

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2006-53724

(P2006-53724A)

(43) 公開日 平成18年2月23日(2006.2.23)

(51) Int. Cl.	F I	テーマコード (参考)
<b>G06F 17/30 (2006.01)</b>	G06F 17/30 330B	5B075
<b>G06F 12/00 (2006.01)</b>	G06F 17/30 140	5B082
	G06F 12/00 513D	
	G06F 12/00 547Z	

審査請求 未請求 請求項の数 8 O L (全 18 頁)

(21) 出願番号	特願2004-234344 (P2004-234344)	(71) 出願人	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成16年8月11日 (2004.8.11)	(74) 代理人	100068504 弁理士 小川 勝男
		(74) 代理人	100086656 弁理士 田中 恭助
		(74) 代理人	100094352 弁理士 佐々木 孝
		(72) 発明者	今木 常之 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内
		(72) 発明者	西澤 格 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内

最終頁に続く

(54) 【発明の名称】 XMLデータ管理方法

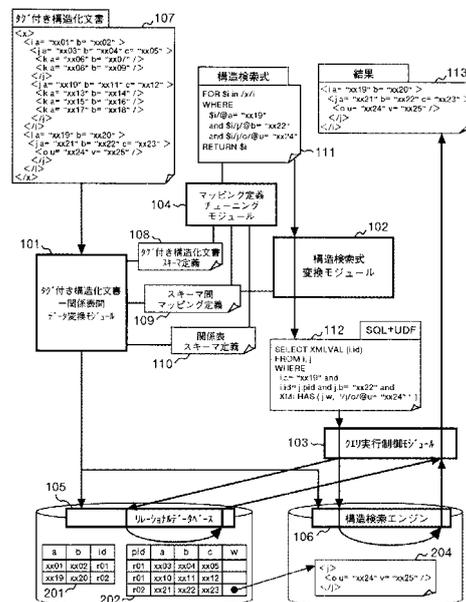
(57) 【要約】

【課題】 XML文書 - 関係表間スキーママッピング定義の最適化および透過的な構造検索

【解決手段】 マッピング定義チューニングモジュール104は、構造検索式の発行履歴を参照して、発行頻度の高い検索処理の効率化を目的に、XML文書がリレーショナルデータベース105と構造検索エンジン106に適切に分解されて格納されるよう、スキーマ間マッピング定義109を変更する。構造検索式変換モジュール102は、スキーマ間マッピング定義109に基づいて構造検索式を変換する。クエリ実行制御モジュール103は、リレーショナルデータベース105と構造検索エンジン106それぞれにクエリを発行して、それぞれの結果から元の構造検索式に対する結果を再構成する。

【選択図】 図1

図 1



## 【特許請求の範囲】

## 【請求項 1】

木構造型のタグ付き構造化文書を、関係データベース、および構造検索専用データベースを用いて管理する方法において、

構造化文書格納定義に従って、単一の構造化文書を、関係データベースに格納する第 1 の構造部分と、構造検索専用データベースに格納する第 2 の構造部分に分解し、

該第 1 の構造部分について、タグを取り除いたデータ自体を抽出して、該格納定義で対応付けられた関係表のカラムに格納し、

該第 2 の構造部分について、タグを含んだまま構造検索専用データベースに格納し、

元の構造化文書に対する構造検索式を、該格納定義に従って、該関係データベース用の第 1 の検索式と、該構造検索専用データベース用の第 2 の検索式に変換し、

該関係データベースに対して該第 1 の検索式を発行して、その結果を受け取り、該構造検索専用データベースに対して該第 2 の検索式を発行して、その結果を受け取り、

両結果から、該構造検索式に対する結果と等価な構造化文書を構築する手順を有することを特徴とする XML データ管理方法。

10

## 【請求項 2】

木構造型のタグ付き構造化文書を、関係データベース、および構造検索専用データベースを用いて管理する方法において、

構造化文書格納定義に従って、単一の構造化文書を、関係データベースに格納する第 1 の構造部分と、構造検索専用データベースに格納する第 2 の構造部分に分解し、

該第 1 の構造部分について、タグを取り除いたデータ自体を抽出して、該格納定義で対応付けられた関係表のカラムに格納し、

該第 2 の構造部分について、タグを含んだまま構造検索専用データベースに格納し、

元の構造化文書に対する構造検索式を、該格納定義に従って、該関係データベース用の第 1 の検索式と、該構造検索専用データベース用の第 2 の検索式に変換し、

該第 2 の検索式をそれと同等の構造検索処理を実行する該第 1 の検索式に埋め込み可能な、関係データベース検索式拡張関数で表現し、

該拡張関数を該第 1 の検索式に埋め込んだ拡張関数付き関係データベース用検索式を生成し、

該拡張関数付き関係データベース用検索式を、該関係データベースに対して発行し、該構造検索

20

30

式に対する結果と等価な構造化文書を取得することを特徴とする XML データ管理方法。

## 【請求項 3】

タグ付き構造化文書に対する構造検索式の発行履歴を記録し、

発行頻度の高い検索式の検索処理効率を指標として、

構造化文書格納定義、および、関係表のスキーマ定義を更新することを特徴とする請求項 1 に記載の XML データ管理方法。

## 【請求項 4】

タグ付き構造化文書の中で、該構造検索専用データベースに格納された部分構造化文書に対する同一の構造検索式の発行頻度が所定数を越える場合に、

該構造検索専用データベースに格納された該部分構造化文書を格納する関係表を新規に作成し、

該構造化文書格納定義に、該部分構造化文書と、該新規に作成した関係表との対応付けを追記し、

該構造検索専用データベースに格納された該部分構造化文書を、該新規に作成した関係表に格納し直すことを特徴とする請求項 3 に記載の XML データ管理方法。

40

## 【請求項 5】

タグ付き構造化文書の中で、あるタグで示される要素が該タグと同型の要素を子要素として持つ自己再帰的な構造を持ち、

かつ該タグのデータが関係データベースに格納されており、

50

かつ該タグに対する階層の深さを問わない同一の構造検索式が出現する場合に、  
該自己再帰的なデータを、構造検索専用データベースに格納し直すことを特徴とする請求項 3 に記載の XML データ管理方法。

【請求項 6】

タグ付き構造化文書の中で、関係表データベースに格納されたデータに対する多段の階層を指定した同一の構造検索式が出現する場合に、

該データを構造検索専用データベースに格納し直すことを特徴とする請求項 3 に記載の XML データ管理方法。

【請求項 7】

タグ付き構造化文書の中で、関係表データベースに格納されたデータに対する多段の階層を指定した同一の構造検索式が出現する場合に、

該構造検索式に登場する全階層のデータを並べたカラムを持つ単一の関係表を新規に作成し、

該データを該新規に作成した関係表に格納し直すことを特徴とする請求項 3 に記載の XML データ管理方法。

【請求項 8】

タグ付き構造化文書の中で、関係表データベースに格納されたデータに対する多段の階層を指定した同一の構造検索式が複数出現する場合に、

該複数の構造検索式に登場する全階層のデータの和集合を並べたカラムを持つ単一の関係表を新規に作成し、

該データを、該新規に作成した関係表に格納し直すことを特徴とする請求項 3 に記載の XML データ管理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データベース管理システムに関する。特にリレーショナルデータベース（あるいは関係データベース、以下、RDB という）を用いる XML（*eXtensible Markup Language*）文書の管理方法に係わり、特に一つの XML 文書を XML 構造検索エンジンと RDB に分解して管理する方法に係わり、特に該文書に対する検索履歴に基づいて分解方法を適宜改善しつつ、ユーザに対してはこの分解方法について透過的な構造検索インタフェースを提供する方法に関する。

