し、また、他のグラフはそのデータセットのインスタンスを含むデータベーステーブルに作用する。このような場合、データファイル及びデータベーステーブルは一般に二つの異なるデータフォーマットを有するであろう。

【 0 0 3 6 】

他のシナリオにおいては、データ管理システムは同一のデータセットの異なるバージョンを各々異なる方法でアクセスできる。一つのグラフは、たとえばデータファイルを標準ファイル入出力メカニズムを用いて読み込むことによってデータセットのインスタンスを直接的にアクセスできる。他のグラフは、ネットワークを介して利用可能なデータ収納庫（repository）のような外部ソースに問い合わせることによってファイルを取得（retrieve）できる。また、あるグラフは、類似した外部問い合わせ（たとえばネットワークされたデータベースに対する問い合わせ）によってデータベーステーブルをアクセスできる。

20

【 0 0 3 7 】

また、データ管理システムは、同一データセットの異なるインスタンスを各々異なる方法で参照することができる。たとえば一つのグラフはパラメータに従って異なるデータ場所をアクセスできる。このようなパラメータは時間と共に多くのデータ場所を指し示すことができる。そのパラメータがグラフの実行間で変化すれば、複数回動作するグラフは異なる機会に異なる場所をアクセスできる。

30

【 0 0 3 8 】

あるシナリオにおいては、一つのグラフ内における一つのデータセットの表現は、単一構成要素ではなく、むしろ構成要素の集まり、たとえば、それ自体が複数の構成要素を有する一つのグラフとして実行されるグラフ内の「サブグラフ」である。この集まりは、一つ以上のデータセットを含むことができ、且つ、一つ以上の実行可能な構成要素を有することができる。

【 0 0 3 9 】

上述のシナリオのすべては、データ管理システムによって取り扱われるデータを視覚化し且つ解析することに対する問題を潜在的に提示する。ユーザが所与の一つのデータセットと相互に作用する構成要素の統合されたビューを必要とするとき、存在し得るデータセットの異なるインスタンスを調整（reconcile）するために種々のアプローチが用いられ得る。

40

【 0 0 4 0 】

一つのアプローチは、同一のデータセットの複数のインスタンスを特定し且つそれらの間にリンケージ（つながり、linkage）を生成する自動的なメカニズムである。しかしながら、幾つかの自動的なメカニズムは欠点、つまり以下に示す3つの欠点を有する。第1に、このメカニズムは、一つのデータセットの各インスタンスが特定な方法で格納されること、たとえば統一名称方式（unified naming scheme）及びディレクトリ構造の下で格納されることを必要とする。これは、データ管理システムに関連するストレージシステム内において各々を特定し且つ場所を突き

50

とめる方法をそのメカニズムに提供する。しかしながら、この構成はデータ管理システムの柔軟性を限定すると共に、このシステムの幾つかの使用方法について過度に制限することになり得る。

【0041】

第2に、動作についての幾つかのシナリオにおいて、そのメカニズムは同一のデータセットのインスタンスを正しく特定せず、また、正しいリンケージを形成しない場合がある。たとえば、このことは、データセットを外部参照エンティティを用いてアクセスし、且つ、その自動的なメカニズムがそのエンティティに対するアクセスを有しない場合に発生する可能性が高い。同様に、このことは、構成要素がデータセットをパラメータリストにおける独立パラメータに従ってアクセスし、且つ、この自動的なメカニズムがそのパラメータリストをアクセスする方法または解釈する方法を有しない場合に発生する可能性が高い。さらに、このことは、一つのデータセットが、一つ以上のデータセット構成要素及び実行可能な構成要素（たとえばサブグラフ）から構成される複雑なエンティティによって表現されているときに発生する可能性が高い。自動的なメカニズムが、構成要素のどの特定の組み合わせが特定のデータセットを表現するかを識別できない場合もある。

10

【0042】

第3に、そのメカニズムは、データセットのインスタンスの間に冗長または不要のリンケージを形成する場合がある。たとえば、データ管理システムによって取り扱われるデータセットのいくつかは、たとえばエラーログのような外部からのデータを表現する場合がある。これらのデータセットの間のいかなるリンケージも不要である。さらに、データ管理システムによって取り扱われるデータセットのインスタンスのいくつかは、冗長な（重複する）インスタンス（たとえばキャッシュされたデータまたはデータの他の一時的なコピー）である場合がある。この形式のデータに接続するリンケージは直ちに陳腐化し、且つ、データ管理システムを分析するユーザを当惑させるであろう。

