

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5238105号
(P5238105)

(45) 発行日 平成25年7月17日(2013.7.17)

(24) 登録日 平成25年4月5日(2013.4.5)

(51) Int.Cl. F I
G O 6 F 17/30 (2006.01) G O 6 F 17/30 3 4 O Z

請求項の数 4 (全 32 頁)

| | | | |
|---------------|-------------------------------|-----------|---|
| (21) 出願番号 | 特願2007-549011 (P2007-549011) | (73) 特許権者 | 000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番1号 |
| (86) (22) 出願日 | 平成17年12月9日(2005.12.9) | (74) 代理人 | 100074099 弁理士 大菅 義之 |
| (86) 国際出願番号 | PCT/JP2005/022699 | (74) 代理人 | 100133570 弁理士 ▲徳▼永 民雄 |
| (87) 国際公開番号 | W02007/066414 | (72) 発明者 | 松浦 正卓 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 |
| (87) 国際公開日 | 平成19年6月14日(2007.6.14) | (72) 発明者 | 林 宏也 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 |
| 審査請求日 | 平成20年6月13日(2008.6.13) | | |
| 審判番号 | 不服2011-16153 (P2011-16153/J1) | | |
| 審判請求日 | 平成23年7月26日(2011.7.26) | | |

最終頁に続く

(54) 【発明の名称】 プログラム、及びデータ抽出方法

(57) 【特許請求の範囲】

【請求項1】

取得可能なデータのなかから、指定された第一の抽出条件を満たすデータを抽出できるデータ抽出装置を実現させるためにコンピュータに実行させるプログラムであって、

前記データを取得する機能と、

前記第一の抽出条件を入力する機能と、

前記入力する機能により、二つ以上、入力された前記第一の抽出条件それぞれを、該第一の抽出条件それぞれに含まれる複数の部分条件に分割し、該分割によって得られる部分条件を複数の前記第一の抽出条件に共通して含まれる共通条件と該共通条件以外の非共通条件との組み合わせで表現する第二の抽出条件に前記第一の抽出条件それぞれを変換し、前記第二の抽出条件の前記部分条件それぞれを前記データの形式に対応する照合用オートマトンにそれぞれ変換し、前記データより前記照合用オートマトンを用いて前記部分条件単位で該部分条件を満たすデータをそれぞれ抽出し、該抽出したデータそれぞれから、前記第二の抽出条件それぞれを満たすデータを抽出することにより、前記第一の抽出条件を満たすデータを抽出する機能と、

を実現させるためのプログラム。

【請求項2】

請求項1記載のプログラムであって、

前記入力する機能は、前記第一の抽出条件それぞれと併せて、該第一の抽出条件それぞれと対応付けたデータの出力先に関する出力条件をそれぞれ入力することができ、

前記出力条件に従って、該出力条件と対応付けられた抽出条件を満たすデータを出力するプログラム。

【請求項 3】

請求項 1 記載のプログラムであって、

前記抽出する機能により前記第二の抽出条件毎に抽出したデータはそれぞれ異なる出力先に出力するプログラム。

【請求項 4】

取得可能なデータのなかから、指定された第一の抽出条件を満たすデータを抽出するためのデータ抽出方法において、

前記第一の抽出条件を二つ以上、取得した場合に、前記第一の抽出条件それぞれを、該第一の抽出条件それぞれに含まれる複数の部分条件に分割し、該分割によって得られる部分条件を複数の前記第一の抽出条件に共通して含まれる共通条件と該共通条件以外の非共通条件との組み合わせで表現する第二の抽出条件に前記第一の抽出条件それぞれを変換し、前記第二の抽出条件の前記部分条件それぞれを前記データの形式に対応する照合用オートマトンにそれぞれ変換し、前記データより前記照合用オートマトンを用いて前記部分条件毎に該部分条件を満たすデータをそれぞれ抽出して該抽出によって得たデータを記憶ユニットに記憶し、

前記抽出によって得たデータを、前記記憶ユニットから読み出し、該読み出したデータそれぞれから、前記第二の抽出条件それぞれを満たすデータを抽出することにより、前記第一の抽出条件を満たすデータを抽出する、

処理を情報処理装置に実行させることを特徴とするデータ抽出方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、取得可能なデータのなかから指定された抽出条件を満たすデータを抽出するための技術に関する。

【背景技術】

【0002】

取得可能なデータのなかから任意のデータを抽出することができるデータ抽出装置は、現在、様々な用途に広く用いられている。インターネットで公開されている情報の検索では、検索エンジンとして用いられている。ユーザはそのデータ抽出装置を用いることにより、大量のデータのなかから所望のデータを迅速に得ることができる。

【0003】

データ抽出装置は、予め定められた単位でデータを抽出する。その単位となるのは、例えばファイル、或いはレコードである。文書、及びインターネット上の Web ページはファイルに相当する。顧客の利用実績 P O S (Point Of Sales) データや H H T (Hand Held Terminal) データなどはレコード単位で管理されるのが普通である。

【0004】

図 1 は、従来のデータ抽出方法を説明する図である。ここで、図 1 を参照して、そのデータ抽出方法について具体的に説明する。

図 1 に示す従来のデータ抽出方法は、例えばクレジットカード会社で行われる場合のものである。表記した「JOURNAL」は、ファクトデータをレコード単位で格納したジャーナルファイルを表している。「MASTER」は、クレジットカードの所有者である顧客のデータをレコード単位で格納したマスタファイルを表している。それにより、図 1 に示すデータ抽出方法は、SQL (Structured Query Language) を用いて、共に複数、存在するジャーナルファイル、及びマスタファイルのなかから所望のものを連結 (JOIN) させ、その連結結果から所望のレコードを抽出する場合の例を表している。

【0005】

連結させるジャーナルファイル、マスタファイルのそれぞれの条件は、FROM 句内の WHERE 句に記述されている。そこに記述された条件により、マスタファイルは現在の

10

20

30

40

50

ものが選択され、ジャーナルファイルは2004年のものが選択される。そのFROM句内のFROM句には、ファイル間におけるレコードの対応関係はクレジットカードナンバーにより特定することが記述されている。連結結果から抽出されるレコードに格納されるデータの項目は、SELECT句に記述されている。そこに記述された項目は、顧客の指名(V.NAME)、その年齢(V.AGE)、利用回数(V.SALES_NUM)、売上額(V.SALES)である。連結結果から抽出するレコードの条件は、WHERE句に記述されている。そこに記述された条件は、カードの種類がゴールドカード、というものである。このようなことから、2004年に利用し、現在もゴールドカードを持つ顧客のレコードが検索結果として抽出される。

【0006】

連結結果から抽出されるレコードを異ならせるには、WHERE句に記述する抽出条件を変更すれば良い。シルバーカードを持つ顧客のレコードを抽出させるのであれば、例えば図2に示すように、「GOLD」の記述を「SILVER」に変更すれば良い。それにより、2004年に利用し、現在もシルバーカードを持つ顧客のレコードが検索結果として抽出される。

【0007】

このように、従来のデータ抽出方法では、所望のデータを得るための抽出条件を決定し、その抽出条件毎に検索を行わせるようになっていた。このため、データを抽出する目的の数、つまり検索に使用する抽出条件の数が多くなるほど、全ての抽出結果を得るまでに要する時間が長くなり、効率的な作業が行えなくなるという問題点があった。

【0008】

現在、デジタルデータで扱う情報の種類、及びその量は非常に増大しつつある。そのため、今後は従来のデータ抽出方法では対応するのが非常に困難となるのが予想される。このこともあって、膨大なデータのなかからでも必要な種類のデータを全てより迅速に得られるようにすることが重要であると考えられる。

【特許文献1】特開2002-222194号公報

【特許文献2】特開2005-70911号公報

【特許文献3】特開平6-319906号公報

【発明の開示】

【0009】

本発明は、膨大なデータのなかからでも必要な種類のデータを全てより迅速に得られるようにする技術を提供することを目的とする。

本発明の第1、及び第2の態様のプログラムは共に、取得可能なデータのなかから指定された抽出条件を満たすデータを抽出できるデータ抽出装置を実現させるためにコンピュータに実行させることを前提とし、それぞれ以下の機能を実現させる。

【0010】

第1の態様のプログラムは、データを取得する機能と、抽出条件を入力する機能と、入力する機能により一つ以上、入力された抽出条件を用いて、該抽出条件毎にデータを抽出する機能と、抽出する機能により抽出条件毎に抽出されたデータをそれぞれ異なる出力先に出力する機能と、を実現させる。

【0011】

第2の態様のプログラムは、データを取得する機能と、抽出条件を入力する機能と、入力する機能により入力された抽出条件を構成する条件式を複数の部分条件式に分割し、該分割によって得られる部分条件式の組み合わせで表現する形式に該抽出条件を変換して、該部分条件式単位で該部分条件式を満たすか否か確認することにより、取得する機能により取得したデータのなかで該抽出条件を満たすデータを抽出する機能と、を実現させる。

【0012】

本発明のデータ抽出方法は、取得可能なデータのなかから指定された抽出条件を満たすデータを抽出するために適用されることが前提であり、対象となるデータが異なる抽出条件を複数、入力可能とさせ、抽出条件が1つ以上、入力された場合に、該抽出条件毎にデ

10

20

30

40

50

ータの抽出を行い、該抽出によって得たデータを、該データが満たす抽出条件に応じた出力先に出力する。

【0013】

本発明では、対象となるデータが異なる抽出条件を複数、入力可能とさせ、抽出条件が1つ以上、入力された場合に、抽出条件毎にデータの抽出を行い、それによって得たデータを、そのデータが満たす抽出条件に応じた出力先にそれぞれ出力する。このため、ユーザは、複数の抽出条件を定義して入力することにより、1度に複数の抽出結果を得ることができる。それにより、必要な全ての抽出結果をより迅速に得ることができる。この結果、高い作業効率も容易に実現させることができる。

【0014】

本発明では、入力された抽出条件は、それを構成する条件式を複数の部分条件式に分割し、その分割によって得られる部分条件式の組み合わせで表現する形式に変換して、部分条件式単位でその部分条件式を満たすか否か確認することにより、データのなかで抽出条件を満たすデータを抽出する。部分条件式の組み合わせで表現する形式に抽出条件を変換することにより、異なる条件式に同じ部分条件式が存在していても、条件式毎に部分条件式をデータが満たすか否かの確認を行う必要性を回避できるようになる。このため、より小さい負荷でデータ抽出を行えることとなる。

【図面の簡単な説明】

【0015】

【図1】従来のデータ抽出方法を説明する図である。

【図2】従来のデータ抽出方法で異なる種類のデータを抽出させるための抽出条件の相違を説明する図である。

【図3】本実施の形態によるデータ抽出装置の昨日構成を説明する図である。

【図4】本実施の形態によるデータ抽出装置100が可能なデータ抽出を説明する図である。

【図5】本実施の形態によるデータ集計装置を実現できるコンピュータのハードウェア構成の一例を示す図である。

【図6】XMLデータの構成例を説明する図である。

【図7】CSVデータの構成例を説明する図である。

【図8】抽出条件群の内容例を説明する図である。

【図9】タグDFA例を説明する図である。

【図10】階層照合NFA例を説明する図である。

【図11】CSV解析DFA例を説明する図である。

【図12】キーワードDFA例を説明する図である。

【図13】論理テーブル例を説明する図である。

【図14】出力バッファの管理方法を説明する図である。

【図15】抽出条件入力部110が実行する処理のフローチャートである。

【図16】データ入力構造検索部120が実行する処理のフローチャートである。

【図17】抽出条件判定部130が実行する処理のフローチャートである。

【図18】データ判定部140が実行する処理のフローチャートである。

【図19】本実施の形態によるデータ抽出装置の適用例を説明する図である(その1)。

【図20】本実施の形態によるデータ抽出装置の適用例を説明する図である(その2)。

【図21】本実施の形態によるデータ抽出装置の適用例を説明する図である(その3)。

【図22】本実施の形態によるデータ抽出装置の適用例を説明する図である(その4)。

【図23】本実施の形態によるデータ抽出装置の適用例を説明する図である(その5)。

【図24】本実施の形態によるデータ抽出装置の適用例を説明する図である(その6)。

【発明を実施するための最良の形態】

【0016】

以下、本発明の実施の形態について、図面を参照しながら詳細に説明する。

図3は、本実施の形態によるデータ抽出装置の機能構成を説明する図である。

10

20

30

40

50

そのデータ抽出装置 100 は、入力装置 210 からデータ 211 としてテキストデータを入力し、そのデータ 211 を指定された抽出条件群 220 により振り分けて出力するものとして実現されている。そのために、抽出条件入力部 110、データ入力構造検索部 120、抽出条件判定部 130、データ判定部 140、外部出力用の出力バッファ 150、及びデータ出力部 160 を備えている。ここでは便宜的に、入力装置 210 から入力するデータ 211 として、図 6 に示すような XML (eXtensible Markup Language) データ、及び図 7 に示すような CSV (Comma Separated Values) データのみを想定する。それらのデータは共にテキストデータである。