【背景技術】

【0002】

現在、XML 文書の管理に特化したネイティブ XML データベース（NXDB）と呼ばれる製品がいくつか存在する。しかし、NXDB は何れも発展途上であり、一般に大量データの管理や集計処理の目的には性能的に不十分であるため、基幹系業務などには適さないとされている。XML の基幹業務応用は、XBRL（*eXtensible Business Reporting Language*）などのビジネス関連の XML 仕様の登場により今後の発展が期待されるため、大量の XML 文書を十分な性能で処理可能な技術が必要とされている。一方、主要 RDB 製品においても XML 文書管理機能が提供されている。RDB は長年にわたる改良により大量データの処理にも十分耐えうる性能を提供するため、XML の基幹系業務応用にも適していると言える。

【0003】

主要 RDB 製品に関する代表的な XML 文書管理方法については、非特許文献 1 および非特許文献 2 に記載されている。

【0004】

非特許文献 1 の方法は、管理対象である XML 文書の文書スキーマと、RDB の関係表スキーマとの間の対応関係に従って、XML 文書に含まれるタグ付けされたデータを構造分解して、複数の関係表に分けて値単位で格納する。このような格納方式を、以下、XML 文書スキーマ - 関係表スキーマ間のマッピング方式と呼ぶ。非特許文献 1 の方法は、X

10

20

30

40

50

M L 文書スキーマの定義を元に、その定義に妥当である X M L 文書を格納するための関係表スキーマの定義、および、該 X M L 文書スキーマと該関係表スキーマとの間の対応関係の定義（以下、スキーママッピング定義という）を、自動的に作成する。また X M L の標準検索仕様 X P a t h 形式の構造検索式を、スキーママッピング定義に従って、関係表検索式（以下、S Q L 式という）に自動変換する。

【 0 0 0 5 】

非特許文献 2 の方法も、基本的にスキーママッピング方式である。ただしスキーママッピング定義は、R D B に格納されたデータから X M L 文書を構築する方向で、ユーザがマニュアルで定義する。また X M L の標準検索仕様 X Q u e r y 形式の構造検索式を、スキーママッピング定義に従って、S Q L 式に自動変換する。

10

【 0 0 0 6 】

【非特許文献 1】店 ML Schemas in Oracle XML DB R. Murthy, S, Banerjee; VLDB2003

【非特許文献 2】徹 uerying XML Views of Relational Data J. Shanmugasundaram, et al., VLDB2001

【発明の開示】

【発明が解決しようとする課題】

【 0 0 0 7 】

一般に、R D B での X M L 文書管理は、X M L 文書に含まれるタグ付けされたデータを構造分解して、複数の関係表に分けて値単位で格納する方式に則っている。このような X M L 文書スキーマ - 関係表スキーマ間のスキーママッピング方式には、X M L 文書を管理するうえで以下のような欠点が存在する：

20

( a ) 検索効率を考慮したマッピング定義

一般に、マッピング方法の違いによって検索性能は異なってくる。最適な検索性能を得るためには、マッピング定義のチューニングが必要であるが、ユーザにとってこの作業は大変な負担となる。

( b ) 非定型データの管理

X M L では、厳密なスキーマ定義に従わない非定型部分データを文書中に含むことが可能であり、これによるデータ表現の柔軟性が X M L 利用拡大の大きな要因となっているが、R D B ではこのようなデータを L O B とよばれる一次元の文字列データとして管理することになるため、その部分に対して高度な検索をかけることができない。

30

( c ) 複雑な構造を持つ文書の管理

X M L では、木構造に基づいたデータモデルにより、複雑なデータ構造を表現することが可能である。一方、関係表は一次元の値の集まりを単位としてデータを管理するため、木構造のような複雑なデータは、複数の関係表間における外部参照関係によって表現しなくてはならない。しかし、X M L 文書スキーマの階層が深い場合は多数の関係表に分けて管理することになるため、検索効率および格納効率の点で望ましくない。このように、R D B と X M L とのデータモデルの違いに基づく関係表での管理が非効率的な X M L 文書が存在する。

( d ) 検索指定方法

関係表にスキーママッピングした X M L 文書に対する検索は、そのマッピング定義に沿って定義される必要があるため、ユーザがマッピング定義を意識して関係表検索式（以下、S Q L 式）を記述する必要がある。また、( a ) の課題にあげたように検索効率性を考慮してマッピング定義を変更した場合は、S Q L 式も記述し直す必要がある。一般に、ユーザにとっては、X M L 文書スキーマのみを意識して構造検索を指定できることが理想であり、X M L 文書の管理においては本来存在しないこれらの必要性は、ユーザにとって大変な負担となる。

40

【 0 0 0 8 】

上記の X M L 文書スキーマ - 関係表スキーマ間スキーママッピング方式における ( a ) ~ ( d ) の欠点を克服するために、本発明ではそれぞれ以下の課題を解決することを目的とする。

50

## 【0009】

第一に、スキーママッピング定義の自動チューニング機能を提供すること。

## 【0010】

第二に、従来LOBで管理していたような非定型部分データに対しても構造検索機能を提供すること。

## 【0011】

第三に、関係表での管理が非効率的なデータを切り分けて、効率的な手段で管理すること。

## 【0012】

第四に、XML文書の関係表への格納方法に関して、透過なXML文書の構造検索機能を提供すること。 10

## 【課題を解決するための手段】

## 【0013】

まず、第二、第三の課題を解決するために、RDBの外部データベース、あるいはRDBのプラグインとして存在するXML構造検索エンジンと連携する。

## 【0014】

従来、関係表のLOBカラムに格納していた非定型部分データを構造検索エンジンに格納することによって、第二の課題を解決する。関係表での管理が非効率的なデータも、代わりに構造検索エンジンで管理することによって、第三の課題を解決する。

## 【0015】

また、第一の課題を解決するために、クエリ発行履歴を参照し、頻出クエリの検索性能効率化を指標として適切なスキーママッピング定義を導出するマッピング定義チューニングモジュールを導入する。第三の課題解決における関係表での管理が不適切なXML部分データの切り分けもこのモジュールで行う。 20

## 【0016】

さらに、第四の課題を解決するために、XML文書に対する構造検索式を、スキーママッピング定義に基づいてSQL式に自動変換するクエリリライト機能を提供する。検索対象が構造検索エンジンで管理されている部分データにも及ぶ場合は、この検索エンジンへの検索式をUDF (User Define Function) として含むSQL式に変換する。このクエリリライトにより、ユーザは、第一～第三の課題解決における、XML文書の関係表および構造検索エンジンへの格納方法の違いに対して、透過的に構造検索指定が可能となる。 30

## 【発明の効果】

## 【0017】

XML文書スキーマ - 関係表スキーマ間のスキーママッピング方式において、

(1) クエリ発行履歴に基づいて、検索処理コストを削減するようにスキーママッピング定義を自動的に改善することが可能である。

(2) 非定型の部分データを構造検索エンジンで管理することによって、該部分データに対する構造検索が可能である。

(3) 関係表での管理が非効率的なデータを切り分けて構造検索エンジンで管理することによって、非効率的な検索処理を回避することが可能である。 40

(4) クエリリライト機能により、XML文書の関係表および構造検索エンジンへの格納方法に関し、透過的にXML文書に対する構造検索を指定することが可能である。

## 【発明を実施するための最良の形態】

## 【0018】

以下、本発明の実施の一形態を、図面を参照しながら説明する。なお簡単のために、本明細書中では以下に述べる発明の実施の形態を単に「本実施例」と呼ぶことにする。

## 【0019】

図1を用いて、本実施例の概略構成について説明する。

## 【0020】

本実施例のシステムは、以下に挙げる4つのモジュールを基本構成要素として成立している：

- ・ タグ付き構造化文書 - 関係表間データ変換モジュール 101
- ・ 構造検索式変換モジュール 102
- ・ クエリ(問合せ)実行制御モジュール 103
- ・ マッピング定義チューニングモジュール 104

