

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第6674904号
(P6674904)

(45) 発行日 令和2年4月1日 (2020. 4. 1)

(24) 登録日 令和2年3月11日 (2020. 3. 11)

(51) Int. Cl.

F I

G O 6 F 9/50 (2006. 01)

G O 6 F 9/48 (2006. 01)

G O 6 F 8/00 (2018. 01)

G O 6 F 9/50 1 2 O A

G O 6 F 9/48 3 7 O

G O 6 F 8/00

請求項の数 37 (全 17 頁)

(21) 出願番号	特願2016-567355 (P2016-567355)	(73) 特許権者	509123208
(86) (22) 出願日	平成27年5月22日 (2015. 5. 22)		アビニシオ テクノロジー エルエルシー
(65) 公表番号	特表2017-522630 (P2017-522630A)		アメリカ合衆国 O 2 4 2 1 マサチュー
(43) 公表日	平成29年8月10日 (2017. 8. 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/US2015/032193		トリート 2 0 1
(87) 国際公開番号	W02015/183738	(74) 代理人	100079108
(87) 国際公開日	平成27年12月3日 (2015. 12. 3)		弁理士 稲葉 良幸
審査請求日	平成30年3月15日 (2018. 3. 15)	(74) 代理人	100109346
(31) 優先権主張番号	62/004, 406		弁理士 大貫 敏史
(32) 優先日	平成26年5月29日 (2014. 5. 29)	(74) 代理人	100117189
(33) 優先権主張国・地域又は機関	米国 (US)		弁理士 江口 昭彦
(31) 優先権主張番号	14/470, 501	(74) 代理人	100134120
(32) 優先日	平成26年8月27日 (2014. 8. 27)		弁理士 内藤 和彦
(33) 優先権主張国・地域又は機関	米国 (US)		

最終頁に続く

(54) 【発明の名称】 ワークロード自動化およびデータシステム分析

(57) 【特許請求の範囲】

【請求項 1】

コンピュータシステムによって、複数のジョブの実行順序を示すジョブ依存性情報と、データ系統情報との間のリンクを識別することにより、前記ジョブ依存性情報と前記データ系統情報との結合に基づく依存性情報を判断することと、

前記依存性情報に基づいて、前記データ系統情報により指定されるデータソースにおける前記ジョブ依存性情報により指定されるジョブの実行予定についての変更の影響を判断することであって、

前記ジョブ依存性情報に含まれる、前記ジョブのジョブスケジューリング情報を識別するクエリを受信することと、

前記ジョブスケジューリング情報とデータ系統情報との間の前記リンクを識別することと、

前記ジョブスケジューリング情報に基づいて、前記データ系統情報により指定される前記データソースへの影響を識別することと、

によって前記変更の影響を判断することと、
を含む、コンピュータにより実施される方法。

【請求項 2】

ワークロード自動化システムに関連付けられたワークロードリポジトリから、スケジューリング情報を取得することと、

前記スケジューリング情報を変換することと、

データ記憶に前記スケジューリング情報を記憶させることであって、前記データ記憶が、前記データ系統情報を記憶していることと、により、ジョブ依存性情報を取得すること、を更に含む、請求項 1 に記載の方法。

【請求項 3】

前記データ系統情報および前記ジョブ依存性情報によって参照されるデータソースを識別することにより、リンクを生成することを更に含む、請求項 1 に記載の方法。

【請求項 4】

前記データソースを識別することが、前記データ系統情報および前記ジョブ依存性情報内の同一の名前によって参照されるデータソースを識別することを含む、請求項 3 に記載の方法。

【請求項 5】

前記データソースを識別することが、ユニフォームリソースロケータを使用してデータソースを識別することを含む、請求項 3 に記載の方法。

【請求項 6】

前記データソースを識別することが、データベース、テーブルスペース、およびテーブル名を使用して、リレーショナルデータベーステーブルを識別することを含む、請求項 3 に記載の方法。

【請求項 7】

前記ジョブ依存性情報および前記データ系統情報によって参照される実行可能プログラムを識別することにより、リンクを生成することを更に含む、請求項 1 に記載の方法。

【請求項 8】

前記実行可能プログラムを識別することが、前記実行可能プログラムに提供されるパラメータに少なくとも部分的に基づいて、実行可能プログラムを識別することを含む、請求項 7 に記載の方法。

【請求項 9】

前記実行可能プログラムを識別することが、永続的なデータ記憶上の前記実行可能プログラムの位置に基づいて、実行可能プログラムを識別することを含む、請求項 7 に記載の方法。

【請求項 10】

前記判断された影響を指定する報告を生成することを更に含む、請求項 1 に記載の方法。

【請求項 11】

前記実行予定についての変更により前記データソースが影響を受ける時に、1 以上の警報を生成することを更に含む、請求項 1 に記載の方法。

【請求項 12】

前記判断された影響が、前記実行予定についての変更により前記データソースの中の特定のデータセットが不正確になりそうなことを明示する時に、1 以上の警報を生成することを更に含む、請求項 1 に記載の方法。

【請求項 13】

1 つまたは複数のコンピュータによる実行時に、前記 1 つまたは複数のコンピュータに動作を行わせるコンピュータプログラム命令を記憶した非一時的コンピュータ記憶媒体であって、前記動作が、

複数のジョブの実行順序を示すジョブ依存性情報と、データ系統情報との間のリンクを識別することにより、前記ジョブ依存性情報と前記データ系統情報との結合に基づく依存性情報を判断することと、

前記依存性情報に基づいて、前記データ系統情報により指定されるデータソースにおける前記ジョブ依存性情報により指定されるジョブの実行予定についての変更の影響を判断することであって、

前記ジョブ依存性情報に含まれる、前記ジョブのジョブスケジューリング情報を識別するクエリを受信することと、

10

20

30

40

50

前記ジョブスケジューリング情報とデータ系統情報との間の前記リンクを識別することと、

前記ジョブスケジューリング情報に基づいて、前記データ系統情報により指定される前記データソースへの影響を識別することと、

によって前記変更の影響を判断することと、
を含む、非一時的コンピュータ記憶媒体。

【請求項 1 4】

ワークロード自動化システムに関連付けられたワークロードリポジトリから、スケジューリング情報を取得することと、

前記スケジューリング情報を変換することと、

データ記憶に前記スケジューリング情報を記憶させることであって、前記データ記憶が、前記データ系統情報を記憶していることと、により、ジョブ依存性情報を取得することを更に含む、請求項 1 3 に記載の媒体。

【請求項 1 5】

前記データ系統情報および前記ジョブ依存性情報によって参照されるデータソースを識別することにより、リンクを生成することを更に含む、請求項 1 3 に記載の媒体。

【請求項 1 6】

前記データソースを識別することが、前記データ系統情報および前記ジョブ依存性情報内の同一の名前によって参照されるデータソースを識別することを含む、請求項 1 5 に記載の媒体。

【請求項 1 7】

前記データソースを識別することが、ユニフォームリソースロケータを使用してデータソースを識別することを含む、請求項 1 5 に記載の媒体。

【請求項 1 8】

前記データソースを識別することが、データベース、テーブルスペース、およびテーブル名を使用して、リレーショナルデータベーステーブルを識別することを含む、請求項 1 5 に記載の媒体。

