

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5340689号  
(P5340689)

(45) 発行日 平成25年11月13日 (2013.11.13)

(24) 登録日 平成25年8月16日 (2013.8.16)

(51) Int.Cl.

F I

G 0 6 Q 10/00 (2012.01)

G 0 6 Q 10/00 1 4 O

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/30 2 3 O Z

G 0 6 F 12/00 (2006.01)

G 0 6 F 12/00 5 1 3 Z

G 0 6 Q 10/10 (2012.01)

G 0 6 Q 10/10 1 3 O Z

請求項の数 9 (全 23 頁)

(21) 出願番号 特願2008-265354 (P2008-265354)  
 (22) 出願日 平成20年10月14日 (2008.10.14)  
 (65) 公開番号 特開2010-97263 (P2010-97263A)  
 (43) 公開日 平成22年4月30日 (2010.4.30)  
 審査請求日 平成23年9月16日 (2011.9.16)

(73) 特許権者 390024350  
 株式会社ジャストシステム  
 徳島県徳島市川内町平石若松 1 〇 8 番地 4  
 (74) 代理人 100095407  
 弁理士 木村 満  
 (72) 発明者 岡 智幸  
 大阪府大阪市東淀川区東中島 1 丁目 3 番 1  
 4 号 株式会社キーエンス内  
 審査官 大野 朋也

最終頁に続く

(54) 【発明の名称】 データベース生成装置、データベース生成方法及びコンピュータプログラム

(57) 【特許請求の範囲】

【請求項 1】

表データが含まれる一又は複数の電子文書ファイル中から抽出したデータに基づいて新たなデータベースを生成するデータベース生成装置において、

一又は複数の前記電子文書ファイルを取得する電子文書ファイル取得手段と、

取得した一又は複数の前記電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出する罫線情報抽出手段と、

抽出した複数の罫線情報に基づいて前記電子文書ファイルの内容を解析する解析手段と

、

前記罫線情報に基づいた前記電子文書ファイルの内容の解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定するデータ抽出規則特定手段と、

特定したデータベース項目及びデータ抽出規則にて、一又は複数の前記電子文書ファイルから前記データベース項目及び対応するデータを抽出するデータ抽出手段と、

抽出されたデータベース項目及び対応するデータを一覧表示する表示手段と、

表示されたデータベース項目及び対応するデータが適正である旨を示す確定情報の入力を受け付ける確定情報受付手段と

を備えることを特徴とするデータベース生成装置。

【請求項 2】

前記データ抽出規則特定手段は、

10

20

一又は複数の前記電子文書ファイルの指定を受け付けるファイル指定受付手段と、  
指定を受け付けた一又は複数の電子文書ファイルの表データの範囲指定を受け付ける範囲指定受付手段と、

受け付けた範囲指定に従って、前記データベース項目及び前記データ抽出規則を特定する特定手段と

を備えることを特徴とする請求項1記載のデータベース生成装置。

【請求項3】

異なる電子文書ファイルから抽出した表データの位置の相違に関する情報、及び/又は異なる電子文書ファイルから抽出したデータベース項目の相違に関する情報を少なくとも含む表データの相違に関するゆらぎ情報を抽出するゆらぎ情報抽出手段と、

10

前記データベース項目、前記データ抽出規則及び抽出されたゆらぎ情報に基づいて、前記データ抽出規則の変更部分が存在するか否かを判断する判断手段と、

該判断手段で変更部分が存在すると判断した場合、同一の前記データベース項目に対しては同一の、異なるデータベース項目に対しては異なるタグ情報を付与するタグ情報付与手段と

を備え、

前記データ抽出手段は、前記データ抽出規則の変更部分を反映して前記データベース項目及び対応するデータを抽出するようにしてあり、

前記表示手段は、前記データベース項目に付与されているタグ情報に従って前記データベース項目を配列して、前記データベース項目及び対応するデータを一覧表示するようにしてあることを特徴とする請求項1又は2記載のデータベース生成装置。

20

【請求項4】

前記解析手段は、

罫線により区切られた区画が複数列又は複数行にわたって同一であるか否かを判断する手段を備え、

該手段で同一であると判断した場合、前記データ抽出手段は、複数列又は複数行にわたって同一である最初の行又は列での区画に相当するデータベース項目にてデータを抽出するようにしてあることを特徴とする請求項1乃至3のいずれか一項に記載のデータベース生成装置。

【請求項5】

30

表データが含まれる一又は複数の電子文書ファイル中から抽出したデータに基づいて新たなデータベースを生成するデータベース生成装置で実行することが可能なデータベース生成方法において、

前記データベース生成装置は、

一又は複数の前記電子文書ファイルを取得し、

取得した一又は複数の前記電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出し、

抽出した複数の罫線情報に基づいて前記電子文書ファイルの内容を解析し、

前記罫線情報に基づいた前記電子文書ファイルの内容の解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定し、

40

特定したデータベース項目及びデータ抽出規則にて、一又は複数の前記電子文書ファイルから前記データベース項目及び対応するデータを抽出し、

抽出されたデータベース項目及び対応するデータを一覧表示し、

表示されたデータベース項目及び対応するデータが適正である旨を示す確定情報の入力を受け付けることを特徴とするデータベース生成方法。

【請求項6】

前記データベース生成装置は、

一又は複数の前記電子文書ファイルの指定を受け付け、

指定を受け付けた一又は複数の電子文書ファイルの表データの範囲指定を受け付け、

受け付けた範囲指定に従って、前記データベース項目及び前記データ抽出規則を特定す

50

ることを特徴とする請求項 5 記載のデータベース生成方法。

【請求項 7】

前記データベース生成装置は、

異なる電子文書ファイルから抽出した表データの位置の相違に関する情報、及び／又は異なる電子文書ファイルから抽出したデータベース項目の相違に関する情報を少なくとも含む表データの相違に関するゆらぎ情報を抽出し、

前記データベース項目、前記データ抽出規則及び抽出されたゆらぎ情報に基づいて、前記データ抽出規則の変更部分が存在するか否かを判断し、

変更部分が存在すると判断した場合、同一の前記データベース項目に対しては同一の、異なるデータベース項目に対しては異なるタグ情報を付与し、

前記データ抽出規則の変更部分を反映して前記データベース項目及び対応するデータを抽出し、

前記データベース項目に付与されているタグ情報に従って前記データベース項目を配列して、前記データベース項目及び対応するデータを一覧表示することを特徴とする請求項 5 又は 6 記載のデータベース生成方法。

【請求項 8】

前記データベース生成装置は、

罫線により区切られた区画が複数列又は複数行にわたって同一であるか否かを判断し、

同一であると判断した場合、複数列又は複数行にわたって同一である最初の行又は列の区画に相当するデータベース項目にてデータを抽出することを特徴とする請求項 5 乃至 7 のいずれか一項に記載のデータベース生成方法。

【請求項 9】

表データが含まれる一又は複数の電子文書ファイル中から抽出したデータに基づいて新たなデータベースを生成するデータベース生成装置で実行することが可能なコンピュータプログラムにおいて、

前記データベース生成装置を、

一又は複数の前記電子文書ファイルを取得する電子文書ファイル取得手段、

取得した一又は複数の前記電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出する罫線情報抽出手段、

抽出した複数の罫線情報に基づいて前記電子文書ファイルの内容を解析する解析手段、

前記罫線情報に基づいた前記電子文書ファイルの内容の解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定するデータ抽出規則特定手段

、  
特定したデータベース項目及びデータ抽出規則にて、一又は複数の前記電子文書ファイルから前記データベース項目及び対応するデータを抽出するデータ抽出手段、

抽出されたデータベース項目及び対応するデータを一覧表示する表示手段、及び

表示されたデータベース項目及び対応するデータが適正である旨を示す確定情報の入力を受け付ける確定情報受付手段

として機能させることを特徴とするコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、異なるデータ形式を有する複数の表データが存在する場合であっても、容易に一のデータベースを生成することができるデータベース生成装置、データベース生成方法及びコンピュータプログラムに関する。

【背景技術】

【0002】

関係データベースを生成する場合、事前に生成されている表データを利用することが多い。表データのデータベース項目が一致している場合には、複数の表データをマージすることにより容易に新たなデータベースを生成することができる。

## 【 0 0 0 3 】

しかし、表データのデータ形式が標準化されていない場合、表データのデータ形式は作成者に依存しており、また同一のアプリケーションであっても使用するソフトウェアプログラムによってデータベース項目の相違、データベース項目の配列の相違等が存在しており、そのままマージすることができない。斯かる問題を解決するべく、従来は中間ファイルフォーマットを用いて、データベース項目が相違している、あるいはデータベース項目の配列順序が相違している複数の表データをマージして1つの表データを生成している。

## 【 0 0 0 4 】

例えば特許文献1では、表データの間接ファイルとして良く用いられているCSVファイルを用い、複数の表データファイルから1つのデータベースを生成しているデータベース管理システムが開示されている。

10

【特許文献1】特開2006-059135号公報

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【 0 0 0 5 】

しかし、特許文献1のように、CSVファイルを介して複数の表データをマージする場合、どのデータベース項目が相違しているか、どのデータベース項目の配列順序が相違しているか等に関する情報を事前に知っておく必要があり、これらの情報に応じて適切な変換手順を設定しておかないと、所望の表データを生成することができないという問題点があった。

