

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4868733号
(P4868733)

(45) 発行日 平成24年2月1日(2012.2.1)

(24) 登録日 平成23年11月25日(2011.11.25)

(51) Int. Cl. F I
G 0 6 F 17/21 (2006.01) G O 6 F 17/21 5 7 0 L
 G O 6 F 17/21 5 0 1 T

請求項の数 6 (全 15 頁)

(21) 出願番号	特願2004-340802 (P2004-340802)	(73) 特許権者	000001007
(22) 出願日	平成16年11月25日(2004.11.25)		キヤノン株式会社
(65) 公開番号	特開2006-154952 (P2006-154952A)		東京都大田区下丸子3丁目30番2号
(43) 公開日	平成18年6月15日(2006.6.15)	(74) 代理人	100076428
審査請求日	平成19年11月14日(2007.11.14)		弁理士 大塚 康德
		(74) 代理人	100112508
			弁理士 高柳 司郎
		(74) 代理人	100115071
			弁理士 大塚 康弘
		(74) 代理人	100116894
			弁理士 木村 秀二
		(74) 代理人	100130409
			弁理士 下山 治
		(74) 代理人	100134175
			弁理士 永川 行光

最終頁に続く

(54) 【発明の名称】 構造化文書処理装置及び構造化文書処理方法、プログラム

(57) 【特許請求の範囲】

【請求項1】

第1構造化文書、第2構造化文書、それぞれが前記第1構造化文書に含まれる要素の要素名を指定する複数の第1ロケーション情報、及び、それぞれが前記第2構造化文書に含まれる要素の要素名を指定する複数の第2ロケーション情報に基づいて、第3構造化文書を構築する構造化文書処理装置であって、

前記第1構造化文書、前記第2構造化文書、前記複数の第1ロケーション情報、及び、前記複数の第2ロケーション情報を、前記複数の第1ロケーション情報のそれぞれと前記複数の第2ロケーション情報のそれぞれを対応付けて入力する入力手段と、

前記入力手段により入力された前記第1構造化文書に含まれる要素のうち、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素のテキスト要素値を抽出する抽出手段と、

前記複数の第1ロケーション情報のそれぞれに対応付けられた前記複数の第2ロケーション情報のそれぞれで指定された要素名と要素名が一致する前記第2構造化文書中の要素のテキスト要素値として、前記抽出手段により前記第1構造化文書から抽出された、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素のテキスト要素値を挿入して、前記第3構造化文書を構築する構造構築手段と

を備えることを特徴とする構造化文書処理装置。

【請求項2】

第1構造化文書、第2構造化文書、それぞれが前記第1構造化文書に含まれる要素の要

10

20

素名を指定する複数の第1ロケーション情報、及び、それぞれが前記第2構造化文書に含まれる要素の要素名を指定する複数の第2ロケーション情報に基づいて、第3構造化文書を構築する構造化文書処理装置における構造化文書処理方法であって、

入力手段が、前記第1構造化文書、前記第2構造化文書、前記複数の第1ロケーション情報、及び、前記複数の第2ロケーション情報を、前記複数の第1ロケーション情報のそれぞれと前記複数の第2ロケーション情報のそれぞれを対応付けて入力する入力工程と、

抽出手段が、前記入力工程で入力された前記第1構造化文書に含まれる要素のうち、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素のテキスト要素値を抽出する抽出工程と、

構造構築手段が、前記複数の第1ロケーション情報のそれぞれに対応付けられた前記複数の第2ロケーション情報のそれぞれで指定された要素名と要素名が一致する前記第2構造化文書中の要素のテキスト要素値として、前記抽出工程で前記第1構造化文書から抽出された、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素のテキスト要素値を挿入して、前記第3構造化文書を構築する構造構築工程とを有することを特徴とする構造化文書処理方法。

【請求項3】

第1構造化文書、第2構造化文書、それぞれが前記第1構造化文書に含まれる要素の要素名を指定する複数の第1ロケーション情報、及び、それぞれが前記第2構造化文書に含まれる要素の要素名を指定する複数の第2ロケーション情報に基づいて、第3構造化文書を構築する構造化文書処理装置における構造化文書処理方法をコンピュータに実行させるプログラムであって、前記プログラムは、

前記第1構造化文書、前記第2構造化文書、前記複数の第1ロケーション情報、及び、前記複数の第2ロケーション情報を、前記複数の第1ロケーション情報のそれぞれと前記複数の第2ロケーション情報のそれぞれを対応付けて入力する入力手順と、

前記入力手順で入力された前記第1構造化文書に含まれる要素のうち、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素のテキスト要素値を抽出する抽出手順と、

前記複数の第1ロケーション情報のそれぞれに対応付けられた前記複数の第2ロケーション情報のそれぞれで指定された要素名と要素名が一致する前記第2構造化文書中の要素のテキスト要素値として、前記抽出手順で前記第1構造化文書から抽出された、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素のテキスト要素値を挿入して、第3構造化文書を構築する構造構築手順と

をコンピュータに実行させることを特徴とするプログラム。

【請求項4】

第1構造化文書、第2構造化文書、それぞれが前記第1構造化文書に含まれる要素の要素名を指定する複数の第1ロケーション情報、及び、それぞれが前記第2構造化文書に含まれる要素の要素名を指定する複数の第2ロケーション情報に基づいて、第3構造化文書を構築する構造化文書処理装置であって、

前記第1構造化文書、前記第2構造化文書、前記複数の第1ロケーション情報、及び、前記複数の第2ロケーション情報を、前記複数の第1ロケーション情報のそれぞれと前記複数の第2ロケーション情報のそれぞれを対応付けて入力する入力手段と、

前記入力手段により入力された前記第1構造化文書に含まれる要素のうち、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素の属性値を抽出する抽出手段と、