【0017】

抽出条件入力部 110 によって入力される抽出条件群 220 は、例えば図 8 に示すような内容のものである。その図 8 では、(1) ~ (3) に分けてそれぞれ抽出条件、及び出力条件を示している。そのように分けて示す抽出条件は全て、ユーザが所望のデータ 211 を抽出するためのものである。抽出条件と併せて示す出力条件は、その抽出条件によって抽出されるデータ 211 の出力先、及びそのファイル名を指定するものである。それにより、抽出条件群 220 は、所望のデータ 211 別に、そのデータ 211 が満たすべき抽出条件、及びその出力先ファイル名を指定するものとなっている。そのようにデータ 211 の出力先を任意に指定できるようにしたのは、データ 211 をより迅速に所望の形で利用するのを可能とさせるためである。以降、(1) に記述された抽出条件は「抽出条件 1」と表記する。これは他でも同様である。

【0018】

図 4 は、本実施の形態によるデータ抽出装置 100 が可能なデータ抽出を説明する図である。ここで図 4 を参照して、そのデータ抽出について具体的に説明する。

図 8 に示す抽出条件群 220 は、データ 211 として XML データを想定したものである。図 4 では、CSV データを想定した抽出条件群 220 を示している。「Query」は抽出条件に相当し、「OutFile」は出力条件に相当する。Query (抽出条件) として表記した「\$X」は、項目名「X」を表し、「\$__」は任意の項目名を表している。それにより、例えば Query 1 で表記した「\$X == 'X1' OR \$X == 'Xa」は、項目名「X」のデータが X1 または Xa であるデータ 211 が抽出の対象であることを示している。その表記が「\$__ == 'Xa」となっている Query では、任意の項目のデータとして Xa が存在しているデータ 211 が抽出の対象であることを示している。そのデータ 211 は XML データ、及び CSV データの何れであっても、ファイルとしてまとめて入力させても良いが、一つずつ順次、入力させても良い。一つずつ入力させる場合、XML データでは図 6 に示すようなものとなり、CSV データでは、図 7 において、先頭に「000001」~「000007」を表記した行のようなものとなる。ここでは便宜的に、それらのデータのまとまりをレコードと呼ぶことにする。また、2つの「'」の間に記述された文字列については「キーワード」と呼ぶことにする。そのキーワードは、図 8 に示す抽出条件群 220 では 2つの「"」の間に記述された文字列が相当する。

【0019】

本実施の形態では、文字列照合方式を用いて、抽出条件群 220 で指定された抽出条件の何れかを満たすデータ 211 を抽出し、満たす抽出条件に対応付けられた出力条件で指定された出力先ファイル名のファイルに出力する。それにより、Query 1 を満たすデータ 211 はファイル名「result1.csv」のファイル 231 として、Query 2 を満たすデータ 211 はファイル名「result2.csv」のファイル 232 として、Query 3 を満たすデータ 211 はファイル名「result3.csv」のファイル 233 として、それぞれ出力される。入力されたデータ 211 とファイル 231 ~ 3 の何れかに出力されるデータ 211 の対応関係は、図中に表記の(1) ~ (6)により示している。

【0020】

各抽出条件はそれぞれ単独で考慮されるため、抽出条件は全て任意に定義することがで

10

20

30

40

50

きる。このため、XMLデータやCSVデータなどのデータ211の種類毎に1つ以上の抽出条件を定義することもでき、また、その構造別に1つ以上の抽出条件を定義することもできるようになっている。従って、対象とするデータ211間でスキーマがどのように相違していても、その相違の影響は確実に回避させることができる。

【0021】

上述したようなことから、抽出条件間は排他関係としなくとも良い。それにより、Query1とQuery2では条件式(論理式)「 $X = 'Xa'$ 」を満たすデータ211をそれぞれ抽出する内容となっている。同様にQuery2とQuery3では条件式「 $X = 'Xb'$ 」を満たすデータをそれぞれ抽出する内容となっている。この結果、ファイル231、232には共に(4)を表記したデータ211が出力され、ファイル232、233には共に(5)を表記したデータ211が出力されている。

10

【0022】

このように、抽出条件群220により複数の抽出条件が指定されると、抽出条件毎にそれを満たすデータ211を振り分けて指定の出力先に出力するようになっている。このため、ユーザは、抽出条件群220として複数の抽出条件、及び出力条件を定義するだけで1度に複数の抽出結果を得ることができる。それにより、必要な全ての抽出結果はより迅速に得ることができる。この結果、高い作業効率も容易に実現させることができる。

【0023】

上述したように、本実施の形態では文字列照合方式を採用している。その文字列照合方式は、抽出条件で指定した文字列と対象のデータ211との照合を、そのデータ211の先頭より後方に向かって逐次、行っていくことにより、その文字列がデータ211中に存在するか否かを調べるものである。その文字列照合方式では、先頭より後方に向かった走査を1回、行うだけで、抽出条件群220で定義された抽出条件の何れをデータ211が満たしているか確認することができる。そのため、定義された抽出条件の数に係わらず、常に迅速に抽出すべきデータ211を抽出することができる。その参考文献としては、例えば特許文献1、及び2が挙げられる。

20

【0024】

図3の説明に戻る。

抽出条件入力部110は、上述したような抽出条件群220を入力し、抽出条件毎に、その抽出条件を解析して対応のオートマトンを生成する。それにより、抽出条件がXMLデータ用のものであればタグDFA(Deterministic Finite state Automaton)170、階層照合NFA(Non-deterministic Finite state Automaton)171、及びキーワードDFA180が生成される。抽出条件がCSVデータ用のものであればCSV解析DFA172、及びキーワードDFA180が生成される。論理テーブル190は、キーワードDFA172と同様に、抽出条件が想定するデータ211の種類に係わらず生成される。

30

【0025】

抽出条件群220の作成は基本的に、ユーザによるデータ入力によって行われる。本実施の形態によるデータ抽出装置100と接続された端末装置で抽出条件群220を作成する場合、例えばユーザは抽出条件群220作成用の画面を表示させ、その画面上に所望の内容の抽出条件群220を入力する。その入力後、データ抽出を指示すると、作成された抽出条件群220がデータ抽出装置100に出力される。

40

【0026】

上記論理テーブル190としては、抽出条件群220が図8に示す内容であった場合、抽出条件入力部110によって図13に示すようなものが生成される。図13に示すように、その論理テーブル190は、A論理テーブル190a、及びZ論理テーブル190bから構成されている。

【0027】

A論理テーブル190aは、抽出条件を構成する条件式(論理式)を関係演算子(図8中では「 $=$ 」及び「 $<$ 」が相当)で分解して、その条件式が表現する論理により細分化し(図8では抽出条件2を構成する条件式「 $/root/Company/code <$

50

「99」は「/root/Company/code」「<99」に分解される)、細分化した条件式(部分条件式)毎に固有の論理番号を付した構成のものである。Z論理テーブル190bは、条件式、或いは抽出条件を部分条件式、或いは条件式に付した論理番号の組み合わせで表現し、表現した組み合わせ毎に固有の論理番号を付した構成のものである。組み合わせる論理番号はA論理テーブル190a、及びZ論理テーブル190bの何れのものであっても良い。その論理番号を用いて条件式、或いは抽出条件を表現することにより、A論理テーブル190a、或いはZ論理テーブル190bで参照すべきレコード(行)を特定できるようにさせている。特に図示していないが、そのZ論理テーブル190bには、論理番号の組み合わせ毎に、その組み合わせで表現される条件式、或いは抽出条件が成立しているか否かを示す符号を格納できるようになっている。以降テーブル190a、及び190bでそれぞれ割り当てる論理番号を区別するために、A論理テーブル190aの論理番号には「A」、Z論理テーブル190bの論理には「Z」をそれぞれ先頭に付して表記する。

10

【0028】

Z論理テーブル190bで論理番号Z1が割り当てられた組み合わせは「A1xA2」である。その組み合わせ「A1xA2」は、論理番号A1の部分条件式(/root/origin)が成立し、且つ論理番号A2の部分条件式("atcg")が成立するデータ211が抽出対象であることを示す形式の論理式となっている。それにより、組み合わせ(論理式)「A1xA2」中の「x」は、論理番号A1、及びA2の部分条件式の論理積を行うことを示す論理演算子となっている。その論理式は、抽出条件1の内容を表している。同様に、論理番号Z4、及びZ5の各論理式はそれぞれ抽出条件3、及び2の内容を表している。抽出条件2は $Z5 = Z2 \times Z3$ になっている。ここで190bのテーブル内で、 $Z2 = A3 \times A4$ により $A3 = /root/Company/code$ 、 $A4 = <99$ に対応する。

20

【0029】

また、 $Z3 = A1 \times A5$ により、 $A1 = /root/origin$ 、 $A5 = "gtac"$ に対応する。したがって、抽出条件2は、Z論理番号Z5と介して、A論理番号A3、A4、A1、A5に対応し、図8で示す抽出条件2の論理積(AND)は、図13で示す論理テーブルとその要素間のリンク状態で示される。図8の抽出条件3は図13の抽出条件3、Z論理番号4、A論理番号A1、A6の論理テーブルとその要素間のリンクで示される。すなわち、抽出条件3は $Z4 = A1 \times A6$ ($A1 = /root/origin$ 、 $A6 = "aacg"$)としてA論理番号に対応している。すなわち、このような論理番号によって各抽出条件で形成される論理テーブルを使って抽出条件毎のデータ判別が可能となる。

30

【0030】

図13に示す検索結果判定情報195は、抽出条件毎に、その抽出条件を表現する論理番号の組み合わせに対して付された論理番号、その抽出条件を満たすデータ211を格納すべき出力バッファ150を示す番号(図中「出力バッファNo.」と表記)、及びファイルディスクリプタ(対応付けられた出力条件)がまとめられたものである。それにより、何れかの抽出条件を満たすデータ211は、検索結果判定情報195を参照して出力すべき出力バッファ150に出力された後、出力すべきファイルに出力される。

40

【0031】

上記オートマトン(タグDFA170、階層照合NFA171、キーワードDFA180、CSV解析DFA172)は検索条件中の文字列をデータ211と照合するための状態遷移テーブルである。状態間は遷移の方向を示す矢印で結んで表現される。先頭を初期状態とし、この初期状態からデータ211中の文字列に応じて順次、状態を遷移させる。遷移させる状態には、検索条件中の文字列の最後に位置する文字に相当する受理状態が1つ以上、含まれている。それによりオートマトンは、データ211中に検出すべき文字列が存在していれば、何れかの受理状態に遷移するように生成される。受理状態に遷移した場合、その受理状態に応じたヒット情報を出力するようになっている。そのヒット情報は

50

、遷移した受理状態に応じた特有のものであり、オートマトンの生成時に併せて生成される。

【0032】

上記タグDFA170は、キーワードと照合すべき文字列（要素内容）が存在する要素までの検索パスを検出するためのものである。抽出条件群220が図8に示す内容であった場合、抽出条件入力部110によって図9に示すようなタグDFA170が最終的に生成される。図8に示す抽出条件群220では、検索パスとして「/root/origin」及び「/root/Company/code」が存在することから、それぞれがタグ名である文字列「root」「origin」「Company」及び「code」をそれぞれ検出できるように生成されている。それらの文字列の最後に位置する文字「t」「n」「y」及び「e」の何れかに相当する受理状態まで遷移することで、その文字に対応する文字列が検出されたことを示すヒット情報170a～dの何れかが出力される。

10

【0033】

階層照合NFA171は、現在、対象とする検索パスを管理するためのものである。抽出条件群220が図8に示す内容であった場合、抽出条件入力部110によって図10に示すような階層照合NFA171が最終的に生成される。そのNFA171は、図10に示すように、何れかの検索パスに記述されたタグ名を単位とした状態遷移が行われるように生成されている。このため、その状態遷移は開始タグ、及び終了タグによって発生する。ここでは、「4」、及び「2」を表記した状態が受理状態に相当する。

【0034】

「4」を表記した受理状態に遷移したことは、検索パス「/root/Company/code」が検出されたことを意味する。それにより、その検索パスで指定されたノードでは、その値が99未満か否か、つまり論理番号A4の部分条件式（論理）が成立するか否かの照合を行うためのヒット情報171aが出力される。そのヒット情報171aは、照合の対象となる部分条件式を示す論理番号（ここではA4）、検索パスの階層の深さを示す階層情報、及びその部分条件式で関係を確認すべき内容を示す比較情報（ここでは<99）を含むものである。同様に「2」を表記した受理状態に遷移したことは、検索パス「/root/origin」が検出されたことを意味するから、その検索パスで指定されたノード、つまりタグ名「origin」のタグでは、その文字列が「atcg」「gtac」或いは「aacg」の何れと一致するか否かの照合を行うためのヒット情報171b～dが出力される。それらのヒット情報171b～dで比較情報を示していないのは、それらに表記した論理番号に対応する部分条件式の照合はキーワードDFA180により行うためである。