【請求項 1 9】

前記ジョブ依存性情報および前記データ系統情報によって参照される実行可能プログラムを識別することにより、リンクを生成することを更に含む、請求項 1 3 に記載の媒体。

【請求項 2 0】

前記実行可能プログラムを識別することが、前記実行可能プログラムに提供されるパラメータに少なくとも部分的に基づいて、実行可能プログラムを識別することを含む、請求項 1 9 に記載の媒体。

【請求項 2 1】

前記実行可能プログラムを識別することが、永続的なデータ記憶上の前記実行可能プログラムの位置に基づいて、実行可能プログラムを識別することを含む、請求項 1 9 に記載の媒体。

【請求項 2 2】

前記判断された影響を指定する報告を生成することを更に含む、請求項 1 3 に記載の媒体。

【請求項 2 3】

前記実行予定についての変更により前記データソースが影響を受ける時に、1 以上の警報を生成することを更に含む、請求項 1 3 に記載の媒体。

【請求項 2 4】

前記判断された影響が、前記実行予定についての変更により前記データソースの中の特定のデータセットが不正確になりそうなことを明示する時に、1 以上の警報を生成することを更に含む、請求項 1 3 に記載の媒体。

【請求項 2 5】

1 つまたは複数のコンピュータと、

10

20

30

40

50

前記 1 つまたは複数のコンピュータによる実行時に、前記 1 つまたは複数のコンピュータに動作を行わせるように実行可能な命令を記憶する、1 つまたは複数の記憶装置と、を備えるシステムであって、前記動作が、

複数のジョブの実行順序を示すジョブ依存性情報と、データ系統情報との間のリンクを識別することにより、前記ジョブ依存性情報と前記データ系統情報との結合に基づく依存性情報を判断することと、

前記依存性情報に基づいて、前記データ系統情報により指定されるデータソースにおける前記ジョブ依存性情報により指定されるジョブの実行予定についての変更の影響を判断することであって、

前記ジョブ依存性情報に含まれる、前記ジョブのジョブスケジューリング情報を識別するクエリを受信することと、

前記ジョブスケジューリング情報とデータ系統情報との間の前記リンクを識別することと、

前記ジョブスケジューリング情報に基づいて、前記データ系統情報により指定される前記データソースへの影響を識別することと、

によって前記変更の影響を判断することと、を含むシステム。

【請求項 2 6】

ワークロード自動化システムに関連付けられたワークロードリポジトリから、スケジューリング情報を取得することと、

前記スケジューリング情報を変換することと、

データ記憶に前記スケジューリング情報を記憶することであって、前記データ記憶が、前記データ系統情報を記憶することと、により、ジョブ依存性情報を取得することを更に含む、請求項 2 5 に記載のシステム。

【請求項 2 7】

前記データ系統情報および前記ジョブ依存性情報によって参照されるデータソースを識別することにより、リンクを生成することを更に含む、請求項 2 5 に記載のシステム。

【請求項 2 8】

前記データソースを識別することが、前記データ系統情報および前記ジョブ依存性情報内の同一の名前によって参照されるデータソースを識別することを含む、請求項 2 7 に記載のシステム。

【請求項 2 9】

前記データソースを識別することが、ユニフォームリソースロケータを使用してデータソースを識別することを含む、請求項 2 7 に記載のシステム。

【請求項 3 0】

前記データソースを識別することが、データベース、テーブルスペース、およびテーブル名を使用して、リレーショナルデータベーステーブルを識別することを含む、請求項 2 7 に記載のシステム。

【請求項 3 1】

前記ジョブ依存性情報および前記データ系統情報によって参照される実行可能プログラムを識別することにより、リンクを生成することを更に含む、請求項 2 5 に記載のシステム。

【請求項 3 2】

前記実行可能プログラムを識別することが、前記実行可能プログラムに提供されるパラメータに少なくとも部分的に基づいて、実行可能プログラムを識別することを含む、請求項 3 1 に記載のシステム。

【請求項 3 3】

前記実行可能プログラムを識別することが、永続的なデータ記憶上の前記実行可能プログラムの位置に基づいて、実行可能プログラムを識別することを含む、請求項 3 1 に記載のシステム。

【請求項 3 4】

10

20

30

40

50

前記判断された影響を指定する報告を生成することを更に含む、請求項 2 5 に記載のシステム。

【請求項 3 5】

前記実行予定についての変更により前記データソースが影響を受ける時に、1 以上の警報を生成することを更に含む、請求項 2 5 に記載のシステム。

【請求項 3 6】

前記判断された影響が、前記実行予定についての変更により前記データソースの中の特定のデータセットが不正確になりそうなことを明示する時に、1 以上の警報を生成することを更に含む、請求項 2 5 に記載のシステム。

【請求項 3 7】

複数のジョブの実行順序を示すジョブ依存性情報と、データ系統情報との間のリンクを識別することにより、前記ジョブ依存性情報と前記データ系統情報との結合に基づく依存性情報を判断する手段と、

前記依存性情報に基づいて、前記データ系統情報により指定されるデータソースにおける前記ジョブ依存性情報により指定されるジョブの実行予定についての変更の影響を判断する手段であって、

前記ジョブ依存性情報に含まれる、前記ジョブのジョブスケジューリング情報を識別するクエリを受信する手段と、

前記ジョブスケジューリング情報とデータ系統情報との間の前記リンクを識別する手段と、

前記ジョブスケジューリング情報に基づいて、前記データ系統情報により指定される前記データソースへの影響を識別する手段と、

によって前記変更の影響を判断する手段と、を備えるシステム。

【発明の詳細な説明】

【技術分野】

【0001】

優先権出願

本出願は、2014年5月29日に出願された「WORKLOAD AUTOMATION AND DATA LINEAGE ANALYSIS」と題する米国仮特許出願第62/004,406号、および、2014年8月27日に出願された「WORKLOAD AUTOMATION AND DATA LINEAGE ANALYSIS」と題する米国特許出願第14/470,501号に基づき、優先権を主張する。両者の内容全体が、参照により本明細書に組み込まれる。

【背景技術】

【0002】

ワークロードの自動化は、概して、一般にジョブをセットアップするプロセスを指し、そのように人間との対話処理なしで完結するように実行されることができる。全ての入力パラメータが、スクリプト、コマンドライン引数、ワークフロー自動化システム、制御ファイル、またはジョブ制御言語によって事前定義されている。ジョブは、利用可能な処理リソースと、予め定義された依存性に基づいてスケジュールされる。

【0003】

データ系統は、データの発生元、およびデータがどこに移動し、データが経時的にどのように変化するかを説明する。この用語は、多様なプロセスを通じて進むにつれて、データがどうなるかを説明することもできる。データ系統は、情報がどのように使用されるかを分析し、特定の目的に合った情報の主要ピットを追跡する取り組みに役立つことができる。