前記複数の第1ロケーション情報のそれぞれに対応付けられた前記複数の第2ロケーション情報のそれぞれで指定された要素名と要素名が一致する前記第2構造化文書中の要素の属性値として、前記抽出手段により前記第1構造化文書から抽出された、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素の属性値を挿入して、前記第3構造化文書を構築する構造構築手段と

を備えることを特徴とする構造化文書処理装置。

10

20

30

40

50

【請求項5】

第1構造化文書、第2構造化文書、それぞれが前記第1構造化文書に含まれる要素の要素名を指定する複数の第1ロケーション情報、及び、それぞれが前記第2構造化文書に含まれる要素の要素名を指定する複数の第2ロケーション情報に基づいて、第3構造化文書を構築する構造化文書処理装置における構造化文書処理方法であって、

入力手段が、前記第1構造化文書、前記第2構造化文書、前記複数の第1ロケーション情報、及び、前記複数の第2ロケーション情報を、前記複数の第1ロケーション情報のそれぞれと前記複数の第2ロケーション情報のそれぞれを対応付けて入力する入力工程と、

抽出手段が、前記入力工程で入力された前記第1構造化文書に含まれる要素のうち、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素の属性値を抽出する抽出工程と、

構造構築手段が、前記複数の第1ロケーション情報のそれぞれに対応付けられた前記複数の第2ロケーション情報のそれぞれで指定された要素名と要素名が一致する前記第2構造化文書中の要素の属性値として、前記抽出工程で前記第1構造化文書から抽出された、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素の属性値を挿入して、前記第3構造化文書を構築する構造構築工程と

を有することを特徴とする構造化文書処理方法。

【請求項6】

第1構造化文書、第2構造化文書、それぞれが前記第1構造化文書に含まれる要素の要素名を指定する複数の第1ロケーション情報、及び、それぞれが前記第2構造化文書に含まれる要素の要素名を指定する複数の第2ロケーション情報に基づいて、第3構造化文書を構築する構造化文書処理装置における構造化文書処理方法をコンピュータに実行させるプログラムであって、前記プログラムは、

前記第1構造化文書、前記第2構造化文書、前記複数の第1ロケーション情報、及び、前記複数の第2ロケーション情報を、前記複数の第1ロケーション情報のそれぞれと前記複数の第2ロケーション情報のそれぞれを対応付けて入力する入力手順と、

前記入力手順で入力された前記第1構造化文書に含まれる要素のうち、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素の属性値を抽出する抽出手順と、

前記複数の第1ロケーション情報のそれぞれに対応付けられた前記複数の第2ロケーション情報のそれぞれで指定された要素名と要素名が一致する前記第2構造化文書中の要素の属性値として、前記抽出手順で前記第1構造化文書から抽出された、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素の属性値を挿入して、第3構造化文書を構築する構造構築手順と

をコンピュータに実行させることを特徴とするプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、第1構造化文書、第2構造化文書に基づいて、第3構造化文書を構築する構造化文書処理装置、構造化文書処理方法、及びプログラムに関する。

【背景技術】

【0002】

電子的な構造化文書情報の効果的な情報利用技術は、インターネットを含む広範な情報の交換・流通にとってきわめて重要な位置を占めている。例えば、XML (eXtensible Markup Language) に代表されるこれらの技術は、WEBによる情報環境に向けて開発されたものであり、構造化言語として標準化されている。

【0003】

しかしながら、構造化文書構造の自動解析及び他の構造化文書への変換に関する体系的な言語処理技術は欠如している。従来、入力される構造化文書情報から必要な情報を抽出し、それらをまとめて、別の構造を持った構造化文書として出力するためには、あらかじ

10

20

30

40

50

め入力される構造化文書の構造と、出力する構造化文書の構造とを把握した上で、XSLT (XML Transformations)の作成や構造化文書から情報を抽出し、新たな別の構造をもった構造化文書として出力するプログラミングを行っていた。

【0004】

上述の従来技術として、例えば、以下の特許文献1及び特許文献2に示されるものがある。

【特許文献1】特開2004-30582号公報

【特許文献2】特開2004-38334号公報

【発明の開示】

【発明が解決しようとする課題】

10

【0005】

しかしながら、入力される構造化文書と、出力する構造化文書の構造があらかじめ分からなければ、入力される構造と出力する構造を意識したXSLTの作成や、プログラミングができず、入力される構造化文書から、必要な情報を抽出し、新たな別の構造をもった構造化文書として出力することは困難であった。

【0006】

そこで、本発明は、上記の問題点を鑑みて、構造化文書を別の構造化文書に変換処理を行う際、予め、構造化文書の属性を把握してプログラミングを行うことなく、ロケーションパスの対応付けに基づいて変換処理を行うことを可能にする構造化文書処理技術を提供することを目的とする。

20

【課題を解決するための手段】

【0007】

上記の目的を達成するべく、本発明にかかる構造化文書処理装置は、第1構造化文書、第2構造化文書、それぞれが前記第1構造化文書に含まれる要素の要素名を指定する複数の第1ロケーション情報、及び、それぞれが前記第2構造化文書に含まれる要素の要素名を指定する複数の第2ロケーション情報に基づいて、第3構造化文書を構築する構造化文書処理装置であって、

前記第1構造化文書、前記第2構造化文書、前記複数の第1ロケーション情報、及び、前記複数の第2ロケーション情報を、前記複数の第1ロケーション情報のそれぞれと前記複数の第2ロケーション情報のそれぞれを対応付けて入力する入力手段と、

30

前記入力手段により入力された前記第1構造化文書に含まれる要素のうち、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素のテキスト要素値を抽出する抽出手段と、

前記複数の第1ロケーション情報のそれぞれに対応付けられた前記複数の第2ロケーション情報のそれぞれで指定された要素名と要素名が一致する前記第2構造化文書中の要素のテキスト要素値として、前記抽出手段により前記第1構造化文書から抽出された、前記複数の第1ロケーション情報のそれぞれで指定された要素名と要素名が一致する要素のテキスト要素値を挿入して、前記第3構造化文書を構築する構造構築手段と