20

【0043】

代替のアプローチは、ユーザが同一データセットのインスタンスをユーザインターフェースを介して手動で統合するシステムである。ユーザが一つのデータセットの複数のインスタンスの間の本質的なリンケージを見逃す可能性は低く、且つ、ユーザが一つのデータセットの複数のインスタンスの間の冗長または不要なリンケージを生成する可能性も低い。しかしながら、データ管理システムが数百・数千の構成要素を有する場合、ユーザが手動で必要なリンケージを生成するために必要な時間はとてつもなく長い。

30

【0044】

部分的に自動化されたアプローチにおいては、データセットマップ（dataset mapper）が使用されることにより、幾つかの自動的な分析を提供され、且つ、大きな及び/または複雑なシステムのユーザにとって手が出せないことがない方法でそのユーザとの所定の交流を可能とする。

【0045】

図2は、例示的なデータセットマップ100の一実施形態のブロック線図であって、関連する主要な要素間の相互関係を示す。データセットマップ100は、一つ以上のグラフ180、180a、180b、180cのセットを解析することができる。各グラフは、一つ以上のデータセット構成要素182、182a、182bに関連し、各データセット構成要素は、データファイル、データベーステーブル、サブグラフまたはデータセットを表現する他の種類の構成要素に対応する可能性がある。マップ100は、同一のデータセット186のインスタンスを含むデータセットの構成要素間のリンケージを形成する目的でグラフを解析する。マップ100は、組込み規則110、ユーザ定義規則120及びヒューリスティクス（経験則、heuristics）130の組合せに従って各データセット構成要素を処理してデータセット構成要素182がデータ管理システム170に知られたデータソース176、176a、176bを表現する幾つかのデータセットの一つのインスタンスを含むか否かを判別する。マップ100はこの情報をユーザインターフェース160に渡し、インターフェース160はデータセット構成要素182に対応する正しいデータ

40

50

セットがあればこの正しいデータセットをユーザ 162 に選択させることを許容する。たとえば、ユーザインターフェース 160 は、単一のデータセットの異なるバージョンまたはインスタンスを特定する一つ以上の基準への合致に基づいて可能性のある候補のマッピングのリストを提示する。組込み規則、ユーザ定義規則及びヒューリスティクスに基づく基準を含むそのような基準の例について以下に詳述する。候補マッピングのリストは一つ以上の基準への合致の数量化に従って順位付け (order) できる (たとえば、所与のデータセットのインスタンスである可能性がより大きい候補はリストにおいてより高位に順序付けられる)。次に、マップ 100 は、データセット構成要素 182 がデータソース 176 を表現するデータセットのインスタンスを含むことを示すデータセットリンケージマッピング 140 を発生する。

10

【0046】

さらに、データセット構成要素 182 は、対応する連結されたデータソース 176 のフォーマット 174 とは異なるデータフォーマット 184 を有することができる。データ管理システム 170 の要求に依り、ユーザはデータセットのすべてのインスタンスに対して単一のデータフォーマットを確立することを選択することができる。このシステムは各データソース 176、176a、176b に対してフォーマット 174、174a、174b を格納する。代替として、ユーザは、データセット構成要素 182 のフォーマット 184 と、対応するデータソース 176 の確立されたフォーマット 174 と、の間に随意 (optional) のマッピング 142 を生成することを選択することができる。その随意のマッピング 142 は、データ管理システム 170 がデータセットの各インスタンスに対するデータの形式についての情報を保持することを許容する。

20

【0047】

また、マップ 100 は、ユーザが、実行可能な構成要素と、他のリンケージは有しないであろう単一のデータセット構成要素と、の間のリンケージを示すことを可能とする。たとえば、データセット構成要素は、単一の読取者を伴うソースデータセット、または、単一の書込者を伴う目標データセットに対応することができる。データセットオブジェクトが、システムに存在し、且つ、正しいレコードフォーマット、文書、データプロフィール等の他の関連するメタデータを有すれば、リンケージはデータセット構成要素を正しいデータセットにマップさせることができる。