20

## 【 0 0 0 6 】

また、表データのセル位置に基づいて自動的に表データをマージする方法も考えられているが、表データの開始セル位置が一定ではなく、開始セル位置の相違に基づいてマージ対象となるセルの位置補正を行わないと、複数の異なるファイル又は異なるシート上の表データを正しくマージすることはできない。

## 【 0 0 0 7 】

本発明は斯かる事情に鑑みてなされたものであり、異なるデータ形式を有する複数の表データが存在する場合であっても、新たなデータベースを正しく生成することができるデータベース生成装置、データベース生成方法及びコンピュータプログラムを提供することを目的とする。

30

## 【課題を解決するための手段】

## 【 0 0 0 8 】

上記目的を達成するために第1発明に係るデータベース生成装置は、表データが含まれる一又は複数の電子文書ファイル中から抽出したデータに基づいて新たなデータベースを生成するデータベース生成装置において、一又は複数の前記電子文書ファイルを取得する電子文書ファイル取得手段と、取得した一又は複数の前記電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出する罫線情報抽出手段と、抽出した複数の罫線情報に基づいて前記電子文書ファイルの内容を解析する解析手段と、前記罫線情報に基づいた前記電子文書ファイルの内容の解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定するデータ抽出規則特定手段と、特定したデータベース項目及びデータ抽出規則にて、一又は複数の前記電子文書ファイルから前記データベース項目及び対応するデータを抽出するデータ抽出手段と、抽出されたデータベース項目及び対応するデータを一覧表示する表示手段と、表示されたデータベース項目及び対応するデータが適正である旨を示す確定情報の入力を受け付ける確定情報受付手段とを備えることを特徴とする。

40

## 【 0 0 0 9 】

また、第2発明に係るデータベース生成装置は、第1発明において、前記データ抽出規則特定手段は、一又は複数の前記電子文書ファイルの指定を受け付けるファイル指定受付手段と、指定を受け付けた一又は複数の電子文書ファイルの表データの範囲指定を受け付ける範囲指定受付手段と、受け付けた範囲指定に従って、前記データベース項目及び前記

50

データ抽出規則を特定する特定手段とを備えることを特徴とする。

【0010】

また、第3発明に係るデータベース生成装置は、第1又は第2発明において、異なる電子文書ファイルから抽出した表データの位置の相違に関する情報、及び/又は異なる電子文書ファイルから抽出したデータベース項目の相違に関する情報を少なくとも含む表データの相違に関するゆらぎ情報を抽出するゆらぎ情報抽出手段と、前記データベース項目、前記データ抽出規則及び抽出されたゆらぎ情報に基づいて、前記データ抽出規則の変更部分が存在するか否かを判断する判断手段と、該判断手段で変更部分が存在すると判断した場合、同一の前記データベース項目に対しては同一の、異なるデータベース項目に対しては異なるタグ情報を付与するタグ情報付与手段とを備え、前記データ抽出手段は、前記データ抽出規則の変更部分を反映して前記データベース項目及び対応するデータを抽出するようにしてあり、前記表示手段は、前記データベース項目に付与されているタグ情報に従って前記データベース項目を配列して、前記データベース項目及び対応するデータを一覧表示するようにしてあることを特徴とする。

10

【0011】

また、第4発明に係るデータベース生成装置は、第1乃至第3発明のいずれか1つにおいて、前記解析手段は、罫線により区切られた区画が複数列又は複数行にわたって同一であるか否かを判断する手段を備え、該手段で同一であると判断した場合、前記データ抽出手段は、複数列又は複数行にわたって同一である最初の行又は列での区画に相当するデータベース項目にてデータを抽出するようにしてあることを特徴とする。

20

【0012】

次に、上記目的を達成するために第5発明に係るデータベース生成方法は、表データが含まれる一又は複数の電子文書ファイル中から抽出したデータに基づいて新たなデータベースを生成するデータベース生成装置で実行することが可能なデータベース生成方法において、前記データベース生成装置は、一又は複数の前記電子文書ファイルを取得し、取得した一又は複数の前記電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出し、抽出した複数の罫線情報に基づいて前記電子文書ファイルの内容を解析し、前記罫線情報に基づいた前記電子文書ファイルの内容の解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定し、特定したデータベース項目及びデータ抽出規則にて、一又は複数の前記電子文書ファイルから前記データベース項目及び対応するデータを抽出し、抽出されたデータベース項目及び対応するデータを一覧表示し、表示されたデータベース項目及び対応するデータが適正である旨を示す確定情報の入力を受け付けることを特徴とする。

30

【0013】

また、第6発明に係るデータベース生成方法は、第5発明において、前記データベース生成装置は、一又は複数の前記電子文書ファイルの指定を受け付け、指定を受け付けた一又は複数の電子文書ファイルの表データの範囲指定を受け付け、受け付けた範囲指定に従って、前記データベース項目及び前記データ抽出規則を特定することを特徴とする。

【0014】

また、第7発明に係るデータベース生成方法は、第5又は第6発明において、前記データベース生成装置は、異なる電子文書ファイルから抽出した表データの位置の相違に関する情報、及び/又は異なる電子文書ファイルから抽出したデータベース項目の相違に関する情報を少なくとも含む表データの相違に関するゆらぎ情報を抽出し、前記データベース項目、前記データ抽出規則及び抽出されたゆらぎ情報に基づいて、前記データ抽出規則の変更部分が存在するか否かを判断し、変更部分が存在すると判断した場合、同一の前記データベース項目に対しては同一の、異なるデータベース項目に対しては異なるタグ情報を付与し、前記データ抽出規則の変更部分を反映して前記データベース項目及び対応するデータを抽出し、前記データベース項目に付与されているタグ情報に従って前記データベース項目を配列して、前記データベース項目及び対応するデータを一覧表示することを特徴とする。

40

50

## 【 0 0 1 5 】

また、第 8 発明に係るデータベース生成方法は、第 5 乃至第 7 発明のいずれか 1 つにおいて、前記データベース生成装置は、罫線により区切られた区画が複数列又は複数行にわたって同一であるか否かを判断し、同一であると判断した場合、複数列又は複数行にわたって同一である最初の行又は列での区画に相当するデータベース項目にてデータを抽出することを特徴とする。

## 【 0 0 1 6 】

次に、上記目的を達成するために第 9 発明に係るコンピュータプログラムは、表データが含まれる一又は複数の電子文書ファイル中から抽出したデータに基づいて新たなデータベースを生成するデータベース生成装置で実行することが可能なコンピュータプログラムにおいて、前記データベース生成装置を、一又は複数の前記電子文書ファイルを取得する電子文書ファイル取得手段、取得した一又は複数の前記電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出する罫線情報抽出手段、抽出した複数の罫線情報に基づいて前記電子文書ファイルの内容を解析する解析手段、前記罫線情報に基づいた前記電子文書ファイルの内容の解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定するデータ抽出規則特定手段、特定したデータベース項目及びデータ抽出規則にて、一又は複数の前記電子文書ファイルから前記データベース項目及び対応するデータを抽出するデータ抽出手段、抽出されたデータベース項目及び対応するデータを一覧表示する表示手段、及び表示されたデータベース項目及び対応するデータが適正である旨を示す確定情報の入力を受け付ける確定情報受付手段として機能させることを特徴とする。

## 【 0 0 2 0 】

第 1 発明、第 5 発明、及び第 9 発明では、一又は複数の電子文書ファイルを取得し、取得した一又は複数の電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出して、抽出した複数の罫線情報に基づいて電子文書ファイルの内容を解析する。罫線情報に基づいた電子文書ファイルの内容の解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定し、特定したデータベース項目及びデータ抽出規則にて、一又は複数の電子文書ファイルからデータベース項目及び対応するデータを抽出する。抽出されたデータベース項目及び対応するデータを一覧表示し、表示されたデータベース項目及び対応するデータが適正である旨を示す確定情報の入力を受け付ける。罫線情報に基づいて電子文書ファイル中の表データの位置を特定することができ、新たに生成するデータベースの基礎となるデータベース項目及びデータを抽出するデータ抽出規則を正しく特定することができる。

## 【 0 0 2 1 】

第 2 発明及び第 6 発明では、一又は複数の電子文書ファイルの指定を受け付け、指定を受け付けた一又は複数の電子文書ファイルの表データの範囲指定を受け付ける。受け付けた範囲指定に従って、データベース項目及びデータ抽出規則を特定する。これにより、新たに生成するデータベースに使用するデータベース項目及びデータ抽出規則に不備が生じた場合であっても、ユーザの範囲指定により適切なデータベース項目及びデータ抽出規則を特定することができ、表データをマージすることを半自動化することができる。

## 【 0 0 2 2 】

第 3 発明及び第 7 発明では、異なる電子文書ファイルから抽出した表データの位置の相違に関する情報、及び / 又は異なる電子文書ファイルから抽出したデータベース項目の相違に関する情報を少なくとも含む表データの相違に関するゆらぎ情報を抽出する。データベース項目、データ抽出規則及び抽出されたゆらぎ情報に基づいて、データ抽出規則の変更部分が存在するか否かを判断し、変更部分が存在すると判断した場合、同一のデータベース項目に対しては同一の、異なるデータベース項目に対しては異なるタグ情報を付与する。データ抽出規則の変更部分を反映してデータベース項目及び対応するデータを抽出し、データベース項目に付与されているタグ情報に従ってデータベース項目を配列して、データベース項目及び対応するデータを一覧表示する。これにより、複数のファイル上で表

