

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4039484号
(P4039484)

(45) 発行日 平成20年1月30日(2008.1.30)

(24) 登録日 平成19年11月16日(2007.11.16)

(51) Int.Cl. F I
G 0 6 F 12/00 (2006.01) G O 6 F 12/00 5 4 6 Z
 G O 6 F 12/00 5 4 7 Z

請求項の数 24 (全 42 頁)

(21) 出願番号	特願2002-267625 (P2002-267625)	(73) 特許権者	390009531
(22) 出願日	平成14年9月13日 (2002.9.13)		インターナショナル・ビジネス・マシー ズ・コーポレーション
(65) 公開番号	特開2003-323332 (P2003-323332A)		I N T E R N A T I O N A L B U S I N E S S M A S C H I N E S C O R P O R A T I O N
(43) 公開日	平成15年11月14日 (2003.11.14)		アメリカ合衆国10504 ニューヨーク 州 アーモンク ニュー オーチャード ロード
審査請求日	平成15年1月8日 (2003.1.8)		
(31) 優先権主張番号	特願2002-53265 (P2002-53265)	(74) 代理人	100086243
(32) 優先日	平成14年2月28日 (2002.2.28)		弁理士 坂口 博
(33) 優先権主張国	日本国(JP)	(74) 代理人	100091568
前置審査			弁理士 市位 嘉宏
		(74) 代理人	100108501
			弁理士 上野 剛史

最終頁に続く

(54) 【発明の名称】 X P a t h評価方法、これを用いたXML文書処理システム及びプログラム

(57) 【特許請求の範囲】

【請求項1】

コンピュータを用いて、所定のデータファイルを対象として当該データファイルが X P a t hで指定される特定部分を持つか否かを評価する X P a t h評価方法において、

メモリから評価を行うべき複数の X P a t h式を読み出し、当該複数の X P a t h式から、当該複数の X P a t h式の共通部分を共通のノードとし、当該複数の X P a t h式の他の部分を個別のノードとする木構造のデータ構造を生成する工程と、

メモリから処理対象のデータファイルを入力し、前記データ構造のノードごとに、当該ノードに対応する前記 X P a t h式の部分で指定される前記特定部分を当該処理対象のデータファイルが持つか否かを評価し、評価結果として得られた当該特定部分を結合する工程と、

10

前記複数の X P a t h式の各 X P a t h式に関して、個々の当該 X P a t h式を構成する部分に対応する前記ノードごとの評価結果に基づいて各 X P a t h式の評価結果を得、得られた評価結果をメモリに格納する工程とを含むことを特徴とする X P a t h評価方法。

【請求項2】

前記データ構造を生成する工程は、

個々の前記 X P a t h式をロケーションステップごとに分解する工程と、

分解して得られた各ロケーションステップにノードを対応させ、複数の X P a t h式に共通するロケーションステップを1つのノードに対応させたデータ構造を生成する工程と

20

を含むことを特徴とする請求項 1 に記載の X P a t h 評価方法。

【請求項 3】

前記データ構造を生成する工程では、

複数の前記ロケーションステップが共通の特定の形で記述されている場合に、当該ロケーションステップにおける述部の式の評価結果をキーとして、X P a t h における残りのロケーションステップを検索するためのハッシュテーブルまたは二分探索木を生成し、

前記ノードごとの評価結果を結合して各 X P a t h 式の評価結果を得る工程では、

前記ハッシュテーブルまたは前記二分探索木を用いた検索により得られた X P a t h の前記残りのロケーションステップに関して、当該残りのロケーションステップで指定される前記特定部分を前記処理対象のデータファイルが持つか否かを評価し、その評価結果と前記ノードごとの評価結果とに基づいて各 X P a t h 式の評価結果を得ることを特徴とする請求項 2 に記載の X P a t h 評価方法。

10

【請求項 4】

前記データ構造を生成する工程は、

前記複数の X P a t h 式が演算式または関数を含む場合に、当該演算式または関数を部分式に分解する工程と、

分解して得られた前記部分式のうち、共通する部分式を 1 つにまとめて有向非循環グラフを生成する工程とを含み、

前記ノードごとの評価結果に基づいて各 X P a t h 式の評価結果を得る工程では、

前記有向非循環グラフを用い、複数の X P a t h 式に共通する部分式については評価結果を当該複数の X P a t h 式で共用して各 X P a t h 式を評価し、その評価結果と前記ノードごとの評価結果とに基づいて各 X P a t h 式の評価結果を得ることを特徴とする請求項 1 に記載の X P a t h 評価方法。

20

【請求項 5】

コンピュータを用いて、所定のデータファイルを対象として当該データファイルが X P a t h で指定される特定部分を持つか否かを評価する X P a t h 評価方法において、

メモリから評価を行うべき複数の X P a t h 式を含む X P a t h セットを読み出し、当該複数の X P a t h セットに共通して含まれる X P a t h 式と所定の X P a t h セットに固有の X P a t h 式とを分類し、当該 X P a t h 式の冗長な要素を省略したデータ構造を生成する工程と、

30

前記 X P a t h セットに固有の X P a t h 式について評価する工程と、

処理対象のデータファイルが前記固有の X P a t h 式で指定される特定部分を持つと評価される場合に、当該 X P a t h セット中の他の X P a t h セットと共通する X P a t h 式について評価する工程と

を含むことを特徴とする X P a t h 評価方法。

【請求項 6】

前記データ構造を生成する工程では、前記 X P a t h セットに含まれる連結 X P a t h 式に関して、当該連結 X P a t h 式を構成する各 X P a t h 式間の依存関係に基づき、評価が不要な X P a t h 式を評価対象から除くことにより、当該連結 X P a t h 式を簡略化することを特徴とする請求項 5 に記載の X P a t h 評価方法。

40

【請求項 7】

前記データ構造を生成する工程では、複数の X P a t h 式間の依存関係に基づき、評価結果が他の X P a t h 式に依存する X P a t h 式を評価対象から除くことを特徴とする請求項 5 に記載の X P a t h 評価方法。

【請求項 8】

前記データ構造を生成する工程では、複数の X P a t h 式の間で共通する部分と各 X P a t h 式に固有の部分とを分割し、各部分を個別の X P a t h 式として評価対象とすることを特徴とする請求項 5 に記載の X P a t h 評価方法。

【請求項 9】

複数の X P a t h 式を格納した X P a t h 格納部と、

50

前記 X P a t h 格納部に格納されている前記複数の X P a t h 式を読み出し、当該複数の X P a t h 式から、当該複数の X P a t h 式の共通部分を共通のノードとし、当該複数の X P a t h 式の他の部分を個別のノードとする木構造のデータ構造を生成するデータ構造生成部と、

処理対象のデータファイルを入力し、前記データ構造のノードごとに、当該ノードに対応する前記 X P a t h 式の部分で指定される特定部分を当該処理対象のデータファイルが持つか否かを評価し、当該ノードごとの評価結果として得られた当該特定部分を結合し、前記複数の X P a t h 式の各 X P a t h 式に関して、個々の当該 X P a t h 式を構成する部分に対応する前記ノードごとの評価結果に基づいて各 X P a t h 式の評価結果を得る評価実行部と

10

を備えることを特徴とする文書処理システム。

【請求項 10】

前記データ構造生成部は、

個々の前記 X P a t h 式をロケーションステップごとに分解するステップ分解手段と、
分解して得られた各ロケーションステップにノードを対応させ、複数の X P a t h 式に共通するロケーションステップを1つのノードに対応させたデータ構造を生成する木生成手段と

を備えることを特徴とする請求項 9 に記載の文書処理システム。

【請求項 11】

前記データ構造生成部は、前記 X P a t h 格納部に格納されている X P a t h 式が変更された場合に、新たに追加された X P a t h 式のうちで前記データ構造に対応するノードがない部分について新たなノードを生成して当該データ構造に加え、削除された X P a t h 式の部分に対応する前記データ構造のノードのうちで当該 X P a t h 式に固有の部分に対応するノードのみを除去することにより、当該データ構造を更新することを特徴とする請求項 9 に記載の文書処理システム。

20

【請求項 12】

前記データ構造生成部は、他の X P a t h 式と重複する X P a t h 式または評価結果が他の X P a t h 式に依存する X P a t h 式を評価対象から除くことを特徴とする請求項 9 に記載の文書処理システム。

【請求項 13】

前記データ構造生成部にて生成されたデータ構造を格納し保存するデータ構造格納部をさらに備え、

前記評価実行部は、前記データ構造格納部に格納されている前記データ構造を用いて前記複数の X P a t h 式の評価を行うことを特徴とする請求項 9 に記載の文書処理システム。

30

【請求項 14】

前記評価実行部は、前記データ構造におけるノードごとに対応する部分の評価を行い、この部分的な評価結果を結合して個々の X P a t h 式全体の評価結果を得ると共に、複数の X P a t h 式間で共通するノードに対応する部分の評価結果を当該複数の X P a t h 式の評価において共用することを特徴とする請求項 9 に記載の文書処理システム。

40

【請求項 15】

複数の X P a t h 式にて構成される X P a t h セットを格納した X P a t h 格納部と、
前記 X P a t h 格納部から評価を行うべき複数の X P a t h 式を含む X P a t h セットを読み出し、当該複数の X P a t h セットに共通して含まれる X P a t h 式と所定の X P a t h セットに固有の X P a t h 式とを分類し、当該 X P a t h 式の冗長な要素を省略したデータ構造を生成するデータ構造生成部と、

処理対象のデータファイルを入力し、前記 X P a t h セットに固有の X P a t h 式について当該データファイルが当該 X P a t h で指定される特定部分を持つか否かを評価し、当該処理対象のデータファイルが当該固有の X P a t h 式で指定される特定部分を持つと評価される場合に、当該 X P a t h セット中の他の X P a t h セットと共通する X P a t

50

h 式について当該データファイルが当該 X P a t h で指定される特定部分を持つか否かを評価し、得られた各 X P a t h 式の評価結果に基づいて前記データファイルに対する X P a t h セットの評価結果を得る評価実行部とを備えることを特徴とする文書処理システム。

【請求項 16】

コンピュータを制御して、所定のデータファイルを対象として当該データファイルが X P a t h で指定される特定部分を持つか否かを評価するプログラムであって、

評価を行うべき複数の X P a t h 式から、当該複数の X P a t h 式の共通部分を共通のノードとし、当該複数の X P a t h 式の他の部分を個別のノードとする木構造のデータ構造を生成する処理と、

前記データ構造のノードごとに、当該ノードに対応する前記 X P a t h 式の部分で指定される前記特定部分を処理対象のデータファイルが持つか否かを評価し、評価結果を結合する処理と、

前記複数の X P a t h 式の各 X P a t h 式に関して、個々の当該 X P a t h 式を構成する部分に対応する前記ノードごとの評価結果に基づいて各 X P a t h 式の評価結果を得る処理とを

前記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 17】

前記プログラムによる前記データ構造を生成する処理は、

個々の前記 X P a t h 式をロケーションステップごとに分解する処理と、

分解して得られた各ロケーションステップにノードを対応させ、複数の X P a t h 式に共通するロケーションステップを 1 つのノードに対応させたデータ構造を生成する処理とを含むことを特徴とする請求項 16 に記載のプログラム。

【請求項 18】

前記プログラムによる前記データ構造を生成する処理では、

複数の前記ロケーションステップが共通の特定の形で記述されている場合に、当該ロケーションステップにおける述部の式の評価結果をキーとして、X P a t h における残りのロケーションステップを検索するためのハッシュテーブルまたは二分探索木を生成し、

前記ノードごとの評価結果に基づいて各 X P a t h 式の評価結果を得る処理では、

前記ハッシュテーブルまたは前記二分探索木を用いた検索により得られた X P a t h の前記残りのロケーションステップに関して、当該残りのロケーションステップで指定される前記特定部分を前記処理対象のデータファイルが持つか否かを評価し、その評価結果と前記ノードごとの評価結果とに基づいて各 X P a t h 式の評価結果を得ることを特徴とする請求項 17 に記載のプログラム。

【請求項 19】

前記プログラムによる前記データ構造を生成する処理は、

前記複数の X P a t h 式が演算式または関数を含む場合に、当該演算式または関数を部分式に分解する処理と、

分解して得られた前記部分式のうち、共通する部分式を 1 つにまとめて有向非循環グラフを生成する処理とを含み、

前記ノードごとの評価結果に基づいて各 X P a t h 式の評価結果を得る処理では、

前記有向非循環グラフを用い、複数の X P a t h 式に共通する部分式については評価結果を当該複数の X P a t h 式で共用して各 X P a t h 式を評価し、その評価結果と前記ノードごとの評価結果とに基づいて各 X P a t h 式の評価結果を得ることを特徴とする請求項 16 に記載のプログラム。

【請求項 20】

前記プログラムは、評価を行うべき複数の X P a t h 式が変更された場合に、新たに追加された X P a t h 式のうちで前記データ構造に対応するノードがない部分について新たなノードを生成して当該データ構造に加え、削除された X P a t h 式の部分に対応する前記データ構造のノードのうちで当該 X P a t h 式に固有の部分に対応するノードのみを除

10

20

30

40

50

去することにより、当該データ構造を更新する処理を、前記コンピュータにさらに実行させることを特徴とする請求項16に記載のプログラム。

【請求項21】

コンピュータを制御して、所定のデータファイルを対象として当該データファイルがXPathで指定される特定部分を持つか否かを評価するプログラムであって、

複数のXPath式を含んで構成された評価を行うべきXPathセットから当該複数のXPathセットに共通して含まれるXPath式と所定のXPathセットに固有のXPath式とを分類し、当該XPath式の冗長な要素を省略したデータ構造を生成する処理と、

前記XPathセットに固有のXPath式について評価する処理と、

処理対象のデータファイルが前記固有のXPath式で指定される特定部分を持つと評価される場合に、当該XPathセット中の他のXPathセットと共通するXPath式について評価する処理と

を前記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22】

前記プログラムによる前記データ構造を生成する処理では、前記XPathセットに含まれる連結XPath式に関して、当該連結XPath式を構成する各XPath式間の依存関係に基づき、評価が不要なXPath式を評価対象から除くことにより、当該連結XPath式を簡略化することを特徴とする請求項21に記載のプログラム。

【請求項23】

前記プログラムによる前記データ構造を生成する処理では、複数のXPath式間の依存関係に基づき、評価結果が他のXPath式に依存するXPath式を評価対象から除くことを特徴とする請求項21に記載のプログラム。

【請求項24】

前記プログラムによる前記データ構造を生成する処理では、複数のXPath式の間で共通する部分と各XPath式に固有の部分とを分割し、各部分を個別のXPath式として評価対象とすることを特徴とする請求項21に記載のプログラム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、XPath(XML Path Language)を用いてXML文書やHTML文書における特定部分を指し示す場合に、XPath式を効率よく評価する技術に関する。

【0002】

【従来の技術】

XML文書の特定部分を指し示すためのパス言語として、W3C(World Wide Web Consortium)から勧告として公開されているXPathがある。XPathは、例えば、XPathPointer、XSLT、XQueryなどの構成要素として使用され、また所定のアプリケーション・プログラムの中でXML文書のDOM(Document Object Model)ツリーをアクセスするのにも使用される。

【0003】

ところで、XML文書を用いた実際の情報処理においては、1つのXML文書に対して複数のXPath式を評価することがよく行われる。例えば、XSLTスタイルシートでは、テンプレート・ルールごとにXPath式がパターンとして指定される。したがって、複雑なXSLTスタイルシートには多くのXPath式が含まれ、処理対象であるXML文書に対して、かかる複数のXPath式を評価することが必要となる。

【0004】

また、インターネット上でのデータ交換等に用いられるウェブページに対してアノテーションを付加することにより、所定のウェブページを様々な形に再利用することが可能であり、新たなアプリケーションを開発可能であることが広く知られている。ウェブページの記述に用いられるHTML文書は、XML文書と同様にXPath式にて所定の部分を指

10

20

30

40

50

し示すことができるため、アノテーションをウェブページ中のエレメントに対応付けるためにも XPath が用いられる。ところで、ウェブページに対するアノテーションの付加を効率的に行うために、特定のアノテーションを多数のウェブページに対して適用し、使い回すことが考えられる。この場合、所定のアノテーションが所定のウェブページに対して適用可能かどうかを判断するために、アノテーション内の複数の XPath 式が対象のウェブページにおける所定のエレメントを正しく指しているかどうかを評価することが必要となる。

【 0 0 0 5 】

XPath による XML 文書の特定部分を指し示す機能は、対象となる XML 文書が XPath にて指定される特定の部分を持つかどうかをチェックする条件とみなすことができる。例えば、米国 BEA 社のサーバシステムである WebLogic Collaborate (<http://www.bea.com/index.html>) は、XML で表現されるメッセージのルーティングとフィルタリングの条件を記述するために XPath を用いる。このような用途では、1 つの XML 文書が到着するたび多数の XPath 式を評価することとなる。

10

【 0 0 0 6 】

このように、1 つの XML 文書に対して複数の XPath 式を評価する必要がある場合、XPath 式の評価を効率的に行うことが要求される。従来、この種の技術としては、XML で記述された文書に対する各利用者の購読条件を XPath で記述しておき、XML 文書が到着するたびに XPath 式とのマッチングを検査して、検査を通過する文書を、当該 XPath を購読条件とする利用者に配信するものがある (例えば、非特許文献 1 参

20

【 0 0 0 7 】**【非特許文献 1】**

Altinel M., Franklin, M., "Efficient Filtering of XML Documents for Selective Dissemination of Information", International Conference on Very Large Data Bases, 2000.