【発明の概要】

【0004】

概略的な態様1では、方法は、ジョブ依存性情報を取得する動作を含み、ジョブ依存性情報が、複数のジョブの実行順序を指定する。方法は、データ記憶と変換との間の依存関

10

20

30

40

50

係を識別するデータ系統情報を取得する動作も含み、少なくとも1つの変換が、第1のデータ記憶からデータを受け取り、第2のデータ記憶に対してデータを生じさせる。方法は、ジョブ依存性情報とデータ系統情報との間にリンクを生成する動作も含む。方法は、ジョブ依存性情報、生成されたリンク、およびデータ系統情報に基づいて、複数のアプリケーションのうちのアプリケーションの実行予定についての変更の影響を判断する動作も含む。

【0005】

本態様の他の実施形態は、対応するコンピュータシステム、装置、および1つまたは複数のコンピュータ記憶装置に記録されたコンピュータプログラムを含み、それぞれが、方法の動作を実行するように構成される。1つまたは複数のコンピュータのシステムは、作
10
動中にシステムに動作を実行させるシステムにインストールされた、ソフトウェア、ファームウェア、ハードウェア、またはそれらの組合せを有することによって、特定の動作を実行するように構成され得る。1つまたは複数のコンピュータプログラムは、データ処理装置による実行時に装置に動作を実行させる命令を含むことによって、特定の動作を実行するように構成され得る。

【0006】

方法は、態様1による態様2を含み、態様2では、ジョブ依存性情報を取得することが、ワークロード自動化システムに関連付けられたワークロードリポジトリから、スケジューリング情報を取得することと、スケジューリング情報を変換することと、データ記憶に
20
スケジューリング情報を記憶させることとであって、データ記憶が、データ系統を記憶していることと、を含む。方法は、態様1または2による態様3を含み、態様3では、リンクを生成することが、データ系統情報およびジョブ依存性情報によって参照されるデータソースを識別することを含む。方法は、態様1、2または3による態様4を含み、態様4では、データソースを識別することが、データ系統情報およびジョブ依存性情報内の同一の名前によって参照されるデータソースを識別することを含む。方法は、態様1、2、3または4による態様5を含み、態様5では、データソースを識別することが、ユニフォームリソースロケータを使用してデータソースを識別することを含む。方法は、態様1、2、3、4または5による態様6を含み、態様6では、データソースを識別することが、データベース、テーブルスペース、およびテーブル名を使用して、リレーショナルデータベ
30
ーステーブルを識別することを含む。方法は、態様1、2、3、4、5または6による態様7を含み、態様7では、リンクを生成することが、ジョブ依存性情報およびデータ系統情報によって参照される実行可能プログラムを識別することを含む。方法は、態様1、2、3、4、5、6または7による態様8を含み、態様8では、実行可能プログラムを識別することが、実行可能プログラムに提供されるパラメータに少なくとも部分的に基づいて、実行可能プログラムを識別することを含む。方法は、態様1、2、3、4、5、6、7または8による態様9を含み、態様9では、実行可能プログラムを識別することが、永続的なデータ記憶上の実行可能プログラムの位置に基づいて、実行可能プログラムを識別することを含む。方法は、態様1、2、3、4、5、6、7、8または9による態様10を含み、態様10では、影響を判断することが、ジョブスケジューリングデータを識別するク
40
エリを受信することと、ジョブスケジューリングデータとデータ系統情報との間のリンクを識別することと、スケジューリングデータに基づいて、データ系統情報への影響を識別することと、を含む。

【発明の効果】

【0007】

本明細書に記載される主題の特定の実施形態は、以下の利点のうちの1つまたは複数を実現するように実装可能である。データ処理システムの全体像を考察することができる。データ依存性を説明するデータ系統情報を、スケジューリング依存性を説明するワークフロー自動化情報と結合することができる。データ系統情報またはジョブスケジューリング情報のいずれかの変更の影響を判断することができる。これにより、技術プロセスの開発者や管理者が、ワークフローをより効率的かつ中断の少ない方法で監視し、調整すること
50

に役立てることができる。全体として、プロセス障害、リソース消費、およびデータ処理期間のそれぞれが、それによって減少され得る。

【図面の簡単な説明】

【0008】

【図1】例示的なスケジューリング図である。

【図2】例示的なデータ系統図である。

【図3】スケジューリング情報とデータ系統情報を統合する、例示的なシステムである。

【図4】データ分析技術が使用され得るデータ処理システムの例を示す。

【図5】スケジューリング情報への変更の影響を識別する、例示的な処理のフローチャートである。

10

【図6】相互に関連付けられた、依存性情報およびデータ系統情報の例を示す。

【発明を実施するための形態】

【0009】

ジョブ依存性情報およびデータ系統情報は、結合されて事業の状態の全体像を提供することができる。従来、ジョブ依存性情報およびデータ系統情報は、異なる情報システムおよびデータベースにまたがって細分化されている。ジョブ依存性情報は、異なるジョブまたはタスクの実行の間に確立されている順序を説明する。データ系統情報は、データソースおよびデータシンクが、事業全体にわたってどのように関係しているかを説明する。ユーザは、特定のジョブが遅延しそうかどうか、または、どの報告もしくはデータシンクが影響を受けそうか、などの疑問に答えたいと望むことがある。本明細書で説明するシステムは、これらの全く異なるデータソースを統合する。

20

【0010】

ジョブ依存性情報は、ワークロード自動化プログラムまたはジョブスケジューリングプログラムから取得されることができる。ワークロード自動化プログラムまたはジョブスケジューリングプログラムは、複雑な依存性を伴うワークロードタイプの様々なセットを調整する。概して、ジョブ依存性情報は、異なるタスクが実行されるべき順序を定義する。アプリケーションのスケジューリングは、典型的には、データの依存性を考慮に入れておらず、それは、本来データを意識したものではない。アプリケーションのスケジューリングは、単に、異なるタスクが実行し得る順序を定めているにすぎない。この順序付けは、データ依存性に基づく可能性があるが、リソース割り当て、全体の実行時間、および他の効率の最適化に基づく可能性もある。ワークロード自動化システムにおけるタスクは、データフローグラフ、Javaプログラム、ファイル転送コマンド、ビジネススイートソフトウェア統合、ウェブサービスアクセス、メッセージング、または任意の他の実行可能なプロセスを含んでもよい。ユーザは、スケジュール変更の影響、例えば、ジョブが遅延するかどうか、を判断したい場合がある。ワークロード自動化システムは、システム内で定義されている詳細を見るように、その機能が制限されている。

30

【0011】

一方、データ系統情報は、データがシステムによって処理される順序を識別する。概して、データ系統情報は、データの発生元、および、データ処理アプリケーションの間にデータがどこに移動するか、またはデータがどのように変換されるかを含み、データのライフサイクルを説明する。データ系統情報は、多様なプロセスによってデータが変換されるにつれて、データに何が起こるかを説明する。概して、データ系統情報の分析は、情報がどのように使用されるかを識別し、特定の目的に合った情報の主要部分を追跡するために使用される。ジョブ依存性情報をデータ系統リポジトリに統合することによって、プロセスおよびデータのより堅牢な視界を展開することができる。