データの位置が相違する場合、データベース項目が相違する場合、データベース項目の配列順序が相違する場合等であっても、斯かる相違に起因して変更されたデータ抽出規則に従ってデータを抽出し、同一のデータベース項目については同一のタグ情報をキー情報として集約することができ、新たな異なるデータベース項目については、異なるタグ情報により新規のデータベース項目として追加することができる。したがって、ユーザがデータベース項目の相違点を事前にすべて知ることができない場合であっても、データベース項目が重複又は欠落することなく新たなデータベースを生成して一覧表示することが可能となる。

【 0 0 2 3 】

第 4 発明及び第 8 発明では、罫線により区切られた区画が複数列又は複数行にわたって同一であるか否かを判断し、同一であると判断した場合、複数列又は複数行にわたって同一である最初の行又は列での区画に相当するデータベース項目にてデータを抽出する。これにより、連続して区画が同一である場合には一の項目がさらに複数の項目に分割される可能性が少ないことから、各データベース項目に対応するデータを漏れなく抽出することが可能となる。

【 発明の効果 】

【 0 0 2 4 】

本発明によれば、罫線情報に基づいて電子文書ファイル中の表データの位置を特定することができ、新たに生成するデータベースの基礎となるデータベース項目及び対応するデータを抽出するデータ抽出規則を正しく特定することができる。

【 発明を実施するための最良の形態 】

【 0 0 2 5 】

以下、本発明の実施の形態に係るデータベース生成装置について、図面に基づいて具体的に説明する。以下の実施の形態は、特許請求の範囲に記載された発明を限定するものではなく、実施の形態の中で説明されている特徴的事項の組み合わせの全てが解決手段の必須事項であるとは限らないことは言うまでもない。

【 0 0 2 6 】

また、本発明は多くの異なる態様にて実施することが可能であり、実施の形態の記載内容に限定して解釈されるべきものではない。実施の形態を通じて同じ要素には同一の符号を付している。

【 0 0 2 7 】

以下の実施の形態では、コンピュータシステムにコンピュータプログラムを導入したデータベース生成装置について説明するが、当業者であれば明らかな通り、本発明はその一部をコンピュータで実行することが可能なコンピュータプログラムとして実施することができる。したがって、本発明は、データベース生成装置というハードウェアとしての実施の形態、ソフトウェアとしての実施の形態、又はソフトウェアとハードウェアとの組み合わせの実施の形態をとることができる。コンピュータプログラムは、ハードディスク、DVD、CD、光記憶装置、磁気記憶装置等の任意のコンピュータで読み取ることが可能な記録媒体に記録することができる。

【 0 0 2 8 】

( 実施の形態 1 )

図 1 は、本発明の実施の形態 1 に係るデータベース生成装置の構成例を示すブロック図である。本発明の実施の形態 1 に係るデータベース生成装置 1 は、少なくとも CPU ( 中央演算装置 ) 1 1、メモリ 1 2、記憶装置 1 3、I/O インタフェース 1 4、ビデオインタフェース 1 5、可搬型ディスクドライブ 1 6、通信インタフェース 1 7 及び上述したハードウェアを接続する内部バス 1 8 で構成されている。

【 0 0 2 9 】

CPU 1 1 は、内部バス 1 8 を介してデータベース生成装置 1 の上述したようなハードウェア各部と接続されており、上述したハードウェア各部の動作を制御するとともに、記憶装置 1 3 に記憶されているコンピュータプログラム 1 0 0 に従って、種々のソフトウェ

10

20

30

40

50

ア的機能を実行する。メモリ 12 は、SRAM、SDRAM等の揮発性メモリで構成され、コンピュータプログラム 100 の実行時にロードモジュールが展開され、コンピュータプログラム 100 の実行時に発生する一時的なデータ等を記憶する。

【0030】

記憶装置 13 は、内蔵される固定型記憶装置（ハードディスク）、SRAM等の揮発性メモリ、ROM等の不揮発性メモリ等で構成されている。記憶装置 13 に記憶されているコンピュータプログラム 100 は、プログラム及びデータ等の情報を記録したDVD、CD-ROM等の可搬型記録媒体 90 から、可搬型ディスクドライブ 16 によりダウンロードされ、実行時には記憶装置 13 からメモリ 12 へ展開して実行される。もちろん、通信インタフェース 17 を介してネットワーク 2 に接続されている外部のコンピュータからダウンロードされたコンピュータプログラムであっても良い。

10

【0031】

また記憶装置 13 は、電子文書ファイル記憶部 131、データ抽出規則記憶部 132、データベース記憶部 133 及びゆらぎ情報記憶部 134 を備えている。電子文書ファイル記憶部 131 には、表データを内容に含み、新たなデータベースを生成するための基礎となる電子文書ファイルを記憶する。

【0032】

データ抽出規則記憶部 132 には、例えば電子文書ファイルに含まれる表データのうち最大のサイズを有する表データを選択する、ファイルの先頭から  $n$  ( $n$  は自然数) 番目の表データを選択する等の、表データからデータベース項目及び対応するデータを抽出するデータ抽出規則を記憶している。

20

【0033】

データベース記憶部 133 には、複数の電子文書ファイルに含まれる表データをマージして新たに生成されたデータベースを記憶する。ゆらぎ情報記憶部 134 には、マージする対象となる表データ間の相違に関する情報、いわゆるゆらぎ情報を記憶する。ゆらぎ情報としては、例えば表データの開始セルの位置の相違に関する表位置ゆらぎ情報、表データの項目の順序が相違する、新規項目の存在、項目の抜けの存在等の項目の相違に関する項目ゆらぎ情報等がある。また、英語表記での大文字と小文字との相違、全角と半角との相違等も含む広い概念である。

【0034】

30

通信インタフェース 17 は内部バス 18 に接続されており、インターネット、LAN、WAN等の外部のネットワーク 2 に接続されることにより、外部のコンピュータ等とデータ送受信を行うことが可能となっている。電子文書ファイル記憶部 131 は、データベース生成装置 1 の記憶装置 13 に備えることに限定されるものではなく、外部のコンピュータの記憶装置に記憶されることによりネットワーク 2 上に点在していても良い。

【0035】

I/Oインタフェース 14 は、キーボード 21、マウス 22 等のデータ入力媒体と接続され、データの入力を受け付ける。また、ビデオインタフェース 15 は、CRTモニタ、LCD等の表示装置 23 と接続され、所定の画像を表示する。

【0036】

40

図 2 は、本発明の実施の形態 1 に係るデータベース生成装置 1 の機能ブロック図である。電子文書ファイル取得部 201 は、一又は複数の表データを含む電子文書ファイルを取得する。電子文書ファイルは、記憶装置 13 内に記憶されている電子文書ファイルを電子文書ファイル記憶部 131 に集約しても良いし、ネットワーク 2 を介して外部のコンピュータから取得しても良い。また、キーボード 21、マウス 22 等の入力装置を介して入力しても良い。

【0037】

罫線情報抽出部 202 は、取得した一又は複数の電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出する。具体的には、罫線で囲まれている部分を表データと認識し、その他の罫線がどのように配置されているかに関する情報を取得する。

50

## 【0038】

解析部203は、抽出した複数の罫線情報に基づいて電子文書ファイルの内容を解析する。具体的には、罫線によりレコード単位で項目がどのように区分けされているかを判断し、見出し部とデータ部とを区別する。

## 【0039】

データ抽出規則特定部204は、解析部203での解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定する。特定されたデータベース項目及びデータ抽出規則は、データ抽出規則記憶部132に記憶される。

## 【0040】

データ抽出規則特定部204は、ファイル指定受付部205、範囲指定受付部206、及び特定部207を備えても良い。解析部203の解析結果だけでは、データベース項目を特定することができない場合もありうるからである。このような場合、手動にてデータベース項目及びデータ抽出規則の特定を受け付ける。

## 【0041】

データ抽出部208は、ファイル指定受付部205にて一又は複数の電子文書ファイルの指定を受け付け、範囲指定受付部206にて複数のシートが存在する場合にはシートの指定、及びシートに含まれる表データ中にて該表データと他の表データとのマージ対象となる範囲指定を受け付ける。特定部207は、受け付けた範囲指定に従って、データベース項目及びデータ抽出規則を特定する。

## 【0042】

データ抽出部208は、特定したデータベース項目及びデータ抽出規則にて、一又は複数の電子文書ファイルからデータベース項目及び対応するデータを抽出する。抽出されたデータベース項目及び対応するデータはデータベース記憶部133に記憶される。

## 【0043】

表示部209は、抽出されたデータベース項目及び対応するデータを表示装置23にて一覧表示し、確定情報受付部210は、表示されたデータベース項目及び対応するデータが適正であるか否かをユーザが判断し、ユーザが適正であると判断した場合、すなわちデータベース項目に重複、抜け等が生じておらず、適正にマージされていると判断した場合には、ユーザによる適正であると判断した旨を示す確定情報の入力を受け付ける。

## 【0044】

図3は、本発明の実施の形態1に係るデータベース生成装置1のCPU11のデータベース生成処理の手順を示すフローチャートである。図3において、データベース生成装置1のCPU11は、一又は複数の表データを含む電子文書ファイルを取得する(ステップS301)。電子文書ファイルは、記憶装置13の電子文書ファイル記憶部131内に記憶されている電子文書ファイルを読み出しても良いし、ネットワーク2を介して外部のコンピュータから読み出しても良い。また、キーボード21、マウス22等の入力装置を介して入力を受け付けても良い。