【 0 0 0 8 】**【発明が解決しようとする課題】**

上述したように、XML 文書や HTML 文書に対する処理において、1 つのデータファイル (文書) に対して複数の XPath 式を評価する場合、効率的に XPath 式の評価を行うことが要求される。

30

しかし、上記文献に開示された方法のように、XPath 式 1 つあたりの評価の実行時間を改善する方法は知られているが、XPath 式の数の増加に比例して、全体の評価に要する時間が直線的に増加してしまうため、全体的な実行時間の短縮には限界があった。これは、従来の XPath の評価方法が、複数の XPath 式を評価する場合に、各 XPath 式を互いに独立なものとして扱っていることに起因する。

【 0 0 0 9 】

ところが、1 つのデータファイルに対して評価されることが想定されている XPath 式が複数ある場合、それらの XPath 式は対象となるデータファイルの構造や要素の値のバリエーションなどから制約を受けるので、複数の XPath 式の中に類似なものが含まれる。したがって、類似する XPath 式から共通の部分を取り出して評価し、類似する XPath 式の間で共通部分の評価結果を共用することにより、複数の XPath 式を個別に評価するよりも高速に実行できると考えられる。

40

【 0 0 1 0 】

また、評価されるべき複数の XPath 式の間に関係がある場合、この依存関係を利用して、XPath 式の評価に要する処理を単純化することができると考えられる。ここで、XPath 式の間に関係とは、例えば、所定の XPath 式が指し示すウェブページ中に 2 つのテーブルコンテンツ (テーブル [1]、テーブル [2]) がある場合、テーブ

50

ル[1]が存在しなければテーブル[2]は存在し得ないというような関係、すなわち、複数のX P a t hの1つまたはその一部を評価することにより、他のX P a t hの評価結果が自明となる関係である。

【0011】

そこで、本発明は、XML文書やHTML文書等のデータファイルに対して複数のX P a t h式を評価する場合に、かかるX P a t h式間の類似性や依存関係を利用して、効率よく複数のX P a t h式の評価結果を求める方法を提供することを目的とする。

【0012】**【課題を解決するための手段】**

上記の目的を達成する本発明は、コンピュータを用いて所定のデータファイル(XML文書やHTML文書)を対象として複数のX P a t hを評価する、次のようなX P a t h評価方法として実現される。すなわち、このX P a t h評価方法は、評価を行うべき複数のX P a t h式における共通部分を木構造における共通のノードとして持つデータ構造を生成するステップと、処理対象のデータファイルに対し、生成されたデータ構造のノードごとに評価を行うステップと、このノードごとの評価結果を結合して個々のX P a t h式の評価結果を得るステップとを含むことを特徴とする。

10

【0013】

ここで、このデータ構造を生成するステップは、詳しくは、個々のX P a t h式をステップごとに分解するステップと、分解して得られた各ステップにノードを対応させ、複数のX P a t h式に共通するステップを1つのノードに対応させたデータ構造を生成するステップとを含む。

20

さらに好ましくは、このデータ構造を生成するステップは、X P a t h式の複数のステップが共通の特定の形で記述されている場合、具体的には、「[(式)(比較演算子)(定数)]」の形になる述部が存在する場合に、特殊なデータ構造を用いることができる。例えば、比較演算子が「=」や「!=」の場合、ステップにおける述部の式の評価結果をキーとしてX P a t hにおける残りのステップを検索するためのハッシュテーブルを生成し、生成されたハッシュテーブルを用いた検索により得られたX P a t hの残りのステップを評価することができる。また、比較演算子が「>」や「>=」、「<」、「<=」の場合、ステップにおける述部の式の評価結果をキーとしてX P a t hにおける残りのステップを検索するための二分探索木を生成し、生成された二分探索木を用いた検索により得られたX P a t hの残りのステップを評価することができる。

30

【0014】

また、このデータ構造を生成するステップは、複数のX P a t h式が演算式または関数を含む場合に、この演算式または関数を部分式に分解するステップと、分解して得られた部分式のうち、共通する部分式を1つにまとめて有向非循環グラフを生成するステップとを含み、X P a t h式の評価を行うステップは、生成された有向非循環グラフを用い、複数のX P a t h式に共通する部分式については評価結果を複数のX P a t h式で共用して、各X P a t h式を評価するステップを含む。

このとき、個々の部分式については、上述したように共通部分を木構造における共通ノードとして持つデータ構造を生成することができ、これを用いた評価を行うことができる。

40

【0015】

本発明による他のX P a t h評価方法は、評価を行うべき複数のX P a t h式の中から共通部分を抽出するステップと、処理対象であるデータファイルに対して、X P a t h式における共通部分と他の部分とを分けて評価し、評価結果を結合して個々のX P a t h式に対する評価結果を得るステップとを含み、このX P a t h式の評価を得るステップは、所定のX P a t h式を評価する際、共通部分に関して、他のX P a t hを評価した際に得られた評価結果がある場合に、この評価結果を当該所定のX P a t h式における共通部分の評価結果として用いるステップを含むことを特徴とする。

【0016】

本発明による、さらに他のX P a t h評価方法は、コンピュータを用いて、所定のデータ

50

ファイルを対象として複数のX P a t hを評価するX P a t h評価方法において、メモリから評価を行うべき複数のX P a t h式を含むX P a t hセットを読み出し、このX P a t hセットに含まれるX P a t h式の冗長な要素を省略したデータ構造を生成するステップと、処理対象のデータファイルを入力し、当該データファイルに対し、このデータ構造を用いてX P a t hセットの評価を行い、得られた評価結果をメモリに格納するステップとを含むことを特徴とする。

【0017】

より詳しくは、このデータ構造を生成するステップでは、複数のX P a t hセットに共通して含まれるX P a t h式と所定のX P a t hセットに固有のX P a t h式とを分類し、X P a t hセットの評価を行う際の評価対象となるX P a t h式を絞り込む。また、X P a t hセットに含まれる連結X P a t h式に関して、この連結X P a t h式を構成する各X P a t h式の間
10
の依存関係に基づき、評価が不要なX P a t h式を評価対象から除くことにより、当該連結X P a t h式を簡略化する。さらに、複数のX P a t h式の間
の依存関係に基づき、評価が不要なX P a t h式を評価対象から除く。さらにまた、複数のX P a t h式の間で共通する部分と各X P a t h式に固有の部分とを分割し、各部分を個別のX P a t h式として評価対象とする。

また、X P a t hセットの評価は、まず、各X P a t hセットに固有のX P a t h式を評価し、固有のX P a t h式が全て処理対象のデータファイルにマッチする場合に、このX P a t hセット中の他のX P a t hセットと共通するX P a t h式を評価する。

【0018】

また、上記の目的を達成する本発明は、次のように構成された文書処理システムとしても実現される。すなわち、この文書処理システムは、複数のX P a t h式を格納したX P a t h格納部と、このX P a t h格納部に格納されている複数のX P a t h式の冗長な要素を省略したデータ構造を生成するデータ構造生成部と、このデータ構造生成部にて生成されたデータ構造を用いて、処理対象であるデータファイルに対し、複数のX P a t h式
20
の評価を行う評価実行部とを備えることを特徴とする。

さらに、この文書処理システムは、データ構造生成部にて生成されたデータ構造を格納し保存するデータ構造格納部を備える構成とすることができる。この場合、評価実行部は、このデータ構造格納部に格納されているデータ構造を繰り返し用いて複数のデータファイル
30
に対してX P a t h式の評価を行うことができる。

【0019】

ここで、このデータ構造生成部は、詳しくは、個々のX P a t h式をステップごとに分解するステップ分解手段と、分解して得られた各ステップにノードを対応させ、複数のX P a t h式に共通するステップを1つのノードに対応させたデータ構造を生成する木生成手段とを備える。

また、このデータ構造生成部は、X P a t h格納部に格納されているX P a t h式が変更された場合に、新たに追加されたX P a t h式のうちでデータ構造に対応するノードがない部分について新たなノードを生成して既存のデータ構造に加える。一方、削除されたX P a t h式の部分に対応するノードのうちでX P a t h式に固有の部分に対応するノードのみを除去することにより、データ構造を更新する。
40

さらに、このデータ構造生成部は、他のX P a t h式と重複するX P a t h式または評価結果が他のX P a t h式に依存するX P a t h式を評価対象から除く。

さらにまた、この評価実行部は、データ構造におけるノードごとに対応する部分の評価を行い、この部分的な評価結果を結合して個々のX P a t h式全体の評価結果を得ると共に、複数のX P a t h式間で共通するノードに対応する部分の評価結果を、複数のX P a t h式の評価において共用する。

【0020】

また、上記の目的を達成する本発明は、上述したX P a t h評価方法の各ステップに対応する処理をコンピュータに実行させ、またはコンピュータを制御して上述した文書処理システムとしての各機能を実現するプログラムとしても実現することができる。このプログ
50

ラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記録媒体に格納して配布したり、ネットワークを介して配信したりすることにより提供することができる。

【0021】

【発明の実施の形態】

以下、添付図面に示す実施の形態に基づいて、この発明を詳細に説明する。

なお、本実施の形態では、XML文書やウェブページの記述に用いられるHTML文書など、XPath式で所定の部分を指し示すことができるデータファイルを対象とする。したがって、以下の説明においてXML文書を対象として記載している部分をHTML文書に対して適用することも、反対にHTML文書を対象として記載している部分をXML文書に対して適用することも可能である。

10

〔第1の実施の形態〕

図1は、第1の実施の形態によるXML文書処理システムを実現するのに好適なコンピュータ装置のハードウェア構成の例を模式的に示した図である。

図1に示すコンピュータ装置は、演算手段であるCPU (Central Processing Unit: 中央処理装置) 101と、M/B (マザーボード) チップセット102及びCPUバスを介してCPU 101に接続されたメインメモリ103と、同じくM/Bチップセット102及びAGP (Accelerated Graphics Port) を介してCPU 101に接続されたビデオカード104と、ビデオカード104にて生成されたグラフィックデータを表示する表示装置110と、PCI (Peripheral Component Interconnect) バスを介してM/Bチップセット102に接続されたハードディスク105及びネットワークインターフェイス106と、さらにこのPCIバスからブリッジ回路107及びISA (Industry Standard Architecture) バスなどの低速なバスを介してM/Bチップセット102に接続されたフロッピーディスクドライブ108及びキーボード/マウス109とを備える。表示装置110としては、例えば液晶ディスプレイ (LCD) やCRTディスプレイを用いることができる。

20

なお、図1は本実施の形態を実現するコンピュータ装置のハードウェア構成を例示するに過ぎず、本実施の形態を適用可能であれば、他の種々の構成を取ることができる。例えば、ビデオカード104を設ける代わりに、ビデオメモリのみを搭載し、CPU 101にてイメージデータを処理する構成としても良いし、音声による入出力を行うためのサウンド機構を設けたり、ATA (AT Attachment) などのインターフェイスを介してCD-ROM (Compact Disc Read Only Memory) やDVD-ROM (Digital Versatile Disc Read Only Memory) のドライブを設けたりしても良い。

30

【0022】

図2は、第1の実施の形態によるXML文書処理システムの機能を説明するブロック図である。

図2を参照すると、本実施の形態は、処理対象であるXML文書の特定部分を指し示す複数のXPath式を格納したXPath格納部10と、XPath格納部10に格納されているXPath式から本実施の形態による効率的な評価を行うためのデータ構造を生成するデータ構造生成部20と、データ構造生成部20にて生成されたデータ構造を格納するデータ構造格納部30と、データ構造格納部30に格納されているデータ構造を用いて処理対象であるXML文書に対してXPath式の評価を行う評価実行部40とを備え、また処理対象であるXML文書を格納するXML文書格納部50と、評価実行部40による評価を経たXML文書を対象として所定の情報処理を実行するXML文書処理部60とを備える。

40

【0023】

図2に示したシステム構成において、データ構造生成部20、評価実行部40及びXML文書処理部60は、図1に示したプログラム制御されたCPU 101にて実現される仮想的なソフトウェアブロックである。CPU 101を制御してこれらの機能を実現するプログラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記録媒体に格納して配布したり、ネットワークを介して配信したりすることにより提供され、メインメモリ103

50

に読み込まれる。また、X P a t h格納部 1 0、データ構造格納部 3 0及びX M L文書格納部 5 0は、メインメモリ 1 0 3にて実現される。なお、メインメモリ 1 0 3に保持されるデータやプログラムは、必要に応じてハードディスク 1 0 5などの記憶装置に退避させることができる。

【 0 0 2 4 】

図 2 に示したシステム構成において、X P a t h格納部 1 0には、X M L文書に対して評価を行うべき、予め用意された複数のX P a t h式が格納されている。データ構造格納部 3 0には、データ構造生成部 2 0にて生成され、評価実行部 4 0にて用いられるデータ構造が格納される。X M L文書格納部 5 0には、図 1 に示したネットワークインターフェイス 1 0 6や、キーボード/マウス 1 0 9等の入力デバイス、処理対象であるX M L文書を格納した記録媒体から当該X M L文書を読み出すドライブ装置などを介して入力されたX M L文書が格納されている。

10

【 0 0 2 5 】

データ構造生成部 2 0は、X P a t h格納部 1 0に格納されているX P a t h式から、評価実行部 4 0による評価に用いるデータ構造を生成する。また、X P a t h格納部 1 0に格納されているX P a t hが変更（追加または削除）されたならば、この変更に応じてデータ構造を更新する。このデータ構造を用いることにより、X M L文書に対してX P a t h格納部 1 0に格納されている複数のX P a t h式を個別に評価する場合に比べて効率の良い評価を行うことが可能となる。このデータ構造についての詳細、このデータ構造を用いた評価方法及びこのデータ構造の生成方法については後述する。

20

【 0 0 2 6 】

評価実行部 4 0は、X M L文書格納部 5 0に格納されているX M L文書に対して、データ構造生成部 2 0にて生成されたデータ構造を用いてX P a t h式の評価（マッチング）を行う。これにより、当該入力文書がX P a t h格納部 1 0に格納されているX P a t h式にて指定される特定部分を持つかどうかを、効率よく評価することができる。評価実行部 4 0による評価結果はメインメモリ 1 0 3の所定の領域に格納される。

X M L文書処理部 6 0は、メインメモリ 1 0 3の所定領域から評価実行部 4 0の評価結果を読み出し、当該評価結果に基づいて、X M L文書格納部 5 0に格納されているX M L文書がX P a t h格納部 1 0に格納されているX P a t h式にて指定される特定部分を持つならば、X M L文書格納部 5 0から当該X M L文書を処理対象として読み出し、目的とする所定の情報処理、すなわちアプリケーションとして用意された情報処理を実行する。

30

【 0 0 2 7 】

次に、本実施の形態におけるX P a t h式の評価を効率化する方法を詳細に説明する。本実施の形態では、複数のX P a t h式における共通部分を共用するためのデータ構造を生成し、このデータ構造を用いてX M L文書に対する評価を行うことにより、評価処理の高速化を実現する。これを実現する方策として、本実施の形態では、次の3つの手法を提示する。

方法 1 . 複数のロケーション・パスの間でのステップの評価の共有。

方法 2 . 複数の演算式や関数の間での部分式や引数の評価の共有。

方法 3 . ロケーション・パスの述部の高速化。

40

【 0 0 2 8 】

以下、各手法について説明する。

方法 1 . 複数のロケーション・パスの間でのステップの評価の共有

X P a t h格納部 1 0に、次の3つのX P a t h式P 1、P 2、P 3が格納されており、X M L文書に対してこれらのX P a t h式を評価する場合を例として説明する。

P 1 : /profile/demographics/age[text() < 20]