20

30

【0035】

階層照合NFA171における状態遷移は、図9に示すタグDFA170を用いて行われる。例えばタグ名である文字列「root」をタグDFA170により検出すると、つまりタグDFA170によりヒット情報170aを出力すると、NFA171では「0」を表記した初期状態から「1」を表記した状態に遷移する。次にタグDFA170により文字列「origin」を検出すると、NFA171では「1」を表記した状態から「2」を表記した状態に遷移する。このとき、タグDFA170により文字列「Company」を検出すると、NFA171では「1」を表記した状態から「3」を表記した状態に遷移する。それらの何れの文字列もタグDFA170により検出できなければ、NFA171では「1」を表記した状態から「0」を表記した初期状態に遷移する。そのように遷移させることにより、階層照合NFA171を用いて検索パスに沿った階層の移動の有無を把握し、対象とする検索パスを管理する。

40

【0036】

CSV解析DFA172は、キーワードと照合すべき文字列（要素内容）が存在する要素までの検索パスを検出するためのものである。その要素が2つのダブルコーテーション間に存在するCSVデータ（図7）では、抽出条件入力部110によって図11に示すようなCSV解析DFA172が生成される。図11中に表記した「0x」はそれに続くシ

50

ンボルが16進数表現であることを表している。

【0037】

キーワードDFA180は、抽出条件により指定されたキーワードと一致する文字列をデータ211中から検出するためのものである。抽出条件群220が図8に示す内容であった場合、抽出条件入力部110によって図12に示すようなキーワードDFA180が最終的に生成される。それに登録された何れかのキーワードの最後に位置する文字に相当する受理状態まで遷移した場合、つまり文字列「aacg」「acgt」及び「gtac」の何れかを検出できた場合、検出された文字列に応じてヒット情報180a~cの何れかが出力される。

【0038】

データ入力構造検索部120は、入力装置210から所定量ずつ連続的にデータ211を入力し、そのデータ211の種類に応じて、照合に用いるオートマトンを決定する。それにより、データ211がXMLデータであれば、タグDFA170、及び階層照合NFA171を用いて抽出条件の何れかに記述された検索パスの検出を行う。データ211がCSVデータであれば、CSV解析DFA172を用いて抽出条件の何れかに記述された項目名の検出を行う。検索パス、或いは項目名を検出すると、その検索パスによって指定されたノード、或いはその項目名のセルが開始する位置を示すデータ位置情報、及び検出された文字列を示すノード・セル情報を抽出条件判定部130に通知する。それらの情報は例えばヒット情報として生成するものか、或いはそれを含むものである。それらの情報の通知は、データ211の終端を検出するまで、検索パス、或いは項目名を検出する度に
20
行う。その終端の検出は、XMLデータではルートタグと組になる終了タグの検出に相当し、CSVデータでは所定個数のセルの検出に相当する。データ入力構造検索部120による検索パス、或いは項目名の検出は、A論理テーブル190aに格納された部分条件式が成立することの確認に相当する。

【0039】

抽出条件判定部130は、データ入力構造検索部120から通知されたデータ位置情報が示すデータ位置より、キーワードDFA180を用いた照合を行う。その照合の結果、そのデータ位置から何れかのキーワードと一致する文字列、或いは関係演算子が示す関係を満たす値(図8に示す抽出条件群220では99未満の値)が存在することを確認すると、Z論理テーブル190bの該当論理番号の箇所にそのことを示す符号(以降「真符号」と表記し、それと異なる符号を「偽符号」と表記する)を格納する。その確認ができる前にデータ211の終端を検出した場合には、その終端の位置を示すデータ位置情報をデータ入力構造検索部120に通知する。それにより、構造検索部120は、データ211の終端を自身が検出したか否かに係わらず、その終端まで走査が終了したことをデータ判定部140に通知する。

【0040】

抽出条件判定部130は、上記通知を行うか、或いは構造検索部120が終端を検出するまで、構造検索部120から情報が通知される度にキーワードDFA180を用いた照合を行う。この結果、データ211が抽出条件2を満たしている場合には、論理番号Z2、及びZ3の符号として真符号が順次、格納され、最後に論理番号Z5の符号として真符号が格納されることになる。そのようにして、対象とするデータ211が論理式を満たす論理番号の箇所にのみ真符号が格納されることから、Z論理テーブル190bを参照することにより、データ211が満たす抽出条件を確認できるようになっている。

【0041】

このようにして本実施の形態では、抽出条件を構成する条件式をそれが表現する論理により細分化し、その細分化によって得られた部分条件式(細分化論理)単位での照合を行うようにしている。それにより、一致する文字列、或いは検索パスの検出、関係演算子で表す関係の確認、及びそのようなことを行うべき箇所の特定、などをそれぞれ個別に実施している。そのようにすると、より柔軟に対応することが可能となり、データ211の種類やその構造などの情報がたとえ不足していたとしても、ユーザは得られている情報から
50

所望のデータ 211 が満たす内容を抽出条件としてより容易に定義できるようになる。このため、ユーザにとっての高い利便性が実現される。

【0042】

部分条件式（細分化論理）は、同じ、或いは他の抽出条件で別に存在する場合がある。図 8 に示す例では、部分条件式「/ root / origin」は抽出条件 1 ~ 3 の何れにも記述されている。しかし、そのような複数の同じ記述は、条件式を細分化することにより、一つの部分条件式として残せば済むようになる。それにより、抽出条件の数や内容に係わらず、成立するか否か確認すべき部分条件式は必要最小限に抑えることができる。条件式、或いは抽出条件は複数の部分条件式の組み合わせで表現される。このため、それらが成立するか否かはより迅速に行えることとなる。

10

【0043】

データ判定部 140 は、Z 論理テーブル 190b を参照して、データ 211 が満たす抽出条件を確認する。その確認により、何れかの抽出条件を満たしていることが判明すると、検索結果判定情報 195（図 13）を参照して、出力すべき出力バッファ 150 にデータ 211 を出力して格納する。

【0044】

図 14 は、出力バッファの管理方法を説明する図である。

データ 211 を対応する出力バッファ 150 への出力は、出力バッファ情報 151、及びバッファ情報 152 により管理している。出力バッファ情報 151 は、抽出条件群 220 により確保した出力バッファ 150 の数を示す取得バッファ数情報、及びバッファ情報 152 にアクセスするためのポインタ情報を備えている。そのバッファ情報 152 は、取得バッファ数情報が示す数のレコードを備えたものであり、各レコードには、対応する出力バッファ 150（ここでは出力バッファ 150a ~ c のうちのの一つ）に関する複数の情報を有する個別バッファ情報 153（ここでは個別バッファ情報 153a ~ c のうちのの一つ）がそれぞれ格納されている。それら出力バッファ情報 151、及びバッファ情報 152 を格納するエリアは出力バッファ 150 と共に、データ抽出装置 100 に搭載、或いは接続された記憶装置 1401 上に確保されている。タグ D F A 170、階層照合 N F A 171、CSV 解析 D F A 172、キーワード D F A 180、及び論理テーブル 190 も例えばその記憶装置 1401 に格納される。

20

【0045】

その個別バッファ情報 153 は、対応する出力バッファ 150 にアクセスするためのポインタ情報、そのデータ 211 を格納可能な全サイズを表す全バッファサイズ、そのサイズのなかでデータ 211 を格納可能な残りのサイズを表す残バッファサイズ、確保した出力バッファ 150 自体のサイズを表す出力バッファサイズ、を有している。各レコードに付した番号の大小関係は抽出条件の番号のそれと同じとさせている。つまり、レコード番号 0 のレコードは抽出条件 1 に対応している。それにより、データ 211 が満たす抽出条件に対応するレコードを特定できるようにさせている。

30

【0046】

上述したようなことから、データ判定部 140 は、Z 論理テーブル 190b を参照してデータ 211 が満たす抽出条件が存在していることを確認すると、検索結果判定情報 195 を参照してその抽出条件を確認し、出力バッファ情報 151、及びバッファ情報 152 を参照する。それにより、確認した抽出条件に対応するレコードをバッファ情報 152 から取り出し、そのレコードに格納された個別バッファ情報 153 により指定される出力バッファ 150 にデータ 211 を出力する。残バッファサイズは、出力するデータ 211 のサイズにより更新する。

40

【0047】

データ出力部 160 は、各出力バッファ 150 の例えば残バッファサイズを監視し、そのサイズが所定値以下になるか、或いは入力装置 210 から入力して処理するデータ 211 が無くなった場合に、検索結果判定情報 195 を参照して、出力バッファ 150 に格納されているデータ 211 を対応するファイルに出力する。それにより、出力条件で指定さ

50

れた出力先ファイル名のファイルに、これまでに抽出したデータ 2 1 1 を保存する。ここでは、3つのファイル 2 3 1 ~ 2 3 3 は共に同じ出力装置 2 3 0 上に保存させている。

【 0 0 4 8 】

図 5 は、データ抽出装置 1 0 0 を実現できるコンピュータのハードウェア構成の一例を示す図である。抽出装置 1 0 0 は複数のコンピュータ（データ処理装置）により実現させても良いが、ここでは図 5 に構成を示す 1 台のコンピュータによって実現されていることを前提として説明することとする。

【 0 0 4 9 】

図 5 に示すコンピュータは、CPU 5 1、メモリ 5 2、入力装置 5 3、出力装置 5 4、外部記憶装置 5 5、媒体駆動装置 5 6、及びネットワーク接続装置 5 7 を有し、これらが 10
バス 5 8 によって互いに接続された構成となっている。同図に示す構成は一例であり、これに限定されるものではない。

【 0 0 5 0 】

メモリ 5 2 は、データを一時的に格納する RAM 等のメモリである。外部記憶装置 5 5、若しくは媒体駆動装置 5 6 がアクセスする可搬記録媒体 MD に記憶されているプログラム、あるいはデータが一時的に格納される。CPU 5 1 は、プログラムをメモリ 5 2 に読み出して実行することにより、全体の制御を行う。そのプログラムは、ネットワーク接続装置 5 7 によりネットワークを介して取得したものであっても良い。

【 0 0 5 1 】

入力装置 5 3 は、例えば、キーボード、マウス等の入力機器と接続されているか、或いはそれらを有するものである。そのような入力機器に対するユーザの操作を検出し、その検出結果を CPU 5 1 に通知する。 20

【 0 0 5 2 】

出力装置 5 4 は、例えばディスプレイと接続されているか、或いはそれを有するものである。CPU 5 1 の制御によって送られてくるデータをディスプレイ上に出力させる。

ネットワーク接続装置 5 7 は、例えばイントラネットやインターネット等のネットワークを介して、他の装置と通信を行うためのものである。外部記憶装置 5 5 は、例えばハードディスク装置である。主に各種データやプログラムの保存に用いられる。

【 0 0 5 3 】

記憶媒体駆動装置 5 6 は、フレキシブル・ディスク、光ディスク（ここでは CD - R O M、CD - R、及び DVD 等を含む）、或いは光磁気ディスク等の可搬型の記録媒体 MD にアクセスするものである。 30

【 0 0 5 4 】

図 3 に示す出力装置 2 3 0 は、図 5 に示す構成では外部記憶装置 5 5、記録媒体 MD が装着された媒体駆動装置 5 6、或いはネットワーク接続装置 5 7 によりアクセス可能な外部装置に相当する。入力装置 2 1 0 は、記録媒体 MD が装着された媒体駆動装置 5 6、或いはネットワーク接続装置 5 7 によりアクセス可能な外部装置に相当する。抽出条件群 2 2 0 の入力は、入力装置 5 3、記録媒体 MD が装着された媒体駆動装置 5 6、或いはネットワーク接続装置 5 7 により行うことができる。図 1 4 に示す記憶装置 1 4 0 1 は、例えば外部記憶装置 5 5、及びメモリ 5 2 の少なくとも一方に相当する。 40

【 0 0 5 5 】

検索条件入力部 1 1 0 は、例えば出力装置 5 4 を除く各部 5 1 ~ 5 3、及び 5 5 ~ 5 8 によって実現される。データ入力構造検索部 1 2 0、及びデータ出力部 1 6 0 は共に、例えば入力装置 5 3、及び出力装置 5 4 を除く各部 5 1、5 2、及び 5 5 ~ 5 7 によって実現される。抽出条件判定部 1 3 0、及びデータ判定部 1 4 0 は共に、例えば入力装置 5 3、出力装置 5 4、及びネットワーク接続装置 5 7 を除く各部 5 1、5 2、5 5、5 6、及び 5 8 によって実現される。

【 0 0 5 6 】

次に、上述した各部 1 1 0、1 2 0、1 3 0、及び 1 4 0 の動作について、図 1 5 ~ 図 1 8 に示す各処理のフローチャートを参照して詳細に説明する。それらの処理は何れも、 50

例えばCPU 51が、外部記憶装置55、若しくは媒体駆動装置56に装着された可搬記録媒体MDに記憶されているプログラムをメモリ52に読み出して実行することにより実現される。

【0057】

図15は、抽出条件入力部110が実行する処理のフローチャートである。始めに図15を参照して、その処理について詳細に説明する。その処理は、例えば抽出条件群220の入力をユーザが入力装置53、或いはネットワークを介して指示することで起動される。その場合、抽出条件群220は入力装置53、或いはネットワーク接続装置57を介して入力される。