以下、それぞれのモジュールについて概説する。

#### 【0021】

タグ付き構造化文書 - 関係表間データ変換モジュール 101 は、タグ付き構造化文書(以下、XML文書という) 107 を構造分解してタグを取り除いたデータ本体を、リレーショナルデータベース(以下、RDBという) 105 の関係表のカラムに対応して属性値を格納する。ただし一部のXML文書については、タグが付いた部分木単位でXML文書専用の構造検索エンジン 106 に格納する。XML文書のうち、RDB 105 に格納する部分、格納先の関係表カラム、および構造検索エンジン 106 にタグごと格納する部分の区別は、タグ付き構造化文書スキーマ定義(以下、XML文書構造定義という) 108、スキーマ間マッピング定義 109、および関係表スキーマ定義 110 に従って決定される。

10

#### 【0022】

XML文書構造定義 108 はXML文書の構造定義を、関係表スキーマ定義 110 は関係表のスキーマ定義をそれぞれ表す。XML文書 107 は、XML文書構造定義 108 に対して妥当である必要があるし、RDB 105 に格納されている関係表 201、202 は、関係表スキーマ定義 110 に従って構成されている。スキーマ間マッピング定義 109 は、XML文書のノード値(タグで修飾された要素値、あるいは属性値)とそれを格納する関係表のカラムの対応付けを定義する。

20

#### 【0023】

構造検索式変換モジュール 102 は、XQuery、XPathなどのXMLの標準検索仕様に従ってユーザが定義した構造検索式 111 を、RDB 105 用の検索仕様であるSQL言語の検索式(以下、SQL式という) 112 に変換するモジュールである。この変換は、スキーマ間マッピング定義 109 に従って行われる。検索範囲が構造検索エンジンに格納した部分XML文書にも及ぶ場合には、SQL式 112 中に構造検索エンジン用の拡張関数(UDF)を埋め込んだ式に変換する。

30

#### 【0024】

クエリ実行制御モジュール 103 は、UDFを含んだSQL式 112 を、SQL部分とUDF部分に分離し、前者をRDB 105 に、後者を構造検索エンジン 106 に対して発行し、その結果を統合して、元の構造検索式 111 に対する結果 113 を構築するモジュールである。このモジュールは、RDB 105 にプラグイン処理機構がある場合は、その機能上で自然に実現される(この場合については、図3を用いて後述する)。

#### 【0025】

マッピング定義チューニングモジュール 104 は、ユーザの構造検索式 111 の発行履歴を参照して、頻出する検索式の処理の効率化を指標として、XML文書構造定義 108 を参照しつつ、スキーマ間マッピング定義 109、および関係表スキーマ定義 110 を適宜更新する。関係表スキーマ定義 110 の更新に伴う関係表の変更は、RDB 105 の機能に任せる。

40

#### 【0026】

以下、図1に示すシステムを実現するためのハードウェア構成について説明する。本システムは、ハードウェア的にはCPU、メモリ、外部記憶装置、入力装置、表示装置などを備える1台又は複数台の計算機によって構成される。XML文書 107、XML文書構造定義 108、スキーマ間マッピング定義 109 および関係表スキーマ定義 110 は、ファイルとして記憶装置上に格納される。構造検索式 111 は、テキストエディタを介して入力装置から入力されるか、図示しないアプリケーションプログラムを介して生成され、

50

メモリに格納される。結果 1 1 3 は、メモリに格納され、表示装置やプリンタに出力されるか、さらに処理のためにアプリケーションに渡されるデータである。構造検索エンジン 1 0 6 は、記憶装置に格納される XML 文書の木構造ファイルを有し、これらファイルを管理するためのデータベース・マネージメント・システムである。タグ付き構造化文書 - 関係表間データ変換モジュール 1 0 1、構造検索式変換モジュール 1 0 2、クエリ実行制御モジュール 1 0 3 およびマッピング定義チューニングモジュール 1 0 4 は、計算機のメモリに格納され、その CPU によって実行されるプログラムである。R D B 1 0 5 は、記憶装置上に格納されるリレーショナルデータベースを有し、このデータベースを管理するためのデータベース・マネージメント・システムである。データ変換モジュール 1 0 1、構造検索式変換モジュール 1 0 2、クエリ実行制御モジュール 1 0 3 およびマッピング定義チューニングモジュール 1 0 4 の一部又は全部が R D B 1 0 5 に組み込まれて実装されてもよい。これらモジュール、R D B 1 0 5 および構造検索エンジン 1 0 6 は、同一の計算機上で実行されてもよいし、その一部又は全部がネットワークを介して異なる計算機上で実行されてもよい。また本システムは、クライアント - サーバ型のシステムで実現されてもよい。

10

#### 【 0 0 2 7 】

以上が、本実施例の概略である。以降、本実施例における、データ変換モジュール 1 0 1 の動作概要を図 2 で、構造検索式変換モジュール 1 0 2 の動作概要を図 3 で、クエリ実行制御モジュール 1 0 3 の動作概要を、実現方法のバリエーション別に図 3、図 4、図 5 を用いて説明する。

20

#### 【 0 0 2 8 】

図 2 を用いて、データ変換モジュール 1 0 1 の動作について説明する。本説明では、XML 文書構造定義 1 0 8 に対して妥当である XML 文書 1 0 7 を R D B 1 0 5 に格納する場合を例にとる。XML 文書構造定義 1 0 8 は、ルート要素  $x$  の下に複数の  $i$  要素が出現し、各  $i$  要素は属性  $a$  ,  $b$  を持ち、さらにその下には複数の  $j$  要素が出現し、各  $j$  要素は属性  $a$  ,  $b$  ,  $c$  を持ち、さらにその下には複数の  $k$  要素が出現し、各  $k$  要素は属性  $a$  ,  $b$  を持つことを表している。 $str$  は文字列を示す。 $x 0 \dots n$  は、対応する要素が 0 個から  $n$  個まで出現可能なことを示す。また、各  $j$  要素の下には、 $k$  要素以外にも任意の要素が登場し得ることを { ANY } で示す。なお、図 2 の XML 文書構造定義 1 0 8 の記法は実施例を限定するものではなく、同様の意味を表現し得る XML 文書構造の定義仕様であれば、どのような記法でも適用可能である。例えば、XML 文書の標準的な文書構造定義仕様である D T D ( Document Type Definition ) では、上記と同様の文書構造定義を以下のように表現する：

30

```
<!ELEMENT x i * >
<!ELEMENT i j * >
<!ATTLIST i
  a CDATA #REQUIRED
  b CDATA #REQUIRED >
<!ELEMENT j ANY >
<!ATTLIST j
  a CDATA #REQUIRED
  b CDATA #REQUIRED
  c CDATA #REQUIRED >
<!ELEMENT k EMPTY >
<!ATTLIST k
  a CDATA #REQUIRED
  b CDATA #REQUIRED >
```

40

一方、XML 文書 1 0 7 を格納する関係表 2 0 1、2 0 2、および 2 0 3 のスキーマは、関係表スキーマ定義 1 1 0 で与えられる。この例では、関係表  $i$  が  $a$  ,  $b$  ,  $id$  の 3 つのカラムを、関係表  $j$  が  $pid$  ,  $a$  ,  $b$  ,  $c$  ,  $w$  ,  $id$  の 6 つのカラムを、関係表  $k$  が  $p$

50

`id`、`a`、`b`、`pid`の4つのカラムを持つことをそれぞれ表現している。ここで`id`と`pid`は親と子のつながりを示す識別子である。`id`は自身を親とする識別子、`pid`は子に設けられる識別子であり、どの親に接続するかを示す識別子である。`id`と`pid`が同一である場合に親子関係の接続があることを示す。なお、図2の関係表スキーマ定義110の記法は実施例を限定するものではなく、同様の意味を表現し得る関係表スキーマの定義仕様であれば、どのような記法でも適用可能である。例えば、一般に関係表のスキーマはSQL式で作成時に定義するため、そのSQL式を関係表スキーマ定義110として利用できる。