P 2 : /profile/interests/sport[text() = 'Soccer']

P 3 : /profile/demographics/age[text() >= 40 and text() < 50]

【 0 0 2 9 】

まず、データ構造生成部 2 0により、上記X P a t h式P 1、P 2、P 3をまとめたデー

50

タ構造が生成される。

データ構造生成部 20 は、このデータ構造を生成するため、上記の X P a t h 式をステップごとに分解する。ステップが軸指定子 (axis) とノードテストの組と述部からなるときには、さらにこれらも分解する。すなわち、X P a t h 式 P 1 は次の P 1 ' に、P 2 は P 2 ' に、P 3 は P 3 ' に、それぞれ分解される。

P 1 ' : profile + demographics + age + [text() < 20]

P 2 ' : profile + interests + sport + [text() = 'Soccer']

P 3 ' : profile + demographics + age + [text() >= 40 and text() < 50]

P 1 '、P 2 '、P 3 ' を比較すると、これら 3 つの X P a t h 式は「profile」を共通ステップとして持つ。また、P 1 ' と P 3 ' とは、さらに続くステップの「demographics」を共通に持つ。そこで、共通のステップをまとめることによって、これらの X P a t h 式を木構造で表現することができる。

すなわち、かかる処理においてデータ構造生成部 20 は、X P a t h 式をステップごとに分解するステップ分解手段と、各ステップに対応するノードを持つ木構造 (データ構造) を生成する木生成手段として機能する。

図 3 は、データ構造生成部 20 にて生成されるデータ構造 (木構造) を示す図である。

【 0 0 3 0 】

次に、評価実行部 40 により、上記のように生成された木データ構造を用いて、XML 文書に対する X P a t h 式 P 1、P 2、P 3 の評価が行われる。

ここでは、図 4 に示す XML 文書に対し、X P a t h 式 P 1、P 2、P 3 を評価する例を考える。

(1) X P a t h 式 P 1 の評価

最初に図 4 に示した XML 文書 D 1 に対し、X P a t h 式 P 1 の先頭のステップ「profile」を評価する。評価結果は、次の 4 要素を持つノードセットになる。

N 1 :

{ <name>Alan</name> ,

<demographics><age>35</age></demographics> ,

<location><city>Osaka</city></location> ,

<interests><sport>Soccer</sport><music>Classical</music><book>History</book></interests> }

【 0 0 3 1 】

ノードセット N 1 の 4 つの要素それぞれに、X P a t h 式 P 1 の次のステップ「demographics」を適用すると、結果は

{ } , { <age>35</age> } , { } , { }

となる。この 4 つのノードセットを結合して、このステップのノードセットを作ると、次のようになる。

N 2 :

{ <age>35</age> }

このノードセット N 2 の要素に「age」を適用すると、次のようになる。

N 3 :

{ 35 }

さらに、このノードセット N 3 に「[text() < 20]」適用すると、結果は空のノードセットになり、これが X P a t h 式 P 1 の評価結果として得られる。

【 0 0 3 2 】

(2) X P a t h 式 P 2 の評価

次に、X P a t h 式 P 2 を評価する。このとき、XML 文書 D 1 に対する X P a t h 式 P 2 の先頭のステップ「profile」の評価結果は、すでにノードセット N 1 として得られている。したがって、この評価のための処理は実際には行わず、すでに得られている評価結果を使用する。ノードセット N 1 の各要素に「interests」を適用すると、結果はそれぞれ

10

20

30

40

50

{ }, { }, { }, { < sport > Soccer < / sport >, < music > Classical < / music >, < book > History < / book > }

となる。この4つのノードセットを結合して、このステップのノードセットを作ると、次のようになる。

N 4 :

{ < sport > Soccer < / sport >, < music > Classical < / music >, < book > History < / book > }

【 0 0 3 3 】

さらに、このノードセットN 4の3つの要素に対して、X P a t h式P 2の残りの表現である「 sport 」を適用し、その結果に「 [text () = ' Soccer '] 」を適用すると、結果は

{ Soccer }, { }, { }

となる。これを結合すると

{ Soccer }

となり、これがX P a t h式P 2の評価結果として得られる。

【 0 0 3 4 】

(3) X P a t h式P 3の評価

最後にX P a t h式P 3を評価する。このとき、X M L文書D 1に対するX P a t h式P 3のロケーション・パスのうち、「 / profile / demographics / age 」の部分の評価値は、すでにノードセットN 3として得られている。したがって、この部分に対する評価のための処理は実際には行わず、すでに得られている評価結果を使用する。ノードセットN 3の各要素に対して、X P a t h式P 3の最後の述部「 [text () >= 40 and text () < 50] 」を適用すると、結果は空のノードセットになり、これがX P a t h式P 3の評価結果として得られる。

【 0 0 3 5 】

以上のようにして、図3に示したデータ構造を用い、複数のX P a t h式に共通するステップについては1度だけ実際の計算を行い、評価結果を共有することによって、複数のX P a t h式の計算時間を短縮することができる。

上記の例では、任意の順番でX P a t h式を評価するように記述したが、図3のデータ構造から導かれる深さ優先または幅優先の順番で評価することにより、評価を行う際に必要なメモリの記憶領域を減らすことができる。例えば、深さ優先で評価した場合、所定のノードを保持するために必要な記憶領域は、当該ノードに対して全ての部分木の評価が終わったならば回収できる。一方、幅優先で評価した場合、上位ノードの評価結果を保持するために必要な記憶領域は、当該評価結果が下位ノードの評価に使用された後に回収できる。

【 0 0 3 6 】

図5は、幅優先で複数のX P a t h式を評価する場合の評価実行部40による処理手順を説明するフローチャートである。

図5を参照すると、評価実行部40は、まず、X M L文書全体を1つの要素とする集合を、現在のノードセットにする(ステップ501)。そして、全てのロケーション・パスの先頭ステップを、ステップ集合の要素にする(ステップ502)。次に、全てのステップ集合の要素に対し、次の処理を行う。すなわち、まず、現在のノードセットの各要素に対して、ステップをそれぞれ評価する(ステップ504)。そして、全ての評価結果を結合する。これを現在のノードセットとする(ステップ505)。

【 0 0 3 7 】

ステップ504、505の処理を、全てのステップ集合の要素に対して行ったならば、次に、全てのステップ集合の要素に対して、次のステップの集合を求める。そして、評価結果を結合し、現在のステップ集合と置き換える(ステップ503、506)。

この後、ステップ集合に要素が残っていればステップ503～506の処理を繰り返し、要素が空になったならば処理を終了する。

すなわち、かかる処理で評価実行部40は、各X P a t h式の評価において、データ構造のノードに対応するステップごとに評価を行う部分評価手段と、ステップごとの評価結果

10

20

30

40

50

を結合して X P a t h 式全体の評価結果を得る評価結果結合手段として機能する。そして、他の X P a t h 式の共通するステップの評価において、当該ステップに関して他の X P a t h 式を評価した際的评价結果が存在するならば、この評価結果を着目中の X P a t h 式の該当ステップの評価結果として用いることにより、当該ステップに対する評価処理を省略する。

【 0 0 3 8 】

この方法の基礎になる X P a t h の性質は、ロケーション・パスの途中までの表現を評価して得たノードセットに対して、そのノードセットの各要素をコンテキストノードとしてパスの残りを評価し、最後に結果を結合すると、パス全体を評価したのと同じものが得られることである。すなわち、上記の手順では、各 X P a t h 式のロケーション・パスのうちで共通部分を抽出し、当該共通部分と他の部分とを分けて評価した上で評価結果を結合する。この過程で、評価結果を共有できる部分（ロケーション・パスの共通部分）に対する評価結果が既に存在するならば、かかる部分に対する評価処理を省略して、その既存の評価結果を使用する。したがって、この性質を満たすようにパスが分解できる場合には、この方法が適用できることとなる。

上記の説明では、軸指定子 (axis) を省略してデフォルトの「child」が選択される場合を例として用いたが、軸指定子が「descendant::age」などのように「child」以外を明示的に書いてある場合も方法 1 で扱うことができる。また、ステップの区切りとして「/」だけでなく「//」の場合も扱うことができる。さらに、ノードテストとして要素名だけでなく、「.」「..」「*」「comment()」などを記述することもできる。例えば、次に示す

```
/*//interest/./demographics/decendant::age[text() < 20]
```

【 0 0 3 9 】

方法 2 . 複数の演算式や関数の間での部分式や引数の評価の共有

上記の方法 1 では、ロケーション・パスの評価を高速化する手法を説明した。

X P a t h においては、ロケーション・パス式が主要な式であるが、パスの述部の条件を記述したり、他のプログラムが使うためのデータを生成したりするために、四則演算式、比較演算式、論理演算式などの演算式や関数も使用される。そこで、以下に説明するロケーション・パス以外の式を高速化する手法を用いる。

【 0 0 4 0 】

X P a t h 格納部 1 0 に、次の 3 つの X P a t h 式 P 7、P 8、P 9 が格納されており、X M L 文書に対してこれらの X P a t h 式を評価する場合を例として説明する。

P 7 :

```
(/CPEXMLv1/person/partyActivities/food[favoriteFood='Hamburger'] or
(not(/CPEXMLv1/person/partyActivities/hobby[typeName='SPORT']) and
 /CPEXMLv1/person/personDemographics/birthDate[year<1960])) and
/CPEXMLv1/person/partyActivities/hobby/startDate[year=1986]
```

P 8 :

```
(/CPEXMLv1/person/partyDemographics[gender = 'MALE'] or
(not(/CPEXMLv1/person/partyActivities/hobby[typeName='SPORT']) and
 /CPEXMLv1/person/personDemographics/birthDate[year<1960])) and
/CPEXMLv1/person/partyActivities/hobby/startDate[year>1990]
```

P 9 :

```
/CPEXMLv1/person/partyActivities/food[favoriteFood = 'Hamburger'] and
/CPEXMLv1/person/partyActivities/hobby/startDate[year>1990]
```

【 0 0 4 1 】

データ構造生成部 2 0 は、まず、上記の X P a t h 式 P 7、P 8、P 9 をロケーション・パスからなる部分式に分解する。すなわち、X P a t h 式 P 7 は P 7 1 ~ P 7 4 に、P 8 は P 8 1 ~ P 8 4 に、P 9 は P 9 1、P 9 2 に、それぞれ分解される。

P 7 1 : /CPEXMLv1/person/partyActivities/food[favoriteFood = 'Hamburger']

P 7 2 : /CPEXMLv1/person/partyActivities/hobby[typeName='SPORT']
 P 7 3 : /CPEXMLv1/person/personDemographics/birthDate[year<1960]
 P 7 4 : /CPEXMLv1/person/partyActivities/hobby/startDate[year=1986]
 P 8 1 : /CPEXMLv1/person/partyDemographics[gender='MALE']
 P 8 2 : /CPEXMLv1/person/partyActivities/hobby[typeName='SPORT']
 P 8 3 : /CPEXMLv1/person/personDemographics/birthDate[year<1960]
 P 8 4 : /CPEXMLv1/person/partyActivities/hobby/startDate[year>1990]
 P 9 1 : /CPEXMLv1/person/partyActivities/food[favoriteFood='Hamburger']
 P 9 2 : /CPEXMLv1/person/partyActivities/hobby/startDate[year>1990]

【 0 0 4 2 】

10

上記の部分式のうちで、P 7 1 と P 9 1、P 7 2 と P 8 2、P 7 3 と P 8 3、P 8 4 と P 9 2 は、それぞれ同じ式である。したがって、X P a t h 式 P 7、P 8、P 9 は、有向非循環グラフ (DAG) にて表現することができる。

図 6 は、X P a t h 式 P 7、P 8、P 9 を表す有向非循環グラフを示す図である。

図 6 における左側の 6 個のロケーション・パスについては、上述した方法 1 によるステップの評価を共有する手法を用いることができるので、この部分を図 3 に示したような木構造のデータ構造に変換する。生成された有向非循環グラフのデータ構造は、方法 1 で生成された木構造のデータ構造と共にデータ構造格納部 3 0 に格納される。

【 0 0 4 3 】

評価実行部 4 0 は、以上のようにして得られたデータ構造を用いて、X P a t h 式 P 7、P 8、P 9 を評価する。まず、ロケーション・パス部分を、方法 1 を用いて評価する。次に、図 6 の有向非循環グラフを使って各 X P a t h 式を評価する。このとき、共通の部分式を持つものについては、当該部分式を 1 度だけ評価し、共通箇所に対して当該評価結果を再使用するようにする。

20

上記の例では、引数がノードセットで、評価した値が真理値になる論理演算式を例として取り上げたが、四則演算式など他のデータ型の演算式でも同様の最適化が可能である。

【 0 0 4 4 】

方法 3 . ロケーション・パスの述部の高速化

上述した方法 1 では、述部が異なるステップは独立のステップとして処理を行うものとして説明した。しかしながら、類似する述部の共通性を利用して述部の実行を高速化することが可能である。

30

例えば「[age < 20]」などのように、「[(式) (比較演算子) (定数)]」の形になる述部が高速化される。表面的にはこの形になっていなくても、置き換えや部分計算によってこの形に変形できればこの手法が適用できる。例えば、ノードの位置による指定「[2]」は、「[position() = 2]」に置き換えられるので、この手法が適用できる。

【 0 0 4 5 】

ここでは、(1) 比較演算子が「=」の場合、(2) 比較演算子が「!=」の場合、(3) 比較演算子が「>」の場合について、具体的について説明する。

(1) 比較演算子が「=」の場合の高速化

ステップにおける述部の式の評価結果をキーとして、残りのロケーション・パスを検索するためのハッシュテーブルを作る。例えば、次に示す X P a t h 式 P 4 1、P 4 2、P 4 3、P 4 4 のように、/CPEXMLv1/person/partyActivities/hobby[typeName='XYZ']/... の「XYZ」部分が様々な値を取るような X P a t h 式がある場合を考える。

40

P 4 1 :

/CPEXMLv1/person/partyActivities/hobby[typeName='SPORT']/startDate/year

P 4 2 :

/CPEXMLv1/person/partyActivities/hobby[typeName='MUSIC']/composer

P 4 3 :

/CPEXMLv1/person/partyActivities/hobby[typeName='SPORT']/hobbyName

P 4 4 :

50

/CPEXMLv1/person/partyActivities/hobby[typeName = 'COMPUTER']/software

【 0 0 4 6 】

データ構造生成部 20 は、上記の X P a t h 式 P 4 1、P 4 2、P 4 3、P 4 4 に対して、次のようなハッシュテーブルを生成する。

- ・キーが 'SPORT' のときに検索される組が { startDate/year, hobbyName }
- ・キーが 'MUSIC' のときに検索される組が { composer }
- ・キーが 'COMPUTER' のときに検索される組が { software }

また、/CPEXMLv1/person/partyActivities/hobbyまでの部分に関しては、方法 1 によりステップ評価の共有が可能であるので、この部分に対するデータ構造を生成する（ただし、この部分については X P a t h 式 P 4 1、P 4 2、P 4 3、P 4 4 が全て同一であるので、木構造としての分岐はない）。生成されたハッシュテーブルは、方法 1 で生成された木構造のデータ構造と共にデータ構造格納部 30 に格納される。

【 0 0 4 7 】

評価実行部 40 は、まず上述した方法 1 により、/CPEXMLv1/person/partyActivity/hobby ステップの評価を行い、ノードセットを得ておく。そして、ノードセットの各ノードに対して上記のハッシュテーブルを用いた処理を行う。

図 7 は、ハッシュテーブルを用いた評価処理の流れを説明するフローチャートである。

図 7 を参照すると、まず、述部の左辺式（この例の場合は「typeName」）を評価する（ステップ 701）。次に、左辺式の値をキーとしてハッシュテーブルを検索し（ステップ 702）、検索結果が得られた場合は、得られたロケーション・パスを当該ノードに対して評価し、評価結果を返す（ステップ 703、704）。一方、検索結果が得られない場合は、空のノードセットを返す（ステップ 703、705）。

【 0 0 4 8 】

以上説明した方法 3 を用いることにより、1つのノードに対して X P a t h の個数と同じ回数の述部を評価する必要がなくなり、1回のテーブル検索で1つのノードに対して全ての述部を評価したのと同じ結果が得られることとなる。例えば、typeNameが 'MUSIC' の場合、X P a t h 式 P 4 2 の述部が真になり、P 4 1、P 4 3、P 4 4 の述部が偽になることが1回のテーブル検索で判断される。