【0048】

マッピング処理

マップ 100 は、複雑なデータ管理システムに生ずる共通のシナリオを取扱うことができる。第 1 のシナリオにおいては、図 3A に示したように、一つのグラフ 210 は出力としてデータセット構成要素 212 を提供し、他のグラフ 220 は入力として異なるデータセット構成要素 222 を受取る。各データセット構成要素は、同一のデータセット 216 のインスタンスを含む。このデータセットは、データ管理システムに知られたデータソース 176 を表現するデータセットと同じであり得る。さらに、第 1 のデータセット構成要素 212 は第 2 のデータセット構成要素 222 に属するフォーマットと同じデータフォーマット 214 を有し、あるいは、代替として、第 2 のデータセット構成要素は異なるフォーマット 224 を有してもよい。マップ 100 は、第 2 のデータセット構成要素 222 を第 1 のデータセット構成要素 212 によって表現されたデータセット 216 のインスタンスであるとして特定し、適切なリンケージマッピング 140 を生成することができる。

30

40

【0049】

第 2 のシナリオにおいては、図 3B に示したように、グラフ 230 は、外部ソース 239 に対する外部参照 238 を用いて外部データセット構成要素 232 に関連付けられる。その外部データセット構成要素 232 は、データフォーマット 234 を有し、且つ、データセット 236 のインスタンスである。第 1 のシナリオと同様、外部データセット構成要素によって表されたデータセット 236 は、データ管理システム 170 に知られたデータソース 176 を表現するデータセットであり得る。マップ 100 は、この外部データセット構成要素 232 を他のデータセットのインスタンスであるとして特定し、適切なリンケ

50

ージマッピング140を生成することができる。

【0050】

第3のシナリオにおいては、図3Cに示したように、グラフ240は、パラメタリスト247内のパラメータ248を用いてデータセット構成要素242に関連付けられる。参照されるデータセット構成要素242は、データフォーマット244を有し、且つ、データセット246のインスタンスである。第1、第2のシナリオと同様、参照されるデータセット構成要素によって表されたデータセット246は、データ管理システム170に知られたデータソース176を表現するデータセットであり得る。マップ100は、この参照されるデータセット構成要素242を他のデータセットのインスタンスであるとして特定し、適切なリンケージマッピング140を生成することができる。

10

【0051】

第4のシナリオにおいては、図3Dに示したように、グラフ250は外部ソース259に対する外部参照258を用いて外部構成要素251に関連付けられる。外部構成要素251は、データセット構成要素ではなく、むしろ実行可能な構成要素のような他の種類の構成要素である。マップ100は、この外部構成要素251を、データセットリンケージマッピング処理に適用できないものとして特定することができる。

【0052】

第5のシナリオにおいては、図3Eに示したように、グラフ260は、それ自体複数の構成要素から構成されるサブグラフ構成要素263に関連している。これらの構成要素は、少なくとも一つのデータセット構成要素262を含み、且つ、この例では一つ以上の実行可能な構成要素261a、261b、261cを含む。このシナリオにおいては、単一エンティティとしてのサブグラフ263は、少なくとも一つのデータセットを表現する。他のサブグラフの例は複数のデータセット構成要素及びゼロを含む多数の実行可能な構成要素を含むことができる。さらに、サブグラフ263は複数の出力265a、265bを有する。各出力は、データセットの異なるインスタンスを、その出力を受け取る構成要素に提供することができる。また、サブグラフの他の例は、多数の入力を有することもできる。さらに別のサブグラフの例は、それぞれのデータセットに対応する入力または出力を有しない場合もある。サブグラフが少なくとも一つのデータセットを表現している場合には、マップ100は、サブグラフ263を少なくとも一つのデータセットのインスタンスであるとして特定し、少なくとも一つの適切なリンケージマッピング140を生成することができる。

20

30

【0053】

マップの動作シーケンスの例を図4に示す。ステップ302において、マップは第1に一つのグラフに関連付けられた要素のうちどの要素がデータセットを表現するかを特定する。一般に、グラフは一つ以上の入力及び出力を有し、また、各入力及び各出力はデータセットのインスタンスであろう。また、各グラフは、ある中間ステップにおいて、データセットのインスタンスを取扱うことができる。この結果、各グラフは、データセットの候補となり得る複数の構成要素に接続され得る。幾つかの場合には、データ管理システムは、ある構成要素がデータセットを表現するか否かについての情報を含む「構成要素の特性」についての情報を有する。このような場合、ステップ304において、マップは可能性のあるデータセット構成要素をデータセットの候補テーブルに加える。幾つかの場合、構成要素は、データセット構成要素及び実行可能な構成要素を含む複数の構成要素から構成されたサブグラフであり得る。サブグラフは一つのデータセットの少なくとも一つのインスタンスを表現することができるであろう。従って、マップは、そのようなサブグラフの総てのリストを編集(コンパイル)し、ステップ304の一部としてそれらをデータセットの候補テーブルに加える。他の幾つかの場合において、構成要素の特質(nature)がデータ管理システムに利用され得ないものであることがある。その構成要素は、外部エンティティへの参照を介してアクセスされ得る。ここで、その外部エンティティへの参照とは、データベーステーブル、インターネットサーバを指し示すユニフォームリソースロケータ、パラメタリスト内のパラメータ、または、他の形式の参照に対する問合せ(クエリ