【0058】

先ず、ステップ11では、抽出条件群220を入力し、例えばメモリ52に保存する。続くステップ12では、保存した抽出条件群220のなかから1抽出条件を選択して読み出し、それを解析して対応するオートマトンの種類を特定する。その次に移行するステップ13では、特定した種類のオートマトンを生成、或いは更新する。その生成、或いは更新により、抽出条件に記述された文字列が必要に応じてタグDFA170、階層照合NFA171、或いはキーワードDFA180に登録される。

【0059】

ステップ13に続くステップ14では、抽出条件群220のなかに選択していない他の抽出条件が有るか否か判定する。そのような抽出条件が残っていた場合、判定はYESとなって上記ステップ12に戻り、他の選択条件を選択する。そうでない場合には、判定はNOとなり、ステップ15で論理テーブル190の生成と併せて検索結果判定情報195(図13)、出力バッファ情報151、及びバッファ情報152の生成を行い、抽出条件数に応じた出力バッファ150(図14)の確保を行った後、一連の処理を終了する。このようにして、抽出条件群220の入力により、必要なオートマトンの生成に併せて、データ211を出力すべき出力先に出力するための準備が行われる。

【0060】

図16は、データ入力構造検索部120が実行する処理のフローチャートである。次に図16を参照して、その処理について詳細に説明する。その処理は、例えばデータ211の入力装置210からの取り込みが指示されている間、実行される。

【0061】

先ず、ステップ21では、入力装置210から入力すべきデータ211が有るか否か判定する。そのようなデータ211が無かった場合、判定はNOとなり、再度、その判定を行う。それにより、そのデータ211が生じるのを待つ。一方、そうでない場合には、判定はYESとなってステップ22に移行する。

【0062】

ステップ22では、入力装置210から所定量のデータ211を入力する。続くステップ23では、入力したデータ211から一つを選択し、抽出条件入力部110によって決定したオートマトンを用いて、それに登録された文字列の何れかと一致する文字列の検索を行う。

【0063】

その検索は1文字単位で行い、その検索が終了するとステップ24に移行して、対象となる文字列(検索パス、項目名、など)を検出できたか否か判定する。そのような文字列を検出できなかった場合、判定はNOとなってステップ27に移行する。そうでない場合には、判定はYESとなってステップ25に移行する。

【0064】

ステップ25では、データ位置情報等を抽出条件判定部130に通知する。その通知により、抽出条件判定部13はキーワードDFA180を用いた照合を行い、その照合によってデータ211の終端を検出すると、そのデータ位置情報を通知する。このことから、次のステップ26では、その通知が有ったか否か判定する。その通知が有った場合、判定はYESとなってステップ28に移行する。そうでない場合には、判定はNOとなって上

10

20

30

40

50

記ステップ 23 に戻り、検索を続行する。

【0065】

上記ステップ 24 の判定が NO となって移行するステップ 27 では、検索によってデータ 211 の終端を検出したか否か判定する。その終端を検出した場合、判定は YES となってステップ 28 に移行する。そうでない場合には、判定は NO となって上記ステップ 23 に戻り、検索を続行する。

【0066】

ステップ 28 では、データ 211 の終端が検出されたことをデータ判定部 140 に通知する。続くステップ 29 では、入力したデータ 211 のなかで未選択のデータ 211 が有るか否か判定する。未選択のデータ 211 が存在する場合、判定は YES となって上記ステップ 23 に戻り、未選択のデータ 211 を選択して検索を開始する。そうでない場合には、判定は NO となって上記ステップ 21 に戻る。それにより、入力装置 210 に入力すべきデータ 211 が有るか否かの確認を行う。

【0067】

図 17 は、抽出条件判定部 130 が実行する処理のフローチャートである。次に図 17 を参照して、その処理について詳細に説明する。

まず、ステップ 41 では、レコードの終了通知が通知されるのを待つ。その通知を受け取ると、判定が NO となってステップ 42 に移行し、通知されたデータ位置情報、及びキーワード D F A 180 を用いた照合を行う。その次に移行するステップ 43 では、キーワード D F A 180 に登録されたキーワードの何れかと一致する文字列をデータ 211 から検出できたか否か判定する。そのような文字列を検出できた場合、判定は YES となり、ステップ 44 で論理テーブル 190 (Z 論理テーブル 190 b) の該当論理番号の箇所に真符号を設定した後、上記ステップ 41 に戻り、通知待ちの状態に移行する。そうでない場合には、判定は NO となってステップ 45 に移行する。

【0068】

ステップ 45 では、データ 211 の終端を検出したか否か判定する。照合によってその終端を検出した場合、判定は YES となり、そのことを通知するためにデータ位置情報をデータ入力構造検索部 120 にステップ 46 で通知した後、上記ステップ 41 に戻る。そうでない場合には、判定は NO となって上記ステップ 42 に戻り、照合を続行する。

【0069】

上述したようにして、データ入力構造検索部 120 と抽出条件判定部 130 の間では必要な情報のやりとりが随時、行われ、その情報によってそれぞれ処理を進行させる。それにより、1 データ 211 毎に、それが成立する抽出条件を確認し、その確認結果に応じた処理を行うようになっている。

【0070】

図 18 は、データ判定部 140 が実行する処理のフローチャートである。最後に図 18 を参照して、その処理について詳細に説明する。

まず、ステップ 51 では、データ入力構造検索部 120 からデータ 211 の終端が通知されるのを待つ。その通知を受け取ると、判定が NO となってステップ 52 に移行し、論理テーブル 190 を参照して、現在、対象としているデータ 211 が満たす抽出条件を判定する。その後はステップ 53 に移行する。

【0071】

ステップ 53 では、データ 211 が満たす抽出条件が有るか否か判定する。そのような抽出条件が存在した場合、判定は YES となってステップ 54 に移行し、検索結果判定情報 195 (図 13)、出力バッファ情報 151、及びバッファ情報 152 (図 14) を参照してデータ 211 を出力すべき出力バッファ 150 に出力し、対応する個別バッファ情報 153 を更新した後、上記ステップ 51 に戻る。それにより、通知待ちの状態に移行する。一方、そうでない場合には、判定は NO となってそのステップ 51 に戻る。

【0072】

図 19 ~ 図 24 は、上記データ抽出装置の適用例を説明する図である。以降は、図 19

10

20

30

40

50

～図24を参照して、その適用可能な利用法について具体的に説明する。図19～図24において、データ抽出装置は「抽出器」と表記している。

【0073】

図19は、複数のデータ抽出装置100を多段階で使用する場合の例を示している。データ1903を入力するデータ抽出装置100は、そのデータ1903を2つの連結器1910に振り分けている。その二つの連結器1910の一方は、マスタファイル1901のデータをデータ1903と連結させて別のデータ抽出装置100に出力し、そのデータ抽出装置100は連結結果を2つの集計器1920に振り分けている。その2つの集計器1920はそれぞれ異なるデータ抽出装置100に集計結果を出力し、その集計結果を入力するデータ抽出装置100はそのデータをそれぞれ3つのファイルに振り分けて出力している。これらは、二つの連結器1910の他方側でも同様である。

10

【0074】

図20は、入力データの振り分けにデータ抽出装置100を使用する場合の例を示している。その入力データは、ジャーナルファイル2000に格納された各レコードのデータである。データ抽出装置100は、抽出条件を満たすデータをジャーナルファイル2001～3のうちの何れかに振り分けて出力するために用いられている。そのように振り分けるのは、例えばマスタX～Zとの連結条件がそれぞれ異なることに対応するためである。そのように振り分けると、データを3系統で並行して処理することが可能となることから、処理の高速化を実現できる。

【0075】

20

図21は、連結結果のデータの振り分けにデータ抽出装置100を使用する場合の例を示している。その連結結果は、マスタとジャーナルのデータを連結させたものである。データ抽出装置100は、抽出条件1～3の何れかを満たすデータを、その抽出条件に応じてファイル2101～3のうちの何れかに出力するために用いられている。

【0076】

図22は、集計結果のデータの振り分けにデータ抽出装置100を使用する場合の例を示している。その集計結果は、マスタとジャーナルのデータの連結結果に対して集計操作を行ったものである。データ抽出装置100は、抽出条件1～3の何れかを満たす集計結果のデータを、その抽出条件に応じてファイル2201～3のうちの何れかに出力するために用いられている。

30

【0077】

図23は、新聞社等で実施されるクリッピングサービスの提供用にデータ抽出装置100を使用する場合の例を示している。その場合、データ抽出装置100にはサービス登録者毎に、その登録者に送るべき記事データが満たす抽出条件を定義する。その抽出装置100には随時、記事データが入力され、その記事データが満たす抽出条件に応じて対応するファイルに出力される。そのファイルに出力された記事データは、定期的にサービス登録者に配信される。サービス登録者の追加、削除、或いは要求の変更などは、抽出条件の追加、削除、或いは内容の変更によって対応することができる。

【0078】

図24は、ハイウェイ利用調査システムにデータ抽出装置100を使用する場合の例を示している。その場合、ハイウェイのモニタシステムから随時、データがデータ抽出装置100に入力される。その抽出装置100には、必要なデータのみを抽出するための抽出条件を定義する。それにより、抽出装置100は、抽出条件に従ってデータを選別する(フィルタリングする)。選別されたデータは、連結器によりマスタデータと照合され、より詳細なデータに展開される。例では、自動車の番号が「k 2104」のデータに対して会社名「通運」が付加されている。マスタデータと照合されたデータは集計器により、例えば会社毎に集計されて出力される。

40

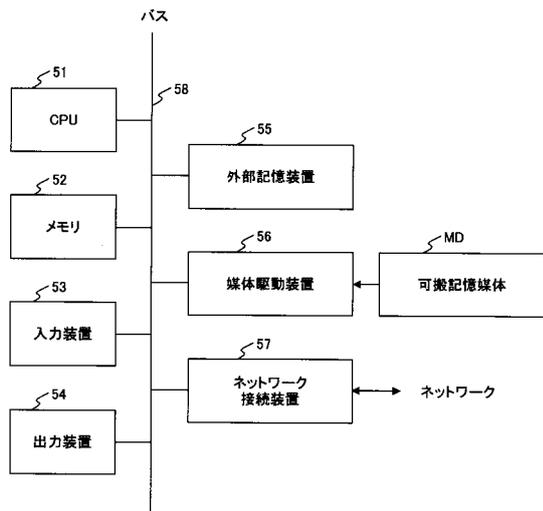
【0079】

なお、本実施の形態では、抽出条件によって出力先を振り分けるデータそのものを外部から入力しているが、そのデータは実際に振り分けるデータの生成用、或いは特定用のも

50

のであっても良い。つまり符号化された圧縮データのようなものであっても良い。そのようなデータの入力、記録媒体MDに記録して行うようにしても良い。

【図5】



【図6】

```

<root>
  <Company>
    <Code>000001</Code>
    <Name>ABCカンパニー</Name>
  </Company>
  <Office>
    <ZipCode>211-0040</ZipCode>
    <Address>神奈川県川崎市上小田中...</Address>
    <Tel>444-999-0001</Tel>
    <Tel>444-999-0002</Tel>
    <Tel>444-999-0003</Tel>
    <Fax>444-999-9999</Fax>
  </Office>
</root>

```

【図7】

```

CompanyCode,CompanyName,ZipCode,Address,Tel,Fax
000001,"AAA株式会社","211-0000","神奈川県川崎市...", "444-999-0001", "444-999-9999"
000002,"BBB商會","164-0000","東京都中野区...", "111-999-0001", "111-999-9999"
000003,"青函会社CCC","230-0000","神奈川県横浜市...", "222-999-0001", "222-999-9999"
000004,"劇団DDD","811-0000","福岡県福岡市...", "333-999-0001", "333-999-9999"
000005,"EEE社","410-0000","静岡県沼津市...", "555-999-0001", "555-999-9999"
000006,"FFFカンパニー","211-0000","神奈川県川崎市...", "444-999-0001", "444-999-9999"
000007,"GGG放送","105-0000","東京都港区...", "666-999-0001", "666-999-9999"
.
.
.
.
.
.
.
.
.
.