X P a t h 式の数 n で、ノードセットに含まれるノード数が m とすると、方法 3 を用いない場合、全てのノードに対して全ての X P a t h を評価するためには、述部のチェックに $O(n \times m)$ の時間を要する（ $O(n \times m)$ は、 $n \times m$ の定数倍以内の時間で計算が可能であることを意味する）。一方、方法 3 を用いた場合、ほぼ定数時間で実行できるハッシュテーブルの検索が m 回行なわれる。したがって、 n が十分大きい場合、方法 3 を用いることにより、評価処理の大幅な高速化を図ることができる。

【 0 0 4 9 】

（ 2 ）比較演算子が「 ! = 」の場合の高速化

上記の比較演算子が「 = 」の場合に、キーが K_i のときに検索される残りのロケーション・パスの集合が P_i であるとする。比較演算子が「 ! = 」の場合には、キーが K_j のときに検索される残りのロケーション・パスの集合が $P_i \setminus P_j$ になるようにハッシュテーブルを構成する。X P a t h 式 P 4 1、P 4 2、P 4 3、P 4 4 の場合、次のような検索結果を返すハッシュテーブルを構成する。

- ・キーが 'SPORT' のとき、{ composer, software }
- ・キーが 'MUSIC' のとき、{ startDate/year, hobbyName, software }
- ・キーが 'COMPUTER' のとき、{ startDate/year, composer, hobbyName }

以上のハッシュテーブルを用いる他、全体の処理については（ 1 ）比較演算子が「 = 」の場合と同様であるため、説明を省略する。

【 0 0 5 0 】

（ 3 ）比較演算子が「 > 」の場合の高速化

比較演算子が「 > 」の場合、データ構造生成部 20 は、ステップの述部の定数をキーとして検索を行う二分探索木を作る。例えば、次の2つの X P a t h 式 P 5、P 6 がある場合

10

20

30

40

50

、
 P 5 : /CPEXMLv1/person/personDemographics/birthDate[year>1990]/gender
 P 6 : /CPEXMLv1/person/personDemographics/birthDate[year>1976]/birthPlace
 次のような計算を行なう二分探索木を構成する。

```

if (key > 1990) then
    {gender, birthPlace}
else
    if (key > 1976) then
        {birthPlace}
    else
        {}
    endif
endif

```

10

生成された二分探索木は、方法 1 で生成された木構造のデータ構造と共にデータ構造格納部 30 に格納される。

20

【 0 0 5 1 】

評価実行部 40 は、まず上述した方法 1 により、/CPEXMLv1/person/personDemographics/birthDateまでのステップの評価を行う。次に、述部の左辺式（この例の場合はyear）を評価し、その値をキーとして二分探索木の検索を行う。この探索の結果、キーが1990よりも大きい値の場合は、「gender」と「birthPlace」の両方が返される。キーが1990以下で1976よりも大きい値の場合には「birthPlace」が返される。そして、どの条件も満たさない場合には空集合が返される。

X P a t h 式の数 n で、ノードセットに含まれるノード数が m とすると、方法 3 を用いない場合、全てのノードに対して全ての X P a t h を評価するためには、述部のチェックに $O(n \times m)$ の時間を要する。一方、方法 3 を用いた場合、 $O(\log n)$ の時間を要する二分探索が m 回行なわれる。したがって、 n が十分大きい場合、方法 3 を用いることにより、評価処理の大幅な高速化を図ることができる。

30

なお、ここでは比較演算子が「>」の場合について説明したが、比較演算子が「>=」、「<」、「<=」の場合も同様に扱うことができるのは言うまでもない。

【 0 0 5 2 】

次に、データ構造生成部 20 が、X P a t h 格納部 10 に格納されている X P a t h 式から評価実行部 40 に用いられるデータ構造（図 3 参照）を生成する方法について説明する。

本実施の形態で用いられるデータ構造の生成には、原則的には木構造を生成する一般的な種々の方法を用いることができ、特にその生成方法は限定されない。しかしながら、本実施の形態では、X P a t h 格納部 10 に格納されている X P a t h が変更（追加または削除）された場合に、効率的にかかる変更をデータ構造に反映させることのできる生成方法を提案し、これについて説明する。

40

【 0 0 5 3 】

X M L 文書に対して X P a t h 式を評価する処理に要する時間を短縮するためには、処理を行うたびに毎回上述したようなデータ構造を生成するのではなく、生成されたデータ構造をデータ構造格納部 30 に保存しておき、複数の X M L 文書に対する処理において使い回すことが考えられる。

さらに、X P a t h 格納部 10 に格納されている X P a t h の集合が変化した場合に、データ構造を最初から生成し直すのではなく、追加された X P a t h 式と削除された X P a

50

th式とを、保存してあるデータ構造に反映するようにすれば、効率的に所望のX Path集合のためのデータ構造を得ることができる。

すなわち、X Path式が新たに追加された場合には、そのX Path式のうちで、既存のデータ構造に対応するノードがないステップについて新たなノードを生成し、このデータ構造に加える。一方、所定のX Path式が削除された場合には、そのX Path式のステップに対応するノードのうちでこのX Path式に固有の部分に対応するノードのみを除去する。これにより、データ構造を効率よく更新する。

【0054】

図8は、データ構造生成部20によるデータ構造の生成方法を説明するフローチャートである。

10

本実施の形態では、X Pathを1つずつ順次追加していくことにより、所望のデータ構造を得ることとする。したがって、最初にデータ構造を生成する際には、X Path格納部10に格納されているX Path集合の中から所定のX Path式を選択し、他のX Path式を順次加えていくこととなり、X Path集合の変更により新たなX Path式が追加された場合には、既に生成されたデータ構造格納部30に格納されているデータ構造に、新たに追加されたX Path式のステップを加えていくこととなる。

【0055】

図8を参照すると、データ構造生成部20は、まず追加するX Path式Pをステップs1、s2、・・・、skに分解する(ステップ801)。次に、パラメータi及びパラメータSを初期化、すなわちi=1、S=「トップレベルのステップの集合」とし(ステップ802)、ステップの集合Sからステップsiを探す(ステップ803)。

20

ステップの集合Sにステップsiが含まれていれば、次に、i+1をパラメータiの新しい値とし(ステップ804、805)、当該新たなiがX Path式Pにおけるステップの数kに達していないか調べる(ステップ806)。そして、iの値がkに達していなければ、パラメータSをS=「データ構造中でステップsiの次のレベルにおけるステップの集合」とし、ステップ803から処理を繰り返す(ステップ807)。

【0056】

一方、ステップ804で、ステップの集合Sにステップsiが含まれていないと判断されたならば、si、・・・、skに対応するステップをデータ構造に追加し(ステップ808)、X Path式Pを当該データ構造の該当ステップに登録して処理を終了する(ステップ809)。

30

また、ステップ806で、新しいiの値がX Path式Pにおけるステップの数kに達した場合も、X Path式Pを当該データ構造の該当ステップに登録して処理を終了する(ステップ809)。

【0057】

ここで、図3に示したデータ構造に、次のX Path式P4を追加する場合を例として、処理を具体的に説明する。

P4 : /profile/location/city[text() = 'Tokyo']

まず、ステップ801で、X Path式P4が、ステップごとに次のように分解される

s1 = profile

s2 = location

s3 = city

s4 = [text() = 'Tokyo']

40

なお、ロケーション・パスが省略形を使って表現してある場合には、省略形を使わない表現に変換してから分解する。例えば、ステップの区切り「//」は「decendant-or-self::node()」に、「.」は「self::node()」に変換しておく。

【0058】

次に、ステップ802で、i=1、S={profile}に初期化される。

ステップs1がステップの集合Sに含まれ、かつ残りのステップが存在するので(ステップ803~806)、ステップ807に進み、データ構造中で「profile」の次のレベル

50

のステップ「{ demographics, interests }」が S に代入される。

すると、今度はステップ 804 の判断において、ステップ s 2 の値「location」はステップの集合 S に含まれないので、ステップ 808 に進み、s 2、s 3、s 4 に対応するステップがデータ構造に追加される。そして、ステップ 809 で、X P a t h 式 P 4 がデータ構造に登録される。

図 9 は、図 3 のデータ構造に X P a t h 式 P 4 が追加された状態のデータ構造を示す図である。

【 0 0 5 9 】

以上で、上述した方法 1 のステップ評価の高速化を実現するためのデータ構造が生成（更新）された。次に、方法 2 による演算式評価の高速化を実現するためのデータ構造の生成（更新）について説明する。

まず、新たに加える演算式からロケーション・パスを抽出し、その中から既存のデータ構造に登録されていないロケーション・パスだけをデータ構造に追加する。ロケーション・パスの追加は、図 8 のフローチャートに示した手順にて行われる。

ロケーション・パス以外の部分式に対しては、共通部分式として既に現れているかどうかを検索する。そして、既に共通部分式として現れている場合には、その共通部分式を評価するデータ構造を再使用するようにする。共通部分式として現れていない場合は、その部分式を評価するためのデータ構造を新たに作り、既存のデータ構造に追加する。

【 0 0 6 0 】

次に、方法 3 による述部評価の高速化を実現するためのデータ構造の生成（更新）について説明する。

図 8 のステップ 808 において、ステップ s i の述部が、[<式> = <定数>]、[<式> != <定数>]、[<式> < <定数>]、[<式> <= <定数>]、[<式> > <定数>]、[<式> >= <定数>]のいずれかである場合、まず、式に対応するハッシュテーブル（または二分木）があるかどうかを調べる。対応するハッシュテーブル（または二分木）がない場合は、これらを生成する。すなわち、定数をキーとした検索項目を既存のハッシュテーブル（または二分木）に追加する。そして、残りのステップ列（s i+1、・・・、s k）を、木構造にして既存のデータ構造に追加する。データ構造への追加は、図 8 のフローチャートに示した手順にて行われる。

【 0 0 6 1 】

図 10 は、データ構造から所定の X P a t h 式を除去する際のデータ構造生成部 20 の処理を説明するフローチャートである。

図 10 を参照すると、データ構造生成部 20 は、まずパラメータ s を s = 「データ構造上で X P a t h 式 P が登録されているステップ」とし（ステップ 1001）、ステップ s から X P a t h 式 P を取り除く（ステップ 1002）。

次に、ステップ s に登録されている X P a t h 式の数 が 0 かどうかを調べ、0 でないならば処理を終了する（ステップ 1003）。

一方、ステップ s に登録されている X P a t h 式の数 が 0 になった場合、次に、パラメータ s p を s p = 「データ構造中でステップ s の 1 つ上位のステップ」とし（ステップ 1004）、ステップ s p からステップ s を取り除く（ステップ 1005）。そして、ステップ s p の 1 つ下位のステップ数 が 0 かどうかを調べ、0 でないならば処理を終了する（ステップ 1006）。

これに対し、ステップ s p の 1 つ下位のステップ数 が 0 になった場合、s = s p としてステップ 1004 に戻り（ステップ 1007）、新たなステップ s に関して処理を繰り返す。

【 0 0 6 2 】

ここで、図 9 に示したデータ構造から X P a t h 式 P 3 を除去する場合を例として、処理を具体的に説明する。

まず、ステップ 1001 で、パラメータ s に、データ構造上の「[text() >= 40 and text < 50]」のステップが割り当てられる。ここから X P a t h 式 P 3 を取り除くと、ステッ

10

20

30

40

50

プ s に登録されている X P a t h 式の数 が 0 になるので、ステップ 1 0 0 4 に進む。そして、ステップ 1 0 0 5 で、ステップ s の上位ステップ「demographics」からステップ s を取り除く。この場合、ステップ 1 0 0 6 の判断において、「demographics」の下位ステップ数が 0 ではないので、ここで処理が終了する。

図 1 1 は、図 9 のデータ構造から X P a t h 式 P 3 を除去した状態のデータ構造を示す図である。

【 0 0 6 3 】

以上で、上述した方法 1 のステップ評価の高速化を実現するためのデータ構造が更新された。次に、方法 2 による演算式評価の高速化を実現するためのデータ構造の更新について説明する。

まず、取り除こうとする演算式からロケーション・パスを抽出し、その中からこの演算式だけが使っているロケーション・パスを求め、当該ロケーション・パスをデータ構造から取り除く。また、所定のロケーション・パスが取り除こうとする演算式だけから使われているかどうかの検査は、ロケーション・パスの参照数を管理するなどの手法を用いて行うことができる。

さらに、ロケーション・パス以外の部分式に対しても同様に、この演算式だけが使っている部分式を求めて、データ構造から取り除く。

【 0 0 6 4 】

次に、方法 3 による述部評価の高速化を実現するためのデータ構造の更新について説明する。

まず、述部の定数に対応する項目をハッシュテーブル（二分木）から取り除く。そして、この項目から検索されるロケーション・パスをデータ構造から削除する。ロケーション・パスの削除は、図 1 0 のフローチャートに示した手順にて行われる。

【 0 0 6 5 】

次に、X P a t h 式及び X M L 文書の具体例を挙げて、本実施の形態の動作についてさらに説明する。

X P a t h 格納部 1 0 に、次の 4 つの X P a t h 式が格納されているものとする。

P 1 1 : /profile/interests[sport/@type = 'Soccer']/music

P 1 2 : /profile/interests[sport/@type = 'Baseball']/book

P 1 3 : /profile/demographics/age

P 1 4 : count(/profile/interests[sport/@type = 'Soccer']) > 1

【 0 0 6 6 】

データ構造生成部 2 0 は、図 8 に示した手順で上記 4 つの X P a t h 式 P 1 1、P 1 2、P 1 3、P 1 4 に対するデータ構造を生成する。なお、以下では、空の木構造（初期状態）に対して、X P a t h 式 P 1 1、P 1 2、P 1 3、P 1 4 を順次追加するものとして説明する。

P 1 1 の追加：

(1) ロケーション・パスを 3 つのステップ「profile」、「interests[sport/@type = 'Soccer']」、「music」に分解する（ステップ 8 0 1）。

(2) ステップ「profile」はデータ構造に登録されていないので、追加し登録する（ステップ 8 0 4、8 0 8、8 0 9）。

(3) ステップ「interests[sport/@type = 'Soccer']」を「interests」と「[sport/@type = 'Soccer']」とに分解する。

(4) ステップ「interests」はデータ構造に登録されていないので、追加し登録する（ステップ 8 0 4、8 0 8、8 0 9）。

(5) 述語のためのハッシュテーブルがないので生成する。

(6) 述語の左辺式であるロケーション・パス「sport/@type」を表現するデータ構造を作り、ハッシュテーブルに付加する。

(7) キーが'Soccer'で値が残りのステップ「music」となるエントリをハッシュテーブルに追加する。

【 0 0 6 7 】

P 1 2 の追加 :

(1) ロケーション・パスを 3 つのステップ「profile」、`「interests[sport/@type = 'Baseball']」`、`「book」` に分解する (ステップ 8 0 1)。

(2) ステップ「profile」と「interests」とは、既にデータ構造に登録されており、また対応するハッシュテーブルに「sport/@type」のためのデータ構造も付加されているので、データ構造への追加は行わず、X P a t h 式 P 1 2 の登録のみを行う (ステップ 8 0 9)。

(3) キーが 'Baseball' で、値が残りのステップ「book」となるエントリをハッシュテーブルに追加する。

10

【 0 0 6 8 】

P 1 3 の追加 :

(1) ロケーション・パスを 3 つのステップ「profile」、「demographics」、「age」に分解する (ステップ 8 0 1)。

(2) ステップ「profile」は登録されているが、「demographics」と「age」は登録されていないので、「demographics」と「age」とをステップとしてデータ構造に追加し、X P a t h 式 P 1 3 を登録する (ステップ 8 0 4、8 0 8、8 0 9)。

【 0 0 6 9 】

P 1 4 の追加 :

(1) 式を部分式に分解して、ロケーション・パス `「/profile/interests[sport/@type = 'Baseball']/」` を抽出する。

(2) ロケーション・パス `「/profile/interests[sport/@type = 'Baseball']」` は、既に登録されているので、その結果を関数 count に渡し、さらに関数 count の結果を比較演算式に渡すようにデータ構造を生成する。

20

【 0 0 7 0 】

以上の結果、4 つの X P a t h 式 P 1 1、P 1 2、P 1 3、P 1 4 を表現するデータ構造が生成される。

図 1 2 は、生成されたデータ構造を示す図である。

生成されたデータ構造はデータ構造格納部 3 0 に格納され、保存される。

【 0 0 7 1 】

次に、図 1 3 に示す X M L 文書が入力され、X M L 文書格納部 5 0 に格納されたものとする。評価実行部 4 0 は、X M L 文書格納部 5 0 からこの X M L 文書を読み出し、データ構造格納部 3 0 に格納されているデータ構造を用いて、X P a t h 式 P 1 1、P 1 2、P 1 3、P 1 4 を評価する。