を備えることを特徴とする。

【発明の効果】

40

【0008】

本発明により、第1構造化文書から第2構造化文書に挿入するテキスト要素値、属性値を簡単に変更することが可能になる。

【発明を実施するための最良の形態】

【0009】

以下、本発明の実施の形態を、具体例を用いて詳細に説明する。

【0010】

(実施形態1)

図1Aは、本発明の実施形態にかかる構造化文書処理装置の構成を示すブロック図である。図1Aにおいて、101は、構造化文書処理装置における解析変換処理部を示す。解

50

析変換処理部101は、不図示のCPU等の制御ユニットの全体的な制御の下、解析変換処理を実行することが可能であり、処理結果を不図示のメモリに格納し、必要に応じてメモリからデータを読み出して処理を実行することが可能である。この解析変換処理部101の内部には、構造化文書(inputA.xml)102、構造化文書(outputB.xml)103やロケーションパスA、B(104、105)を入力処理するための入力部106、構造自動解析部107、ロケーションパス対応付け及びデータ変換処理部113、構造構築部114、出力部115が含まれる。

【0011】

102は、入力部106に入力される構造化文書(inputA.xml)であり、103は、102の構造化文書を変換して出力するための文書の構造を記述する出力用の構造化文書(outputB.xml)(共にXMLデータであり、以下、「XMLデータ」と表現する)である。104、105は、それぞれの構造化文書102、103の内部データを指し示すためのロケーションパスを示す。

10

【0012】

入力部106は、インターネットを含んだネットワークと接続し、ネットワークを介して構造化文書(例えば、XML、SGML、HTML等)及びロケーション情報(ロケーションパス)を受信することができるものとする。入力部106が、構造化文書(XMLデータ)102、103とロケーションパスA、B(104、105)を受信すると、構造自動解析部107は処理を開始する。

【0013】

20

107は、構造自動解析部であり、入力部106が受信した構造化文書(102、103)及びロケーションパスA、B(104、105)の構造解析を行い、その構造解析に基づいて、データの選択、構造の再構築を行う。

【0014】

図1Bは、構造自動解析部107の内部を具体的に示すブロック図である。図1B中、109は、構造解析・分解部であり、例えば、SAX(The Simple API for XML)エンジンを利用して、入力されたXMLデータをパースしていき、不図示のメモリ上にパースによって順番に検出されるXMLデータ構造に含まれる要素ごとにリストを作成し、要素の名前、要素タグの接続関係(親子関係、兄弟関係)、要素タグと要素タグで囲まれる"値"、属性の名前、属性の値を1つの要素に対して付加情報として対応づけ、XMLデータ解析リストとしてXMLデータの構造を分解した情報を格納する。

30

【0015】

108はロケーションパス解析部であり、入力されたロケーションパスA、B(104、105)をそれぞれ解析し、不図示のメモリ上に作成したXMLデータ解析リストの内容とリンクして指し示すことが可能な形式でXMLデータ解析リストに格納する。

【0016】

111はデータ選択・抽出部であり、XMLデータ解析リストから、格納してあるロケーションパスが指し示すデータのリスト番号と、そのリスト番号と親子関係あるいは兄弟関係など関連するリスト番号をすべて選択し、抽出する処理を行う。

【0017】

40

112は構造構築部であり、XMLデータ解析リストから、リスト番号に基づいて、一度分解したXMLデータのみを抽出し、構築しなおして、例えば、XMLデータとして出力する処理を行う。

【0018】

説明を図1Aに戻し、113は、ロケーションパス対応付け・データ変換処理部であり、入力部106に入力されたロケーションパスA104とロケーションパスB105を1対1で対応付け、入力用XMLデータ(inputA.xml)102と、出力用XMLデータ(outputB.xml)103とのそれぞれのXMLデータ解析リストに格納されているデータを、ロケーションパスの対応付けによって変換し、変換した出力用XMLデータのXMLデータ解析リストを出力する。

50

【 0 0 1 9 】

1 1 4 は構造構築部であり、XMLデータ解析リストから再びXMLデータとして構築するという1 1 2と同様の処理を行い、出力部1 1 5から、構築したXMLデータ(output.xml) 1 1 6を出力する。

【 0 0 2 0 】

出力部1 1 5は、インターネットを含んだネットワークと接続し、ネットワークを介して構築した構造化文書(例えば、XML、SGML、HTML等)を他の装置に配信することができるものとする。

【 0 0 2 1 】

図2は、構造解析・分解部1 0 9の処理を説明する図である。図2(a)において、2 0 6は、入力用XMLデータ(inputA.xml) 1 0 2の具体例を示す図である。このXMLデータ2 0 6が構造解析・分解部1 0 9に入力されると、図2(b)に示すフローチャートが構造解析・分解部1 0 9により実行される。

10

【 0 0 2 2 】

図2(b)のステップS 2 0 1でXMLデータを入力し、ステップS 2 0 2でSAXエンジンによって、XMLデータ2 0 6を先頭のタグからパースしていく。ステップS 2 0 3で、メモリ上に、パースによって順番に出現するXMLデータの構造における要素ごとに番号を振ってリストを作成し、ステップS 2 0 4で、要素の階層番号、要素の名前、namespace名、要素タグと要素タグで囲まれる"値"であるテキスト要素値、要素タグの親を示す親要素番号、要素タグの兄弟を表す1つ前兄弟番号、同じ階層で同じ要素名がある場合に判別できるよう、要素が出現した順番を示す要素出現順番号、要素が持っている属性の数、要素が持っている属性の名前と属性の値を全て、1つの要素に対して付加情報として対応づけを行い、ステップS 2 0 5で、全ての要素に対し分解処理が終了したか否かを判断し、終了しない場合は、ステップS 2 0 4の処理を続行し、分解処理が終わり次第パースを終了する(S 2 0 5)。