#### 【0029】

上記のXML文書構造定義108および関係表スキーマ定義110に基づいて、両定義間の値の対応付けを定義するのが、スキーマ間マッピング定義109である。1行目の「`/x/i/@a i . a`」は、XML文書107の*i*要素の属性*a*の値を、関係表*i*(201)のカラム*a*に格納することを表現している。2、4、5、6、8、9行目も同様である。3行目の「`/x/i/j/. . . j . pid = i . id`」は、*j*要素の親を関係表*j*(202)のカラム*pid*で示し、関係表*i*(201)のカラム*id*を外部参照していることを表現している。7行目も同様に、関係表*k*(203)と関係表*j*(202)の間の外部参照を表現している。10行目は、XML文書構造定義108において定義されていない*j*要素の部分内容を関係表*j*(202)のカラム*w*に格納することを示している。ただし実際にはその部分構造は、タグごと構造検索エンジン106に格納され、その格納イメージ204に対して構造検索エンジン106上で付されたファイルのID(この例では*x i - a*)のみが関係表*j*(202)のカラム*w*に格納される。

#### 【0030】

データ変換モジュール101は、まず関係表スキーマ定義110に基づきRDB105を介してその記憶領域内に関係表*i*、*j*、*k*の各領域を確保する。次にデータ変換モジュール101は、XML文書107から*i*、*j*、*k*又は*o*要素の1つを取り出し、XML文書構造定義108を参照して取り出した要素の形式がそのスキーマ定義に合致するか否かチェックする。定義に合致すれば、データ変換モジュール101は、スキーマ間マッピング定義109を参照して取り出した要素の各属性の属性値をRDBの該当する関係表の1レコードとして格納し、そのレコードの*id*と*pid*を設定する。*id*はその関係表のレコード位置に応じた識別子を生成して設定する。*pid*にはメモリに保存された親要素の*id*があればその*id*を格納する。次にデータ変換モジュール101は、当該要素の要素名とその*id*をメモリに一時保存する。データ変換モジュール101は、XML文書107から*o*要素を取り出したとき、XML文書構造定義108を参照して取り出した要素がANYに相当することを認識し、関係表*j*の該当するレコードのカラムに構造検索エンジン106で指定されたファイルIDを設定し、取り出した*o*要素に*j*のタグを付けた部分構造を構造検索エンジン106に送る。構造検索エンジン106は、受け取った部分構造をその記憶領域に格納イメージ204として格納する。データ変換モジュール101は、XML文書107のすべての要素を取り出し終わるまでXML文書107から次の要素を取り出すステップに戻って上記処理を繰り返す。

#### 【0031】

図3を用いて、構造検索式変換モジュール102の動作について説明する。本説明では、図2の例で取り上げたXML文書構造定義108、スキーマ間マッピング定義109、および関係表スキーマ定義110に従って、データ変換モジュール101によりRDB105の関係表201、202、203、および構造検索エンジン106に格納イメージ204として格納されたXML文書107に対して、XQuery標準に従って記述された構造検索式111を処理する場合を例にとる。

#### 【0032】

構造検索式111は、FOR句により変数*\$i*に*i*要素名を代入する。従って、WHERE句は*i*要素の属性*a*が“`xx19`”であり、その*i*要素は、属性*b*が“`xx22`”であるような*j*要素を子要素として持ち、さらにその*j*要素は、属性*u*が“`xx24`”であ

るような  $\circ$  要素を子要素として持つことを条件としている。さらに、上記の条件で抽出した  $i$  要素全体を結果として取得することを要求している。

【0033】

構造検索式変換モジュール102は、スキーマ間マッピング定義109を参照して、上記の意味を持つ構造検索式111を構造指定UDFを含むSQL式112に変換する。マッピング定義109に従うと、上記の検索式の意味は、関係表  $i$  (201) のカラム  $a$  , 関係表  $j$  (202) のカラム  $b$  , および、関係表  $j$  (202) のカラム  $w$  に対して条件を指定していることと等価である。 $\circ$  要素は  $j$  要素の子要素としては定義されていないため、マッピング定義109の10行目を適用して、関係表  $j$  (202) のカラム  $w$  に対する条件となる。但し、このカラムは構造検索エンジン106に格納されたイメージ204への参照であるため、この条件は構造指定UDFで表現される。

10

【0034】

以上から、構造検索式111を変換したクエリ112は、関係表  $i$  (201) から、カラム  $a$  が “  $x x 1 9$  ” であるようなレコードを抽出し、関係表  $j$  (202) のレコードから、カラム  $b$  が “  $x x 2 2$  ”、かつ、カラム  $w$  で参照される構造検索エンジン106の格納イメージ(204)が、属性  $u$  の値 = “  $x x 2 4$  ” であるような  $\circ$  要素を含んでいるようなレコードを抽出し、さらに両レコードの間に外部参照関係が成り立っていることを条件として指定するSQL式として生成される。関係表  $j$  (202) のカラム  $w$  に対する構造指定は、XMLHASというUDFで表現される。これは、第一引数で指定したカラムの値が指し示す構造検索エンジン106上の格納イメージが、第二引数で指定したXPath構造式にマッチする部分データを含むか否かを判定するブール関数である。

20

【0035】

また構造検索式111は、上記の条件を満たす  $i$  要素全体を結果として取得することを要求しているため、SQL式112のSELECT句には、XMLVALというXML文書構築UDFを指定する。これは、引数でID指定された要素について、全ての子孫要素をRDB105および構造検索エンジン106から抽出して、XML文書を再構築して返すスカラ関数である。

【0036】

図3～図5を用いて、クエリ実行制御モジュール103の動作を説明する。

【0037】

図3は、RDB105がプラグイン処理機構301を持つ場合を示している。この場合、クエリ実行制御モジュール103が実現すべき機能はRDB105に組み込まれていることになる。ここでは、説明のためにこの機能を単独で実現するプログラムをクエリ実行制御モジュールと呼び、RDB105自身と区別する。クエリ実行制御モジュール103は、UDFを含むSQL式112をネイティブなSQL部とUDF部に分離し、RDB105がSQL部を処理し、プラグイン処理機構301がUDF部を処理する。構造指定UDFであるXMLHASは、実際には構造検索エンジン106で処理されるため、ほとんどの構造検索エンジンが対応している構造検索仕様XPathのクエリ302として、該エンジンに対して発行する。但し、XMLHASがRDB105に組み込みのプラグインとして実現されている場合は、RDB105上で直接この構造指定UDFを実行する。XML

30

40

【0038】

クエリ実行制御モジュール103は、クエリを実行する際に、先に構造検索エンジン106に対する条件でデータを絞るか、あるいは関係表に対する条件でデータを絞るか、クエリの処理効率を指標にして決定する。

【0039】

SQL式112を例にとると、前者の場合は、まずXMLHASの条件判定に適合する構造検索エンジン106上の格納イメージ204を抽出し、そのID(この例では  $x i - a$ ) と関係表  $j$  (202) のカラム  $w$  の値が一致することも条件に含めて、関係表201, 202からデータを抽出することになる。

50

## 【 0 0 4 0 】

一方、後者の場合は、まず関係表  $i$  ( 2 0 1 ) のカラム  $a$  と関係表  $j$  ( 2 0 2 ) のカラム  $b$ 、および関係表  $i$  ( 2 0 1 ) のカラム  $i d$  と関係表  $j$  ( 2 0 2 ) のカラム  $p i d$  の間の外部参照関係を条件にデータを絞り、抽出した関係表  $j$  ( 2 0 2 ) のレコードのカラム  $w$  の値が指し示す、構造検索エンジン 1 0 6 上の格納イメージ 2 0 4 に対して X M L H A S による条件判定を行う。

## 【 0 0 4 1 】

上記のようなクエリ実行手順の決定は、一般的な R D B 1 0 5 が備える実行計画決定処理により最適化される。従って、プラグイン処理機構 3 0 1 を備える R D B 1 0 5 を利用する場合は、クエリ実行制御モジュール 1 0 3 を新たに設ける必要はない。