30

まず、最初のステップ「profile」を実行し、2 つの要素からなるノードセット N 1 を得る。

N 1 :

```
{ <demographics>
  <age>19</age>
</demographics>,
<interests>
  <sport type@='Baseball' />
  <book>History</book/>
</interests>
}
```

40

50

【 0 0 7 2 】

次に、ノードセット N 1 の各要素にステップ「interests」を適用し、次の 2 つのノードセットを得る。

{ }

{ <sport type@='Baseball' />, <book>History</book> }

そして、この 2 つの併合し、「interests」の結果であるノードセット N 2 を得る。

N 2 :

{ <sport type@='Baseball' />, <book>History</book> }

【 0 0 7 3 】

次に、ノードセット N 2 の 2 つの要素に対して、左辺式「sport/@type」を評価すると、それぞれ、{ } と { 'Baseball' } とが得られる。

そして、「Baseball」をキーとしてハッシュテーブルを検索すると、ステップ「book」と関数countが得られる。「book」をノードセット N 2 の 2 つの要素に対して評価すると、{ } と { History } とが得られる。この 2 つを併合して

X P a t h 式 P 2 の評価値：{ History }

を得る。

【 0 0 7 4 】

次に、{ <sport type@='Baseball' />, } を引数として関数countを評価すると値は 1 となり、この値 1 との比較演算を評価するとfalseになるので

X P a t h 式 P 4 の評価値：false

を得る。

そして、ノードセット N 1 の各要素にステップ「demographics」を適用し、結果を併合すると「{ <age>19</age> }」になり、さらにこの要素にステップ「age」を適用して、

X P a t h 式 P 3 の評価値：{ 19 }

を得る。

さらに、データ構造を訪問し終わっても X P a t h 式 P 1 にたどり着くノードがないので、

X P a t h 式 P 1 の評価値：空のノードセット

を得る。

【 0 0 7 5 】

次に、本実施の形態による XML 文書処理システムによる X P a t h 式の評価処理の実施例を示す。

本実施の形態では、上述したように X P a t h 式をステップごとに細分化し、細分化によってできた X P a t h の断片の実行に関して、実行順序を代えたり実行回数を減らしたりすることにより、全体の実行時間を短縮している。

X P a t h 式の断片もまた X P a t h 式であり、この断片の実行には、Apache XML Projectにて提供される X S L T プロセッサXalan-Java 2 (以下、単にXalanと記す)に含まれる X P a t h プロセッサのパッケージorg.apache.xpathを使用した。例えば、ロケーション・パスの 1 つのステップの実行や算術式の評価のためにXalanの機能をそのまま使っている。このため、Xalanのみを使用して実行した場合と、本実施の形態による XML 文書処理システムでXalanを使用し実行したものの差分によって、本実施の形態による X P a t h の評価における効率化の度合いを直接知ることができる。

【 0 0 7 6 】

本実施例では、CPExchange (IDEAlliance, CPExchange Specification Version 1.0, 2000. (<http://www.cpexchange.org/>)) で定義される XML 文書に対して、多くの X P a t h 式が用意されており、それらの X P a t h 式の中から与えられた XML 文書に適合するものを選ぶことを行なう。X P a t h 式としては、次のようなパターンで X Y Z が異なるものを生成してシステムに登録する。

/CPEXMLv1/person/partyActivities/hobby[typeName='XYZ']

/CPEXMLv1/person/partyActivities/hobby[hobbyName='XYZ']

10

20

30

40

50

/CPEXMLv1/person/partyActivities/food[favoriteFood='XYZ']
 /CPEXMLv1/person/personDemographics/gender[@enumtype='XYZ']
 /CPEXMLv1/personName[firstName='XYZ']
 /CPEXMLv1/person/partyActivities/hobby/startDate[year!='XYZ']
 /CPEXMLv1/person/personDemographics/birthDate[year='XYZ']
 /CPEXMLv1/person/partyActivities/newspaper/startDate[year='XYZ']
 /CPEXMLv1/person/partyActivities/hobby/startDate[year='XYZ']
 /CPEXMLv1/person/partyActivities/magazine/startDate[year='XYZ']

また、システムの動作環境としては、C P U 1 0 1 に米国Intel社のモバイル Pentium IIIの800MHzを用い、メインメモリ103を128MBのRAM (Random Access Memory)、オペレーティングシステムを米国Microsoft社のWindows 2000とした。

【0077】

図14は、上記の条件で、本実施例によるXPath式の評価に要した実行時間を示す図である。図14には、(1)XalanだけでXPath式を評価するシステムと(2)本発明の手法とXalanの組合わせで評価するシステムの2つの実行時間を測定した結果が示されている。

図14を参照すると、Xalanだけを使用したシステムでは、XPath式の数の増加に対して直線的に実行時間が増大している。一方、本実施の形態においてXalanを使用したシステムでは、XPath式の数の増加量に対して実行時間の増大がほとんど見られない。したがって、XPath式の数が増加するにしたがって両社の性能の違いが大きくなって

【0078】

〔第2の実施の形態〕

上記第1の実施の形態では、評価されるべき複数のXPath式に関し、その共通部分を共用するデータ構造を用いてXML文書に対する評価を行うことにより、評価結果を共用し、評価処理の高速化を実現した。これに対し、第2の実施の形態では、評価されるべき複数のXPath式の間依存関係を解析し、得られた依存関係に基づいて当該複数のXPath式の冗長部分を省略したデータ構造を用いて評価を行うことにより、評価処理を単純化し、評価処理の高速化を実現する。なお、本実施の形態では、ウェブページにアノテーションを付加する目的で、アノテーション内のXPath式が対象のウェブページの

所定のエレメントを正しく指しているかどうかを評価する場合を例として説明する。第2の実施の形態によるウェブページ処理システムは、第1の実施の形態と同様に、例えば図1に示すようなハードウェア構成を有するコンピュータ装置にて実現される。

【0079】

図15は、第2の実施の形態によるウェブページ処理システムの機能を説明するブロック図である。

図15を参照すると、本実施の形態は、処理対象であるウェブページ(HTML文書)の特定部分を指し示す複数のXPath式を格納したXPath格納部1510と、XPath格納部1510に格納されているXPath式から本実施の形態による効率的な評価を行うためのデータ構造を生成するデータ構造生成部1520と、データ構造生成部1520にて生成されたデータ構造を格納するデータ構造格納部1530と、データ構造格納部1530に格納されているデータ構造を用いて処理対象であるウェブページに対してXPath式の評価を行う評価実行部1540とを備え、また処理対象であるウェブページを格納する文書格納部1550と、評価実行部1540による評価を経たウェブページを対象として所定の情報処理を実行する文書処理部1560とを備える。

【0080】

図15に示したシステム構成において、データ構造生成部1520、評価実行部1540及び文書処理部1560は、図1に示したプログラム制御されたCPU101にて実現される仮想的なソフトウェアブロックである。CPU101を制御してこれらの機能を実現するプログラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記録媒体に格納

して配布したり、ネットワークを介して配信したりすることにより提供され、メインメモリ103に読み込まれる。また、XPath格納部1510、データ構造格納部1530及び文書格納部1550は、メインメモリ103にて実現される。なお、メインメモリ103に保持されるデータやプログラムは、必要に応じてハードディスク105などの記憶装置に退避させることができる。

【0081】

図15に示したシステム構成において、XPath格納部1510には、ウェブページに対して評価を行うべき、予め用意された複数のXPath式が格納されている。データ構造格納部1530には、データ構造生成部1520にて生成され、評価実行部1540にて用いられるデータ構造が格納される。文書格納部1550には、図1に示したネットワークインターフェイス106や、キーボード/マウス109等の入力デバイス、処理対象であるウェブページを格納した記録媒体から当該ウェブページを読み出すドライブ装置などを介して入力されたウェブページが格納されている。

10

【0082】

データ構造生成部1520は、XPath格納部1510に格納されているXPath式から、評価実行部1540による評価に用いるデータ構造を生成する。このデータ構造を用いることにより、ウェブページに対してXPath格納部1510に格納されている複数のXPath式を個別に評価する場合に比べて効率の良い評価を行うことが可能となる。このデータ構造についての詳細、このデータ構造を用いた評価方法及びこのデータ構造の生成方法については後述する。

20

【0083】

評価実行部1540は、XPathエンジンであり、文書格納部1550に格納されているウェブページに対して、データ構造生成部1520にて生成されたデータ構造を用いてXPath式の評価(マッチング)を行う。これにより、当該入力文書がXPath格納部1510に格納されているXPath式にて指定される特定部分を持つかどうかを、効率よく評価することができる。評価実行部1540による評価結果はメインメモリ103の所定の領域に格納される。

文書処理部1560は、メインメモリ103の所定領域から評価実行部1540の評価結果を読み出し、当該評価結果に基づいて、文書格納部1550に格納されているウェブページがXPath格納部1510に格納されているXPath式にて指定される特定部分を持つならば、文書格納部1550から当該ウェブページを処理対象として読み出し、目的とする所定の情報処理、すなわちウェブページにアノテーションを付加する処理を実行する。

30

【0084】

次に、本実施の形態におけるXPath式の評価を効率化する方法を詳細に説明する。ウェブページにアノテーションを付加する作業において、所定のアノテーションが所定のウェブページに対して適用可能かどうかを判断するには、当該アノテーション内の複数のXPath式が対象のウェブページの所定のエレメントを正しく指しているかどうかを評価することが必要である。これは、複数のXPath式の集合(以下、XPathセットと称す)が当該ウェブページにマッチしているかどうかを評価する問題と言い換えることができる。XPathセットがウェブページにマッチするとは、XPathセットに含まれる全てのXPath式が当該ウェブページの構造にマッチすることを意味する。本実施の形態では、特に、複数のアノテーションの中で所定のウェブページに適用できるものを探索する場合のように、複数のXPathセットのうちでウェブページにマッチするものを判定する場合の評価処理の効率化を対象とする。

40

【0085】

本実施の形態は、複数のXPath式の間依存関係に基づいてXPath式の評価を簡単化するためのデータ構造を生成し、このデータ構造を用いてウェブページに対する評価を行うことにより、上記のような場合の評価処理の高速化を実現する。本実施の形態では、評価実行部1540の種類(下記のタイプ1、タイプ2)に応じて異なるアルゴリズム

50

が適用される。

タイプ1：XPath評価関数の呼び出しオーバーヘッドが大きく、XPath式の評価の回数が少ない方がシステム全体のパフォーマンス（処理効率）を向上させることができる。

タイプ2：XPath評価関数の呼び出しオーバーヘッドが小さく、XPath式を分割して評価を行ってもシステム全体のパフォーマンス（処理効率）は大きく低下しない。

【0086】

まず、評価実行部1540の種類がタイプ1である場合について説明する。

この場合、XPath式の評価回数が少ない方がシステムのパフォーマンスを向上させられるため、データ構造生成部1520は、生成するデータ構造においてXPath式の分割を行わない。具体的には、

1. XPath式の共通性に基づくXPath式の絞込み
2. XPath式の依存関係に基づく連結XPath式の簡略化
3. XPath式の依存関係に基づくXPath式の省略
4. XPath式の統合
5. XPath式のツリー長とDOMツリー上での位置に基づく評価処理の優先順位付け
6. 実行時判断のためのデータ構造の構築

という処理が行われる。各処理について、以下に説明する。なお、以下の説明において、「XPath式が成立する」「XPath式が真になる」とは、「評価対象であるウェブページ中に、XPath式に対応するエレメントが存在している」ことを意味する。

【0087】

1. XPath式の共通性に基づくXPath式の絞込み

複数のXPathセットが共通に含んでいるXPath式は、ウェブページにマッチするXPathセットの判定に関与しない。このようなXPath式の評価結果は、各XPathセットにおいて共通となる（差別化されない）ためである。したがって、複数のXPathセットに含まれないXPath（specific XPath）に対する評価を先に行う。

図16は、XPathセットとこれに含まれるXPath式とを一覧表示した図表である。図16に示す例では、4種類のXPathセット（A、B、C、D）にそれぞれ複数のXPath式が含まれている。図16に示すXPath式には、重複するものが多く含まれている。例えば、番号5（XPathセットA）のXPath式と番号14（XPathセットB）のXPath式とは、全く同一のXPathである。

図17は、図16に示したXPath式のうち重複するものを整理して識別情報（XPathID）を割り当てた様子を示す図である。図17に示す例では、図16に示した31個のXPath式が18個に整理されている。

【0088】

図18は、図17のように整理されたXPath式をさらに解析し、複数のXPathセットに共通するXPath式を抽出した様子をグラフ構造で示した図である。

図17を参照すると、XPath式X7はXPathセットC、Dに共通して含まれ、またXPath式X1、X2、X3、X18は全てのXPathセットに共通して含まれる。図18の例では、これら複数のXPathセットに共通して含まれるXPath式を表すノードから個々のXPathセットを分岐させて、各XPathセットに固有のXPath式を記述している。したがって、図18の例において、評価対象は、各XPathセットを表すノードに付加されている13種類のXPath式（X4～X6、X8～X17）に絞込まれることとなる。

【0089】

2. XPath式の依存関係に基づく連結XPath式の簡略化

連結XPath式は、複数のXPath式をOR記号（|）で結合したものである。例えば、所定のウェブページ（html[1]）の所定の位置（body[1]）に配置された2つのテーブルコンテンツ（table[1]、table[2]）のいずれかを指す連結XPath式、
/html[1]/body[1]/table[1]|/html[1]/body[1]/table[2]

10

20

30

40

50

を考える。これは、所定のウェブページに、/html[1]/body[1]/table[1]または/html[1]/body[1]/table[2]のいずれか一方が存在すれば当該X P a t h式が成立することを示している。このとき、X P a t h式の表記規則から、table[1]が存在しなければ、table[2]も存在しないことは自明である。したがって、この連結X P a t h式に関しては、前半の/html[1]/body[1]/table[1]

のみを評価すれば良いことになる。このように所定のX P a t h式(X P a t h 1)が成立するとき、別のX P a t h式(X P a t h 2)も成立する場合、X P a t h 2はX P a t h 1に依存していると定義する。この性質を利用して連結X P a t hを簡略化する。

【0090】

X P a t h式の間依存関係には、主として次の2種類がある。

- ・同一タグ兄弟(sibling)関係

X P a t h 1 /html[1]/body[1]/table[2]

X P a t h 2 /html[1]/body[1]/table[1]

であるとき、X P a t h 1とX P a t h 2とは兄弟の関係にあり、

X P a t h 1が成立すればX P a t h 2も成立し、

X P a t h 2が不成立ならばX P a t h 1も不成立である。

したがって、「X P a t h 1はX P a t h 2に依存」する。

- ・親ノード関係

X P a t h 1 /html[1]/body[1]/table[1]/tr[1]/td[1]

X P a t h 2 /html[1]/body[1]/table[1]

であるとき、X P a t h 1はX P a t h 2の子孫であり、

X P a t h 1が成立すればX P a t h 2も成立し、

X P a t h 2が不成立ならばX P a t h 1も不成立である。

したがって、「X P a t h 1はX P a t h 2に依存」する。

【0091】

連結X P a t h式の簡略化は、以下の手順で行う。

まず、連結X P a t h式を展開して新たにX P a t h I Dを割り当てる。例えば図17に示したX P a t h式X 1は、

/html[1]/body[1]/table[1]/tbody[1]/tr[1]/html[1]/body[1]/table[1]/tbody[1]/tr[2]

であるから、次の2つのX P a t h式に展開することができる。

X 1 - 1 : /html[1]/body[1]/table[1]/tbody[1]/tr[1]

X 1 - 2 : /html[1]/body[1]/table[1]/tbody[1]/tr[2]

同様に、図18において評価対象とした13種類のX P a t h式を解析し展開する。そして、展開された各X P a t h式の依存関係を調べる。

【0092】

図19は、得られた展開後のX P a t h式の依存関係をツリー構造で示す図である。なお、図中、円で囲まれたX P a t h式X 6 - 5とX P a t h式X 7とは、同一のX P a t h式である。

図19を参照すると、X 1 1 - 3、X 1 1 - 2、X 1 1 - 1のように、で本来1つの連結X P a t h式(X 1 1)に含まれている複数のX P a t h式が依存関係を持つ場合がある。この場合、次の理由により、下位のX P a t h式であるX 1 1 - 2、X 1 1 - 3は評価が不要となる。