40

【0012】

ジョブ依存性情報をデータ系統リポジトリに統合することによって、プロセスおよびデータのより堅牢な視界を展開することができる。

【0013】

ジョブ依存性情報は、ワークロード自動化ツールから抽出されることができ、データ系

50

統情報は、データ系統ツールから抽出されることができる。情報は、共に結合され、以降のアクセスのために、共通のリポジトリに記憶されることができる。

【0014】

図1は、ジョブ「P s i」についての例示的なスケジューリング図100である。スケジューリング図100は、ワークロード自動化システムにおける、スケジューリング動作の例である。スケジューリング図は、ジョブ間のジョブ依存性を示している。多くの実装において、コンポーネント間により大きな相互関係を有する、より複雑な図が存在する。現在のスケジューリング図100は、例示の目的で使用される。ジョブは、後続ジョブの開始前に、先行ジョブが完了している必要がある、階層的順序で定義される。この図では、ジョブは、有向矢印で接続されている。矢印は、先行ジョブから後続ジョブに向かって
10

例えば、「スクリプト コマンド 1」ジョブ102は、「データベース」ジョブ106または「ファイル監視」ジョブ110が開始し得る前に完了しなければならない。同様に、「データベース」ジョブ106および「ファイル監視」ジョブ110は、「FTP」ジョブ112が開始し得る前に完了しなければならない。「実行」ジョブ114は、「ファイル監視」ジョブ110が完了した後で、開始し得る。最後に、「モニタ完了」ジョブ116は、「FTP」ジョブ112および「実行」ジョブ114が完了した後にのみ、実行することができる。

【0015】

ワークロード自動化システムは、様々なジョブについての情報を収集する。例えば、「スクリプト コマンド 1」ジョブ102は、ジョブを定義し、説明する属性104を有
20

する。この例では、属性104は、実行されるジョブの種類を示すジョブタイプと、ジョブの名前を示すジョブ名と、スクリプトの位置を定義するファイルパスと、実行されるスクリプトの名前を示すファイル名と、スクリプトを実行すべきユーザ名を示す実行者と、ジョブの現在の状態（例えば、保留中、実行中、完了済み、失敗）を示す状態と、現在のジョブが完了後にのみ実行できるジョブを示す後続と、実行されるステップを定義するスクリプトと、を含む。

【0016】

他の種類のジョブは、異なる属性を含んでもよい。例えば、「データベース」ジョブ106は、属性108を有する。これらの属性は、SQLコマンド（ここでは、「select
us.order, us.order_amount from ne_production」）、先行ジョブのリスト（ここで
30

は、スクリプト コマンド 1）、および後続ジョブのリスト（ここでは、FTP）を含むが、限定はされない。

【0017】

同様に、「実行」ジョブ114は、実行されるべきプログラムの名前、例えば、「T r a n s f o r m A . e x e」を識別するパラメータ118を含んでもよい。「FTP」ジョブ112は、ファイルおよびファイル転送動作の宛先を識別するパラメータ120を含
40

んでもよい。例えば、パラメータ120は、B r a z i l F e e d . d a t ファイルが、s e r v e r . c o m に転送予定であることを識別する。本明細書において、識別されるパラメータは、単なる例示にすぎない。他のパラメータは、ジョブスケジューリング情報によって定義され、含まれ得る。

【0018】

ワークロード自動化システムは、図示していない、他のジョブに関連してジョブP s i をスケジューリングしてもよい。例えば、ジョブP s i は、ジョブZ e t a（または、図示していない、何らかの他のジョブ）の後に行われるようにスケジューリングされてもよい。ワークロード自動化システムは、ジョブ間のスケジュールを、リソース管理、依存性の報告、利用可能時間、優先度、または他の制約に基づいて判断してもよい。

【0019】

図2は、例示的なデータ系統図200である。データ系統は、概して、データの発生元、ならびにデータがどこへ移動するか、およびデータがどのように変換され処理されるかを含む、データライフサイクルとして定義される。この用語は、多様なプロセスを通じて
50

進むにつれて、データがどうなるかを説明することもできる。データ系統は、情報がどのように使用されるかを分析し、特定の目的に合った情報の主要ビットを追跡する取り組みに役立つことができる。概して、データ系統図は、データソース、データシンクおよび変換の間の関係性を示す図である。各変換は、1つまたは複数のデータソース（例えば、入力データ）を含み、1つまたは複数のデータシンクにデータ（例えば、出力データ）を生じさせることができる。データ系統情報内のそれぞれのデータソース、データシンク、および変換は、本明細書では、まとめてデータ系統要素と呼ばれるものとする。

【0020】

この例では、データソース「U . S . F e e d」202は、「変換A」204の変換にデータを提供する。変換A 204は、「U . S . F e e d」202によって提供されたデータに対して動作を実行し、「中間データセット1」206のデータ記憶に結果を記憶させる。データソース「M e x i c o F e e d」208および「B r a z i l F e e d」214は、「変換C」210の変換にデータを提供する。「変換C」210の変換は、「M e x i c o F e e d」208および「B r a z i l F e e d」214によって提供されたデータに対して動作を実行し、「中間データセット2」212のデータ記憶に結果を記憶させる。データソースは、例えば、単層ファイル、リレーショナルデータベース、オブジェクトデータベース、または、コンピュータシステムにデータを記憶させるための任意の他の機構であってもよい。例えば、「B r a z i l F e e d」214は、「B r a z i l F e e d . d a t」などのファイルであってもよい。変換は、データを操作することが可能な実行可能プログラムであってもよい。例えば、仮想マシン内で実行されるJavaプログラム、実行ファイル、データフローグラフなどであってもよい。例えば、「変換A」204の変換は、「T r a n s f o r m A . e x e」という名前の実行ファイルであってもよい。

【0021】

「中間データセット1」206のデータ記憶および「中間データセット2」212のデータ記憶は、「変換B」216の変換にデータを提供する。「変換B」216の変換は、「中間データセット1」206から提供されたデータを使用し、「中間データセット2」212は、「出力報告」218のデータ記憶に結果を記憶させる。

【0022】

データ系統に記憶されている情報は、データの異なる部分が、データの他の部分にどのように影響を与えるかを識別することができる。例えば、「U . S . F e e d」データソースは、注文、および注文毎の額を含んでもよい。「変換A」204は、地域に基づいて、例えば、ニューイングランド、東部諸州、南部、中西部、平原州で発生した注文などによって、データを集約してもよい。データ系統情報は、「U . S . F e e d」202からの額のフィールドが、「中間データセット1」206の「地域合計」フィールドに集約されると識別してもよい。

【0023】

スケジューリングデータをデータ系統データと結合することによってのみ得られる、いくつかの情報がある。例えば、図1の「データベース ジョブ」106が、図2の「U . S . F e e d」202を生成する場合において、その後ジョブ106が遅延し、または実行に失敗する場合、出力報告218は、遅延するか、または不正確なものとなる。データ系統情報およびジョブスケジューリング情報の両方を検討することなしに、これらの関係性が生じることはない。上記の例を参照すると、「ジョブZ e t a」（ジョブP s iに先行するものとして上述された）が遅延する場合、出力報告218が遅延するか、または不正確なものとなる恐れがあるため、さらに関係性がより複雑になる可能性がある。

【0024】

図3は、スケジューリング情報およびデータ系統情報を統合する例示的システムである。スケジューリングリポジトリ302a~bからのデータは、データ系統リポジトリ306にインポートされることができる。スケジューリングリポジトリ302a~bは、ワークロード自動化システム、例えば、C O N T R O L - M、T I V O L I、T W S A U T O