20

【 0 0 2 3 】

図2(c)の2 0 9は、XMLデータ解析リストの例を示す図であり、不図示のメモリ上にXMLデータの構造を分解した情報がXMLデータ解析リストとして格納される。

【 0 0 2 4 】

図3は、ロケーションパス解析部1 0 8の処理を説明する図であり、3 1 0は入力用XMLデータ1 0 2中のデータを指し示すロケーションパス1 0 4の具体例を示す図である。このロケーションパス3 1 0がロケーションパス解析部1 0 8に入力されると、図3(b)に示すフローチャートがロケーションパス解析部1 0 8により実行される。

30

【 0 0 2 5 】

図3(b)のステップS 3 0 1で、ロケーションパス3 1 0を入力し、ステップS 3 0 2で、ロケーションパス解析部1 0 8は、与えられたロケーションパスを1から順に解析して分解し、不図示のメモリ上にリストとして保存する。

【 0 0 2 6 】

図3(c)の3 0 4は、分解されて、リストとして保存されているロケーションパスの例を示す図である。ロケーションパス3 1 0中において、例えば、"student[position()=2]"という表現があれば、これは、同一階層番号上に同じ要素名が存在し、それを区別するために、その中でも2番目の要素という指定がされていることを示し、この場合、"student#2"という形でメモリ上にロケーションパスが格納される。また、例えば、"class/number"という表現があれば、class@numberという形でロケーションパスが格納される。

40

【 0 0 2 7 】

図4A、Bは、データ選択・抽出部1 1 1の処理内容を説明する図である。

【 0 0 2 8 】

図4Aにおいて、4 0 1は構造解析・分解部1 0 9から出力されるXMLデータ解析リストを例示しており、また、図4Bはデータ選択・抽出部1 1 1の処理の流れを説明するフローチャートである。

50

【 0 0 2 9 】

まず、図 4 B のステップ S 4 0 1 で、構造解析・分解部 1 0 9 は、先に説明した図 3 (c) の 3 0 4 のロケーションパスのリストを順に取得する。

【 0 0 3 0 】

次いで、ステップ S 4 0 2 で、XML データ解析リストから、要素番号順にリスト (図 4 A の 4 0 0) を取得し、取得したリストから要素名と、取得したロケーションパスリストの最下層の要素名 address (図 3 (c) では、リストの 1 番目なら address、リストの 2 番目なら tel、リストの 3 番目なら name) を比較していき、XML データ解析リスト 4 0 1 内の要素名がロケーションパスで指定された要素名と同じ要素名であるか検索する。

【 0 0 3 1 】

ステップ S 4 0 3 において、要素名が同じでなければ (S 4 0 3 - N o)、次の要素番号のリストを取得するステップ S 4 0 2 の処理まで戻る。同じであれば (S 4 0 3 - Y e s)、処理をステップ S 4 0 4 に進め、その要素番号を保存する。

【 0 0 3 2 】

ステップ S 4 0 5 で、その要素番号に対するリスト内から親要素番号を検索する。

【 0 0 3 3 】

ステップ S 4 0 6 において、親要素番号が「-1」であれば (S 4 0 6 - Y e s)、ステップ S 4 1 9 に処理を進め、ヒットしたリスト番号と、それに関係する保存したリスト番号を出力する。1 番最初から親要素番号が「-1」であれば、ルート of の要素名であり、関係する保存したリスト番号は無いことになる。

【 0 0 3 4 】

一方、ステップ S 4 0 6 において、親要素番号が-1 でなければ (S 4 0 6 - N o)、処理をステップ S 4 0 7 に進め、親要素番号が示すリストを取得する。

【 0 0 3 5 】

ステップ S 4 0 8 で、取得してあるロケーションパスに次の要素名、例えば、最下層の要素名の親の要素名 (リスト 1 であれば address の前の security)、が指定されているか検索する。

【 0 0 3 6 】

ステップ S 4 0 9 において、指定の有無を判定し、指定されていなければ (S 4 0 9 - N o)、その親要素番号の示すリストは必要なデータであるので、処理をステップ S 4 1 7 に進め、取得したリスト番号を保存する。

【 0 0 3 7 】

ステップ S 4 0 9 の判定において、指定されていれば (S 4 0 9 - Y e s)、処理をステップ S 4 1 0 に進め、親要素番号の示すリストの要素名と、そのロケーションパスで指定されている要素名を順 (例 security teacher) に比較していく。

【 0 0 3 8 】

ステップ S 4 1 1 で、ロケーションパスで指定された要素名と、XML データ解析リスト内の親要素番号を辿っていき、ひとつでも要素名の不一致があれば (S 4 1 1 - N o)、その親要素番号で指定したリストデータは必要が無いと判断し、次の要素番号のリストを取得する処理に戻る。要素名が一致していれば (S 4 1 1 - Y e s)、ステップ S 4 1 2 でロケーションパスで指定されている要素名に、データとして「#N(N: 数値)」が付加されているか確認する (S 4 1 2)。

【 0 0 3 9 】

ステップ S 4 1 3 で、「#N」が付加されているか判断し、付加されていれば (S 4 1 3 - Y e s)、#N の N で指定された番号と、ステップ S 4 1 4 で要素出現順番号を検索した番号をステップ S 4 1 5 で比較する。

【 0 0 4 0 】

ステップ S 4 1 6 で番号が一致しているか判断し、一致していなければ (S 4 1 6 - N o)、その親要素番号で指定したリストデータは必要が無いと判断し、次の要素番号のリストを取得する処理に戻る。