すなわち、X 1 1 - 1が成立した場合、元の連結X P a t h式X 1 1も成立する。これは結合関係にあるX P a t h式が1つでも成立すれば、かかるX P a t h式を含む連結X P a t h式は成立するという性質を持つからである。一方、依存関係により、X 1 1 - 1が不成立の場合、X 1 1 - 2、X 1 1 - 3も不成立となる。したがって、X 1 1 - 2、X 1 1 - 3はX P a t h式X 1 1の評価結果に影響を与えない。

以上のように、他のX P a t h式に依存するX P a t h式を削除することによって、連結X P a t h式が簡略化される。

10

20

30

40

50

図20は、図19に示したXPath式の依存関係に基づいて連結XPath式を簡略化した様子を示す図である。

【0093】

3. XPath式の依存関係に基づくXPath式の省略

複数のXPath式において、ウェブページの構造やエレメントのバリエーションなどから受ける制約のために、1つのXPath式について評価結果が得られた場合に、他のXPath式の評価結果が自明となる場合がある。例えば、XPathセット中に所定のウェブページ(html[1])の所定の位置(body[1])に配置された2つのテーブルコンテンツ(table[1]、table[2])を指す2つのXPath式、

XPath1 /html[1]/body[1]/table[2]

XPath2 /html[1]/body[1]/table[1]

がある場合を考える。この場合、table[2]を指すXPath1が存在するならば、table[1]を指すXPath2も必ず存在する。したがって、XPath2は評価しなくても良いこととなる。この性質を利用して評価するXPath式を削減する。

【0094】

これらXPathの依存関係についてさらに詳しく説明する。

所定のXPathセットが成立するためには、当該XPathセットに「含まれる全てのXPath式の成立」が条件となる。上記のXPath1、XPath2は、どちらか一方でも不成立の場合、XPathセット全体が成立しない。すなわち、以下のような関係が成り立つ。

XPath1が成立すればXPath2も成立する。

XPath1が不成立ならばXPathセット全体が成立しない。

これにより、XPath2の成立不成立を無視することができる。言い換えると、次のようになる。

「所定のXPath式(上の例でXPath2)が含まれる全てのXPathセットが依存関係ツリー(図19、20に示したようなツリー構造)の末端側に出現するとき、そのXPath式(上の例でXPath2)は成立不成立に関与しないため、省略できる。」

図21は、図20に示したXPath式の依存関係に基づいてXPath式を省略した様子を示す図である。図21(A)は、省略されるXPath式(X2、X7、X9)をマークした状態、図21(B)は、図21(A)でマークされたXPath式を依存関係ツリーから消去した状態を示す。図中のXPath式に付された「A」、「B」、「ABCD」等の符号は、当該XPath式が含まれるXPathセットを示す。なお、図21中では、X11及びX1の添え字表記(X1-1、X11-1の「-1」)をはずしてある。これは、2.の処理によって、これらのXPath式が連結XPath式から単独XPath式に変わったためである。

【0095】

4. XPath式の統合

XPath評価関数の呼び出しオーバーヘッドが大きいタイプ1の評価実行部1540を用いる場合、データ構造生成部1520は、評価を行うべきXPath式の本数を減らすため、2.の処理においてXPath式の依存関係を解析した際にばらばらにした連結XPath式を、再び結合する。

図22は、図21に示したXPath式の依存関係ツリー中のXPath式に対して再結合可能なものを再結合する様子を示す図である。

結合処理は、次の規則に基づいて行う。

(1) 結合後のXPath式的位置は、所属するXPath式の全てが共有する親ノードの直下とする。

(2) 共有する親ノードがない場合、結合後のXPath式自体がルートノードとなる。

図21、22に示す例では、例えば、X5-*というXPath式は、X5-4、X5-1、X5-2の3つが存在する(図21(B)参照)。これらのXPath式を、それぞれルートに向かって辿っていくと、X4で合流する。そこで、X5-*の各XPath式

10

20

30

40

50

をX4の直下のノードとして1つにまとめ(図22(A)参照)、再結合してXPath式X5とする(図22(B)参照)。同様にして、X6、X17が再結合される。

ここで、仮にX5-2がX18の子ノードであったとする。この場合、共有する親ノードは存在しないため、X5-*は1つのルートノードとして、X1、X3、X4、X18と同じレベルのノードとして結合されることとなる。

【0096】

5. XPath式のツリー長とDOMツリー上での位置に基づく評価処理の優先順位付け評価対象であるウェブページにマッチしないXPathセットを効率よく決定しふるい落とすためには、「存在しにくいエレメント」を対象にしたXPath式を先に処理することが好ましい。そこで、次のような基準を用いて評価するXPath式の順番(優先順位)を決定する。

(1) XPath式のツリー長が長いXPath式

(2) 評価対象であるウェブページのDOMツリー上で、より後ろに位置するXPath式

(3) XPath式の依存関係に基づき、成立したときにより多くのXPath式に影響を与えるものを優先する。

を先に評価する。例えば、/table[1]よりもより末端まで含んだ/table[1]/tr[1]/td[1]の評価を先に行う(規則(1))。また、/table[1]と/table[2]とでは、DOMツリー上で巡回したときに、より「後ろに」出現する、言い換えればHTMLのタグ順序で後ろに出現するエレメントを指しているため、/table[2]の評価を先に行う(規則(2))。これは、DOMツリー上で後ろに位置するXPath式(/table[2])の出現率は前に位置するXPath式(/table[1])の出現確率以下だからである。

【0097】

図23は、図22に示した依存関係ツリーの各ノード(XPath式)に、上記の規則に従って優先順位(丸数字)を付加した様子を示す図、図24は、図23に示した優先順位に従ってXPath式を並べた図表である。図24には、さらに各XPath式が依存している他のXPath式(依存XPath)、依存レベル、XPath式のツリー長が記載されている。ここで、依存レベルはルートノードを0としたときのツリーの深さを意味する。また、ツリー長は、XPath式自身のツリーの長さである。連結XPath式のツリー長は、当該連結XPath式を構成するXPath式のうちでツリー長が最も長いXPath式のツリー長となっている。

図23において、例えばXPath式X10、X11は、共に依存レベルが3であり、依存レベルが最も深い。すなわち成立した場合に、より多くのXPath式に影響を与える。X10とX11とでは、X10の方がツリー長が長い。そのため、X10の優先順位が1、X11の優先順位が2となっている。また、依存レベル及びツリー長がいずれも等しいXPath式X12、X13では、/td[2]/table[1]を指すX13の方が評価対象のウェブページのDOMツリー上で後方に位置するので、X13の方が優先順位が高くなっている。

【0098】

6. 実行時判断のためのデータ構造の構築

以上のようにしてXPath式の依存関係に基づいて整理されたXPathセットを、ウェブページとのマッチングの処理において参照可能なデータ構造(例えば有効グラフ構造)に変換し、所定の記憶装置(例えば図1に示したメインメモリ103)に保存する。

【0099】

以上のようにして、XPathセットの評価を効率的に行うためのデータ構造が生成された後、評価実行部1540が、このデータ構造を用いて、文書格納部1550に格納されているウェブページに対する評価を行う。

上述したように、この評価実行部1540は、XPath式の評価の回数が少ない方がシステム全体のパフォーマンスを向上させることができるタイプ(タイプ1)であるため、優先順位の高いXPath式から順に、次の手順で評価を行う。

1. 未処理の X P a t h 式の中で優先順位の最も高い X P a t h 式を選択し、ウェブページに対する評価を行う。

2. 評価対象の X P a t h 式がマッチした場合、当該 X P a t h 式に依存している他の X P a t h 式（優先順位の低い X P a t h 式）はマッチすることが保証されるため評価から除外する。

一方、評価対象の X P a t h 式がマッチしなかった場合、当該 X P a t h 式を含む X P a t h セットは、当該ウェブページにマッチしないため、当該 X P a t h セットに含まれる全ての X P a t h 式を評価対象から除外する。

3. 各 X P a t h セットに固有の X P a t h 式（specific XPath）の全てについて処理がおわった段階で、全ての X P a t h 式のマッチしている X P a t h セットは、当該ウェブページに完全にマッチする可能性がある。そこで、そのような X P a t h セットが存在するかどうかを調べる。

4. そのような X P a t h セットが存在するならば、当該 X P a t h セットに関して、他の X P a t h セットと共有する X P a t h 式の評価を行う。

一方、そのような X P a t h セットが存在しない場合は、当該ウェブページに完全にマッチする X P a t h セットは存在しないため、他の X P a t h セットと共有する X P a t h 式の評価を行わずに処理を終了する。

本実施の形態による X P a t h セットの評価処理を、複数のアノテーションの中からウェブページに付加するアノテーションを決定する処理に用いる場合、上記の手順で全ての X P a t h 式がマッチすると判断された X P a t h セットを持つアノテーションが、ウェブページに付加できるアノテーションとして決定されることとなる。

【 0 1 0 0 】

次に、評価実行部 1 5 4 0 の種類がタイプ 2 である場合について説明する。

この場合、X P a t h 式を分割して評価を行ってもシステム全体のパフォーマンスが大きく低下しないため、データ構造生成部 1 5 2 0 は、生成するデータ構造において X P a t h 式を分割し、個々の X P a t h 式を単純化する。具体的には、

1. X P a t h 式の共通性に基づく X P a t h 式の絞込み
2. X P a t h 式の依存関係に基づく連結 X P a t h 式の簡略化
3. X P a t h 式の依存関係に基づく X P a t h 式の省略
4. X P a t h 式の分割
5. X P a t h 式のツリー長と D O M ツリー上での位置に基づく評価処理の優先順位付け
6. 実行時判断のためのデータ構造の構築

という処理が行われる。上記の処理のうち、1 ~ 3 の処理は、評価実行部 1 5 4 0 の種類がタイプ 1 の場合と同様であるので、説明を省略する。

【 0 1 0 1 】

4. X P a t h 式の分割

X P a t h 評価関数の呼び出しオーバーヘッドが小さいタイプ 2 の評価実行部 1 5 4 0 を用いる場合、データ構造生成部 1 5 2 0 は、複数の X P a t h 式で共通する部分と固有の部分とを細分化してそれぞれを評価することにより、処理の高速化を図ることができる。例えば、次の 2 つの X P a t h 式（X P a t h 1、X P a t h 2）は、/html[1]/body[1]/table[1]/tr[1]/td[1]の部分（ノード）を共有している。

X P a t h 1 /html[1]/body[1]/table[1]/tr[1]/td[1]/font[1]

X P a t h 2 /html[1]/body[1]/table[1]/tr[1]/td[1]/b[1]

したがって、この共通ノードの評価を行い、そのノードから相対パスを font[1] 及び b[1] だけ評価することにより、高速な処理を実現できる。ここで、/html[1]/body[1]/table[1]/tr[1]/td[1] と /html[1]/body[1]/table[1]/tr[1]/td[1]/font[1] との関係性を「包含関係」と呼び、「/html[1]/body[1]/table[1]/tr[1]/td[1]/font[1] が /html[1]/body[1]/table[1]/tr[1]/td[1] を包含している」とする。この包含関係に基づいて、X P a t h 式を共通部分と固有の部分とに分割し、それぞれを部分 X P a t h 式とする。

図 2 5 は、依存関係に基づく簡略化、省略の後、この分割処理を経て得られた X P a t h

10

20

30

40

50

式のリストを示す図表である。図中のXPath式PX1、PX2、PX3が、分割処理に伴って追加された部分XPath式である。この3つのXPath式を分割処理することにより、各XPath式が、図17等に列挙したXPath式に比して非常に簡潔になっていることがわかる。また、PX1、PX2、PX3を含めた依存関係ツリーを図26に示す。

【0102】

5. XPath式のツリー長とDOMツリー上での位置に基づく評価処理の優先順位付け評価対象であるウェブページにマッチしないXPathセットを効率よく決定しふるい落とすためには、「存在しにくいエレメント」を対象にしたXPath式を先に処理することが好ましい。そこで、次のような基準を用いて評価するXPath式の順番(優先順位)を決定する。

10

優先順位を決定する処理は、上述した評価実行部1540の種類がタイプ1である場合の処理と同様であるが、ここでは3.の処理で新たに追加されたXPath式PX1、PX2、PX3を含めて優先順位が決定される。

図27は、図26に示した依存関係ツリーの各ノード(XPath式)に、上記の規則に従って優先順位を付加した様子を示す図、図28は、図27に示した優先順位にしたがってXPath式を並べた図表である。なお、上の4.の処理で複数のXPath式における共通部分が分割されて個別のXPath式とされていることから、XPath式自身のツリー長は優先順位を決定するために参酌しない。

【0103】

20

6. 実行時判断のためのデータ構造の構築

以上のようにしてXPath式の依存関係に基づいて整理されたXPathセットを、ウェブページとのマッチングの処理において参照可能なデータ構造(例えば有効グラフ構造)に変換し、所定の記憶装置(例えば図1に示したメインメモリ103)に保存する。この処理は、上述した評価実行部1540の種類がタイプ1である場合の処理と同様である。

【0104】

以上のようにして、XPathセットの評価を効率的に行うためのデータ構造が生成された後、評価実行部1540が、このデータ構造を用いて、文書格納部1550に格納されているウェブページに対する評価を行う。

30

上述したように、この評価実行部1540は、XPath式の評価の回数が多くてもシステム全体のパフォーマンスが大きく低下しないタイプ(タイプ2)であるため、優先順位の高いXPath式から順に、次の手順で評価を行う。

1. 未処理のXPath式の中で優先順位の最も高いXPath式を選択し、ウェブページに対する評価を行う。このとき、評価対象のXPath式が包含しているXPath式の評価も行う。なお、包含しているXPath式の評価が既になされている場合(より優先順位の高いXPath式の評価に伴って評価されている場合)、その評価結果を利用する。

2. 1.で評価したXPath式がマッチした場合、当該XPath式に依存している他のXPath式(優先順位の低いXPath式)はマッチすることが保証されるため評価から除外する。

40

一方、評価対象のXPath式がマッチしなかった場合、当該XPath式を含むXPathセットは、当該ウェブページにマッチしないため、当該XPathセットに含まれる全てのXPath式を評価対象から除外する。

また、評価対象のXPath式が包含しているXPath式がマッチしなかった場合、当該XPath式を包含する他の全てのXPath式を原則として評価対象から除外する。ただし、連結XPath式については、当該連結XPath式を構成する全てのXPath式が評価されるまでは、当該連結XPath式がマッチするかどうかを判断できないので、評価対象からは除外しない。

3. 各XPathセットに固有のXPath式(specific XPath)の全てについて処理が

50

おわった段階で、全てのX P a t h式のマッチしているX P a t hセットは、当該ウェブページに完全にマッチする可能性がある。そこで、そのようなX P a t hセットが存在するかどうかを調べる。

4. そのようなX P a t hセットが存在するならば、当該X P a t hセットに関して、他のX P a t hセットと共有するX P a t h式の評価を行う。

一方、そのようなX P a t hセットが存在しない場合は、当該ウェブページに完全にマッチするX P a t hセットは存在しないため、他のX P a t hセットと共有するX P a t h式の評価を行わずに処理を終了する。

本実施の形態によるX P a t hセットの評価処理を、複数のアノテーションの中からウェブページに付加するアノテーションを決定する処理に用いる場合、上記の手順で全てのX P a t h式がマッチすると判断されたX P a t hセットを持つアノテーションが、ウェブページに付加できるアノテーションとして決定されることとなる。

【0105】

以上説明した本実施の形態は、X P a t h式的全機能に対して効率的な評価を行うためのデータ構造を生成(最適化)するものではなく、頻繁に利用される一部の機能について最適化する。

図29は、本実施の形態による最適化が可能なX P a t h式の構造を説明する図である。図29に示すX P a t h式において、E xはx番目のエレメントである。また、C xはx番目のsibling内でのポジションを示す数値であり、predicateで表現すると、[position() = Cx]と等しい。すなわち、X P a t h式を2つの部分に分けることができ、前半部分(図29のパート1)が次の条件に当てはまる場合に、本実施の形態の最適化を行うことができる。

- ・ n段目まで子孫方向に1段ずつツリーを移動している。
 - ・ 軸(axis)は子(child)のみであり、他の軸(axis)は使用されていない。
- また、descendant-or-self::node()の省略である「//」指定も使用されていない。
- ・ predicateは[position() = 番号] (省略形は[番号])のみが使用可能。
 - ・ predicateの省略(/tag/)は無し。

後半部分(図29のパート2)には、どのようなX P a t h式が出現してもかまわないが、これらの部分は最適化されない。

言い換えれば、本実施の形態は、ウェブページを記述したHTMLドキュメントのようにツリー長深くかつ不定形であるツリーのノードを指すX P a t h式には、上記のような構造(パート1、パート2に分けることができるという構造)を持ったものが多いという性質を利用している。