```

【図8】

```

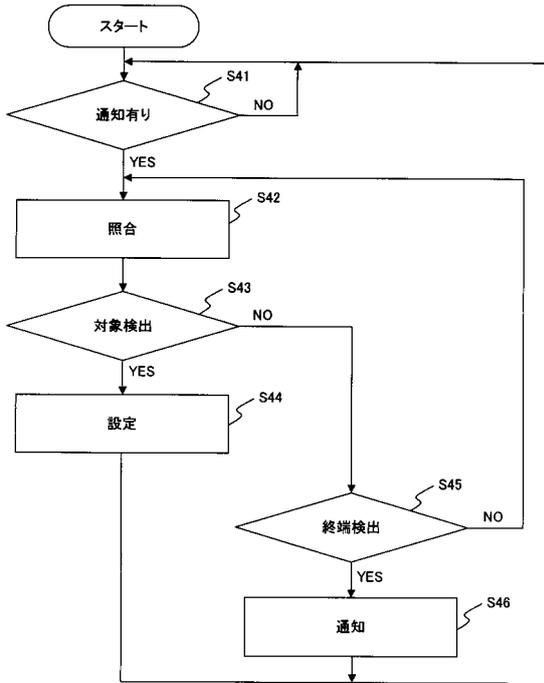
(1)抽出条件: /root/origin = "atcg"
出力条件: "D:\Selection\Output\Result01.dat"

(2)抽出条件: /root/Company/code < 99 AND /root/origin = "gtac"
出力条件: "D:\Selection\Output\Result02.dat"

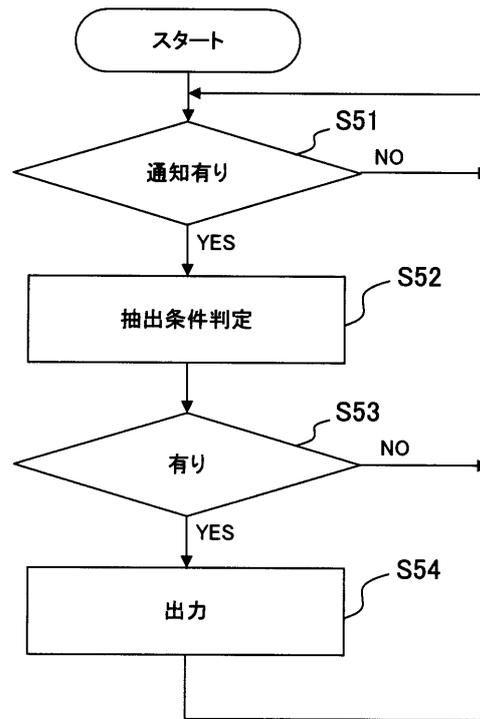
(3)抽出条件: /root/origin = "aacg"
出力条件: "D:\Selection\Output\Result03.dat"

```

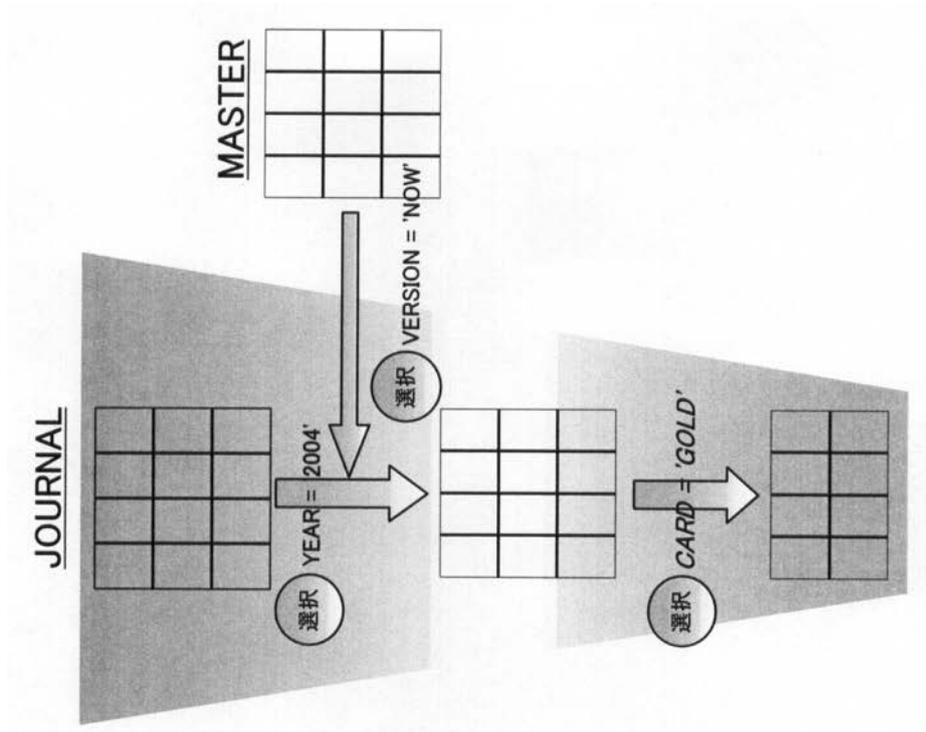

【図17】



【図18】



【 図 1 】



```

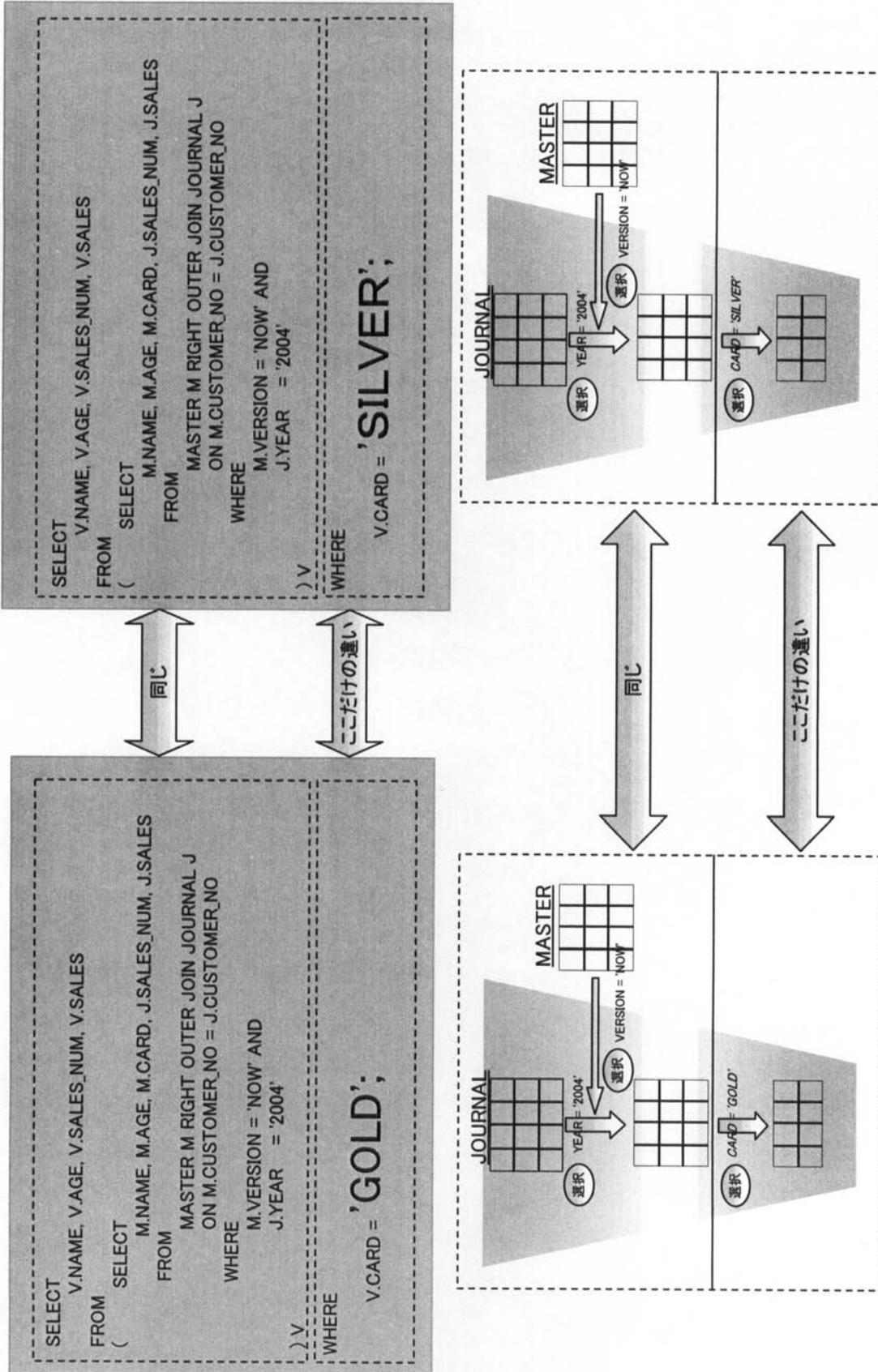
SELECT
V.NAME, V.AGE, V.SALES_NUM, V.SALES
FROM
(
  SELECT
  M.NAME, M.AGE, M.CARD, J.SALES_NUM, J.SALES
  FROM
  MASTER M RIGHT OUTER JOIN JOURNAL J
  ON M.CUSTOMER_NO = J.CUSTOMER_NO
  WHERE
  M.VERSION = 'NOW' AND
  J.YEAR = '2004'
) V
WHERE
V.CARD = 'GOLD';