10

20

30

40

50

S Y S、C A - 7 などに関連付けられたデータリポジトリであってもよい。ワークロード自動化システムのそれぞれについてのデータは、異なるフォーマットで記憶されてもよい。変換コンポーネント 3 0 4 a ~ b は、スケジューリングリポジトリ 3 0 2 a ~ b に記憶されたデータを、結合されたりポジトリ 3 0 6 の記憶用の共通データフォーマットに変換するために使用され得る。いくつかの実装では、変換コンポーネントは、例えば、データフローグラフの計算環境で実行するデータフローグラフであってもよい。

【 0 0 2 5 】

データ系統リポジトリ 3 1 4 からのデータは、結合されたりポジトリに記憶されることもできる。データ系統情報は、リポジトリ内に挿入される前に、変換コンポーネント 3 1 6 によって変換されてもよい。例えば、データのデータタイプは、リポジトリの予定されるデータフォーマットに従うために、あるタイプから別のタイプに変換されてもよい。さらに、データ構造は、例えば、データをジョブスケジューリング情報と効率的に統合するために、データ系統データ構造を単純化することを含み、変更されてもよい。

【 0 0 2 6 】

ジョブスケジューリング情報は、収集され、結合されたりポジトリに統合されることができる。ジョブスケジューリング情報は、結合されたりポジトリで受け入れ可能なフォーマットに修正されることができる。例えば、データの特定のフィールドのフォーマットが、変更されてもよい。異なるデータオブジェクト間の関係性は、機能的に同一または異なる形式に変更されてもよい。ジョブスケジューリング情報を統合することは、以前の大量のワークロードスケジューリングデータから、古いまたは期限を超過した情報を識別すること、および上書きすること、または保管することを含むことができる。ジョブ依存性情報とデータ系統情報とが、共に結合されリンクされる。ジョブ依存性情報は、情報に関連付けられた属性またはパラメータに基づいて、データ系統情報にリンクされてもよい。例えば、ジョブおよびデータ系統要素は、同一の実行ファイル（例えば、上述の「TransformA.exe」）を参照してもよい。実行ファイルは、完全修飾識別子に基づいて識別されることができる。完全修飾識別子は、完全パスを含んでもよく、即ち、コンピュータ、およびハードドライブなどの永続的な記憶装置上の位置が、識別されてもよい。完全修飾識別子は、実行ファイルに提供される任意のパラメータを含んでもよい。同様に、ジョブおよびデータ系統要素は、同一のデータ記憶を参照してもよい。例えば、上記の図では、図 1 の F T P ジョブ 1 1 2 および B r a z i l F e e d データソース 2 1 4 は、「B r a z i l F e e d . d a t」ファイルを参照する。データソースは、完全修飾識別子に基づいて識別され得る。例えば、完全修飾識別子は、単層ファイルを識別する完全パスもしくははユニフォームリソースロケータ（URL）、または、サーバ、データベース、テーブルスペース、リレーショナルデータベース内のテーブル名を識別する情報であってもよい。これらの、または他の共通の要素が識別される時、プロセスは、ジョブ依存性情報とデータ系統情報との間にリンクを生成することができる。

【 0 0 2 7 】

いくつかの実装では、ジョブスケジューリング情報は、定期的に（例えば、週 1 回、1 日 1 回、1 時間 1 回、など）、結合されたりポジトリ 3 0 6 に統合されることができる。いくつかの実装では、ジョブスケジューリング情報に対する変更が、統合プロセスのトリガとなって、ほぼリアルタイムで情報を統合してもよい。例えば、データベースのトリガは、変更が検出された時に統合プロセスを開始してもよい。あるいは、ジョブスケジューリングシステムにおけるコールバック機構が、統合プロセスを開始させてもよい。

【 0 0 2 8 】

結合されたりポジトリ 3 0 6 は、グラフベースのアプリケーションの開発および実行、ならびにグラフベースのアプリケーションと他のシステム（例えば、他のオペレーティングシステム）との間のメタデータの交換を支援するように設計された、拡張性のあるオブジェクト指向データベースシステムであることが好ましい。結合されたりポジトリ 3 0 6 は、ドキュメンテーション、レコードフォーマット（例えば、テーブル内のレコードのフィールドおよびデータタイプ）、変換機能、グラフ、ジョブ、ならびにモニタリング情報

10

20

30

40

50

を含む、あらゆる種類のメタデータのための記憶システムである。

【0029】

結合されたりポジトリ306は、コンピューティングシステムによって処理されるべき実際のデータを表す、データオブジェクトを記憶することもできる。

【0030】

結合されたりポジトリ306に記憶されたデータ系統情報およびジョブ依存性情報の結合は、そうでなければ利用できない、報告および情報を生成するために使用されることができる。これら2つのデータのソースを共に結合することによって、そうでなければ利用できない、ジョブの全体像を見ることが可能となる。例えば、データの結合は、「ジョブが遅延する場合に、任意の所与のデータセットに対して何を意味するか」という問いに対する答えを提供するために使用されることができる。ジョブは、データセット1に直接影響を与えないかもしれないが、ワークロード自動化システムにおけるスケジューリング指示子のために、間接的に影響を与えるかもしれない。

【0031】

監査および報告システム308は、特定のデータセットが、影響を受けようとしている時に、警報を出すことができる。例えば、業務では、特定のデータセットが不正確になりそうな時に警報を出すことを望むかもしれない。

【0032】

情報処理システム310は、ユーザ312にグラフィカルユーザインターフェースを提示することができ、ユーザが、上述のリンクに基づいて、ジョブスケジューリング情報とデータ系統情報との間をナビゲートすることを含み、スケジューリング情報および/またはデータ系統の詳細をドリルダウンし、検討することができるようにする。

【0033】

図4は、データ分析技術が使用され得る、データ処理システム400の例を示す。システム400は、データソース402を含み、データソース402は、例えば、ワークロード自動化システムのデータリポジトリを含む、記憶装置、またはオンラインデータストリームへの接続などの1つまたは複数のデータのソースを含むことができる。各データ記憶は、様々なフォーマットのうちのいずれか（例えば、データベーステーブル、スプレッドシートファイル、単層テキストファイル、またはメインフレームによって使用されるネイティブフォーマット）で、データを記憶し、または提供することができる。実行環境404は、前処理モジュール406と、実行モジュール412とを含む。実行環境404は、UNIXオペレーティングシステムのあるバージョンなどの、適切なオペレーティングシステムの制御下で、例えば、1つまたは複数の汎用コンピュータ上に提供されてもよい。例えば、実行環境404は、局所的（例えば、対称型マルチプロセッシング（SMP）コンピュータなどの、マルチプロセッサシステム）、もしくは局所分散（例えば、クラスタとして連結された複数のプロセッサ、もしくは大規模並列処理（MPP）システム）のいずれか、またはリモートもしくはリモート分散（例えば、ローカルエリアネットワーク（LAN）および/または広域ネットワーク（WAN）を介して連結された複数のプロセッサ）、またはそれらの任意の組合せで、複数の中央処理装置（CPU）またはプロセッサコアを使用するコンピュータシステムの構成を含む、マルチノード並列コンピューティング環境を含むことができる。