【0106】

なお、本実施の形態では、X P a t hセット内の全てのX P a t h式が成立することをX P a t hセット成立の条件として説明した。しかし、このような厳密なマッチングだけではなく、「最もマッチするものを必ず1つ」選択するという課題も存在する。例えば、ウェブページにアノテーションを付加するシステムにおいて所定のウェブページに付加すべきアノテーションを選択する場合、何らかの評価基準を用いてX P a t h式の適合度を算出し、最も適しているものを選択することによって、全くトランスコードできないという状況を回避することができると考えられる。このような用途のために、本実施の形態を拡張し、X P a t hセットの成立を判断する際に曖昧性を許容するアルゴリズムを導入することができる。

【0107】

このアルゴリズムの具体的な手順について説明する。

1. 評価実行部1540によるマッチングの実行時に、所定のX P a t h式がマッチしなかった場合、そのX P a t h式に対してマッチしなかったことを示すフラグ(UNMATCHフラグ)を立てる。

2. マッチしないX P a t h式と同じX P a t hセットに含まれることによって評価され

10

20

30

40

50

なかったX P a t h式に対して、評価が行われていないこと示すフラグ (UNEVALUATEDフラグ) を立てる。

3. 全てのX P a t hセットがマッチしないと評価された段階で、UNEVALUATEDフラグのたっている全てのX P a t h式を評価する。(ただし、既に依存関係や包含関係によって評価結果の判明しているものは除外する) マッチしなかったX P a t h式にはUNMATCHフラグを立てる。

【0108】

X P a t hセットの成立を判断する際に曖昧性を許容するアルゴリズムを導入する場合、適合度の評価手法に関わらず、ここまでの処理は共通である。しかし、これ以降の処理は、適合度の評価手法に応じて異なる処理が行われる。ここでは一例として、「最もマッチ

10

しないX P a t h式の数の少ないX P a t hセット」であり、かつ「DOMツリー上の距離で定義されるX P a t h式の適合度が最も大きいもの」を選択するアルゴリズムを示す。

4. UNMATCHフラグの立っているX P a t h式の最も少ないX P a t hセットを選択する。

5. UNMATCHフラグの立っているX P a t h式の数が同数であるX P a t hセットが複数存在する場合、各X P a t hセットに含まれているX P a t h式の適合度を全て算出する。例えば、適合度Aは次のように算出することができる

$$A X P a t h = k 1 * P + k 2 * S - k 3 * L$$

ここで、各パラメータは以下の通りである。

k 1、k 2、k 3：定数係数

20

P：着目中のX P a t h式の親シーケンス (parent sequence) を辿ったときに、存在するノードに至るまでのステップ数。最悪でもbodyノードは共通であるため、必ず値を持つ。

S：着目中のX P a t h式の兄弟シーケンス (sibling sequence) の距離。(親ノード) / t r [4] が存在せず (親ノード) / t r [3] が存在する場合距離を1、(親ノード) / t r [2] が存在する場合2と定義する。兄弟ノードが存在しない場合、Pの対象となる存在する親方向のノードにおける子ノード内での兄弟ノードを算出。対象となるノードが元々ファーストノード (ノタグ [1]) である場合、この値には一定のデフォルト値を割り当てる。

L：着目中のX P a t h式が指しているノードのツリー長。この距離が長い場合、Pの値などで不利になることが考えられるため、負の係数によって相殺する。

30

以上のような計算によって各X P a t h式の適合度が算出され、最終的に得られた各X P a t h式の適合度の総和をX P a t hセットの適合度とする。

$$A X P a t h \text{セット} = A X p a t h$$

最後に得られたAの値が最も小さいX P a t hセットを選択することで最も適合するX P a t hセットを選択することができる。

【0109】

ただし、このような拡張を行った場合、X P a t hセットが成立する場合の計算量に変化はないが、成立するX P a t hセットが存在しなかった場合の計算量が大幅に増加する可能性がある。そのため、全体のパフォーマンスは、成立するX P a t hセットが存在しない

40

入力の場合に応じて低下してしまう。このような事態を防ぐためには、適合度算出のアルゴリズムを単純化し、高速に算出可能な「曖昧性」にとどめることが必要と考えられる。

【0110】

以上のように、第1、第2の実施の形態 (以下、本実施の形態) では、多数のX P a t h式の評価結果を効率よく求めることによって、XML文書に対する処理の速度を向上させることができる。例えば、X S L Tプロセッサに本実施の形態を組み込むことによって、スタイルシートに多数のX P a t h式が含まれる場合の処理速度の向上が可能になる。

また、X P a t hを、XML文書に対して特定の部分を持つかどうかのチェックをするために使うことによって、利用者毎に指定した配信条件でのニュース送信や、コンテンツ毎

50

に指定した条件でのコンテンツと利用者とのマッチングを、XMLをベースとして行うことができる。この場合、従来はXPath式の数の増加に伴って処理時間が増大するので、大規模なサービスを行うことができなかつたが、本実施の形態によって、処理時間の増大を抑えることができるため、かかるサービスを実現することが可能になる。

【0111】

さらにまた、ネットワーク上のウェブページ（HTML文書）を所定のアノテーションに基づいてトランスコーディングするシステムにおいて、予め用意された多くのアノテーションパターンの中からトランスコーディングしようとするウェブページに適用可能なアノテーション（当該ウェブページ中のエレメントを正しく指すXPathセットを持ったアノテーション）を検索する場合に、本実施の形態を用いることにより、XPathセットの10
 の評価に要する時間を短縮し、高速にアノテーションの検索処理を実行することが可能となる。

【0112】

なお、本実施の形態によるXML文書処理システムは、複数のXPath式があるときに、それらの中に類似なものや重複するもの、あるいは依存関係によって他のXPath式の評価結果によっては評価を要しないものといった冗長な部分が含まれることを利用して効率が向上させている。したがって、そのような冗長な部分がない場合には効率が向上しないこととなる。

しかし、この場合でも、ロケーション・パスの計算はノードセットから別のノードセットを再帰的に求めることで実現されるのでXalanなどをそのまま使った従来の方式と同じ手20
 続きで処理されることになる。また、1つの比較演算子はエントリ数1のテーブル検索で実現されるが、両者ともにノードセットの再構成に比べると無視できるぐらいに高速に実行できる。さらに、述部が決められたパターンでない場合には最適化の対象にならないが、その場合のテストはXalanなどをそのまま使った従来の方式と同じ手続きで計算を行えばよい。

したがって、冗長な部分がなく最適化できないXPath式の集合でも、従来の方式とほぼ同等の効率でXPath式の評価を行うことができる。そして、最適化可能なもの（相互に類似するXPath式）が増えるにしたがって処理の実行効率が向上することとなる。

【0113】

【発明の効果】

以上説明したように、本発明によれば、XML文書に対して複数のXPath式を評価する場合の実行効率を向上させることができるという効果がある。この効果は、当該複数のXPath式の中に冗長な部分が増えれば増えるほど顕著となる。

【図面の簡単な説明】

【図1】 第1の実施の形態によるXML文書処理システムを実現するのに好適なコンピュータ装置のハードウェア構成の例を模式的に示した図である。

【図2】 第1の実施の形態によるXML文書処理システムの機能を説明する図である。

【図3】 第1の実施の形態のデータ構造生成部にて生成されるデータ構造を示す図である。40

【図4】 第1の実施の形態の処理対象であるXML文書の例を示す図である。

【図5】 幅優先で複数のXPath式を評価する場合の評価実行部による処理手順を説明するフローチャートである。

【図6】 第1の実施の形態で評価されるXPath式P7、P8、P9を示す有向非循環グラフを示す図である。

【図7】 第1の実施の形態においてハッシュテーブルを用いた評価処理の流れを説明するフローチャートである。

【図8】 第1の実施の形態のデータ構造生成部によるデータ構造の生成方法を説明するフローチャートである。

【図9】 図3のデータ構造にXPath式P4が追加された状態のデータ構造を示す図 50

である。

【図10】 データ構造から所定のXPath式を除去する際のデータ構造生成部の処理を説明するフローチャートである。

【図11】 図9のデータ構造からXPath式P3を除去した状態のデータ構造を示す図である。

【図12】 XPath式P11、P12、P13、P14を表現するデータ構造を示す図である。

【図13】 処理対象であるXML文書を示す図である。

【図14】 実施例におけるXPath式の評価に要した実行時間を示す図である。

【図15】 第2の実施の形態によるウェブページ処理システムの機能を説明するブロック図である。

10

【図16】 XPathセットとこれに含まれるXPath式とを一覧表示した図表である。

【図17】 図16に示したXPath式のうち重複するものを整理して識別情報(XPathID)を割り当てた様子を示す図である。

【図18】 図17に示したXPath式のうち複数のXPathセットに共通するXPath式を抽出した様子をグラフ構造で示した図である。

【図19】 連結XPath式を展開後のXPath式の依存関係をツリー構造で示す図である。

【図20】 図19に示したXPath式の依存関係に基づいて連結XPath式を簡略化した様子を示す図である。

20

【図21】 図20に示したXPath式の依存関係に基づいてXPath式を省略した様子を示す図である。

【図22】 図21に示したXPath式の依存関係ツリー中のXPath式に対して再結合可能なものを再結合する様子を示す図である。

【図23】 図22に示した依存関係ツリーの各ノード(XPath式)に優先順位を付加した様子を示す図である。

【図24】 図23に示した優先順位に従ってXPath式を並べた図表である。

【図25】 依存関係に基づく簡略化、省略の後、この分割処理を経て得られたXPath式のリストを示す図表である。

30

【図26】 図25に示したPX1、PX2、PX3を含めた依存関係ツリーを示す図である。

【図27】 図26に示した依存関係ツリーの各ノード(XPath式)に、上記の規則に従って優先順位を付加した様子を示す図である。

【図28】 図27に示した優先順位にしたがってXPath式を並べた図表である。

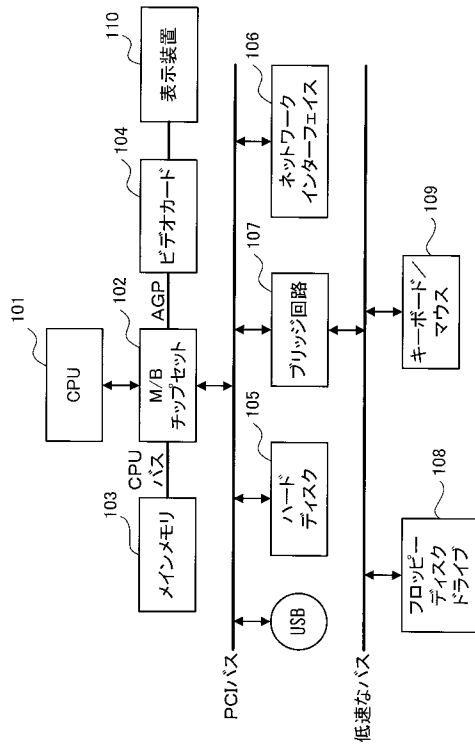
【図29】 第2の実施の形態による最適化が可能なXPath式の構造を説明する図である。

【符号の説明】

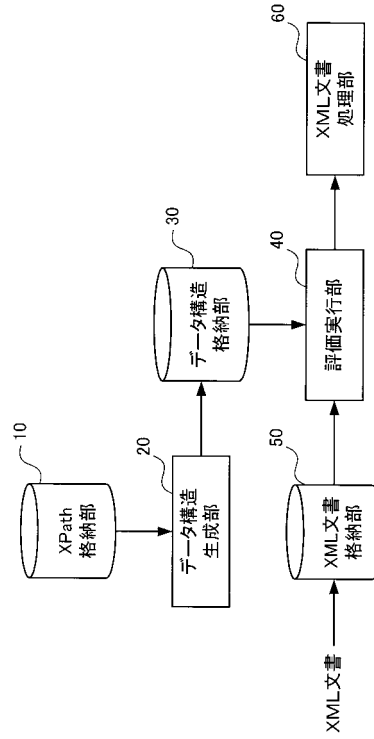
10、1510...XPath格納部、20、1520...データ構造生成部、30、1530...データ構造格納部、40、1540...評価実行部、50...XML文書格納部、60...XML文書処理部、101...CPU、103...メインメモリ、105...ハードディスク、106...ネットワークインターフェイス、109...キーボード/マウス、1550...文書格納部、1560...文書処理部

40

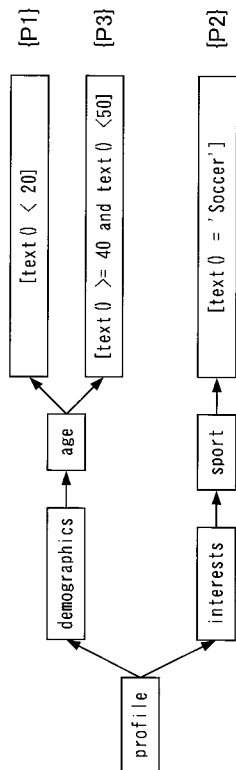
【 図 1 】



【 図 2 】



【 図 3 】



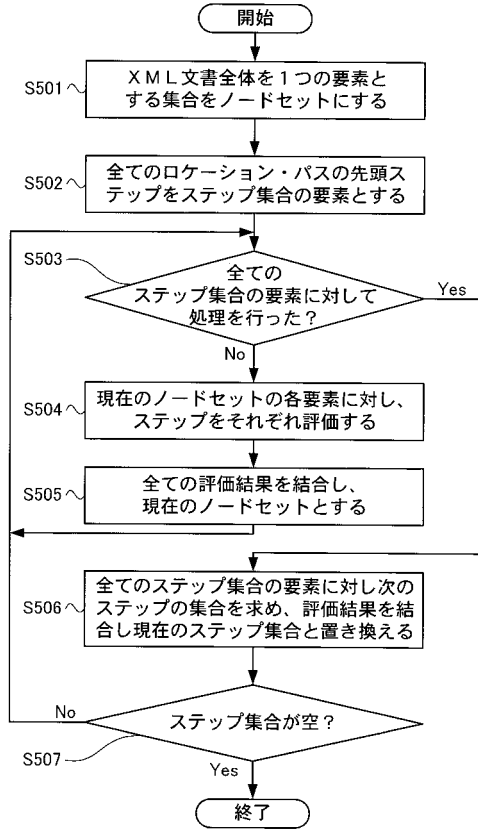
【 図 4 】

```

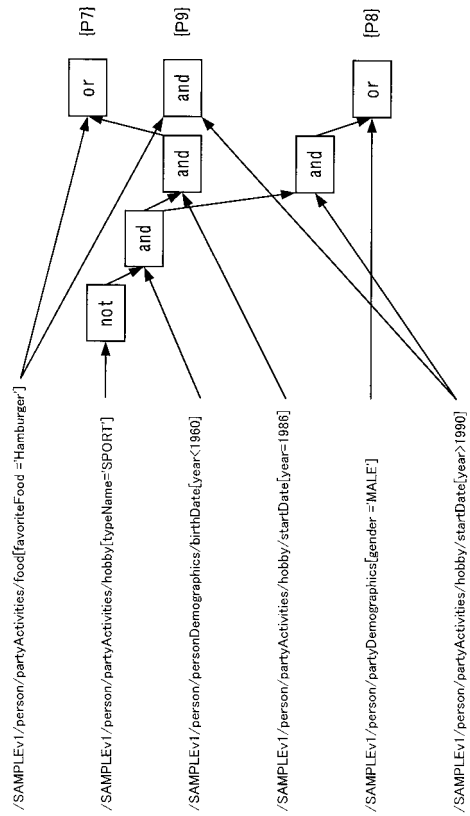
D1:
<profile>
  <name>Alan</name>
  <demographics>
    <age>35</age>
  </demographics>
  <location>
    <city>Osaka</city>
  </location>
  <interests>
    <sport>Soccer</sport>
    <music>Classical</music>
    <book>History</book>
  </interests>
</profile>