10

## 【 0 0 4 2 】

一方、R D B 1 0 5 にプラグイン処理機構 3 0 1 が備わっておらず、クエリ実行制御モジュール 1 0 3 を R D B 1 0 5 の外部に設ける必要がある場合の動作概要について、図 4、図 5 を用いて説明する。クエリ実行制御モジュール 1 0 3 を R D B 1 0 5 の外に新たに設ける場合、上記のようなクエリ実行手順も独自に決定する必要がある。

## 【 0 0 4 3 】

図 4 を用いて、先に構造検索エンジン 1 0 6 に対する条件でデータを絞る場合について説明する。クエリ実行制御モジュール 1 0 3 が U D F を含む S Q L 式 1 1 2 を受けると、該モジュール内の U D F 分離処理 4 0 1 がネイティブな S Q L 式 4 0 3 と U D F 部に分離する。次にクエリ実行制御モジュール 1 0 3 は、U D F 部を構造検索エンジン 1 0 6 に対する X P a t h 検索式 3 0 2 として発行し ( 丸付き数字 1 )、この式にマッチするデータを含む格納イメージ 2 0 4 の I D ( この例では  $x i - a$  ) を獲得し、R D B 1 0 5 に一時表  $x$  ( 4 0 4 ) として格納する。次にクエリ実行制御モジュール 1 0 3 は、関係表  $j$  ( 2 0 2 ) のカラム  $w$  に格納されている I D が、この一時表  $x$  に含まれることも条件にして、S Q L 式 4 0 3 により関係表 2 0 1 , 2 0 2 からデータを抽出する ( 丸付き数字 2 )。ただし S Q L 式 4 0 3 の  $x$  は、一時表  $x$  を意味し、 $x . i d$  は一時表  $x$  の  $i d$  カラムを意味する。S Q L 式 4 0 3 による検索の結果として、クエリ実行制御モジュール 1 0 3 には  $i . i d$  として “  $r 0 2$  ” というデータが返る。

20

## 【 0 0 4 4 】

このようにしてタグ付き構造化文書再構成処理 4 0 2 は、抽出した I D を持つ  $i$  要素を関係表 2 0 1 ~ 2 0 3 および構造検索エンジン 1 0 6 の格納イメージ 2 0 4 のデータから再構成する。このため、タグ付き構造化文書再構成処理 4 0 2 は、スキーマ間マッピング定義 1 0 9 を参照して関係表よりデータを抽出する S Q L 式 4 0 5 ~ 4 0 7 を作成し、これらの S Q L 式を R D B 1 0 5 に対して発行する ( 丸付き数字 3 )。これらは、関係表間の外部参照関係に基づき、抽出した I D “  $r 0 2$  ” を持つ  $i$  要素の全子孫要素を抽出するものである。ここで S Q L 式 4 0 5 の  $i i d$  には “  $r 0 2$  ” が代入される。 $j i d$  には何も代入されず、結果的には S Q L 式 4 0 7 の結果は返らない。本例の場合には S Q L 式 4 0 7 がなくても構わない。一方、 $i$  要素は一部に構造検索エンジン 1 0 6 に保存された格納イメージ 2 0 4 のデータも含むため、タグ付き構造化文書再構成処理 4 0 2 は、それを取得するためのクエリ 4 0 8 を、構造検索エンジン 1 0 6 に対して発行し ( 丸付き数字 3 ) し、格納イメージ 2 0 4 を取得する。タグ付き構造化文書再構成処理 4 0 2 は、X M L 文書構造定義 1 0 8、スキーマ間マッピング定義 1 0 9 および関係表スキーマ定義 1 1 0 を参照し、抽出したデータと格納イメージ 2 0 4 から X M L 文書を再構成し、結果 1 1 3 を得る ( 丸付き数字 4 )。

30

40

## 【 0 0 4 5 】

図 5 を用いて、先に関係表 2 0 1 , 2 0 2 に対する条件でデータを絞る場合について説明する。この場合、U D F 分離処理 4 0 1 は、U D F を含む S Q L 式 1 1 2 を、ネイティブ S Q L 式 5 0 2 のようなクエリに分離する。クエリ実行制御モジュール 1 0 3 は、この S Q L 式を R D B 1 0 5 に発行し ( 丸付き数字 1 )、関係表 2 0 1 , 2 0 2 に対する条件でデータを絞る。その際に、関係表  $j$  ( 2 0 2 ) のカラム  $w$  の値も同時に抽出する。クエ

50

り実行制御モジュール103は、i.idとして“r02”、j.wとして“xi-a”という値を受け取る。構造判定処理501は、SQL式502の結果を受け取り、格納イメージID“xi-a”とXPath式302を構造検索エンジン106に送り(丸付き数字2)、これらの条件に合う格納イメージが構造検索エンジン106に登録されているか否かを判定する。構造検索エンジン106に該当するデータがあれば、格納イメージID“xi-a”をタグ付き構造化文書再構成処理402に渡す。以降の処理は、図4の場合と同一である。

#### 【0046】

なお、以上の説明では、構造検索式変換モジュール102とクエリ実行制御モジュール103を区別して説明したが、これらは一つのモジュールとして実現されていても構わない。その場合は、UDFを含むSQL式112を生成せずに、構造検索式111から直接SQL式403、502に変換するような実施例もあり得ることは自明である。

#### 【0047】

以降、本実施例におけるマッピング定義チューニングモジュール104が行う具体的なマッピング定義改善の処理手順について、図6、図7(a)、図7(b)、図8を用いて説明する。

#### 【0048】

図6を用いて、XML文書構造定義で明示的に定義されていない部分データについての検索頻度が所定数を越える場合の、マッピング定義改善処理について説明する。

#### 【0049】

システムは、タグ付き構造化文書に対する構造検索式の発行履歴を図示しない検索履歴データベースに記録する。マッピング定義チューニングモジュール104は、検索履歴データベースを参照し、同一のUDFについての構造検索式の発行頻度を計数する。図2～図5の例における構造検索式111中のo要素についての条件指定のように、XML文書構造定義に登場せず、従って明示的にスキーマ間マッピングを定義していない部分に対して検索が頻出する場合、マッピング定義チューニングモジュール104は、この部分を格納する関係表とマッピング定義を自動的に生成する。

#### 【0050】

RDBの検索処理性能は、長年の改良の結果、一般的な構造検索エンジンに比べ高速であり、また他の関係表データに対するのと同時に条件指定することを考慮した場合、データは、RDB外部の構造検索エンジンではなく、可能な限り関係表で管理した方が効率的に優れるため、このようなマッピングの変更は性能改善に繋がる。

#### 【0051】

本例では、マッピング定義チューニングモジュール104は、RDB105上に関係表o(601)を新規に作成し、関係表j(202)との間に外部参照関係を規定する。この時、関係表スキーマ定義110は、関係表スキーマ定義603に変更される。関係表oは、カラムpid, u, v, idの4つのカラムを持つと定義される。同時に、スキーマ間マッピング定義109は、スキーマ間マッピング定義604に変更される。マッピング定義チューニングモジュール104は、マッピング定義109の10行目にあった未定義部分を構造検索エンジンにマッピングすることを表す記述を削除し、新たに10行目に関係表jと関係表oの外部参照関係を表す記述、および11, 12行目に、o要素の各属性と関係表oの各カラムとの対応を表す記述を追加する。またマッピング定義チューニングモジュール104は、XML文書構造定義108の{ANY}を<o u="str" b="str"/>x0...nに変更する。

#### 【0052】

マッピング定義チューニングモジュール104が実行する処理手順の詳細は次の通りである。マッピング定義チューニングモジュール104は、スキーマ間マッピング定義109の各定義レコードをたどり、XML文書構造定義108に定義されていない要素の部分内容を見つける。次にマッピング定義チューニングモジュール104は、関係表スキーマ定義110を参照してその部分内容に定義された関係表とカラムの識別子を取得する。次