10

20

30

40

50

【 0 0 4 1 】

一致していれば (S 4 1 6 - Y e s)、その親要素番号を保持し (S 4 1 7)、そのリストデータの親要素番号を取得 (S 4 1 8) する。

【 0 0 4 2 】

そして、再びステップ S 4 0 6 に処理を戻し、親要素番号が -1 でないかの判断の処理からまた繰り返す。

【 0 0 4 3 】

最終的にロケーションパスで指定されたパスで、XMLデータ解析リスト内でヒットした要素番号とそれに関係する要素番号を保持する。具体的には、ロケーションパス (図 3 (c) の 3 0 4) で 1 番目のリストのロケーションパス (teacher/security/address) にヒットする番号は、全部で、5 4 2 1 0 の順でヒットし、これらが抽出される番号となる。1つのロケーションパスに対して処理が終われば、2番目のロケーションパス (class/student[position="2"]/security/tel) を取得し、同様の処理を繰り返す。2番目のリストのロケーションパスにヒットする番号は、1 6 1 4 1 2 1 0 の順でヒットする。3番目のリストのロケーションパス (name) にヒットする番号は、name という指定だけなので、3 2 1 0、8 7 1 0、1 3 1 2 1 0 という順番でヒットする。

10

【 0 0 4 4 】

すべての処理が終わった時点で、抽出した番号で重なるものがあれば、重複した番号は1つだけ残し、その他はすべて削除し、最終的に残った番号を出力する。その結果、データ選択・抽出部 1 1 1 が選択する要素番号のリストは (1) 式のようになる。

20

【 0 0 4 5 】

要素番号のリスト : 0, 1, 2, 3, 4, 5, 7, 8, 12, 13, 14, 16 . . . (1)

これが、すべてのロケーションパスで指定されたデータに関連する要素番号となる。

【 0 0 4 6 】

次に、構造構築部 1 1 2 の処理を図 5 A、B を参照しつつ説明する。

【 0 0 4 7 】

図 5 A は、構造構築部 1 1 2 の処理内容を説明するフローチャートである。まず、ステップ S 5 0 1 において、XMLデータ解析リスト (図 4 A を参照) を取得する。

【 0 0 4 8 】

そして、ステップ S 5 0 2 で、XMLデータ解析リストから抽出する要素番号を保持したリスト ((1) 式を参照) を取得する。

30

【 0 0 4 9 】

ステップ S 5 0 3 で、XMLデータ解析リスト 4 0 1 (図 4 A) から要素番号のリストデータ ((1) 式) を順に取得する。

【 0 0 5 0 】

ステップ S 5 0 4 で、抽出したリストデータの階層番号を取得し、保持する。

【 0 0 5 1 】

ステップ S 5 0 5 で、前回保持した階層番号と今回保持した階層番号を比較する。ステップ S 5 0 6 で、初回の処理でない、もしくは、前回保持した階層番号 > = 今回取得した階層番号の判断が偽ならば (S 5 0 6 - N O)、処理をステップ S 5 0 7 に進め、リストデータから要素名を取得し、「 < 要素名 > 」の開始タグの形にして (' < ' > ' 記号を付加して) 文字列 (図 5 B の 5 2 0) として格納する。

40

【 0 0 5 2 】

ステップ S 5 0 8 で、リストデータ ((1) 式) に基づいて、XMLデータ解析リスト中に属性名、属性値で該当するものがあれば順に取得し、開始タグの中に追加する形で、文字列 (図 5 B の 5 2 0) に追加する。

【 0 0 5 3 】

ステップ S 5 0 9 で、リストデータ ((1) 式) に基づいて、XMLデータ解析リスト中にテキスト要素値で該当するものがあれば取得し、開始タグの後に、文字列 (図 5 中 5

50

20)として追加する。

【0054】

ステップS510で、取得した要素名を「</要素名>」の終了タグの形にして、スタック(図5Bの521)に格納(push)する。

【0055】

ステップS506の判断で真ならば(S506-Yes)、処理をステップS515に進め、スタック(図5Bの521)に格納してある終了タグを1回取り出して(pop)、文字列(図5Bの520)に格納する。

【0056】

ステップS511で、XML解析リスト(図4A)から要素番号のリスト((1)式)に記述されている要素番号のデータ分、データの抽出をしていなければ(S511-Yes)、ステップS503に処理を戻し、抽出していないデータに対して同様の処理を繰り返す。

10

【0057】

一方、リストデータ((1)式)に記述された要素番号のデータをすべて抽出し終わったら(S511-No)、処理をステップS512に進め、図5Bのスタック521に終了タグがまだ残っているか確認する。

【0058】

ステップS513で、終了タグがスタックに残っていれば(S513-Yes)、スタックから終了タグを取得(pop)し、文字列(図5Bの520)に追加する。

20

【0059】

終了タグが残っていなければ(S513-No)、処理をステップS514に進め、格納した文字列(図5Bの520)を出力する。

【0060】

データ選択・抽出部111は、XML解析リスト(図4Aの401)から、要素番号のリスト((1)式)に記述されている要素番号のデータだけを抽出して、文字列を組み立て最終的に出力することができ、その結果、データ選択・抽出部111は、入力されたXMLデータの構造を崩さず、ロケーションパスによって選択されたデータだけを抽出して、図5BにおけるXMLデータ522を出力することができる。

【0061】

30

図6は、ロケーションパス対応付け・データ変換部113の処理内容を説明する図である。これまでは、入力されたXMLデータが1つで、そのデータからロケーションパスで選択したデータだけを抽出してXMLデータとして構造を崩さずに構築しなおして出力するところまで述べたが、前記ロケーションパス対応付け&データ変換処理部を、前記構築構築部の処理が行われる前に加えることで、前記入力データの入力用としたXMLデータ(以下入力用XMLデータとする)、前記入力データの出力用としたXMLデータ(以下出力用XMLデータとする)をそれぞれ読み込み、入力用XMLデータからロケーションパスで選択したデータを、出力用XMLデータのロケーションパスで選択した要素に挿入することで、異なる構造をしたXMLデータによるデータの変換が可能となる。以下その流れを説明する。

40

【0062】

図6の601は、入力用XMLデータであり、これまでの説明の中で使用してきた入力用XMLデータ(inputA.xml)206と同様の構造とデータ内容を有する。602は、出力用XMLデータであり、既にテキスト要素値が一部格納してあり、入力用XMLデータとは構造が異なっている。