40

50

ー)である。これらの場合、一般に、マップは、その(外部エンティティへの)参照によって指し示されたエンティティを独立にアクセスするための手段を有しない。従って、マップは、そのような参照の総てのリストを編集(コンパイル)し、ステップ304の一部としてそれらをデータセット候補テーブルに加える。

【0054】

次に、ステップ306において、所与のデータセットの候補に対し、マップは、データセットの候補がマップする可能性のある既知のデータセットのリストを発生する。マップは、ユーザ定義規則、組込み規則及びヒューリスティックスの組合せを用いて、既知のデータセットの何れがデータセットに候補にマップされるかを評価する。

【0055】

次に、ステップ308において、ユーザは、そのデータセットの候補に対応する既知のデータセットを選択する。また、ユーザは、提案された既知のデータセットの何れもが正しい合致でない場合、すべての既知のデータセットのリスト全体をアクセスするかもしれない。加えて、ユーザは、データセットの候補がデータセットでないことを示すことができる。たとえば、遠隔サーバへの参照は、遠隔実行可能手続(データエンティティでない)への要求であろう。他の例として、データセットの候補はデータを表現できるが、そのデータは、エラーログのようなデータ管理システムに関係しない種類のデータである場合がある。この場合、ユーザは、ユーザインターフェースに対してこのデータはマッピング処理において無視すべきであることを指示することができる。

【0056】

次に、ステップ310において、ユーザは新しくマップされたデータセットのデータフォーマットを特定する。システムは、データフォーマットテンプレートの集合を有し、その一つが選択され得る。代替として、ユーザはユーザインターフェース内に新しいデータフォーマットを生成することができる。

【0057】

次に、ステップ312において、マップは、この情報を用いてデータセットの候補に対するリンケージマッピング、及び、随意に、フォーマットマッピングを発生する。

【0058】

次に、マップがすべてのデータセットの候補を処理していなければ、マップは、ステップ308、310及び312の繰り返しにおけるリンケージ発生のために、次のデータセットの候補をユーザに提供する。

【0059】

次に、ステップ314において、ユーザはデータ管理システムに関連した構成要素を見て、グラフとデータセット構成要素との間の関連の視覚化が正確であることを構成要素間の新しいリンケージに基づいて保証する。ステップ316において、ユーザはリンケージ及びフォーマットマッピングに対して調整するオプションを有する。

【0060】

最後に、ステップ318において、マップは、リンケージ及びフォーマットマッピングをデータ管理システムに引き渡す。そのマッピングは、一つ以上のグラフと一緒に格納され、またはデータ管理システムに関連する別のストレージエンティティ内に格納され、もしくは、他の手段によって格納されることができる。

【0061】

データセットマッピングのメンテナンス

マップ100は、データセットリンケージの完全性(インテグリティ、integrity)に影響する複数のシナリオを取り扱うことができる。

【0062】

第1のシナリオは、新構成要素がデータ管理システム170に追加されたときに新しいデータセットの候補を特定することを含む。このシナリオの下では、マップ100は各構成要素を分析し且つ可能性のあるリンケージをユーザに提示する。マップ100は、どのような新しい構成要素にも作用して適切なリンケージを必要に応じて発生する。

10

20

30

40

50

【 0 0 6 3 】

第2のシナリオは、データ管理システム170が時間と共に変化したときに、既存のリンクージをメンテナンスすることを含む。たとえば、データ管理システムが関連するグラフの通常の動作中にデータセットの新インスタンスが出現することがある。他の一つの例として、データセットがそのアイデンティティ、たとえばデータ管理システムにおけるその名称及び場所、を変更する場合がある。さらなる例として、データセット全体が削除されている場合もある。さらなる一つの例として、データセット候補は、リンクージ生成の先の回において見落とされていたかも知れず、その結果、リンクージの集まりが不完全な場合がある。マッピングシステムのユーザインターフェース160は、ユーザ162に、既存のリンクージを修正させて不完全または期限切れのマッピングを修復させることを許容する。