【0034】

変換モジュール406は、データソース402からデータを読み出し、データを正準形式に変換し、データ記憶416に情報を記憶させる。データソース402を提供する記憶装置は、実行環境404に対して局所的（例えば、実行環境404を提供するコンピュータに接続された記憶媒体（例えば、ハードドライブ408）に記憶されている）であつてもよく、または、実行環境404に対してリモート（例えば、リモート接続（例えば、クラウドコンピューティングインフラストラクチャによって提供される）を介して、実行環境404を提供するコンピュータと通信関係にある、リモートシステム（例えば、メインフレーム410）上で提供されている）であつてもよい。

【 0 0 3 5 】

分析モジュール 4 1 2 は、データ系統情報と結合された、変換モジュール 4 0 6 によって生成された、記憶された情報を使用して、結合されたデータの分析を、結合されていない方法で実行する。例えば、ジョブのスケジュール変更は、ジョブによって直接影響を受けるもの以外のデータ記憶に強い影響を与える。ジョブは、他のジョブに影響を与える可能性があり、それらのジョブのそれぞれが、データソースに影響を与える可能性がある。いくつかのシナリオでは、データソースに対する変更は、同様に、追加のジョブに影響を与える可能性がある。記憶された情報は、データ記憶システム 4 1 6 に記憶されてもよい。データ記憶システム 4 1 6 は、ユーザ 4 2 0 と対話処理を行う分析システム 4 1 8 にもアクセス可能である。ユーザ 4 2 0 は、結合されたデータのドリルダウン分析を実行することが可能である。

10

【 0 0 3 6 】

分析システム 4 1 8 および実行環境 4 0 4 は、いくつかの実装において、有向リンク（作業要素、即ち、データの流れを表す）によって頂点間が接続された、頂点（データ処理コンポーネントまたはデータセットを表す）を含むデータフローグラフとして、計算アプリケーションを実行するシステムを使用して設計されている。例えば、そのような環境は、参照により本明細書に組み込まれる、「Managing Parameters for Graph-Based Applications」と題した米国特許出願公開第 2 0 0 7 / 0 0 1 1 6 6 8 号に、より詳細に記載されている。そのようなグラフベースの計算を実行するシステムは、参照により本明細書に組み込まれる、「EXECUTING COMPUTATIONS EXPRESSED AS GRAPHS」と題する米国特許第 5 , 9 6 6 , 0 7 2 号に記載されている。このシステムに従って作られるデータフローグラフは、グラフコンポーネントによって表される個々のプロセスに出入りする情報を取得し、プロセス間で情報を移動し、およびプロセスに対する実行順序を定義する方法を提供する。このシステムは、プロセス間の通信方法を任意の利用可能な方法（例えば、グラフのリンクに従う通信パスが、TCP/IPもしくはUNIXドメインソケットを使用することができ、共有メモリを使用して、プロセス間でデータを渡すことができる）から選択するアルゴリズムを含む。

20

【 0 0 3 7 】

変換モジュール 4 0 6 は、異なる形式のデータベースシステムを含む、データソース 4 0 2 を具体化し得る様々なタイプのシステムからデータを受信することができる。データは、空値を含む可能性がある、それぞれのフィールド（「属性」または「カラム」とも呼ばれる）についての値を有するレコードとして編成されてもよい。データソースからデータを読み出す時、変換モジュール 4 0 6 は、典型的には、そのデータソース内のレコードを説明する、いくつかの初期フォーマット情報で開始される。場合によっては、データソースのレコード構造は、最初は既知でなくてもよく、その代わりに、データソースまたはデータの分析後に判断されてもよい。レコードについての初期情報は、例えば、固有値を表すビットの数字、レコード内のフィールドの順序、および、ビットによって表される値の型（例えば、文字列、符号付き / 符号なし整数）を含むことができる。

30

【 0 0 3 8 】

図 5 は、スケジューリング情報への変更の影響を識別する、例示的なプロセス 5 0 0 のフローチャートである。プロセスは、プロセスを実行するコンピュータシステムによって実行されてもよい。

40

【 0 0 3 9 】

データ系統情報が、取得 5 0 2 されることができる。データ系統情報は、上述したデータ記憶から取得されてもよい。データ系統情報は、データ記憶と変換との間の依存関係を識別することができる。変換は、1つのデータ記憶からデータを受け取り、別のデータ記憶にデータを生じさせることができる。

【 0 0 4 0 】

ジョブ依存性情報が、取得 5 0 4 されることができる。ジョブ依存性情報は、上記で論

50

じたプロセスを通じて、取得され得る。ジョブ依存性情報は、複数のジョブの実行順序を指定することができる。

【0041】

ジョブ依存性情報およびデータ系統情報の要素のうちの少なくともいくつかの間のリンクが、識別されることができる。リンクは、直接であってもよい（例えば、ジョブが、変換を実行506させてもよい）。リンクは、間接であってもよい（例えば、ジョブが、データフローグラフを実行させてもよく、その場合に、データフローグラフが変換を含む）。リンクは、ジョブスケジューリング情報およびデータ系統情報によって参照されるファイルおよびデータ記憶に基づいて、判断されてもよい。

【0042】

データ記憶上の複数のアプリケーションのうちのアプリケーションの実行予定についての変更の影響が判断508され得る。影響は、ジョブ依存性情報、リンク、およびデータ系統情報に基づいて判断されてもよい。例えば、ユーザは、少なくとも1つのジョブ、実行可能プログラムまたはデータ記憶を識別するクエリを提示してもよい。例えば、特定のジョブ、実行可能プログラム、もしくはデータソースが利用可能でない場合、または特定のジョブが、失敗もしくは、時間通りに実行されない場合に、ユーザは、影響を判断することを希望してもよい。代替的に、または追加的に、ジョブが失敗したこと、または時間通りの完了に失敗したことを、ワークロード自動化システムが識別してもよい。例えば、ジョブ自体によってハンドリングできないエラーが処理中に発生する時に、ジョブが失敗することがある。例えば、期限が経過する時にも、ジョブが失敗することがある。

【0043】

プロセスは、識別されたジョブ、実行可能プログラムまたはデータ記憶に従属しているジョブの全てを識別することができる。プロセスは、識別されたジョブと、従属関係にあるジョブとの間のリンク、およびデータ系統要素を識別することができる。データ系統要素は、データソース、データシンク、およびデータ変換の関係性、または関係性の一部を説明する、データ要素であってもよい。識別されたジョブおよび従属関係にあるジョブにリンクされているデータ系統要素は、ジョブに従属しているデータ系統要素を判断するために使用されることができる。即ち、データ系統要素に続いてアクセスされるデータ系統要素の全てが、リンクによって識別されることができる。

【0044】

プロセスは、再帰的に適用されることができる。例えば、データ系統要素が一度識別されると、追加のリンクが、データ系統要素を再度追加のジョブに関係付ける。追加のジョブは、同様に、再度追加のデータ系統要素にリンクしてもよい。