【0063】

603は、入力用XMLデータに対するロケーションパスを示しており、604は、出力用XMLデータに対するロケーションパスである。

【0064】

ロケーションパス対応付け・データ変換部113では、この両者のロケーションパスを

50

605のようにリストの一番目同士、二番目同士という形で、1対1で対応付ける。このロケーションパスの対応付けの情報を使って、ロケーションパス対応付け・データ変換部113におけるデータ変換処理は、不図示のメモリ上に保持してある入力用XMLデータ601を読み込み、構造解析の結果得られるXMLデータ解析リスト606から、それぞれ、ロケーションパスで選択されたデータを示す要素番号を検索し、出力用XMLデータ602の解析から得られるXMLデータ解析リスト607へコピーする(608)。その結果、ロケーションパス対応付け・データ変換部113は、XMLデータ解析リスト609を求め、出力する。

【0065】

そして、ロケーションパス対応付け・データ変換部113は、このデータの要素番号をすべてリストデータ610として、構造構築部112に入力する。構造構築部112は、XMLデータ解析リスト609と、要素番号のリストデータ610を受け取り、図5Aで説明したリストデータに対応する抽出処理を行うことで、XMLデータ612(図6(e))を出力する。具体的には、入力用XMLデータ601内のデータである、「Aさん」、「東京都」、「045-3333-3333」というテキスト要素値が、出力用XMLデータ602のそれぞれ、ロケーションパスで選択されたテキスト要素値へ格納される。

10

【0066】

この結果、構造の異なるXMLデータ601からロケーションパスで指定したデータを、別の構造を持ったXMLデータ602のロケーションパスで指定した領域にデータを挿入することができ、結果的に、XMLの構造を変換して出力することが可能になる。従って、予め入力用XMLデータと出力用XMLデータの構造を把握していなくても、XMLデータの解析から自動的に変換処理を行うことが可能になるので、動的なデータの交換、動的なXMLの構造変換が可能になる。

20

【0067】

尚、本実施形態では、テキスト要素値の変換しか行っていないが、属性値に対してもロケーションパスの指定で、構造の異なるデータ間で同様の変換処理を行うことが可能になる。

【0068】

また、ロケーションパスの関連付けを示すデータをネットワークを介して、他の装置と交換することも可能であり、異なる構造を有する構造化文書のやり取りが、ネットワーク上でも可能になる。

30

【0069】

以上説明したように本実施形態によれば、構造化文書を別の構造化文書に変換処理を行う際、予め、構造化文書の属性を把握してプログラミングを行うことなく、ロケーションパスの対応付けに基づいて変換処理を行うことが可能になる。

【0070】

(他の実施形態)

本発明の目的は、前述した実施形態の機能を実現するソフトウェアのプログラムコードを記録した記憶媒体を、上述の構造化文書処理装置に供給し、その装置のコンピュータ(またはCPUやMPU)が記憶媒体に格納されたプログラムコードを読み出し実行することによっても、達成されることは言うまでもない。プログラムコードの格納は、クライアントコンピュータに限定されるものではなく、例えば、サーバとして機能するコンピュータに記憶されておくことも可能である。

40

【0071】

この場合、記憶媒体から読出されたプログラムコード自体が前述した実施形態の機能を実現することになり、そのプログラムコードを記憶した記憶媒体は本発明を構成することになる。

【0072】

プログラムコードを供給するための記憶媒体としては、例えば、フレキシブルディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、DVD、磁

50

気テープ、不揮発性のメモリカード、ROMなどを用いることができる。

【0073】

また、コンピュータが読出したプログラムコードを実行することにより、前述した実施形態の機能が実現されるだけでなく、そのプログラムコードの指示に基づき、コンピュータ上で稼働しているOS（オペレーティングシステム）などが実際の処理の一部または全部を行い、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

【0074】

さらに、記憶媒体から読出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPUなどが実際の処理の一部または全部を行い、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

10

【図面の簡単な説明】

【0075】

【図1A】本発明の実施形態にかかる構造化文書処理装置の構成を示すブロック図である。

【図1B】構造自動解析部107の内部を具体的に示すブロック図である。

【図2】構造解析・分解部109の処理を説明する図である。

【図3】ロケーションパス解析部108の処理を説明する図である。

20

【図4A】構造解析・分解部109から出力されるXMLデータ解析リストを例示する図である。

【図4B】データ選択・抽出部111の処理の流れを説明するフローチャートである。

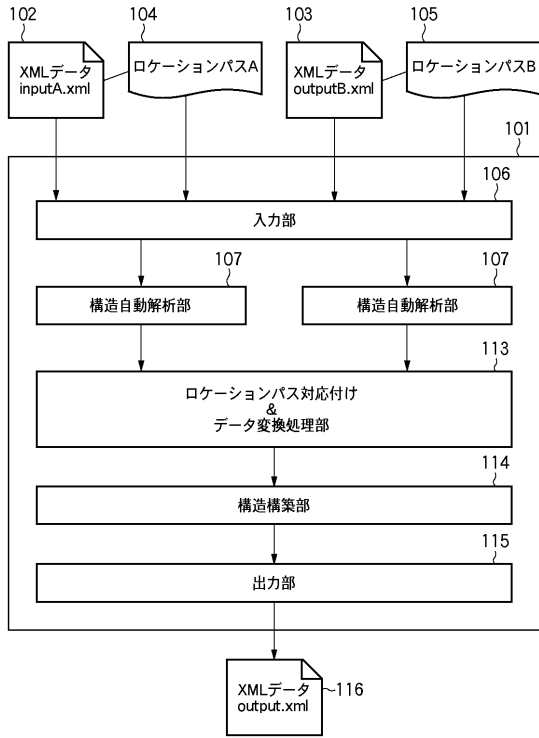
【図5A】構造構築部112の処理内容を説明するフローチャートである。

【図5B】構造構築部112の処理の内容を説明する図である。

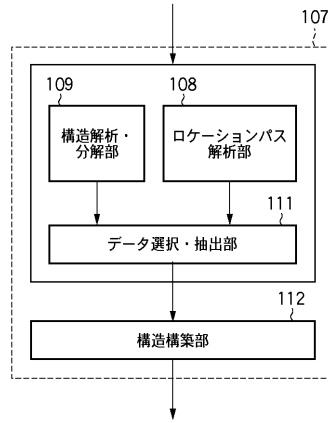
【図6】ロケーションパス対応付け・データ変換部113の処理内容を説明する図である。