10

【 0 0 6 4 】

第3のシナリオは、既知のパターンに判で押したように従うデータセットの参照のリンクージを自動的に更新することを含む。たとえば、グラフは、パラメータリスト247において参照されるデータセットを取り扱うことができる。このようなパラメータリストは時間と共に変化する。パラメータリストがデータ管理システムに知られた標準フォーマットに従うのであれば、マップはパラメータリスト内の変化を特定し且つ既存のリンクージをそれに応じて更新することができる。

【 0 0 6 5 】

データセットリンクージのマッピング

20

図5に示したように、データセットリンクージマッピング140は、構成要素名402、データセット名404、データセット形式406、フォーマット408、マスタデータセット場所410及びフラグ412を含む。構成要素名402はデータセット構成要素またはそのデータセットのインスタンスを表現するサブグラフである。データセット名404はこの構成要素によって表現されたデータセットを指し示す識別子である。データセット形式406は、データファイル、データベーステーブルまたは他の形式等のこのデータセットのインスタンスが分類されるカテゴリを示す。フォーマット408は、このデータセットのインスタンスがそのデータを表現するのに用いるフォーマットまたは配列である。マスタデータセット場所410は、このデータセットを追跡するデータ管理システム内のレコードを指し示す識別子である。最後に、フラグ412は、たとえばユーザがこのデータセットのインスタンスをデータ管理システムに適用できないものとして特定し、このインスタンスはリンクージの集合から排除されるべきであると特定した場合に、このデータセットのインスタンスは無視されるべきか否かを示す。

30

【 0 0 6 6 】

組込み規則

マップ100は、データ管理システムの標準的な慣例 (standard conventions) に従って動作する組込み規則100の集合を有する。データセット構成要素が組込み規則110に従う場合、マップはデータセット構成要素に対応するデータセットを最高位精度で特定することができる。規則の一つの実施態様において、データセットの候補を含む外部参照データベースは、データ管理システムによって用いられる標準化されたディレクトリ構造下で永続記憶装置内に配置されなければならない。さらに、パラメータに従って外部参照データセット構成要素をアクセスするグラフは、データ管理システムがまたアクセスし且つ分解できるパラメータを用いなければならない。さらに、データセット構成要素のフォーマットは、永続記憶装置に内にて利用可能であり、且つ、データ管理システムによってアクセスできるものでなければならない。データ管理システムに依り、他の組込み規則でもよい。

40

【 0 0 6 7 】

ユーザ定義規則

データセットの候補を特定するためにマップが使用する組込み規則に加えて、マップ100は、随意のユーザ定義規則120の集まりを有する。これらのユーザ定義規則120

50

は、どれがユーザの特定のデータ管理システムに適用可能であるかに応じて、ユーザによって使用可能または使用不能に設定され得る。一実施態様において、マップは六個の選択が自由なユーザ定義規則を有する。データベーステーブルという名目の一部の情報（たとえば、そのテーブルを定義したユーザについての情報等）がテーブルのアイデンティティを不明瞭にする場合、マップはデータベーステーブルという名目の一部の情報を無視することができる。さらに、マップは、データベーステーブルの名称からこの情報を除去することができる。さらに、マップは、データ管理システムに関連するデータセットに関係のないデータを含むことが知られている特定のカテゴリのデータファイルを、無視することができる。そのようなカテゴリは、データファイル形式またはデータファイルエクステンションであろう。さらに、マップは、参照をパラメータリスト内の特定のパラメータへと変換し、且つ、この参照をパラメータ自身の名称に置換できる。さらに、マップは、パラメータに対する参照を完全に除去することができる。ユーザは、また、マップが従う他の規則を生成することもできる。