## 【0045】

CPU11は、取得した一又は複数の電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出する(ステップS302)。具体的には、罫線で囲まれている部分を表データと認識し、その他の罫線がどのように配置されているかに関する情報を取得する。

## 【0046】

CPU11は、抽出した複数の罫線情報に基づいて電子文書ファイルの内容を解析する(ステップS303)。具体的には、罫線によりレコード単位で項目がどのように区分けされているかを判断し、見出し部とデータ部とを区別する。

## 【0047】

図4は、罫線情報に基づいて表データの抽出を行う処理の例示図である。具体的には、電子文書ファイルのデータの走査方向につきユーザの指定を受け付け、項目が階層化されているか否かを1行ずつ判定する。図4(a)では、表データを下方向42へ走査する場合を示しており、項目領域41の1行目には項目「材料名」、「重量」、「物質」、「比

10

20

30

40

50

率」が存在することを検出することができる。2行目では、項目「比率」が「平均重量」、「最大重量」に分割され、項目数が増加していることを検出することができる。

【0048】

3行目では、項目名を検出することはできないものの、2行目の項目とセル位置及び項目数が同一であることを検出することができる。したがって、2行目まで見出し部であり、3行目以降がデータ部であることを自動認識することができ、新たなデータベース生成のためのデータ抽出は、3行目以降のデータ部から行うことができる。

【0049】

図4(b)では、表データを右方向44へ走査する場合を示しており、項目領域43の1列目には項目「材料名」、「重量」、「物質」、「比率」が存在することを検出することができ、2列目では、項目「比率」が「平均重量」、「最大重量」に分割され、項目数が増加していることを検出することができる。

10

【0050】

3列目では、項目名を検出することはできないものの、2列目の項目とセル位置及び項目数が同一であることを検出することができる。したがって、2列目まで見出し部であり、3列目以降がデータ部であることを自動認識することができ、新たなデータベース生成のためのデータ抽出は、3列目以降のデータ部から行うことができる。

【0051】

このように走査方向によらず、罫線情報に基づいて、データベース生成時に抽出すべきデータベース項目及び対応するデータのセル位置を正確に検出することができるので、表データの行と列とが反転している場合であっても一のデータベースとしてマージすることが可能となる。

20

【0052】

図3に戻って、データベース生成装置1のCPU11は、解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定する(ステップS304)。CPU11は、一又は複数の電子文書ファイル又は該電子文書ファイル中の表データから、新たに生成するデータベースのデータベース項目及び対応するデータを抽出する(ステップS305)。データを抽出する規則は、データ抽出規則記憶部132に記憶されているデータ抽出規則に従う。

【0053】

30

なお、罫線情報の解析結果だけでは正しくデータベース項目等が特定できない場合も生じうる。この場合、手動にてデータベース項目及びデータ抽出規則の特定を受け付ける。図5は、本発明の実施の形態1に係るデータベース生成装置1のCPU11の手動特定処理の手順を示すフローチャートである。

【0054】

図5において、データベース生成装置1のCPU11は、図3のステップS303の処理の実行終了後、一又は複数の電子文書ファイルの指定を受け付け(ステップS501)、複数のシートが存在する場合にはシートの指定、及びシートに含まれる表データ中にて該表データと他の表データとのマージ対象となる範囲指定を受け付ける(ステップS502)。CPU11は、受け付けた範囲指定に従って、データベース項目及びデータ抽出規則を特定し(ステップS503)、処理を図3のステップS305へ進める。

40

【0055】

図6は、範囲指定が必要となる場合の例示図である。図6(a)は、表データの構造が特段の規則性を有していない場合の例示図である。この場合、キーボード21、マウス22等の入力装置により、表データとして使用する領域61のみを範囲指定として受け付ける。指定を受け付けた範囲に、例えば他の表データのデータベース項目とリンクするようなタグ情報を付加することにより、新たなデータベースのデータとして用いることができる。

【0056】

図6(b)は、表データとして認識できない領域区分となっている場合の例示図である

50

。図6(b)の例では、見出し部として認識すべき領域62が表として認識できる領域、すなわち矩形領域となっていない。この場合、キーボード21、マウス22等の入力装置により、領域62を含めて列ごとの領域63の範囲指定を受け付け、見出し部「材料」、「重量」、「比率1」、「比率2」に対して、他の表データのデータベース項目とリンクするようにタグ情報を付加する。これにより、新たなデータベースのデータとして用いることができる。

#### 【0057】

図3に戻って、データベース生成装置1のCPU11は、抽出されたデータベース項目及び対応するデータを表示装置23にて一覧表示し(ステップS306)、表示されたデータベース項目及び対応するデータが適正である旨の確認情報を受け付けたか否かを判断する(ステップS307)。CPU11が、確認情報を受け付けていないと判断した場合(ステップS307:NO)、再度新たなデータベースを生成するべく処理を終了し、確認情報を受け付けたと判断した場合(ステップS307:YES)、ユーザがデータベース項目に重複、抜け等が生じておらず、適正にマージされていると判断したので、生成された新たなデータベースを検索処理の対象等として記憶装置13に記憶する(ステップS308)。

#### 【0058】

このように、罫線情報に基づいて抽出対象となる表データの存在位置を特定することができ、表データからデータベース項目及び対応するデータを抽出するためのデータ抽出規則を正しく特定することができるので、新たなデータベースを適正に生成することが可能となる。

#### 【0059】

また、マージ対象となる表データ間に、いわゆるゆらぎ情報が存在する場合がある。ここで、「ゆらぎ情報」とは、表データ間の相違に関する情報の総称である。例えば表データの開始セルの位置の相違に関する表位置ゆらぎ情報、表データの項目の順序が相違する、新規項目の存在、項目の抜けの存在等の項目の相違に関する項目ゆらぎ情報等がある。

#### 【0060】

図7は、表データの位置に相違が存在する「表位置ゆらぎ情報」の説明図である。図7(a)から図7(c)に示すように、表データを示す罫線が存在する領域の左上のセル71、72、73のサイズが相違することにより、電子文書ファイル中の表データのセル位置がそれぞれ相違している。表位置ゆらぎ情報が存在する場合、例えばデータ抽出規則を「上から1番目の表」等に特定しておく、又は記憶してあるデータ抽出規則から選択することにより、図7(a)から図7(c)に示すすべての表データを抽出の対象とすることができる。つまり、表データの開始位置(開始セル位置)をあらかじめ定めるのではなく、罫線情報に基づいて表データを抽出し、一の電子文書ファイル中又は一の電子文書ファイルに含まれる一のシート中に複数の表データが抽出された場合には、「上から1番目の表」、「最大の表」、「特定の項目を含む表」等の直接的な位置ではない表の特徴をデータ抽出規則として設定して記憶しておくことにより、表位置のゆらぎに対応して表データを抽出することができる。

#### 【0061】

図8は、表データの項目に相違が存在する「項目ゆらぎ情報」の説明図である。図8(a)を基準とした場合、図8(b)は項目Cと項目Bとの順序が入れ替わっている。従来のCSVファイルを用いて表データをマージする場合には、項目Cと項目Bとの順序が入れ替わっていることを事前にユーザが知っている状態で、入れ替え指示を出す必要があった。

#### 【0062】

本実施の形態1では、項目が入れ替わっていることを検出して、項目名にリンクしたタグ情報を付与する。すなわち図8(a)の見出し部81では、例えば項目Aに対してタグ情報「a」を、項目Bに対してタグ情報「b」を、項目Cに対してタグ情報「c」を、それぞれ付与する。図8(b)の見出し部82では、項目Bと項目Cとの順序が入れ替わっ

10

20

30

40

50

ているが、タグ情報は図8(a)と同様の対応関係で付与しておく。データ抽出時にはタグ情報‘a’、タグ情報‘b’及びタグ情報‘c’を基礎としてデータを集約するので、基礎となる表データで項目がどのように配置されていても、新たなデータベースでは、タグ情報の順に集約することができる。したがって、項目Cと項目Bとの順序が入れ替わっていることを事前にユーザが知らなくても、表データを適正にマージすることが可能となる。

【0063】

また、図8(c)では、見出し部83に新たな項目である項目D、項目Eが存在するのに対し、項目Cが欠落している。この場合も、項目Dに対してタグ情報‘d’を、項目Eに対してタグ情報‘e’を、それぞれ付与することにより、データ抽出時にタグ情報を基礎としてデータを抽出する限り、誤った項目を集約するおそれはない。すなわち、新規に追加された項目は独立して集約することができるし、欠落している項目については本表データからはデータを抽出することがない。

【0064】

以上のように、表データ間にゆらぎ情報が存在する場合であっても、タグ情報を基礎として同一項目についてはデータを集約することができ、ユーザが項目の入れ替わりに関する情報等を正確に把握することなく、新たなデータベースを生成することが可能となる。つまり、抽出すべき項目、項目の順序等をあらかじめ定めるのではなく、レコードを抽出する場合に項目名も抽出し、項目をタグ情報として例えばXMLデータベース形式で記憶し、タグ情報(項目名)によって表を再構成することで、項目のゆらぎに対応して表データを抽出することができる。また、何らかの原因により本来は同じ項目として取り扱うべき項目を異なる項目として認識した場合であっても、表を再構成するときに、同じ項目として取り扱うべき項目を同一化する指示を受け付けることが可能な構成とすることにより、事後的に項目ゆらぎを修正することも可能となる。