10

20

30

40

50

にマッピング定義チューニングモジュール104は、RDB105に対してSQL検索式を送付し、その関係表とカラム位置の属性値を取得する。その属性値が構造検索エンジン106のファイルIDを示しているので、マッピング定義チューニングモジュール104は、構造検索エンジン106からその格納イメージ204を取得する。次にマッピング定義チューニングモジュール104は、RDB105を介してその記憶領域内に関係表oの記憶領域を確保する。次にマッピング定義チューニングモジュール104は、上記のデータ変換モジュール101の処理手順に従って格納イメージ204からo要素を取り出し、RDB105を介して関係表oを作成する。次にマッピング定義チューニングモジュール104は、関係表スキーマ定義110に関係表oの定義を追加し、関係表スキーマ定義110を関係表スキーマ定義603に更新する。次にマッピング定義チューニングモジュール104は、スキーマ間マッピング定義109に関係表oについてのマッピング定義を追加し、スキーマ間マッピング定義109をスキーマ間マッピング定義604に更新する。次にマッピング定義チューニングモジュール104は、XML文書構造定義108の定義文をたどり、未定義の要素を見つけ、o要素の定義に置き換える。次にマッピング定義チューニングモジュール104は、構造検索エンジン106から格納イメージ204を削除する。

10

#### 【0053】

以上の変更が加えられたマッピング定義においては、構造検索式111は、構造検索式変換モジュール102によって、SQL式602に変換されることになる。このSQL式は、構造指定UDFを含まないため、RDB105で処理するのに望ましい形となっている。

20

#### 【0054】

図7(a)及び図7(b)を用いて、再帰的な構造を持つXMLデータ管理の改善を実現する処理手順について説明する。図7(a)に示すように、XML文書701は、XML文書構造定義702に妥当である、自己再帰的な構造を持つ。すなわちj要素の子要素としてj要素自身が複数出現する。このようなXML文書の関係表への格納方法は、スキーマ間マッピング定義703、および関係表スキーマ定義704によって定義される。マッピング定義703の3行目は、j要素の親はi要素がj要素であり、その区別を関係表jのカラムprlの値("i"または"j")で表現することを意味している。XML文書701の格納先となる関係表は、関係表i(705)および関係表j(706)の2つで、関係表iと関係表jの間の外部参照関係、および関係表j内部での自己参照関係が規定されている。

30

#### 【0055】

一方、構造検索式707は、属性aの値が"xx01"であるi要素の子孫要素として任意の階層に出現する、属性aの値が"xx18"であるようなj要素を抽出することを要求している。このことをSQL式で表現するには、再帰クエリを利用する必要がある。構造検索式707は、構造検索式変換モジュール102によって、SQL式708に変換される。このSQL式は、再帰的に関係表j(706)の自己参照関係を辿って、一時表tmpに、i要素の全ての子孫を抽出して行く再帰クエリである。

#### 【0056】

しかし、一般的にRDBの再帰クエリは効率の悪い処理であり、このような構造検索式が頻出する場合には、上記のようなマッピング定義は好ましくない。

40

#### 【0057】

これに対し、再帰構造を持つXML部分データを、敢えて構造検索エンジン106に格納することで改善を図る。一般的に構造検索エンジンは、階層の深いデータに対しても妥当な性能で検索処理が可能であるように設計されているため、関係表で管理するよりも効率が良い場合がある。

#### 【0058】

図7(b)に示すスキーマ間マッピング定義709は、上記のスキーマ間マッピング定義703における3~6行目のj要素を関係表j(706)に対応付けている記述を削除

50

し、新たに3行目に、i要素の子孫を全て構造検索エンジン106に格納する記述を追加している。関係表スキーマ定義710は、関係表i(705)に構造検索エンジン106での格納イメージのIDを格納するカラムwを追加している。

#### 【0059】

以上のマッピング定義においては、XML文書701は、関係表705および構造検索エンジン106の格納イメージ711, 712に分解して格納される。また構造検索式707は、構造検索式変換モジュール102によって、UDFを含むSQL式713に変換される。SQL式713は、SQL式708と比較して再帰を含まないシンプルなクエリとなっており、RDB105と構造検索エンジン106の適切な使い分けが成される。

#### 【0060】

なお、以上のようなRDB105での管理が非効率的であるXML文書を、敢えて構造検索エンジンに格納するように変更する改善手法は、再帰構造を持つXML文書以外でも適用可能である。例えば、階層の深いXML文書を関係表に格納する場合は、多数の関係表を定義してその間の外部参照関係を規定することになるが、このような関係表に対して構造検索をかける場合は、外部参照関係の条件を全てSQL式に加えなくてはならない。このような条件は、RDBにおいては検索コストの高いジョイン操作として処理されるため効率が悪い。このような場合に対しても、図7(b)のようなマッピング定義チューニング手法を適用することによって、検索効率を改善することが可能である。

#### 【0061】

階層の深いXML文書のマッピング定義の改善には、構造検索エンジンを用いない別の手法もある。図8を用いてこれを説明する。構造検索式801は、i要素、その子要素であるj要素、さらにその子要素であるk要素に関する条件を指定するクエリである。スキーマ間マッピング定義109を用いる場合には、この構造検索式は構造変換モジュール102によってSQL式803に変換されることになる。このSQL式には、二つのジョイン操作、“i.id = j.pid”、および“j.id = k.pid”の条件が含まれることになる。これに対し、関係表k(203)を関係表k(802)のように、関係表i(201)のカラムaと関係表j(202)のカラムcの値もレコードに含むように更新することによって、同じ構造検索式を関係表k(802)のみに対するクエリとして実行することが可能となる。

#### 【0062】

関係表スキーマ定義110、およびスキーマ間マッピング定義109は、それぞれ関係表スキーマ定義805、スキーマ間マッピング定義806に更新されることになる。スキーマ間マッピング定義806の1行目は、i要素の属性aの値を関係表k(802)のカラムiaにも格納することを表現している。6行目も同様である。以上のマッピング定義においては、構造検索式801は、構造検索式変換モジュール102によって、SQL式804に変換される。該検索式はジョイン操作を含まないため検索コストが低い。

#### 【0063】

複数の構造検索式の効率化を目的とする場合は、全ての構造検索式のパスの和を取って、上記と同様のマッピング定義改善手法を適用することが可能である。例えば、以下の構造検索式全てに関して効率化を図る場合：

```

・      / x / i [ @ a = “ . . ” ] / / k [ @ a = “ . . ” ]
・      / / j [ @ c = “ . . ” ] / k [ @ b = “ . . ” ]
・      / / i [ @ a = “ . . ” and @ b = “ . . ” ] / j [ @ c = “ . . ” ] / k

```

i要素の属性a, b、およびj要素の属性cの値を含むように関係表kを更新する。

#### 【0064】

このようなマッピング変更は、関係表の正規化を崩すことにあたり、一つの値を複数のカラムで管理することになるため、データの更新時にはオーバヘッドとなる。マッピング定義チューニングモジュール104は、更新クエリの発行履歴も併せて参照し、参照系クエリと更新系クエリの発行頻度の兼ね合いに応じて、このマッピング定義改善手法を適用

10

20

30

40

50

するか否かを決定する。

【0065】

なお、以上の説明で用いたスキーマ間マッピング定義の記法は実施例を限定するものではなく、同様の意味を表現し得る定義仕様であれば、どのような記法でも適用可能である。また、以上は、XML文書の管理方法として説明したが、本実施例における方法は、SGML、HTMLに代表されるタグ付き構造化文書一般の管理方法としても適用可能であることは自明である。