```

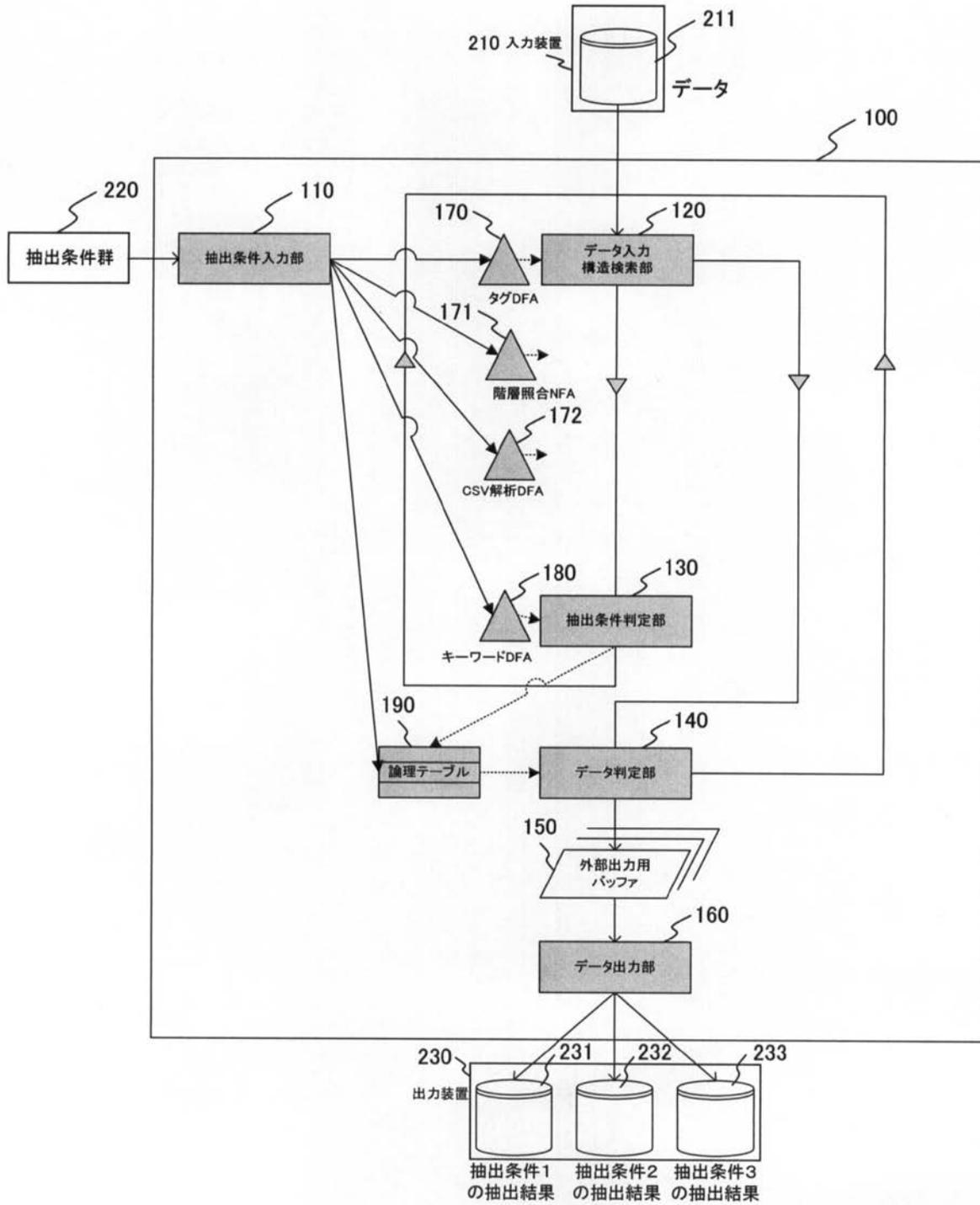
連結時の条件
マスタ/ジャーナルの制約

連結結果に対する条件

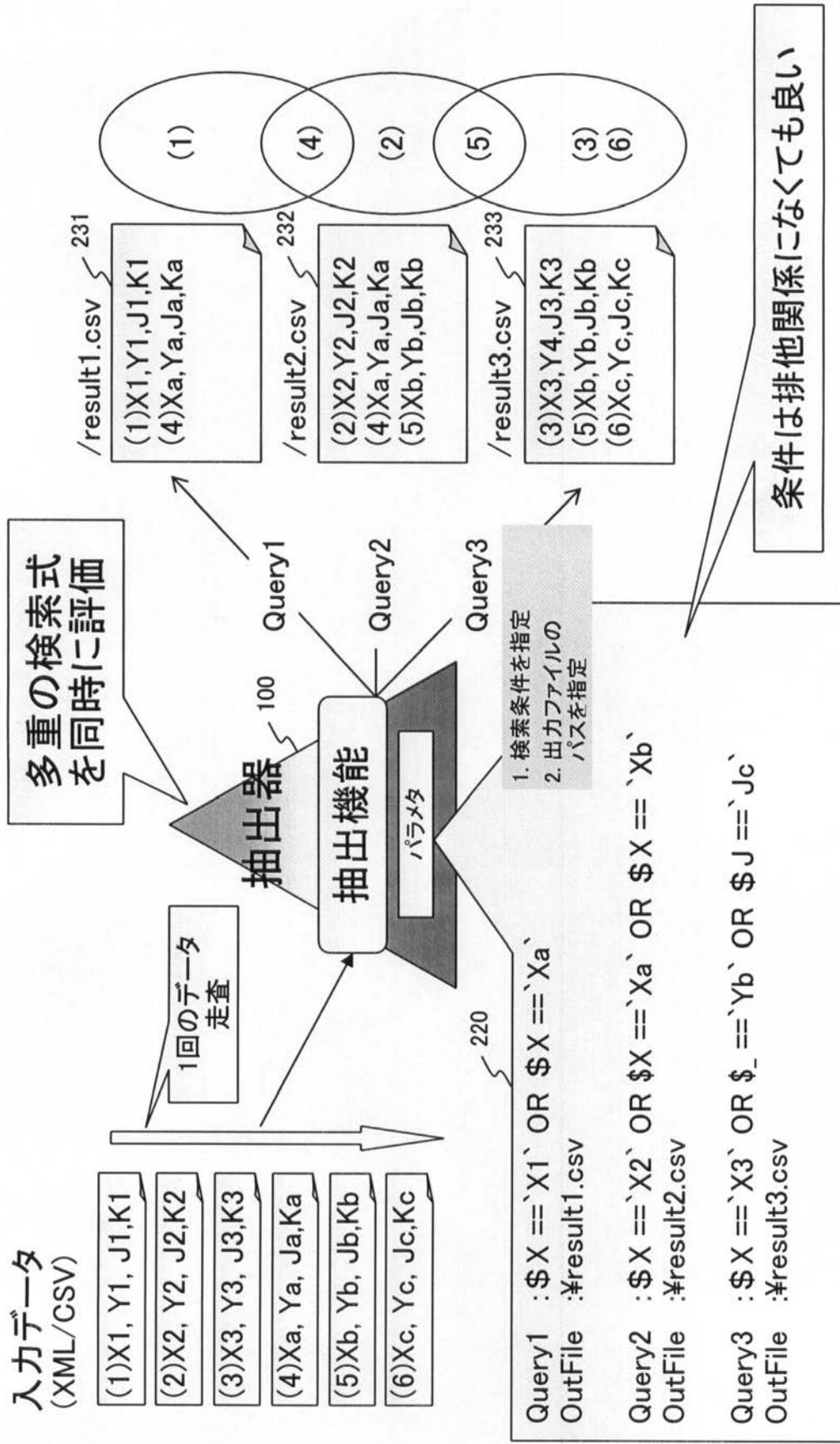
【 図 2 】



【図3】

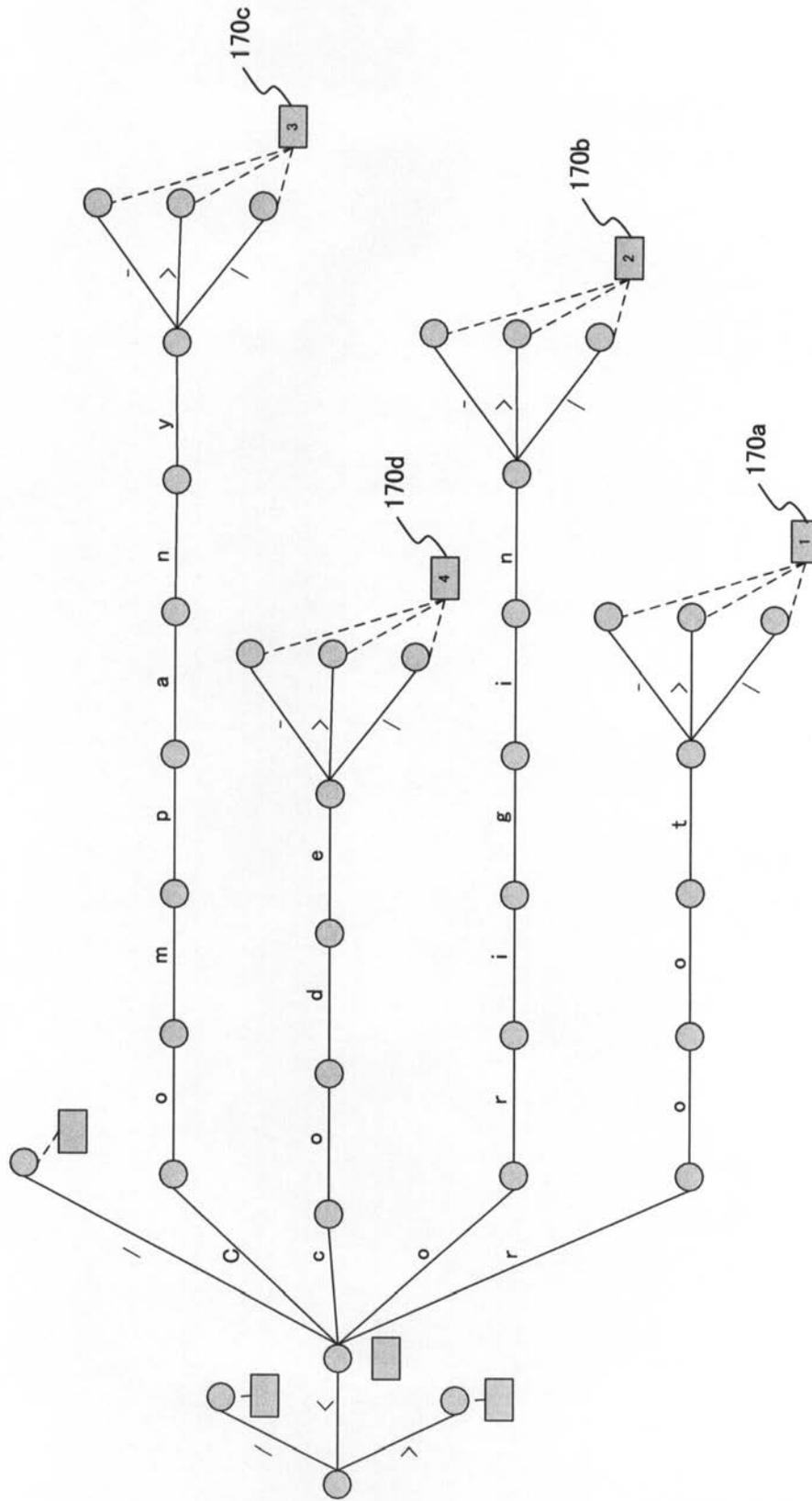


【 図 4 】



【図9】

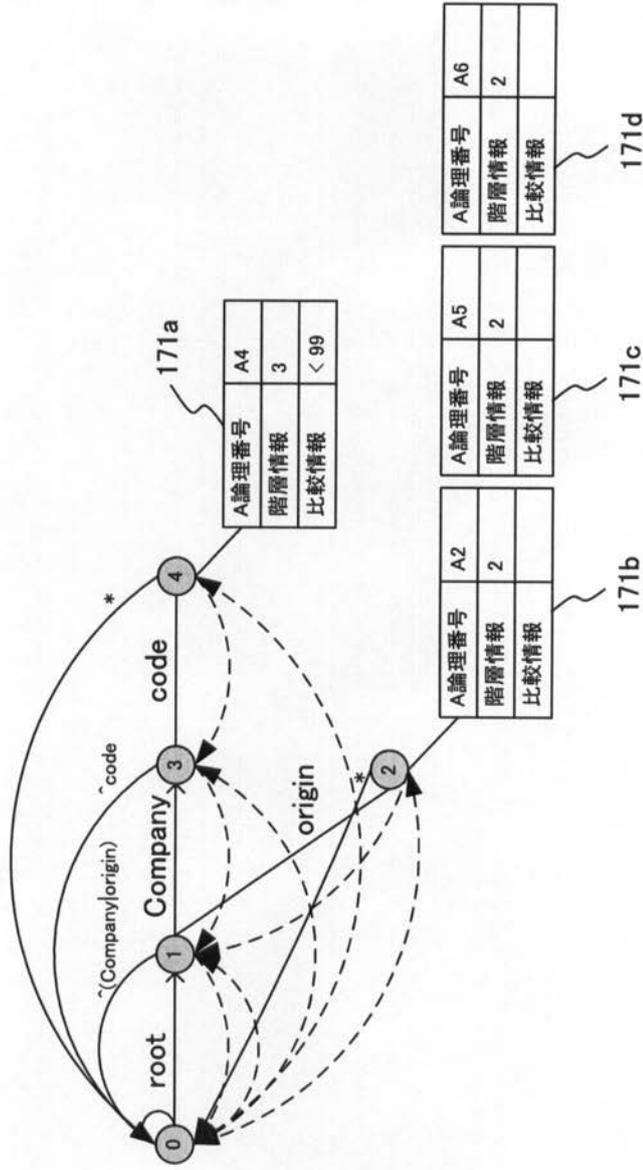
•/root/originノードと/root/Company/code を検出するためのDFA



【 図 1 0 】

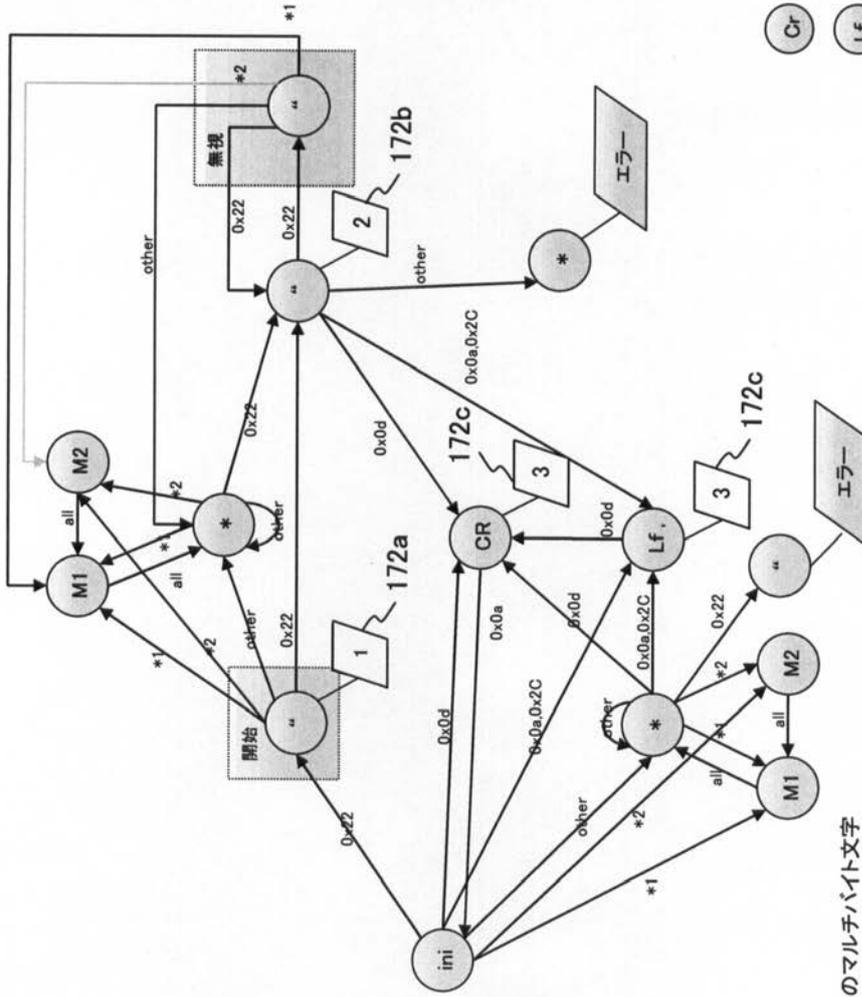
・ /root/origin と /root/Company/code の階層照合のためのNFA

| | |
|---|---------|
| 0 | * |
| 1 | root |
| 2 | origin |
| 3 | Company |
| 4 | code |



—→ 開始タグによる状態遷移を示す
 - -→ 終了タグによる状態遷移を示す

【 図 1 1 】



ヒット情報

- 172a: データ開始位置取得
(ダブルコーテーション時)
- 172b: データ終了位置取得
(ダブルコーテーション時)
- 172c: 項目終了
エラー: CSVエラー