【0045】

例えば、図6は、ジョブ依存性情報およびデータ系統情報の結合に基づいて判断され得る依存性情報の単純化した例を示す。ジョブ600は、2つのサブジョブ、(1)生成daily_sales.datジョブ602、および(2)FTP_daily_sales.datジョブ604を含む。daily_sales_file.datは、データ系統情報によって識別されるものとして、データフローグラフ606によって使用される。この例では、FTP_daily_sales.datジョブは、破線610で示されるように、daily_sales.datデータソース608にリンクされている。input_file.dat608は、集約変換612によって、図示されていない他のデータと集約される。集約変換612は、データソース、quarterly.dat614を生成する。別のジョブ616は、quarterly.datのファイルの生成を監視する、ファイル監視quarterly.datジョブ618を含む。この関係性に基づいて、quarterly.datのデータソースと、ファイル監視quarterly.datジョブ618とは、破線620で表すように、結合されたりポジトリにおいてリンクされる。生成10-K情報ジョブ622は、quarterly.datのファイルを使用して、SEC用の10-K情報を生成する。

【0046】

結合されたスケジュール依存性情報およびデータ系統情報を使用すること、ならびに再度追加のスケジュール依存性情報にリンクすることによって、生成 `daily sales` ジョブと生成 10 - K 情報ジョブがジョブ依存性情報によってリンクされていなくても、システムは、`daily sales.dat` ファイル生成の問題が、10 - K の生成に伴う遅延という結果を引き起こし得るということを判断することができる。

【0047】

上述のデータ統合および分析アプローチは、適切なソフトウェアを実行するコンピューティングシステムを使用して実行され得る。例えば、ソフトウェアは、1つまたは複数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ティングシステム（分散、クライアント/サーバ、もしくはグリッドなどの、様々なアーキテクチャのものであり得る）上で実行する、1つまたは複数のコンピュータプログラム内の手続きを含んでもよい。それぞれのコンピューティングシステムは、少なくとも1つのプロセッサ、少なくとも1つのデータ記憶システム（揮発性および/または不揮発性メモリおよび/または記憶素子を含む）、少なくとも1つのユーザインターフェース（少なくとも1つの入力デバイスまたはポートを用いて入力を受信するための、および、少なくとも1つの出力デバイスまたはポートを用いて出力を提供するための）を含む。ソフトウェアは、例えば、データフローグラフの設計、構成、および実行に関するサービスを提供する、より大きなプログラムの1つまたは複数のモジュールを含んでもよい。プログラムのモジュール（例えば、データフローグラフの要素）は、データリポジトリ内に記憶されたデータモデルに従ったデータ構造または他の編成済みデータとして実装されてもよい。

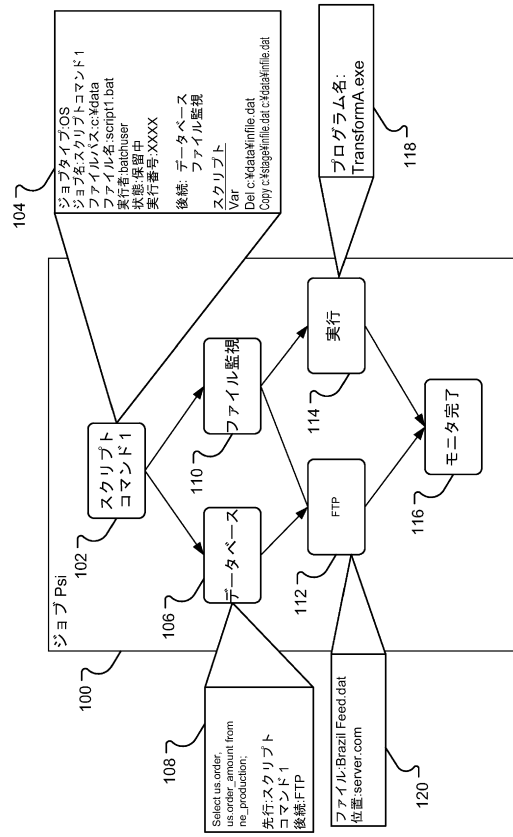
【0048】

ソフトウェアは、CD-ROMもしくは他のコンピュータ可読媒体（例えば、汎用もしくは専用コンピューティングシステムもしくはデバイスによって可読である）などの、有形の非一時的媒体において提供されてもよく、またはネットワークの通信媒体を介して、ソフトウェアが実行されるコンピューティングシステムの有形の非一時的媒体に（例えば、伝播信号において符号化されて）配信されてもよい。処理のいくつかもしくは全ては、専用コンピュータ上で実行されてもよく、または、コプロセッサもしくはフィールドプログラマブルゲートアレイ（FPGA）もしくは専用の特定用途向け集積回路（ASIC）などの、専用ハードウェアを使用してもよい。処理は、ソフトウェアによって指定された計算の異なる部分が、異なる計算要素によって実行される、分散方式で実装されてもよい。そのようなコンピュータプログラムのそれぞれは、記憶デバイス媒体が、本明細書で説明した処理を実行するためにコンピュータによって読み出される時にコンピュータを構成および動作させるための、汎用または専用プログラム可能なコンピュータによってアクセス可能な、記憶デバイスのコンピュータ可読記憶媒体（例えば、ソリッドステートメモリもしくは媒体、または磁気媒体もしくは光学式媒体）上に記憶され、またはダウンロードされるのが好ましい。本発明のシステムは、有形の非一時的媒体として実装され、そのように構成された媒体が、コンピュータを、特定の予め定義された方法で動作させて、本明細書で説明した1つまたは複数の処理ステップを実行させる、コンピュータプログラムで構成されるものと考えられてもよい。

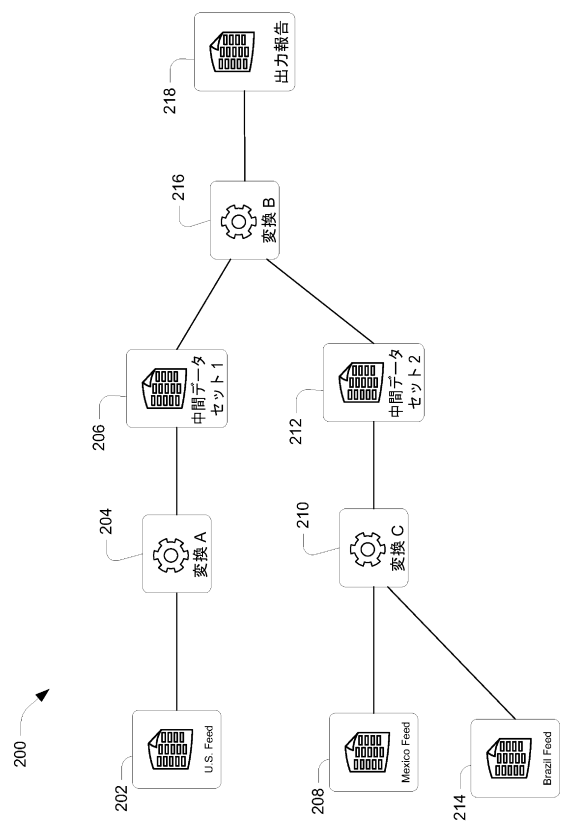
【0049】

いくつかの本発明の実施形態について説明した。それにも関わらず、前述の説明は、以下の特許請求の範囲によって定義される本発明の範囲を例証することを意図するものであり、限定することを意図するものではないと理解されるべきである。したがって、他の実施形態もまた、以下の特許請求の範囲内である。例えば、本発明の範囲から逸脱することなく、様々な修正が行われ得る。さらに、上述したステップのうちのいくつかは、独立した順序であってもよく、したがって、説明した順序とは異なる順序で実行されることができる。