【図面の簡単な説明】

【0066】

【図1】実施例の全体構成図である。

10

【図2】実施例のXML文書のデータ変換機能に関する部分の構成図である。

【図3】実施例のクエリリライト機能に関する部分の構成図である。

【図4】実施例のクエリ実行機能に関する部分の構成図である。

【図5】実施例のクエリ実行機能に関する部分の構成図（続き）である。

【図6】実施例のスキーママッピング改善例を説明する図である。

【図7(a)】実施例のスキーママッピング改善例を説明する図（続き）である。

【図7(b)】実施例のスキーママッピング改善例を説明する図（続き）である。

【図8】実施例のスキーママッピング改善例を説明する図（続き）である。

【符号の説明】

【0067】

20

101 . . . タグ付き構造化文書 - 関係表間データ変換モジュール, 102 . . . 構造検索式変換モジュール, 103 . . . クエリ実行制御モジュール, 104 . . . マッピング定義チューニングモジュール, 105 . . . リレーショナルデータベース, 106 . . . 構造検索エンジン, 107 / 701 . . . タグ付き構造化文書, 108 / 702 . . . タグ付き構造化文書スキーマ定義, 109 / 604 / 703 / 709 / 806 . . . スキーマ間マッピング定義, 110 / 603 / 704 / 710 / 805 . . . 関係表スキーマ定義, 111 / 707 / 801 . . . 構造検索式, 112 / 602 / 708 / 713 / 803 / 804 . . . リライト結果のクエリ, 113 . . . (構造検索式の)結果, 201 ~ 203 / 601 / 705 / 706 / 802 . . . 関係表, 204 / 711 / 712 . . . (構造検索エンジンに対する部分XML文書の)格納イメージ