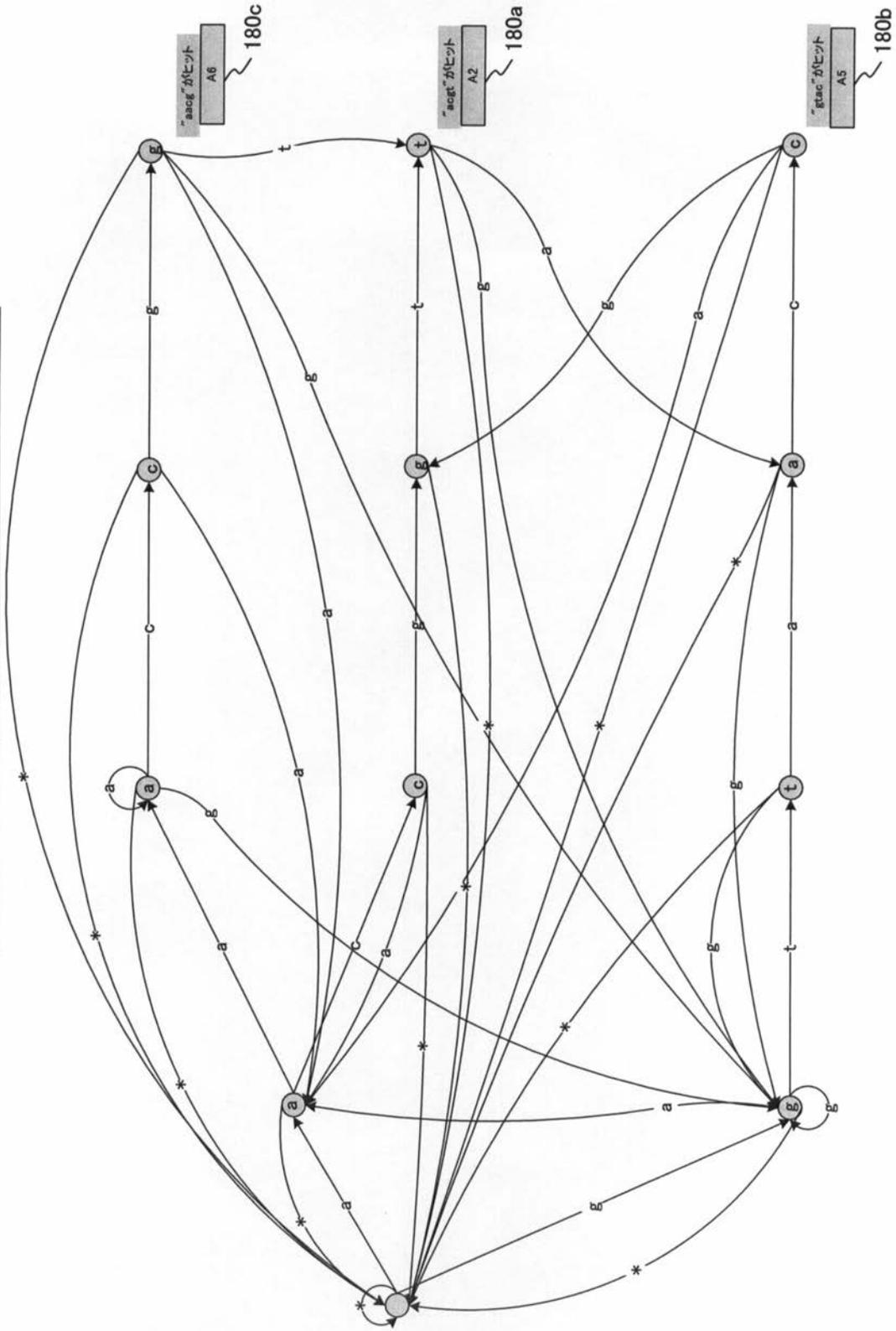
- Cr 0x0aはiniへ、0x0a以外はiniと同じリンク設定を行う。
- Lf iniと同じリンク設定を行う。
- M1 マルチ文字1バイトスキップさせるための節
- M2 マルチ文字2バイトスキップさせるための節