【0065】

図9は、本発明の実施の形態1に係るデータベース生成装置1のCPU11のゆらぎ補正処理の手順を示すフローチャートである。

【0066】

図9において、データベース生成装置1のCPU11は、図3のステップS304の処理実行後、異なる電子文書ファイルから抽出した表データの位置の相違に関する情報、及び/又は異なる電子文書ファイルから抽出したデータベース項目の相違に関する情報を少なくとも含む表データの相違に関するゆらぎ情報を抽出する(ステップS901)。抽出するゆらぎ情報は、上述した2つに限定されるものではない。

【0067】

CPU11は、データベース項目、データ抽出規則及び抽出されたゆらぎ情報に基づいて、データ抽出規則の変更部分が存在するか否かを判断する(ステップS902)。CPU11が、変更部分が存在しないと判断した場合(ステップS902:NO)、CPU11は、ゆらぎ情報に起因するデータ抽出規則に対する何らかの補正処理(以下、ゆらぎ補正)が実行されていないと判断して、処理を図3のステップS307へ進める。CPU11が、変更部分が存在すると判断した場合(ステップS902:YES)、CPU11は、同一のデータベース項目に対しては同一の、異なるデータベース項目に対しては異なるタグ情報を付与する(ステップS903)。

【0068】

CPU11は、データ抽出規則の変更部分を反映してデータベース項目及び対応するデータを抽出し(ステップS904)、データベース項目に付与されているタグ情報に従ってデータベース項目を配列して、データベース項目及び対応するデータを一覧表示し(ステップS905)、処理を図3のステップS307へ進める。

【0069】

このようにすることで、マージ対象となる表データ間にゆらぎ情報が存在する場合であっても、ユーザは事前にその存在を把握しておく必要がなく、適正なゆらぎ補正を実行し

10

20

30

40

50

た状態で新たなデータベース項目及び対応するデータを視認することができる。したがって、ユーザに過剰な負荷がかかることなく、容易に新たなデータベースを生成することが可能となる。

【 0 0 7 0 】

( 実施の形態 2 )

本発明の実施の形態 2 に係るデータベース生成装置の構成は、実施の形態 1 と同様であることから、同一の符号を付することにより、詳細な説明を省略する。本実施の形態 2 では、異なるデータ形式を有する複数の表データが存在する場合であっても、容易に新たなデータベースを生成することができ、ゆらぎ補正の内容を視認することができる点で実施の形態 1 と相違する。

10

【 0 0 7 1 】

図 1 0 は、本発明の実施の形態 2 に係るデータベース生成装置 1 の機能ブロック図である。図 1 0 では、実施の形態 1 と同様の機能ブロックについては同一の符号を付している。電子文書ファイル取得部 2 0 1 は、一又は複数の表データを含む電子文書ファイルを取得する。電子文書ファイルは、記憶装置 1 3 内に記憶されている電子文書ファイルを電子文書ファイル記憶部 1 3 1 に集約しても良いし、ネットワーク 2 を介して外部のコンピュータから取得しても良い。また、キーボード 2 1、マウス 2 2 等の入力装置を介して入力しても良い。

【 0 0 7 2 】

罫線情報抽出部 2 0 2 は、取得した一又は複数の電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出する。具体的には、罫線で囲まれている部分を表データと認識し、その他の罫線がどのように配置されているかに関する情報を取得する。

20

【 0 0 7 3 】

解析部 2 0 3 は、抽出した複数の罫線情報に基づいて電子文書ファイルの内容を解析する。具体的には、罫線によりレコード単位で項目がどのように分けられているかを判断し、見出し部とデータ部とを区別する。

【 0 0 7 4 】

データ抽出規則特定部 2 0 4 は、解析部 2 0 3 での解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定する。特定されたデータベース項目及びデータ抽出規則は、データ抽出規則記憶部 1 3 2 に記憶される。

30

【 0 0 7 5 】

データ抽出規則特定部 2 0 4 は、ファイル指定受付部 2 0 5、範囲指定受付部 2 0 6、及び特定部 2 0 7 を備えても良い。解析部 2 0 3 の解析結果だけでは、データベース項目を特定することができない場合もありうるからである。このような場合、手動にてデータベース項目及びデータ抽出規則の特定を受け付ける。

【 0 0 7 6 】

ゆらぎ情報抽出部 1 0 0 1 は、解析部 2 0 3 での解析結果に基づいて、いわゆるゆらぎ情報が存在する場合には、存在するゆらぎ情報を抽出する。例えば、抽出対象となる表データ間において項目ゆらぎ情報が存在する場合、項目の順序の相違、項目の相違等に関する情報を抽出する。

40

【 0 0 7 7 】

判断部 1 0 0 2 は、データベース項目、データ抽出規則及び抽出されたゆらぎ情報に基づいて、データ抽出規則の変更部分が存在するか否かを判断する。すなわち、ゆらぎ情報の存在によって、タグ情報に基づく表データの抽出規則が変更されるので、変更部分が存在すると判断した場合にはデータ抽出規則に対して何らかのゆらぎ補正が実行されていると判断することができる。

【 0 0 7 8 】

タグ情報付与部 1 0 0 3 は、判断部 1 0 0 2 で変更部分が存在すると判断した場合、同一のデータベース項目に対しては同一の、異なるデータベース項目に対しては異なるタグ情報を付与する。すなわち項目の配列、順序、種類等が異なる表データであっても、タグ

50

情報が同一である場合には同一の項目であることを担保する。これにより、タグ情報に基づいてデータベースを構成することにより、ゆらぎ情報が存在する場合であっても適切に新たなデータベースを生成することができる。

【0079】

データ抽出部1004は、特定したデータベース項目及びデータ抽出規則にて、一又は複数の電子文書ファイルからデータベース項目及び対応するデータを抽出する。データ抽出規則にはタグ情報に関する項目が含まれ、タグ情報に応じてデータベース項目及び対応するデータが抽出され、データベース記憶部133に記憶される。

【0080】

表示部209は、抽出されたデータベース項目及び対応するデータを表示装置23にて一覧表示する。一方、ゆらぎ情報表示部1005は、データ抽出の対象となる電子文書ファイルに関するゆらぎ情報の種類に関する情報、及びゆらぎ情報が存在するシートを特定する情報を表示する。

10

【0081】

図11は、本発明の実施の形態2に係るデータベース生成装置1のCPU11のデータベース生成処理の手順を示すフローチャートである。図11において、データベース生成装置1のCPU11は、一又は複数の表データを含む電子文書ファイルを取得する(ステップS1101)。電子文書ファイルは、記憶装置13内の電子文書ファイル記憶部131に記憶されている電子文書ファイルを読み出しても良いし、ネットワーク2を介して外部のコンピュータから読み出しても良い。また、キーボード21、マウス22等の入力装置を介して入力を受け付けても良い。

20

【0082】

CPU11は、取得した一又は複数の電子文書ファイルに含まれる罫線に関する罫線情報をそれぞれ抽出する(ステップS1102)。具体的には、罫線で囲まれている部分を表データと認識し、その他の罫線がどのように配置されているかに関する情報を取得する。

【0083】

CPU11は、抽出した複数の罫線情報に基づいて電子文書ファイルの内容を解析する(ステップS1103)。具体的には、罫線によりレコード単位で項目がどのように区分けされているかを判断し、見出し部とデータ部とを区別する。区別する方法は実施の形態1と同様であることから、詳細な説明は省略する。

30

【0084】

CPU11は、解析結果に基づいて、生成するデータベースのデータベース項目及びデータ抽出規則を特定する(ステップS1104)。CPU11は、解析結果に基づいて、いわゆるゆらぎ情報が存在する場合には、存在するゆらぎ情報を抽出する(ステップS1105)。すなわち、抽出対象となる表データ間において、項目の順序の相違、項目の相違等に関する情報を抽出する。

【0085】

CPU11は、データベース項目、データ抽出規則及び抽出されたゆらぎ情報に基づいて、データ抽出規則の変更部分が存在するか否かを判断する(ステップS1106)。すなわち、ゆらぎ情報の存在によって、タグ情報に基づく表データの抽出規則が変更されるので、変更部分が存在すると判断した場合にはデータ抽出規則に対して何らかのゆらぎ補正が実行されていると判断することができる。

40

【0086】

CPU11が、データ抽出規則の変更部分が存在すると判断した場合(ステップS1106: YES)、CPU11は、同一のデータベース項目に対しては同一の、異なるデータベース項目に対しては異なるタグ情報を付与する(ステップS1107)。すなわち項目の配列、順序、種類等が異なる表データであっても、タグ情報が同一である場合には同一の項目であることを担保する。これにより、タグ情報に基づいてデータベースを構成することで、ゆらぎ情報が存在する場合であっても適切に新たなデータベースを生成するこ

50

とができる。

【 0 0 8 7 】

C P U 1 1 が、データ抽出規則の変更部分が存在しないと判断した場合（ステップ S 1 1 0 6 : N O ）、C P U 1 1 は、ステップ S 1 1 0 7 をスキップし、特定したデータベース項目及びデータ抽出規則にて、一又は複数の電子文書ファイルからデータベース項目及び対応するデータを抽出する（ステップ S 1 1 0 8 ）。データ抽出規則にはタグ情報に関する項目が含まれ、タグ情報に応じてデータベース項目及び対応するデータが抽出され、データベース記憶部 1 3 3 に記憶される。