10

【0068】

ヒューリスティックス

データセット候補を評価するための組込み規則及びユーザ定義規則に続いて、さらにマップ100はヒューリスティックス130の集合を用いる。ヒューリスティックス130は、マップに、所与のデータセット構成要素の特徴を解析させることを許容し、これらの特徴を既知のデータセットと比較させる。既知のデータセットに類似した特徴を有するデータセット構成要素はそのデータセットの一つのインスタンスである可能性が高い。一つの実施態様においては、マップは二つのヒューリスティックスを用いる。一つのヒューリスティックスは所与のデータセット構成要素のデータの特徴である。たとえば、データセット構成要素に関連するデータが、既知のデータセットに関連するデータと同一量のバイト及びレコードを有するときには、そのデータセット構成要素はそのデータセットのインスタンスの可能性が高い。さらに、データセット構成要素が、既知のデータセットと類似の名前または日付のクリエーションを有するならば、そのデータセット構成要素はそのデータセットのインスタンスである可能性が高い。第2のヒューリスティックスはデータセット構成要素のデータフォーマットである。データセット構成要素が既知のデータセットとデータフォーマットを共有していれば、そのデータセット構成要素はそのデータセットのインスタンスの可能性が高い。このヒューリスティックスは、複数の別個のデータセットが同一のデータフォーマットを用いている状況では、信頼性が低い。

20

30

【0069】

データセットのフォーマット及びマッピング

データソースを表現する各データセットは、データセットの各要素に対して、その要素がどのデータ形式を表しているかを示す関連データフォーマットを有する。たとえば、データベーステーブルのデータフォーマットは、所与のレコード内の各フィールドのデータ形式を示す。データ管理システム170は、データソース176、176a、176bを表現する各データセットに対する単一のデータフォーマット174、174a、174bを保持する。

【0070】

マップ100が新たなデータセット186を表現するデータセット構成要素182に遭遇した場合、マップ100は、データセット構成要素182のデータフォーマット184に基づいてデータ管理システムによって格納されるべき対応するデータフォーマットを生成する。

40

【0071】

データセット構成要素182がデータソース176を表現する既知のデータセットを表している幾つの場合には、データセット構成要素182は、データソース176を表現する既知のデータセットのデータフォーマット174とは異なるデータフォーマット184を有する。データ管理システム170は、データソース176を表すデータセットを、存在するかもしれないデータセットのインスタンスの数とは関係なく、単一のエンティテ

50

ィとして取扱う。その結果、このような状況が発生したときには、データ管理システム 170 はマップ 100 を頼りにして異なるフォーマット 174、184 を統合させる。一実施態様においては、マップは、ユーザ及びデータ管理システムの要求に依存して四つの異なる方法のうちの一つの方法にて各状況に対処することができる。ユーザ 162 は、各状況に対する統合のための四つの方法のうちのも一つを選択することができる。

【0072】

統合のための第 1 の方法の下では、マップ 100 は、データセット構成要素 182 のデータフォーマット 184 を、そのデータセットのマスタデータフォーマットとして用い、それに応じてデータ管理システム 170 を更新する。

【0073】

統合のための第 2 の方法の下では、マップ 100 は、現存のデータセットのデータフォーマット 174 を、そのデータセットのマスタデータフォーマットとして用い、それに応じてデータ管理システム 170 を更新する。

【0074】

統合のための第 3 の方法の下では、マップ 100 は、両方のデータフォーマットを保持し、各データフォーマットのフィールド間のマッピング 142 を発生する。図 6 に示したように、データフォーマットマッピング 142 は、データセットフォーマット 510 のフィールド 512 a、512 b、512 c のどれがデータセットインスタンス（たとえば、データセットの構成要素）のフォーマットのフィールド 522 a、522 b、522 c のどれに対応するかを示す。

【0075】

統合のための第 4 の方法の下では、マップは、いずれか一方のデータフォーマットとして作用できる新しい結合データフォーマットを発生する。

【0076】

一般的なコンピュータの実施形態

上述したデータセットマッピングのアプローチは、コンピュータ上で実行されるソフトウェアを用いて実現できる。たとえば、そのソフトウェアは、一つ以上のプログラムされたまたはプログラム可能なコンピュータシステム（分散型、クライアント/サーバ型、または、グリッド型等の種々アーキテクチャ型である）上で実行される一つ以上のコンピュータプログラムにおける手続を形成する。その各コンピュータシステムは、少なくとも一つのプロセッサ、少なくとも一つの記憶システム（揮発性メモリ、及び不揮発性メモリ、及びまたは記憶素子を含む。）、少なくとも一つの入力装置またはポート、及び、少なくとも一つの出力装置またはポートを含む。そのソフトウェアは、たとえば、データフローグラフの設計及び構成に関係した他のサービスを提供するより大きなプログラムの一つ以上のモジュールを形成できる。グラフのノード及び要素は、コンピュータ読み出し可能媒体に格納されたデータ構造、または、データ貯蔵庫（保管所）に格納されたデータモデルに合致する他の系統的なデータとして実現され得る。