。

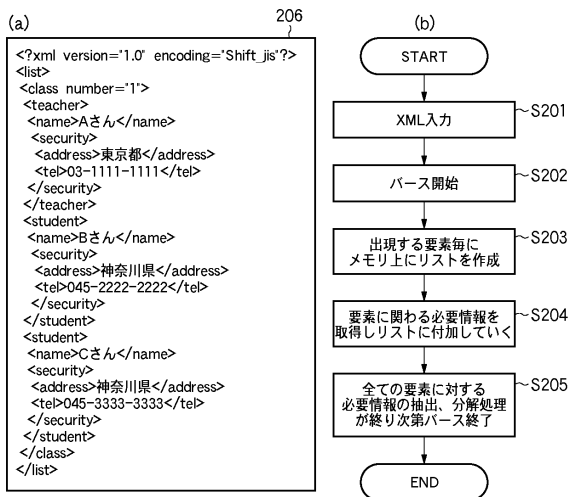
【図1A】



【図1B】



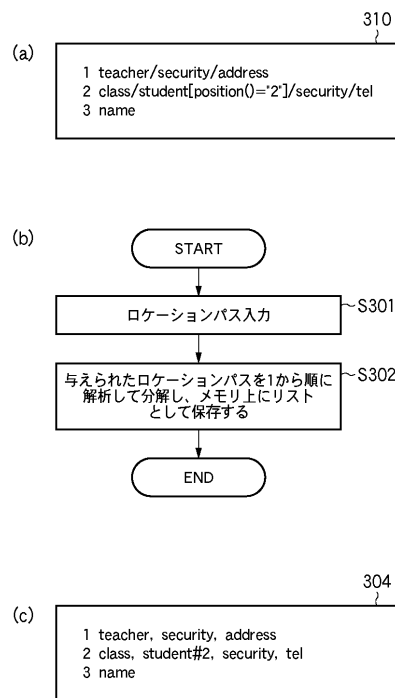
【図2】



(c)

要素番号	階層番号	要素名	name space名	テキスト要素値	親要素番号	1つ前兄弟番号	要素出現順番号	属性数	属性名	属性値
0	0	list	null		-1	0	1	0		
1	1	class	null		0	0	1	1	number	1
2	2	teacher	null		1	0	1	0		
3	3	name	null	Aさん	2	0	1	0		
4	3	security	null		2	3	1	0		
5	4	address	null	東京都	4	0	1	0		
6	4	tel	null	03-1111-1111	4	5	1	0		
7	2	student	null		1	2	1	0		
8	3	name	null	Bさん	7	0	1	0		
9	3	security	null		7	8	1	0		
10	4	address	null	神奈川県	9	0	1	0		
11	4	tel	null	045-2222-2222	9	10	1	0		
12	2	student	null		1	7	2	0		
13	3	name	null	Cさん	12	0	1	0		
14	3	security	null		12	13	1	0		
15	4	address	null	神奈川県	14	0	1	0		
16	4	tel	null	045-3333-3333	14	15	1	0		

【図3】



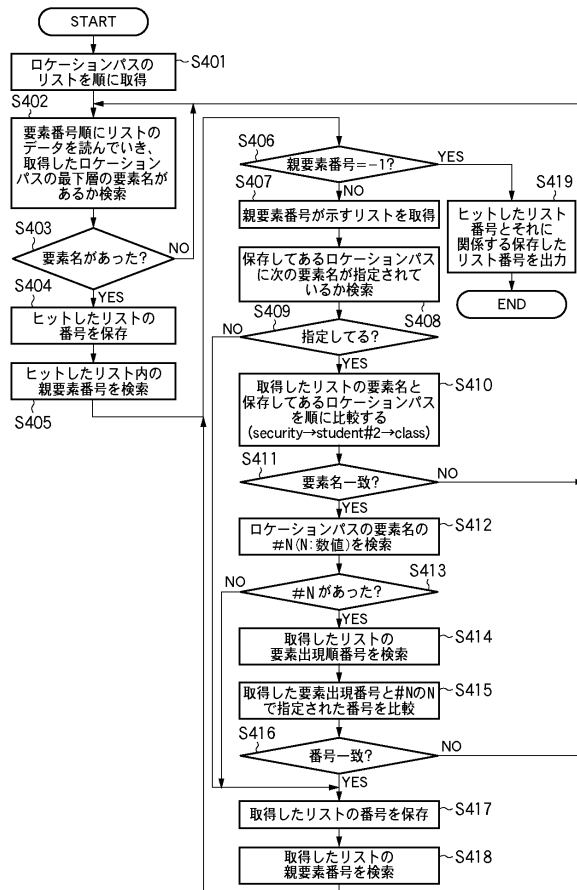
(c)