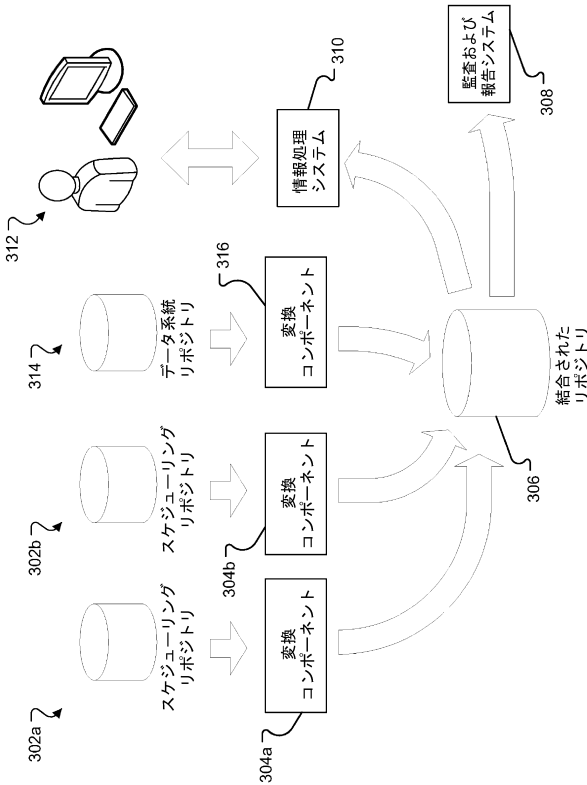
【図 1】



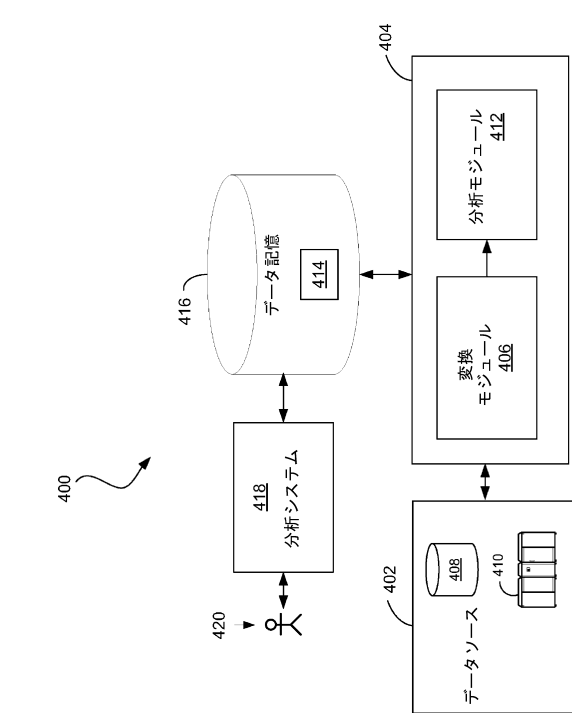
【図 2】



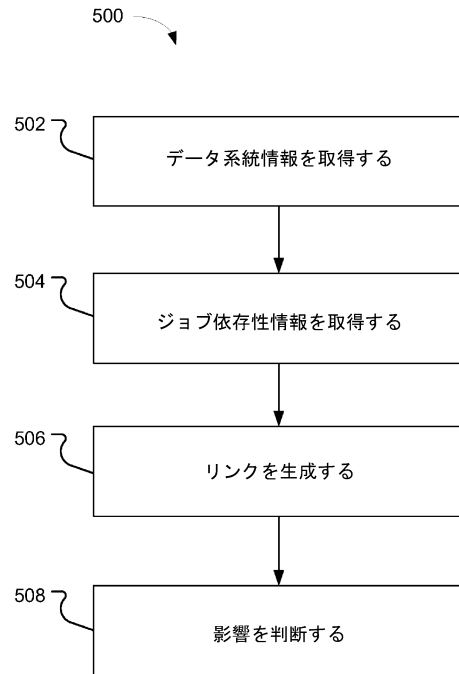
【図 3】



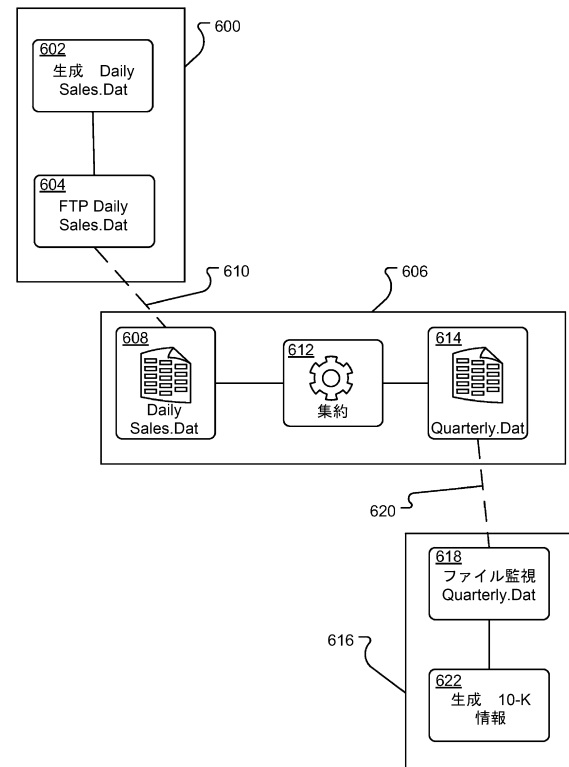
【図 4】



【図 5】



【図 6】



フロントページの続き

- (72)発明者 ウォルフソン, ハリー マイケル
アメリカ合衆国, マサチューセッツ州 02474, アーリントン, メルローズ ストリート 3
6
- (72)発明者 ゴウルド, ジョエル
アメリカ合衆国, マサチューセッツ州 02474, アーリントン, リー テラス 27
- (72)発明者 イエラカリス, アンソニー
アメリカ合衆国, マサチューセッツ州 02460, ニュートン, ワイルドウッド アベニュー
67
- (72)発明者 ウェイクリング, ティム
アメリカ合衆国, マサチューセッツ州 01810, アンドーバー, アボット ストリート 11

審査官 加藤 優一

- (56)参考文献 米国特許出願公開第2009/0241117 (US, A1)

特開2009-163566 (JP, A)
特開2001-022695 (JP, A)
特開2006-268509 (JP, A)
特表2008-507008 (JP, A)
特開2000-066931 (JP, A)
特開平06-083598 (JP, A)
特開2002-342317 (JP, A)
特開2009-268103 (JP, A)
特開2006-120021 (JP, A)

- (58)調査した分野(Int.Cl., DB名)

G06F 9/455 - 9/54
G06F 12/00
G06F 8/00