【0077】

そのソフトウェアは、汎用または専用のプログラマブルコンピュータによって読み出し可能な CD-ROM 等の記憶媒体上に提供され、或いは、そのソフトウェアが実行されるコンピュータにネットワークの通信媒体を用いて（伝播信号に符号化されて）配信される。すべての機能は、専用のコンピュータ上で、または、コプロセッサのような専用ハードウェアを用いて実行され得る。ソフトウェアは、ソフトウェアによって特定される異なる計算の部分が異なるコンピュータによって実行されるという分散方式によって実行され得る。好ましくは、上述の各コンピュータプログラムは、汎用または専用プログラマブルコンピュータによって読み出し可能な記憶媒体またはデバイス（たとえば、固体メモリまたは媒体、もしくは磁氣的または光学的媒体）に格納され或いはダウンロードされ、記憶媒体またはデバイスがコンピュータシステムによって読み出されて上述した手順を実行するときに、コンピュータを構築及び動作させる。また、本発明のシステムは、コンピュータプログラムとともに構築されたコンピュータ読み出し可能記憶媒体として実現されること

10

20

30

40

50

が考えられる。ここで、そのように構築された記憶媒体は、上述の機能を実行するために特別且つ予め定義された様式でコンピュータシステムを作動させる。

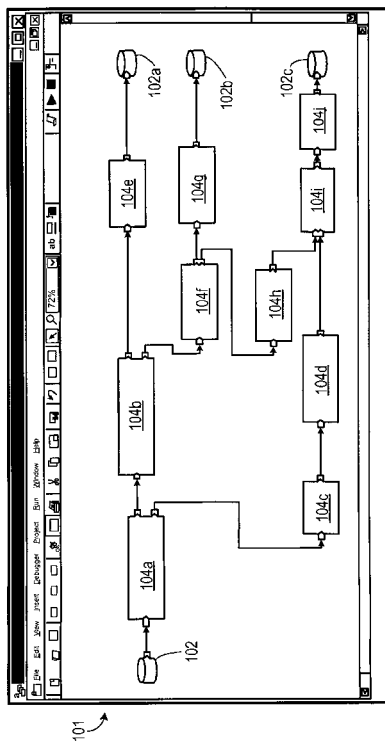
【0078】

多くの本発明の実施例を記載した。しかしながら、本発明の精神及び範囲から逸脱することなく多種の変更がなされ得ることが理解されよう。たとえば上述のいくつかのステップは、順序が独立であり、従って上述の順序と異なる順序で実行され得る。

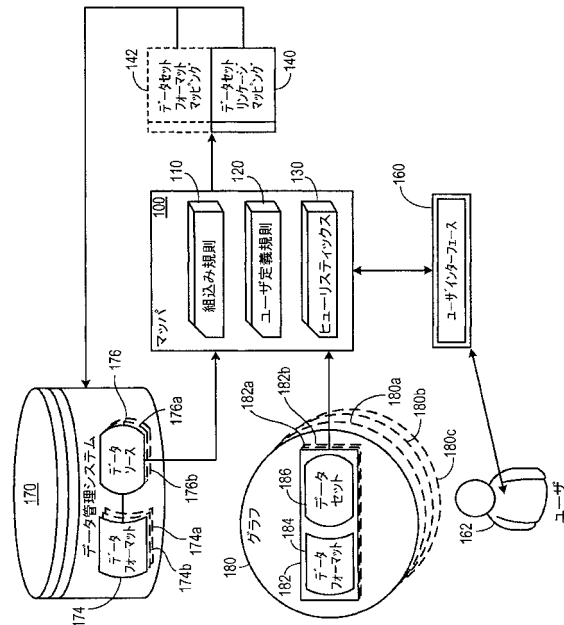
【0079】

上述の記述は説明のためであり、添付された本発明の特許請求の範囲によって定義される発明の範囲を限定する意図はないことが理解されるべきである。たとえば、上述の機能ステップの多くは、全体の処理に実質的に影響を及ぼすことなく異なる順序で実行され得る。他の実施例は特許請求の範囲の請求項の範囲内にある。