【 0 0 8 8 】

C P U 1 1 は、抽出されたデータベース項目及び対応するデータを表示装置 2 3 にて一  
覧表示するとともに（ステップ S 1 1 0 9 ）、データ抽出の対象となる電子文書ファイル  
に関するゆらぎ情報の種類に関する情報、及びゆらぎ情報が存在する位置に関する情報、  
例えばゆらぎ情報が存在する電子文書ファイル名、シート名、セル位置等を表示する（ス  
テップ S 1 1 1 0 ）。

【 0 0 8 9 】

このように、ゆらぎ情報が存在する場合であっても、ゆらぎ情報に応じてデータ抽出規則を変更する、すなわちタグ情報を基礎としたデータ抽出規則を変更することにより、正しくデータベース項目及び対応するデータを集約することが可能となる。また、ゆらぎ情報が生じた電子文書ファイル、シート、セル位置等を視認することができ、ゆらぎ補正をどのように行うかを判断することができる。

【 0 0 9 0 】

図 1 2 は、抽出されたゆらぎ情報の一覧を表示する表示画面の例示図である。図 1 2 の例では、「表位置ゆらぎ」、「項目ゆらぎ」、「文字ゆらぎ」、「エラー」の各項目について、その総数、電子文書ファイル名、シート名等の一覧を、それぞれ領域 1 2 1、1 2 2、1 2 3、1 2 4 に表示している。

【 0 0 9 1 】

また、ゆらぎ情報の一覧表示画面から、ゆらぎ補正がどのように実行されたか確認することが好ましい。図 1 3 は、ゆらぎ補正の内容を確認するための表示画面の例示図である。図 1 3 ( a ) は、図 1 2 と同様、抽出されたゆらぎ情報の一覧を表示する表示画面の例示図であり、図 1 3 ( b ) は、ゆらぎ情報の内容を示す表示画面の例示図であり、図 1 3 ( c ) は、ゆらぎ情報抽出の対象となった表データを並列に表示する表示画面の例示図である。

【 0 0 9 2 】

図 1 3 ( a ) の表示画面において、ゆらぎ補正がどのように行われたのか知りたい項目 1 3 5 について、マウス等の入力装置による選択を受け付ける。項目 1 3 5 の選択を受け付けた場合、図 1 3 ( b ) に示すように、対応する表データのシート自体を表示する表データ表示領域 1 3 6 とゆらぎ情報の内容を表示するゆらぎ情報表示領域 1 3 8 を含む表示画面が表示される。表データ表示領域 1 3 6 には、ゆらぎ情報の対象となった表データのシートが表示され、ゆらぎ情報表示領域 1 3 8 には、表位置ゆらぎ情報に対するゆらぎ補正の内容が表示されている。

【 0 0 9 3 】

図 1 3 の例では、表データが開始するセル位置にゆらぎが生じており、図 1 3 ( b ) の表示画面のようにデータ抽出規則における開始するセル位置 1 3 7 を変更した旨を示す情報が表示されている。すなわち、データ抽出規則では、セル位置「C : 3」からデータを抽出するよう指示されていたのに対し、ゆらぎ情報により、データ抽出規則を、セル位置「D : 5」からデータを抽出するよう変更した旨を示している。

【 0 0 9 4 】

また、シート自体を比較することが可能なように並列して表示しても良い。この場合、図 1 3 ( c ) に示すように、表データ表示領域 1 3 6 の横に並列して、ゆらぎ元表データ表示領域 1 3 9 を表示しておき、ゆらぎの対象となるセル位置 1 3 7、1 4 0 をそれぞれ

10

20

30

40

50

強調表示することにより、どのようにゆらぎ補正が実行されているのか、明確に視認することができる。

【 0 0 9 5 】

また、ゆらぎ情報の一覧表示画面において、エラー表示がなされている表データは、事前に想定しているゆらぎ補正では対応できなかった表データの存在を示している。図 1 4 は、エラー表示の内容を示す表示画面の例示図である。図 1 4 ( a ) は、図 1 2 と同様、抽出されたゆらぎ情報の一覧を表示する表示画面の例示図であり、図 1 4 ( b ) は、エラー情報の内容を示す表示画面の例示図であり、図 1 4 ( c ) は、エラー情報抽出の対象となった表データを並列に表示する表示画面の例示図である。

【 0 0 9 6 】

10

図 1 4 ( a ) の表示画面において、エラー表示がなされている項目 1 4 1 について、マウス等の入力装置による選択を受け付ける。項目 1 4 1 の選択を受け付けた場合、図 1 4 ( b ) に示すように、対応する表データのシート自体を表示する表データ表示領域 1 4 2 とエラー情報の内容を表示するエラー情報表示領域 1 4 4 を含む表示画面が表示される。表データ表示領域 1 4 2 には、表データのシート自体が表示され、エラー情報表示領域 1 4 4 には、エラーが生じた原因を示す原因情報、例えば表データ内の項目がすべて一致していない旨を示す情報が表示されている。

【 0 0 9 7 】

図 1 4 の例では、表データの項目が、データ抽出規則に規定されている新たなデータベースのデータベース項目とすべて一致していないことから、ゆらぎ情報の補正処理を行うことができない。したがって、項目をデータベース項目にあわせるよう補正するには、シート自体を比較することが可能なように並列して表示することが好ましい。

20

【 0 0 9 8 】

図 1 4 ( c ) では、表データ表示領域 1 4 2 の横に並列して、ゆらぎ情報を抽出できなかった元表データ表示領域 1 4 5 を表示している。両者の項目領域 1 4 3 と項目領域 1 4 6 とを比較できるよう、例えば強調表示することにより、項目名がどのように相違しているのかを、明確に視認することができる。

【 0 0 9 9 】

なお、原因情報としては、少なくともエラー情報が生じている電子文書ファイルを特定する情報、エラー情報の内容に関する情報、及びエラー情報が生じているシートを特定する情報を含むことが好ましい。また、原因情報は、データ抽出規則の変更部分に関する情報を含むことも望ましい。これは、ゆらぎ情報を表示する場合にゆらぎ情報が生じている電子文書ファイルを特定する情報、ゆらぎ情報に関する情報、ゆらぎ情報が生じているシートを特定する情報を含むことに準ずる。

30

【 0 1 0 0 】

エラー情報が生じている電子文書ファイルを特定する情報、エラー情報の内容に関する情報、及びエラー情報が生じているシートを特定する情報を含むことにより、ゆらぎ情報を補正しきれなかった電子文書ファイル及びシートを特定することができ、補正ができなかった原因を視認することが可能となる。また、データ抽出規則の相違を視認することができ、データ抽出規則を変更することにより新たなデータベース生成に用いることができるか否かを判断することができる。

40

【 0 1 0 1 】

さらに、原因情報を解析して、エラー情報表示領域 1 4 4 に、エラー情報が生じている表データを抽出することが可能な他のデータ抽出規則を検索して表示することが好ましい。データ抽出規則を変更することが可能である場合には、エラーを解消することができる可能性もあり、エラー情報が生じている表データであってもエラーを解消することで新たなデータベースに取り込むことができるからである。

【 0 1 0 2 】

図 1 5 は、他のデータ抽出規則を示唆する旨の情報を表示する表示画面の例示図である。図 1 5 の例では、表データ表示領域 1 5 2 の横に並列して、ゆらぎ情報を抽出できな

50

った元表データ表示領域 155 を表示している。両方の表データ領域 153 と表データ領域 156 とを比較し、いずれの表も表のサイズが最も大きい表であることを検出した場合、エラー情報表示領域 154 には、いずれの表も 1 番大きな表である旨のメッセージを表示する。斯かる情報が表示されることにより、データ抽出規則として「最も大きいサイズの表データ」を採用した場合には、表データ領域 156 に表示されている表データについても新たなデータベースに取り込むことが可能となる。

#### 【0103】

図 16 は、原因情報を表示する場合の本発明の実施の形態 2 に係るデータベース生成装置 1 の CPU 11 の処理手順を示すフローチャートである。図 16 において、データベース生成装置 1 の CPU 11 は、図 11 のステップ S1110 の処理の後、エラー情報が生じた電子文書ファイルが存在するか否かを判断する（ステップ S1601）。

10

#### 【0104】

CPU 11 が、エラー情報が生じた電子文書ファイルが存在しないと判断した場合（ステップ S1601：NO）、CPU 11 は、処理を終了する。CPU 11 が、エラー情報が生じた電子文書ファイルが存在すると判断した場合（ステップ S1601：YES）、CPU 11 は、エラー情報が生じた原因に関する原因情報を表示装置 23 に表示し（ステップ S1602）、他のデータ抽出規則が記憶装置 13 のデータ抽出規則記憶部 132 に記憶されているか否かを判断する（ステップ S1603）。

#### 【0105】

CPU 11 が、他のデータ抽出規則が記憶されていないと判断した場合（ステップ S1603：NO）、CPU 11 は、処理を終了する。CPU 11 が、他のデータ抽出規則が記憶されていると判断した場合（ステップ S1603：YES）、CPU 11 は、記憶されている他のデータ抽出規則から一のデータ抽出規則を選択し（ステップ S1604）、エラー情報が生じた表データから、データベース項目及び対応するデータを抽出する（ステップ S1605）。

20

#### 【0106】

CPU 11 は、再度エラー情報が生じたか否かを判断し（ステップ S1606）、CPU 11 が、エラー情報が生じたと判断した場合（ステップ S1606：YES）、CPU 11 は、記憶されているすべての他のデータ抽出規則を選択したか否かを判断する（ステップ S1607）。CPU 11 が、すべてのデータ抽出規則を選択したと判断した場合（ステップ S1607：YES）、CPU 11 は、処理を終了する。