30

【図1】

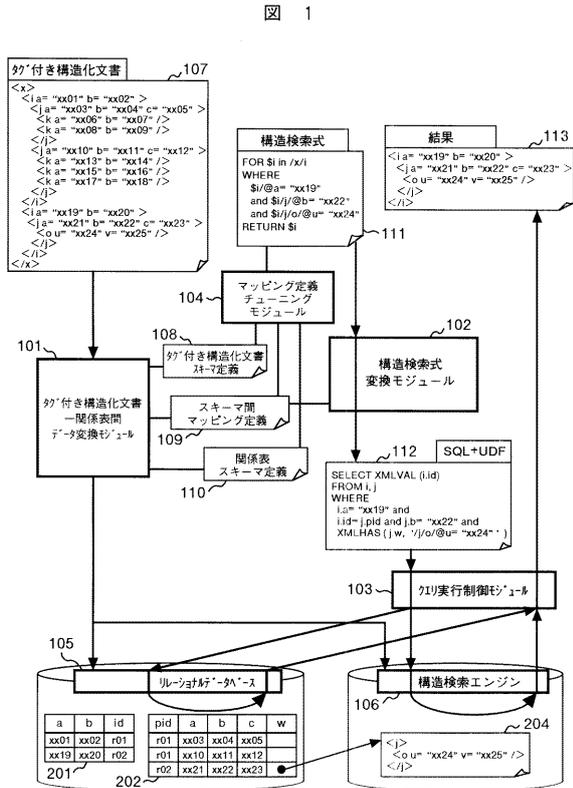


図 1

【図2】

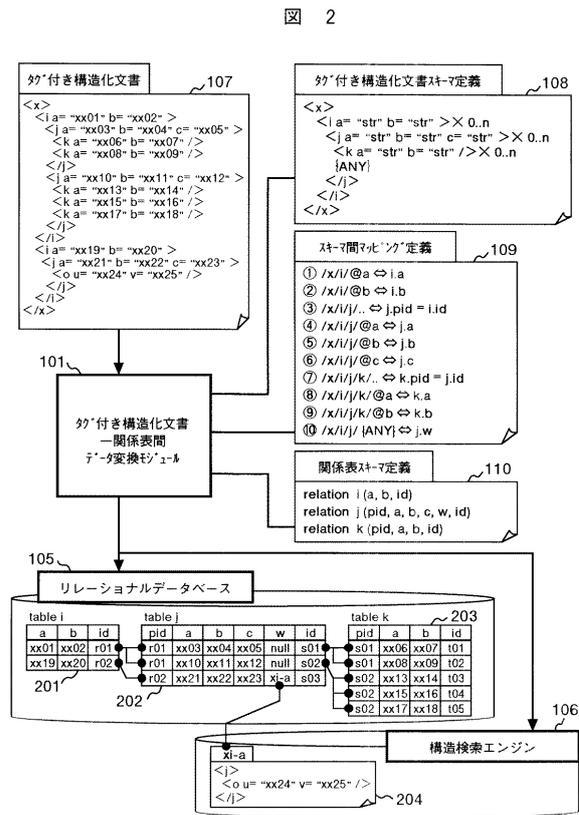


図 2

【図3】

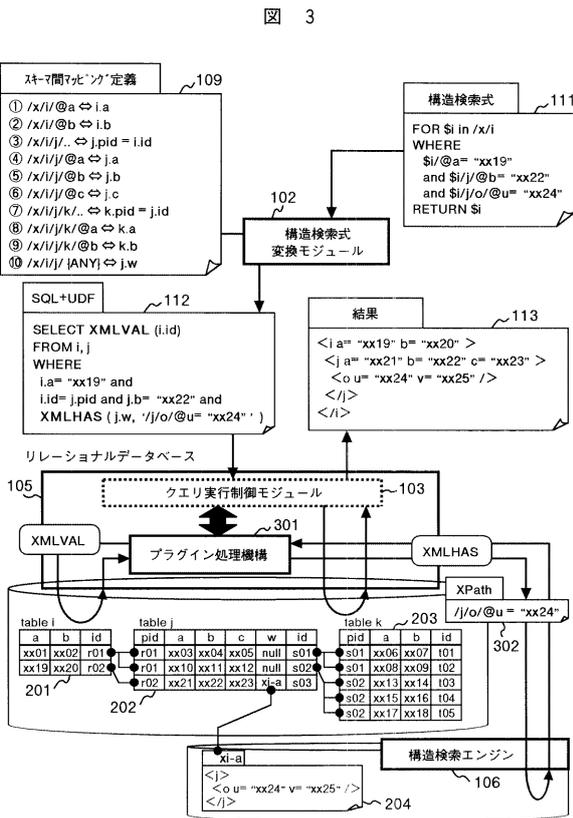


図 3

【図4】

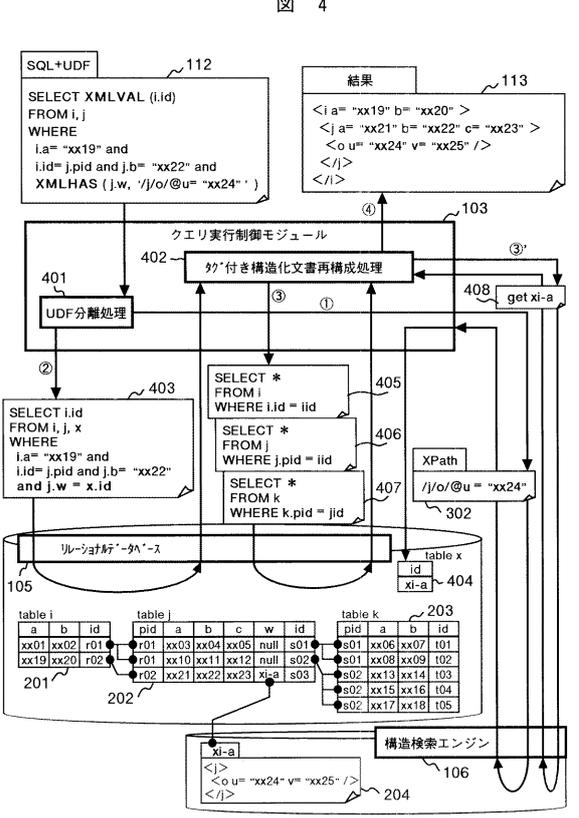
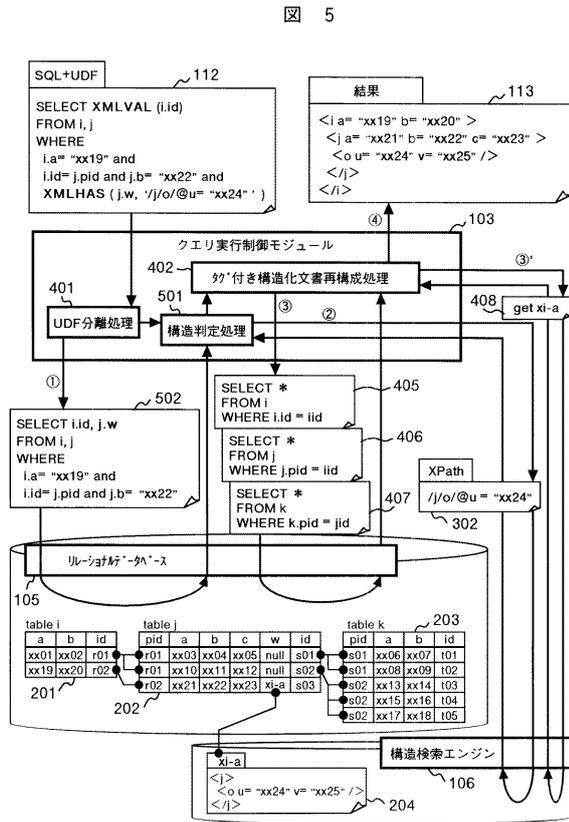
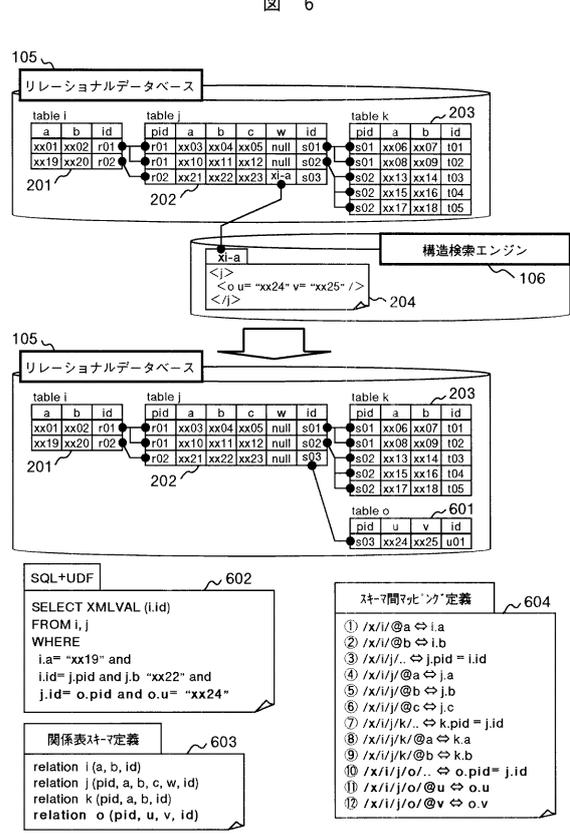


図 4

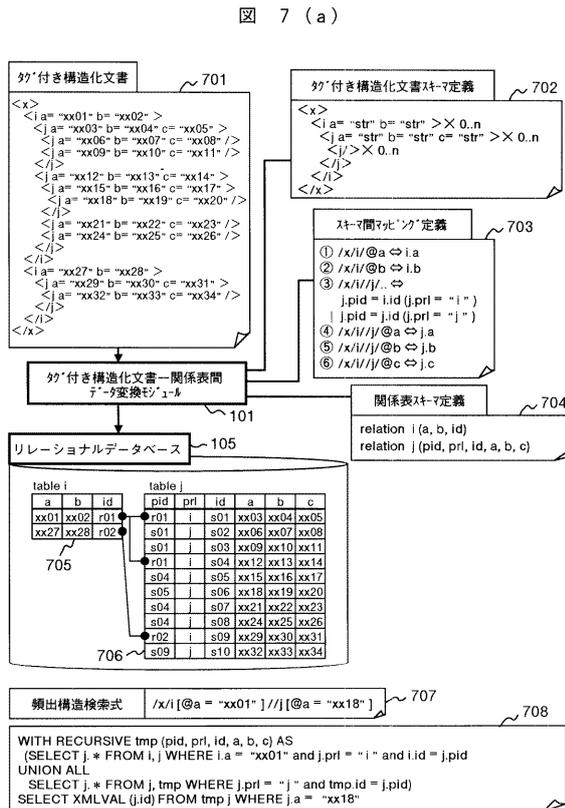
【図5】



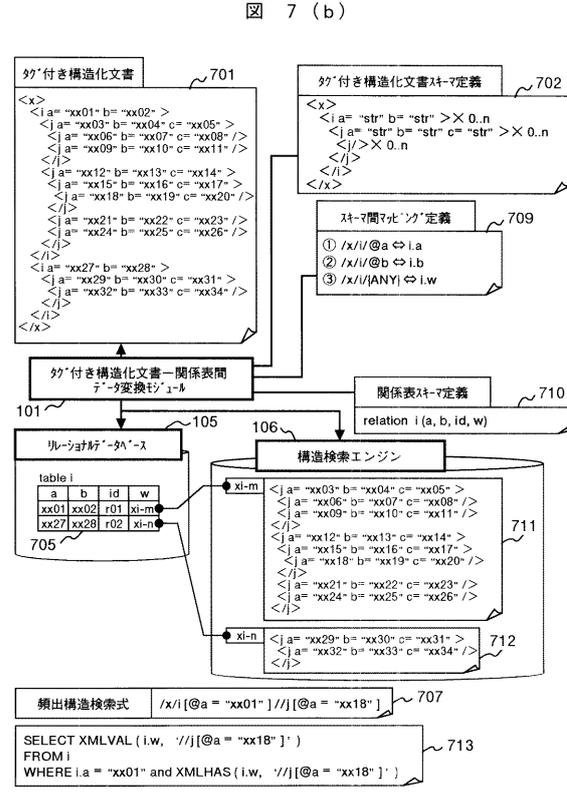
【図6】



【図7(a)】

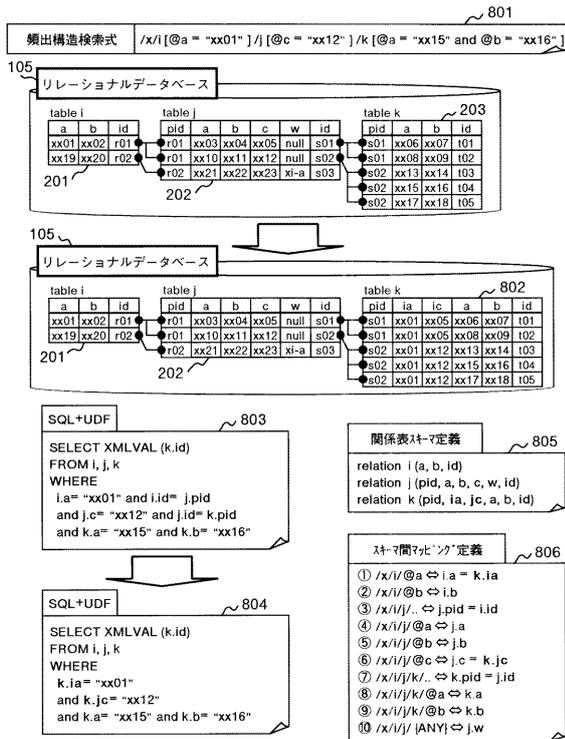


【図7(b)】



【 図 8 】

図 8



フロントページの続き

Fターム(参考) 5B075 ND35 PP23 QT06  
5B082 CA11 FA11 GA03 GA06 GA07 GA08