※1 2byteのマルチバイト文字
EUCの場合: 0xA1~0xFE, 0x8E

※2 3byteのマルチバイト文字
EUCの場合: 0x8F

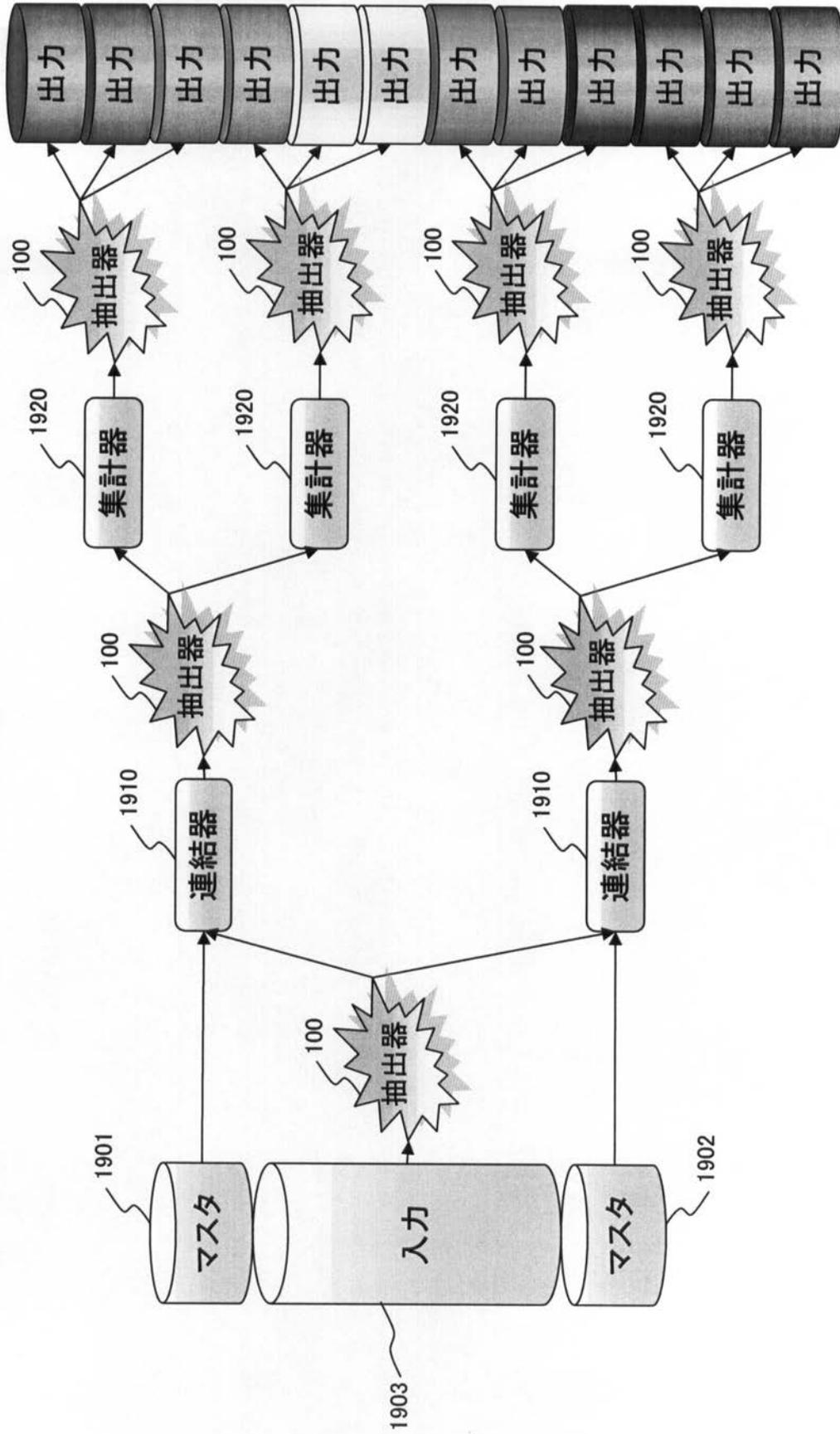
【 図 1 2 】

キーワード1 A2 : "acgt"
キーワード2 A5 : "gtac"
キーワード3 A6 : "aacg" を照合する場合

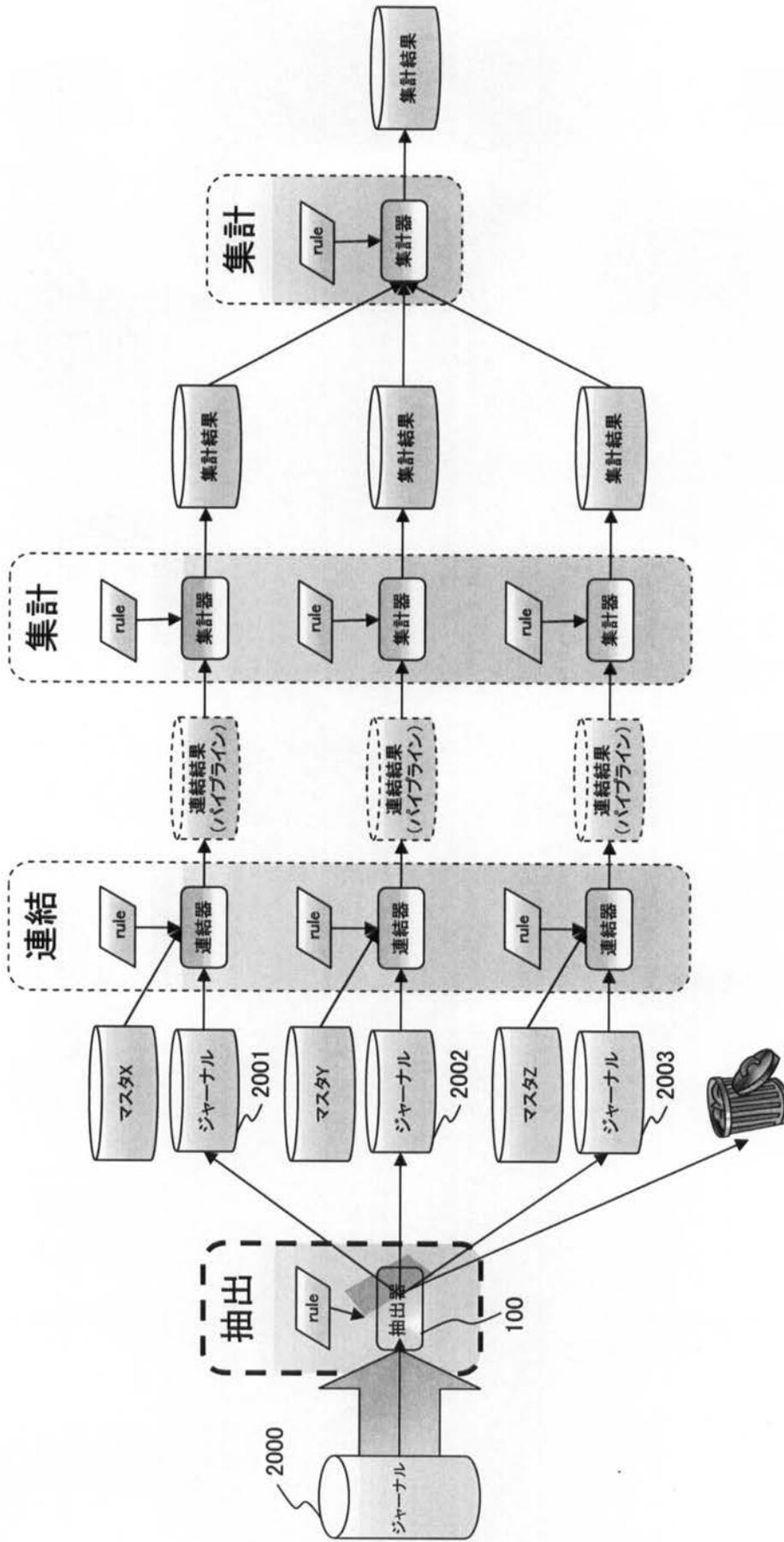


【図19】

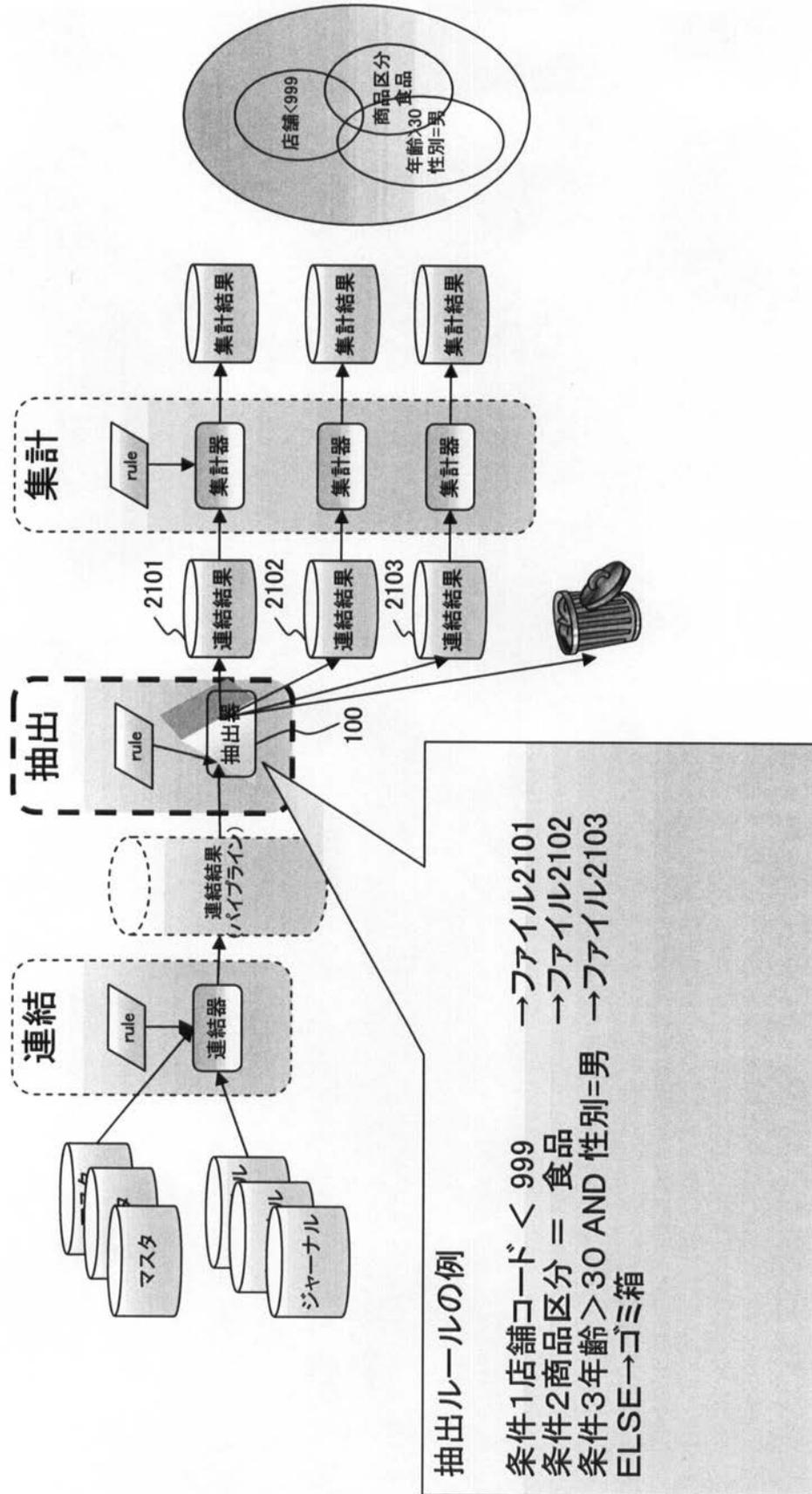
連結・集計と組み合わせてあらゆるところに配置可能



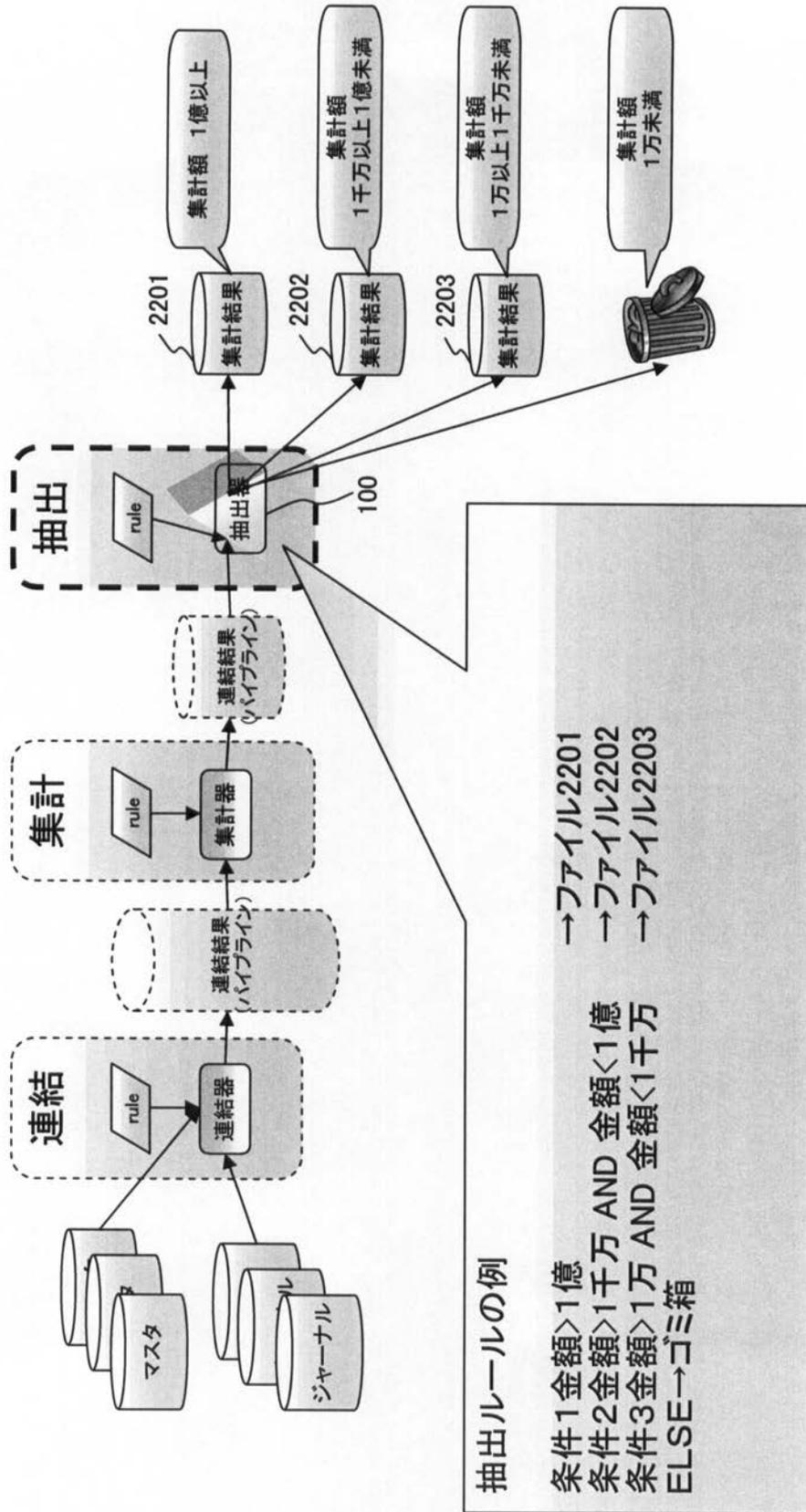
【図20】



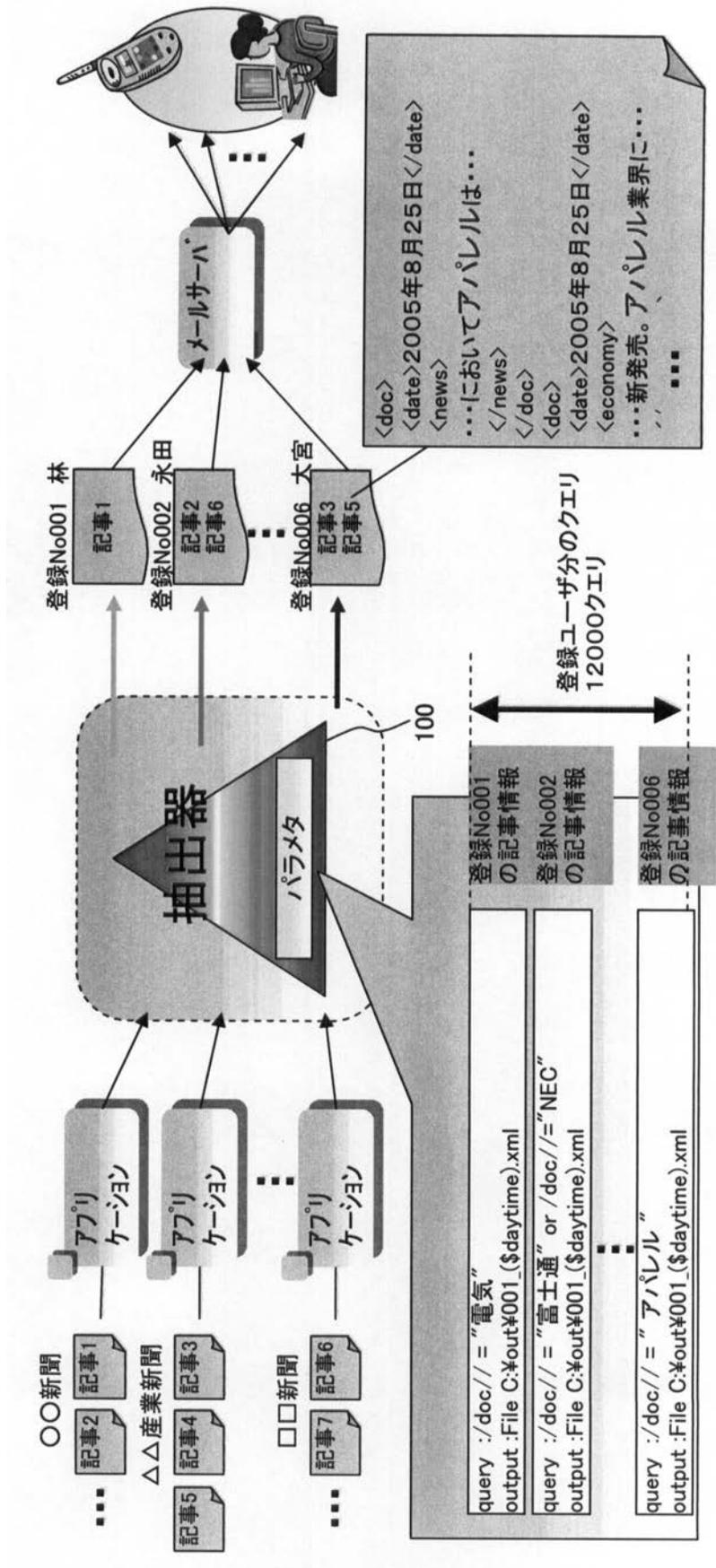
【 図 2 1 】



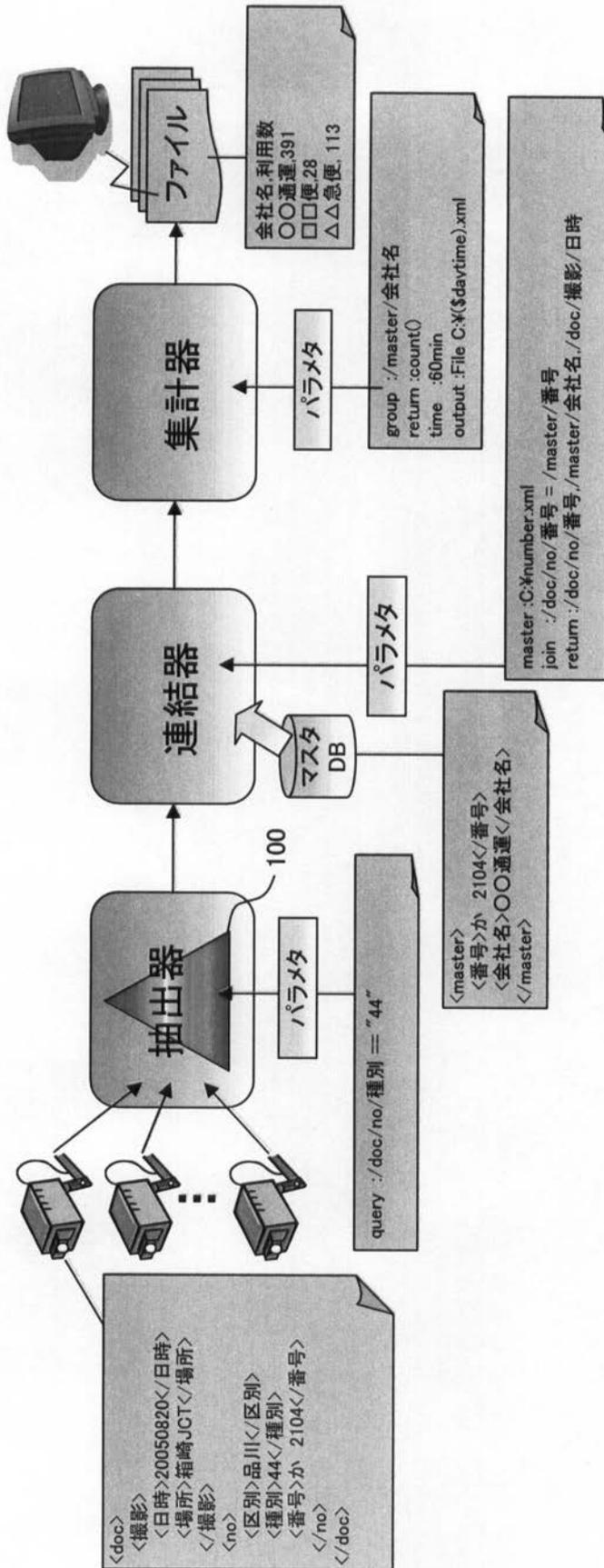
【 図 2 2 】



【 図 2 3 】



【 図 2 4 】



フロントページの続き

- (72)発明者 永田 真彦
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
- (72)発明者 大宮 清英
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

合議体

- 審判長 長島 孝志
審判官 田中 秀人
審判官 仲間 晃

- (56)参考文献 特開2000-339346(JP,A)
特開2005-101928(JP,A)
特開平6-139291(JP,A)
特開2001-344282(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G06F17/30,12/00