30

#### 【0107】

CPU 11 が、まだ選択されていない他のデータ抽出規則が存在すると判断した場合（ステップ S1607：NO）、CPU 11 は、次の他のデータ抽出規則を選択し（ステップ S1608）、処理をステップ S1605 に戻して、上述した処理を繰り返す。CPU 11 が、エラー情報が生じていないと判断した場合（ステップ S1606：NO）、CPU 11 は、選択した他のデータ抽出規則に関するメッセージを表示装置 23 に表示する（ステップ S1609）。

#### 【0108】

このようにすることで、複数の表データ間にゆらぎ情報が存在する場合であっても、ゆらぎ情報を補正して新たなデータベースを生成することができ、どのようにゆらぎ補正を実行したか、あるいはゆらぎ補正を実行することができなかった原因に関する原因情報を表示することにより、元の表データ又はデータ抽出規則を変更するための情報を得ることが可能となる。また、エラー情報が生じた電子文書ファイルを特定することができ、エラー情報が生じた原因を視認することができるとともに、データ抽出規則を変更することによりエラー情報が生じることなく新たなデータベース生成を行うことができるか否かを判断することが可能となる。

40

#### 【0109】

なお、本発明は上記実施例に限定されるものではなく、本発明の趣旨の範囲内であれば多種の変更、改良等が可能である。例えば変更すべきデータ抽出規則が見つかった場合、

50

自動的にデータ抽出規則を変更しても良いし、ユーザによる変更指示の入力を受け付けても良い。また、データ抽出規則が変更された場合、自動的に再度データベース生成処理を実行するようにしても良いし、ユーザによるデータベース再生成指示の入力を受け付けても良い。