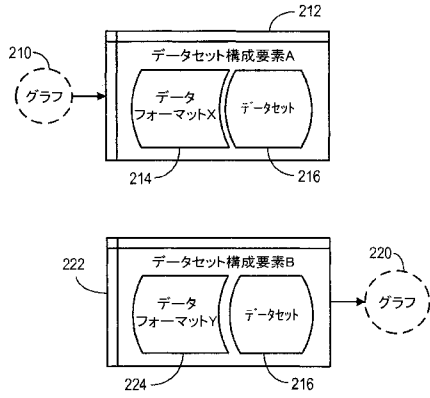
【図1】



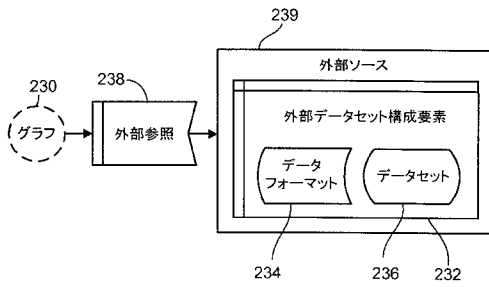
【図2】



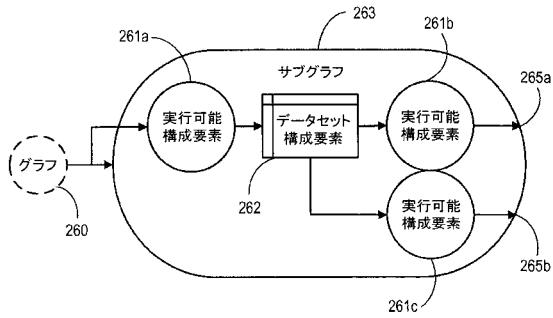
【図3A】



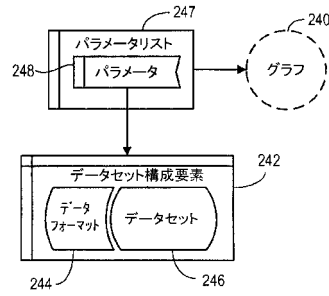
【図3B】



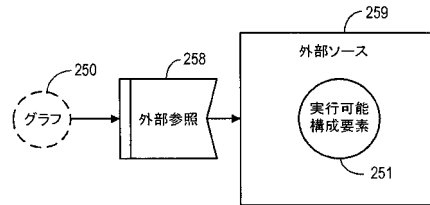
【図3E】



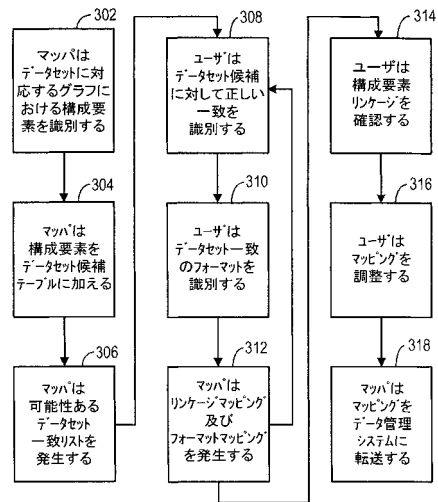
【図3C】



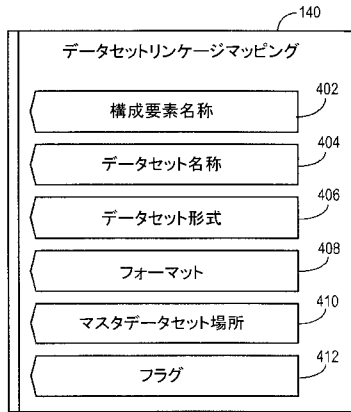
【図3D】



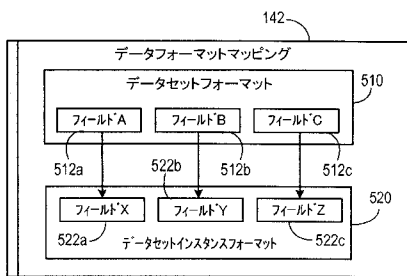
【図4】



【図5】



【図6】



フロントページの続き

(72)発明者 ワイス アダム

アメリカ合衆国 02420 マサチューセッツ州 レキシントン ローソン アベニュー 15

審査官 岡本 俊威

(56)参考文献 特開平11-143755(JP,A)

米国特許第07080088(US,B1)

(58)調査した分野(Int.Cl., DB名)

G06T 1/00

G06F 12/00