1	teacher, security, address
2	class, student#2, security, tel
3	name

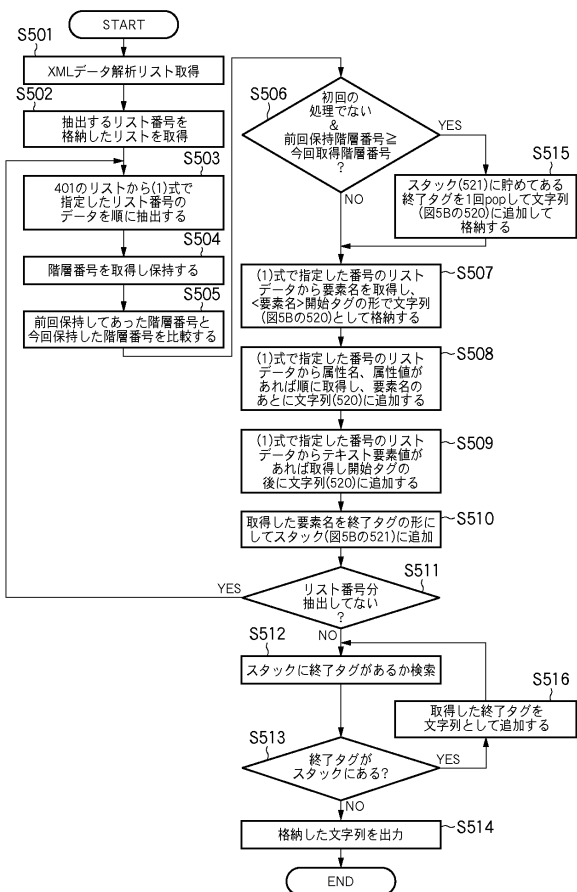
【図4A】

要素番号	階層番号	要素番号	要素名	name_space名	テキスト要素値	親要素番号	1つ前兄弟番号	要素出現番号	属性数	属性名	属性値
0	0	0	list	null		-1	0	1	0		
1	1	1	class	null		0	0	1	1	number	1
2	2	2	teacher	null		1	0	1	0		
3	3	3	name	null	Aさん	2	0	1	0		
4	3	3	security	null		2	3	1	0		
5	4	4	address	null	東京都	4	0	1	0		
6	4	4	tel	null	03-1111-1111	4	5	1	0		
7	2	2	student	null		1	2	1	0		
8	3	3	name	null	Bさん	7	0	1	0		
9	3	3	security	null		7	8	1	0		
10	4	4	address	null	神奈川県	9	0	1	0		
11	4	4	tel	null	045-2222-2222	9	10	1	0		
12	2	2	student	null		1	7	2	0		
13	3	3	name	null	Cさん	12	0	1	0		
14	3	3	security	null		12	13	1	0		
15	4	4	address	null	神奈川県	14	0	1	0		
16	4	4	tel	null	045-3333-3333	14	15	1	0		

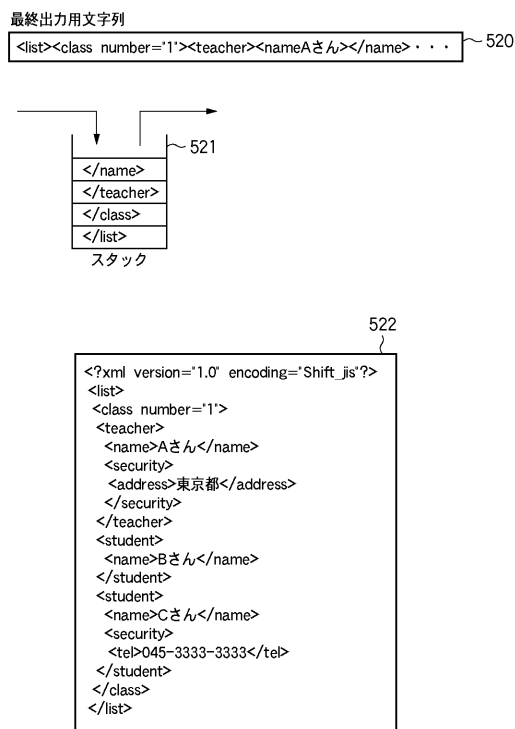
【図4B】



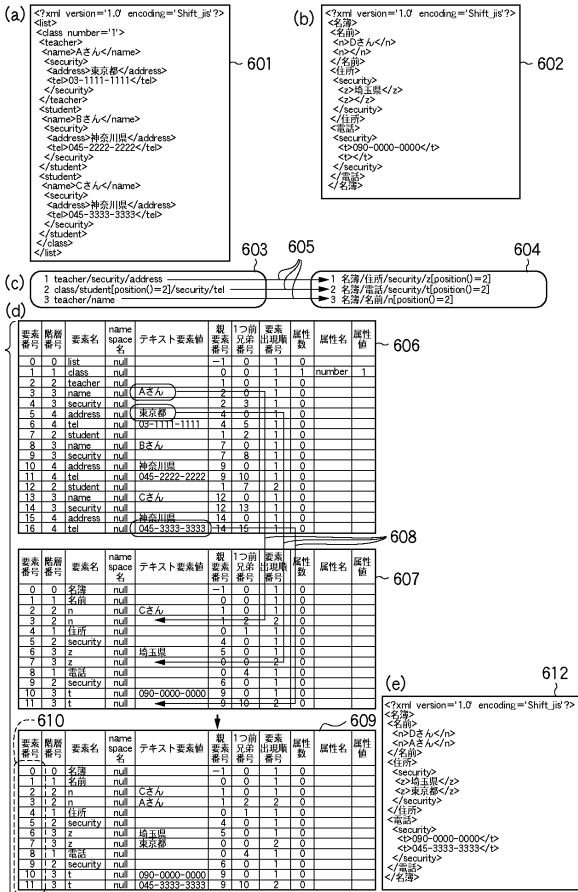
【図5A】



【図5B】



【図6】



フロントページの続き

(72)発明者 岩崎 晋吾
東京都大田区下丸子3丁目30番2号 キヤノン株式会社内

審査官 岩田 淳

(56)参考文献 特開2002-183116(JP,A)
吉田 晃伸, Blogサイト構築で理解するXMLデータベース, 日経Linux, 日本, 日経BP社, 2004年11月8日, 第6巻, 第11号, P.133-138
薬師寺 国安, 外1名, 動いてなんぼのXML, OPEN DESIGN, 日本, CQ出版株式会社, 2001年12月1日, 第8巻, 第12号, P.144-159

(58)調査した分野(Int.Cl., DB名)
G06F 17/21
JSTPlus(JDreamII)