【図面の簡単な説明】

【0110】

【図1】本発明の実施の形態1に係るデータベース生成装置の構成例を示すブロック図である。

【図2】本発明の実施の形態1に係るデータベース生成装置の機能ブロック図である。

【図3】本発明の実施の形態1に係るデータベース生成装置のCPUのデータベース生成処理の手順を示すフローチャートである。

10

【図4】罫線情報に基づいて表データの抽出を行う処理の例示図である。

【図5】本発明の実施の形態1に係るデータベース生成装置のCPUの手動特定処理の手順を示すフローチャートである。

【図6】範囲指定が必要となる場合の例示図である。

【図7】表データの位置に相違が存在する「表位置ゆらぎ情報」の説明図である。

【図8】表データの項目に相違が存在する「項目ゆらぎ情報」の説明図である。

【図9】本発明の実施の形態1に係るデータベース生成装置のCPUのゆらぎ補正処理の手順を示すフローチャートである。

【図10】本発明の実施の形態2に係るデータベース生成装置の機能ブロック図である。

20

【図11】本発明の実施の形態2に係るデータベース生成装置のCPUのデータベース生成処理の手順を示すフローチャートである。

【図12】抽出されたゆらぎ情報の一覧を表示する表示画面の例示図である。

【図13】ゆらぎ補正の内容を確認するための表示画面の例示図である。

【図14】エラー表示の内容を示す表示画面の例示図である。

【図15】他のデータ抽出規則を示唆する旨の情報を表示する表示画面の例示図である。

【図16】原因情報を表示する場合の本発明の実施の形態2に係るデータベース生成装置のCPUの処理手順を示すフローチャートである。

【符号の説明】

【0111】

30

1 データベース生成装置

2 ネットワーク

11 CPU

12 RAM

13 記憶装置

14 I/Oインタフェース

15 ビデオインタフェース

16 可搬型ディスクドライブ

17 通信インタフェース

18 内部バス

40

23 表示装置

90 可搬型記録媒体

100 コンピュータプログラム

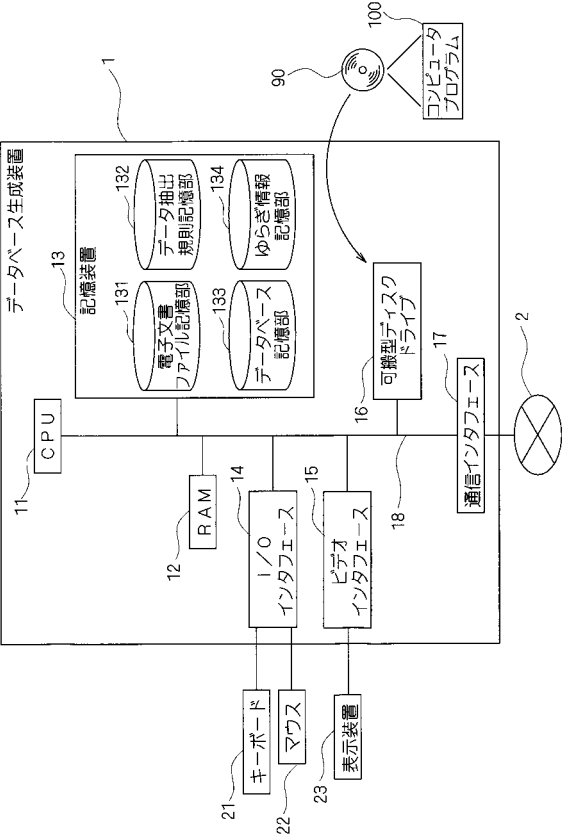
131 電子文書ファイル記憶部

132 データ抽出規則記憶部

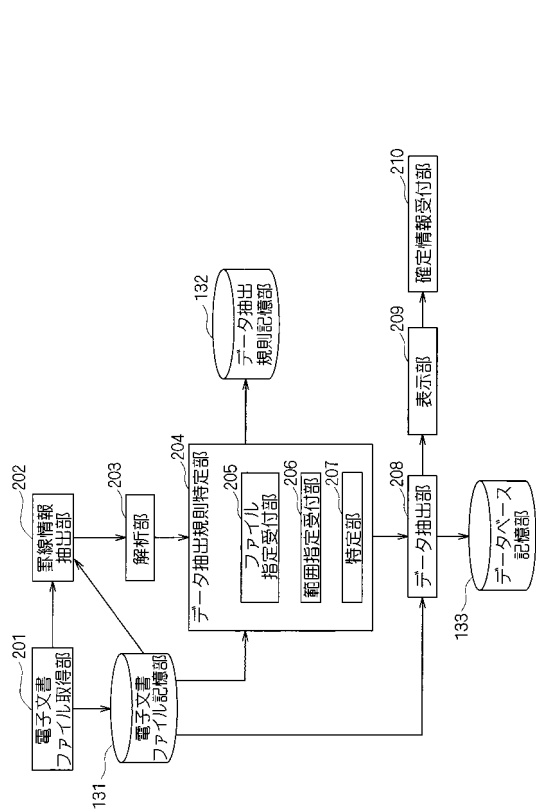
133 データベース記憶部

134 ゆらぎ情報記憶部

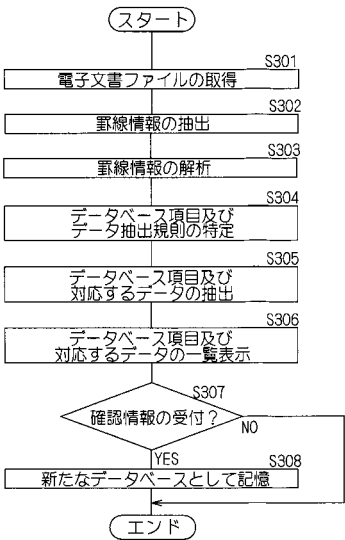
【図 1】



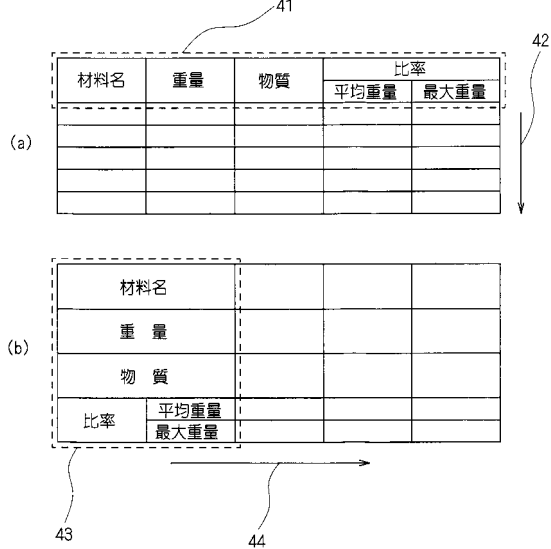
【図 2】



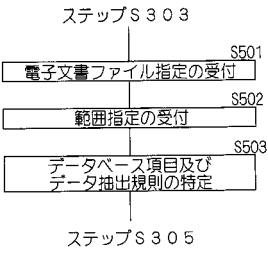
【図 3】



【図 4】



【図 5】



【図 6】

会社名	キーエンス			
	部署	MEQI		
	材料	重量	比率1	比率2

61

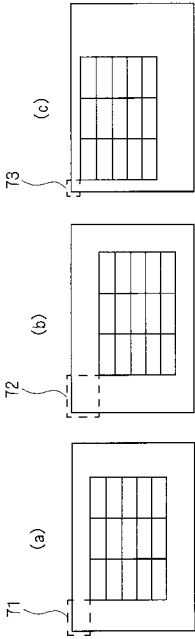
  

材料		比率1	比率2
	重量		

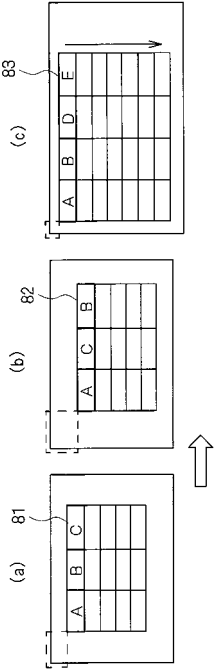
62

63 63 63 63

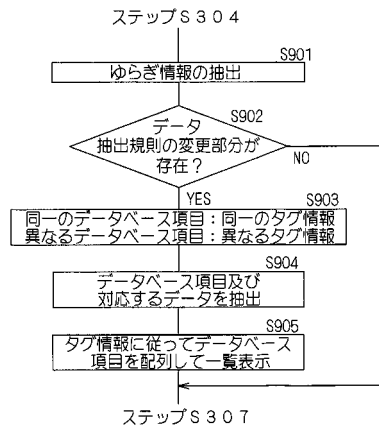
【図 7】



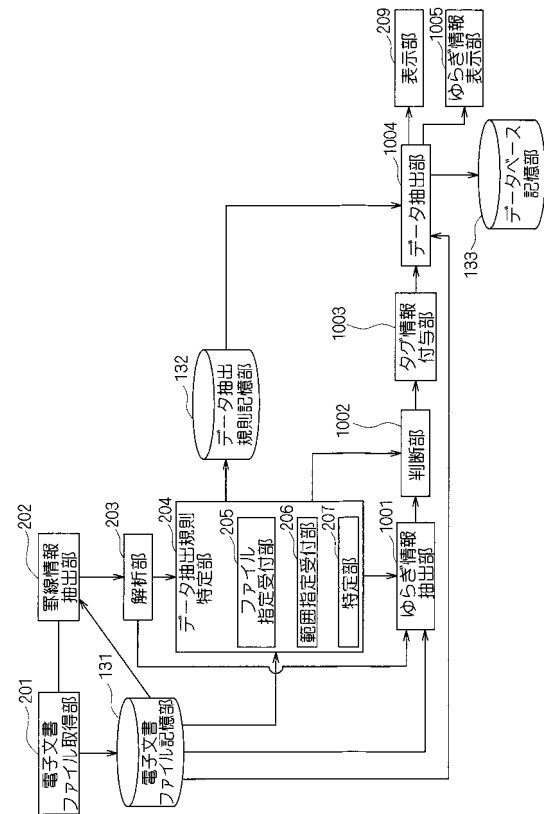
【図 8】



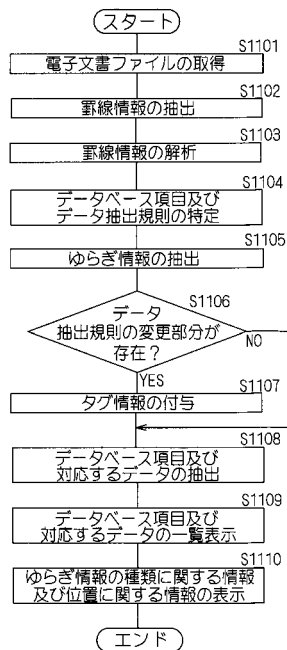
【図 9】



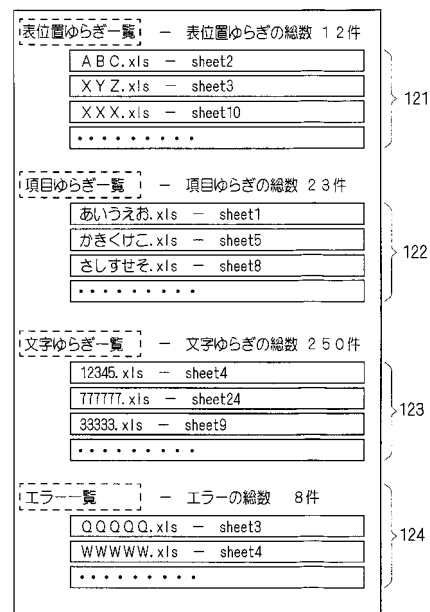
【図 10】



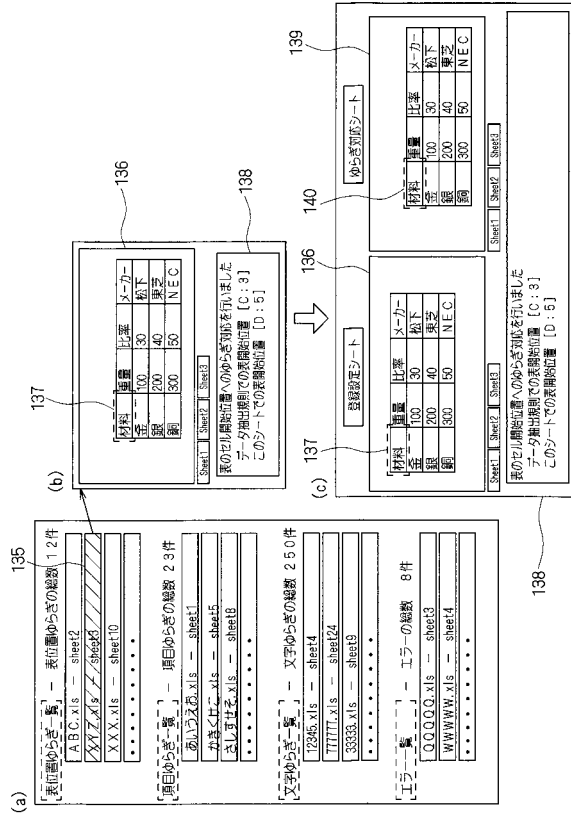
【図 11】



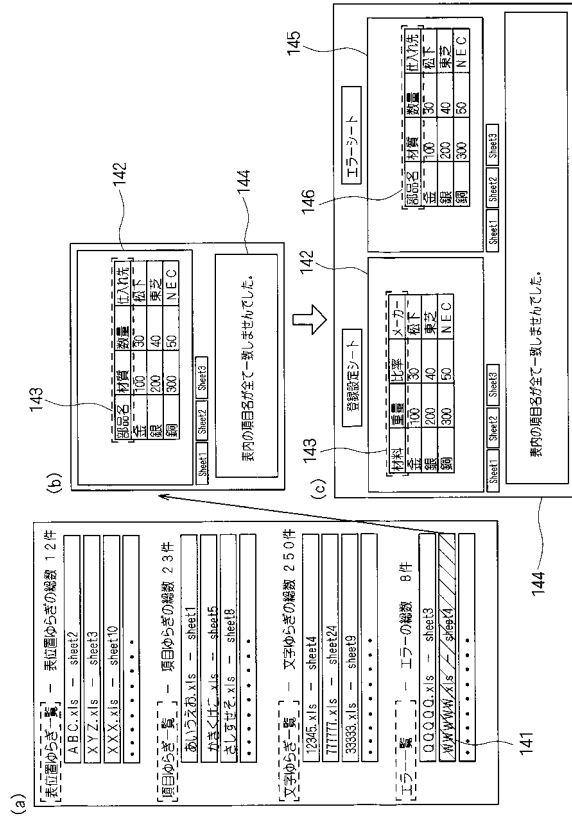
【図 12】



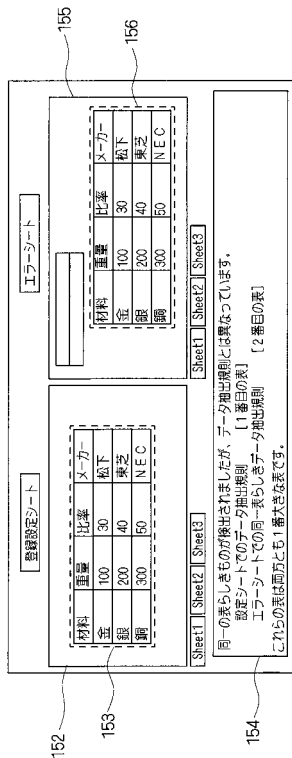
【図 13】



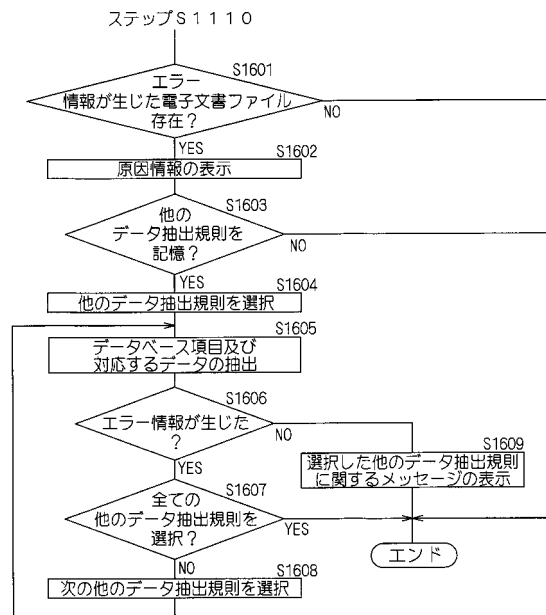
【図 14】



【図 15】



【図 16】



---

フロントページの続き

(56)参考文献 特開2004-252509(JP,A)  
特開2001-331764(JP,A)  
特開平09-050527(JP,A)  
特開2005-242587(JP,A)  
特開平06-243130(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06Q 10/00-50/34  
G06F 12/00  
G06F 17/30