```

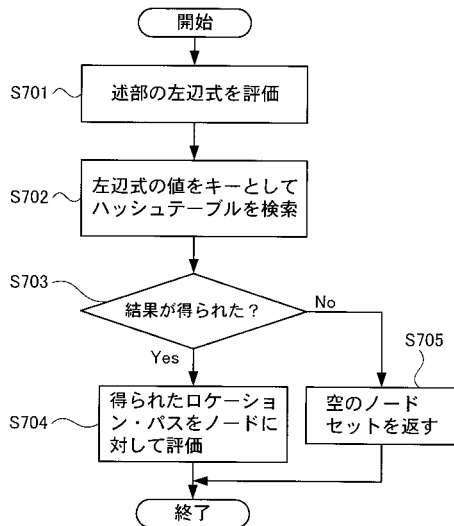
【図5】



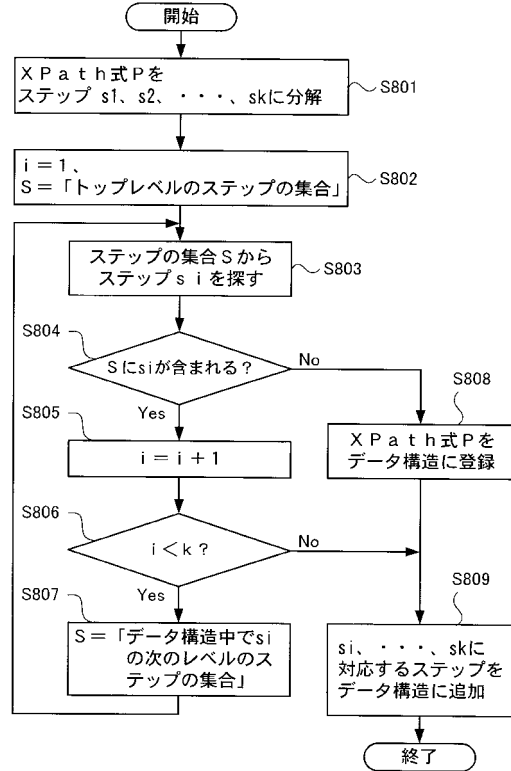
【図6】



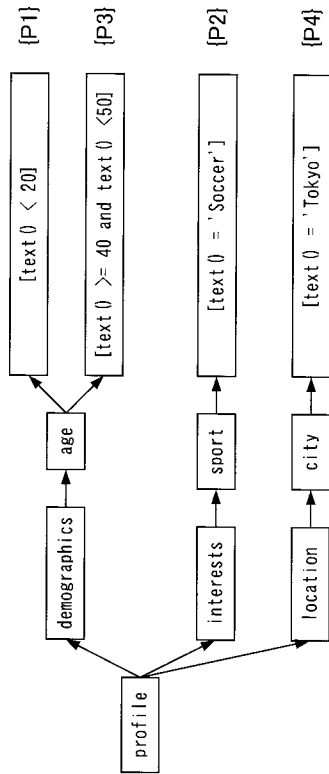
【図7】



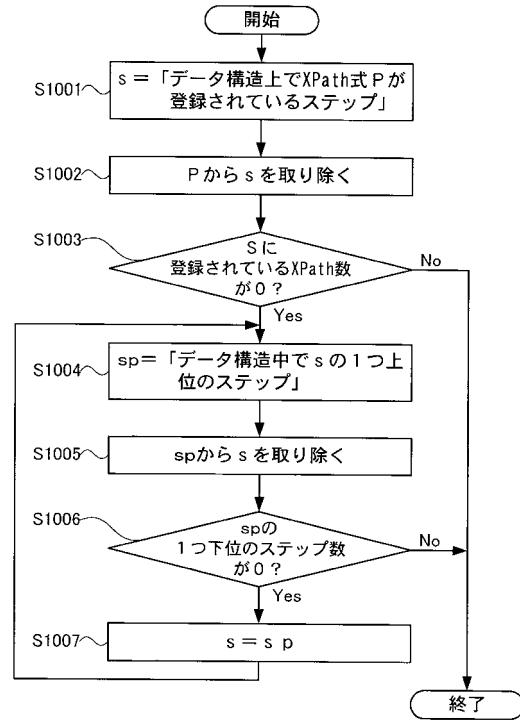
【図8】



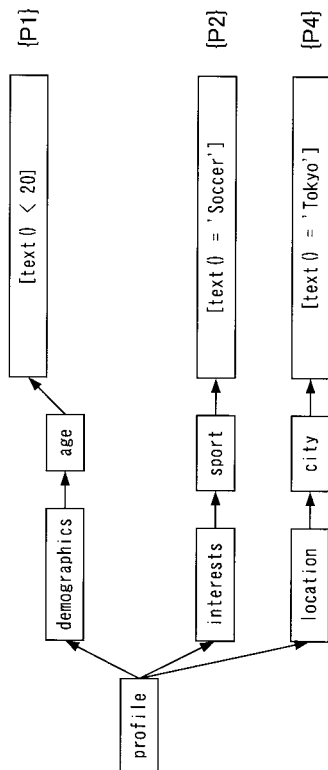
【図9】



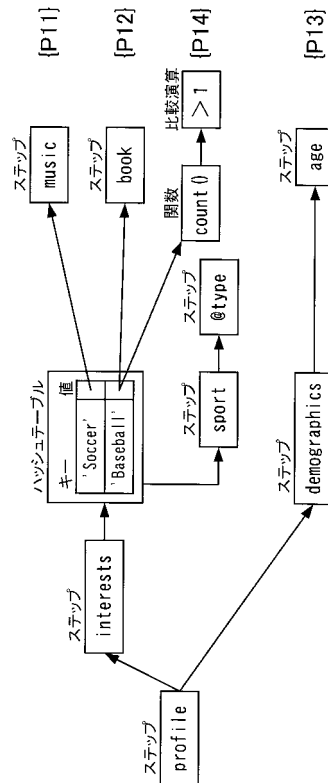
【図10】



【図11】



【図12】



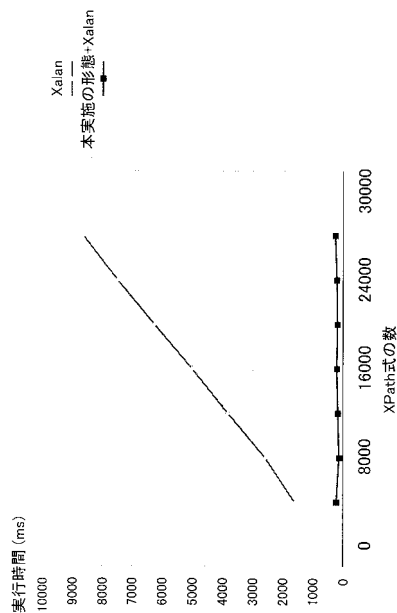
【 図 1 3 】

```

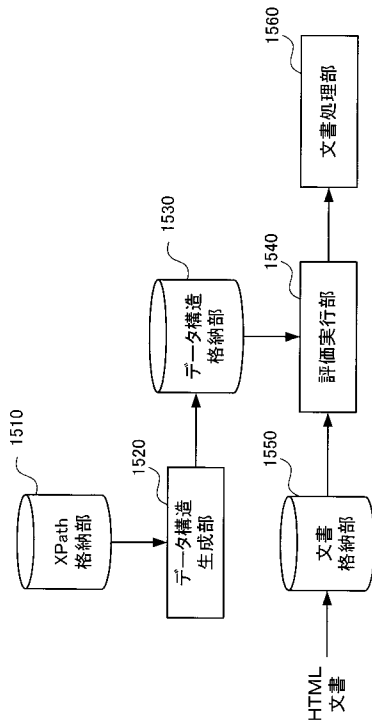
<profile>
  <demographics>
    <age>19</age>
  </demographics>
  <interests>
    <sport type@=' Baseball' />
    <book>History</book>
  </interests>
</profile>

```

【 図 1 4 】



【 図 1 5 】



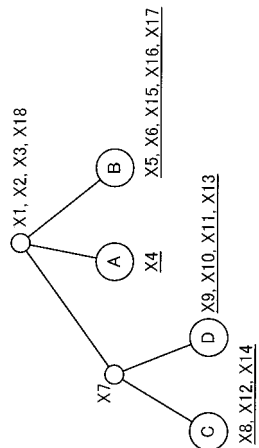
【 図 1 6 】

XPath ID	XPath セット	XPath
1	A	/html/body/table/body/tr/html/body/table/body/tr
2		/html/body/table/body/tr/html/body/table/body/tr
3		/html/body/table/body/tr/html/body/table/body/tr
4		/html/body/table/body/tr/html/body/table/body/tr
5		/html/body/table/body/tr/html/body/table/body/tr
6	B	/html/body/table/body/tr/html/body/table/body/tr
7		/html/body/table/body/tr/html/body/table/body/tr
8		/html/body/table/body/tr/html/body/table/body/tr
9		/html/body/table/body/tr/html/body/table/body/tr
10		/html/body/table/body/tr/html/body/table/body/tr
11		/html/body/table/body/tr/html/body/table/body/tr
12		/html/body/table/body/tr/html/body/table/body/tr
13		/html/body/table/body/tr/html/body/table/body/tr
14		/html/body/table/body/tr/html/body/table/body/tr
15		C
16	/html/body/table/body/tr/html/body/table/body/tr	
17	/html/body/table/body/tr/html/body/table/body/tr	
18	/html/body/table/body/tr/html/body/table/body/tr	
19	/html/body/table/body/tr/html/body/table/body/tr	
20	/html/body/table/body/tr/html/body/table/body/tr	
21	/html/body/table/body/tr/html/body/table/body/tr	
22	/html/body/table/body/tr/html/body/table/body/tr	
23	D	/html/body/table/body/tr/html/body/table/body/tr
24		/html/body/table/body/tr/html/body/table/body/tr
25		/html/body/table/body/tr/html/body/table/body/tr
26		/html/body/table/body/tr/html/body/table/body/tr
27		/html/body/table/body/tr/html/body/table/body/tr
28		/html/body/table/body/tr/html/body/table/body/tr
29		/html/body/table/body/tr/html/body/table/body/tr
30		/html/body/table/body/tr/html/body/table/body/tr
31		/html/body/table/body/tr/html/body/table/body/tr

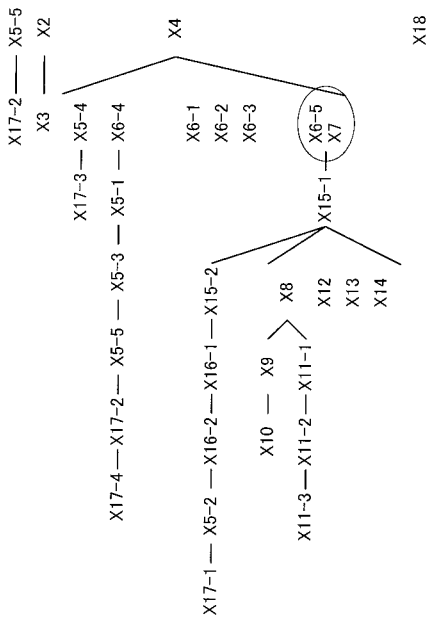
【 図 17 】

XPath ID	XPath
X1	/html[1]/body[1]/table[1]/tbody[1]/tr[1]/td[1]/table[1]/tbody[1]/tr[2]
X2	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[1]/table[1]
X3	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[1]/table[1]
X4	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[1]
X5	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/font[6]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/font[1]/tbody[1]/tr[1]/td[3]/font[8]
X6	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/script[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/font[1]/tbody[1]/tr[1]/td[3]/font[1]
X7	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[1]
X8	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[1]
X9	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[1]/tbody[1]/tr[1]/td[1]/table[2]
X10	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[1]/tbody[1]/tr[1]/td[1]/table[2]
X11	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[1]/tbody[1]/tr[1]/td[1]/table[2]/tbody[1]/tr[1]/td[1]
X12	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[1]
X13	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[1]
X14	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[3]
X15	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[3]
X16	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[5]
X17	/html[1]/body[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/font[1]/tbody[1]/tr[1]/td[3]/font[2]
X18	/html[1]/body[1]/table[3]

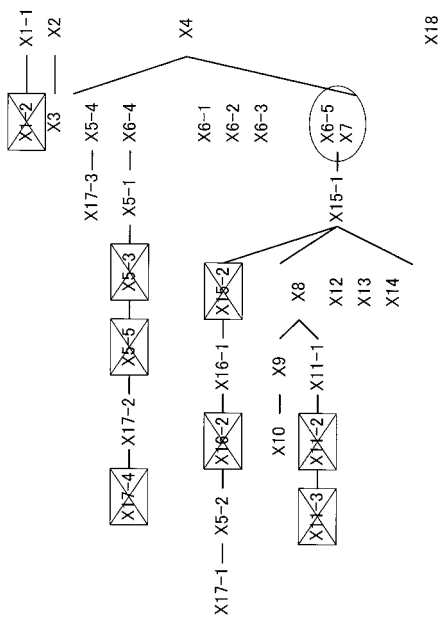
【 図 18 】



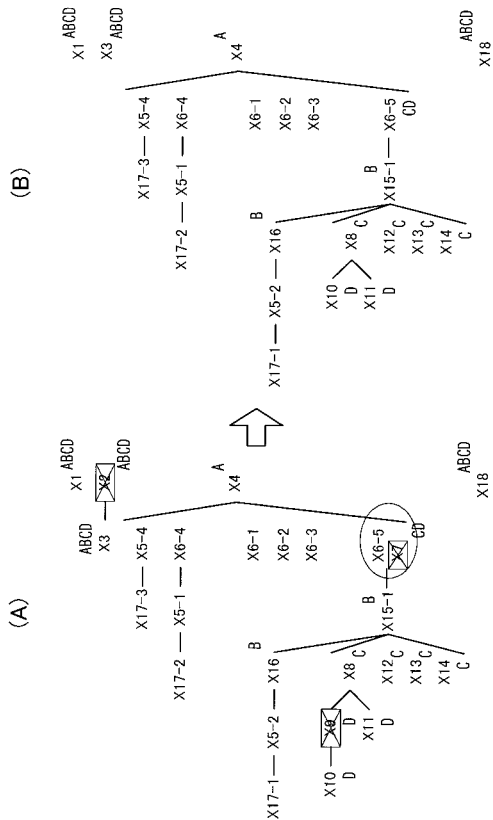
【 図 19 】



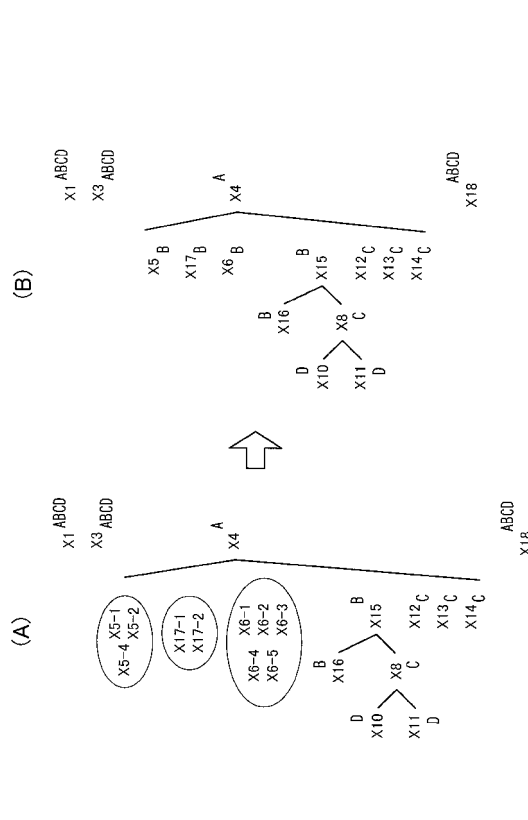
【 図 20 】



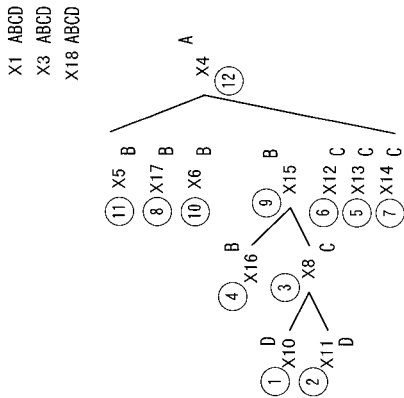
【 図 2 1 】



【 図 2 2 】



【 図 2 3 】



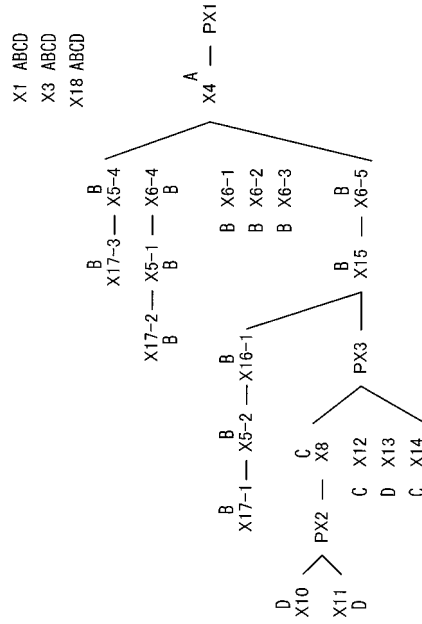
【 図 2 4 】

優先順位	XPath ID	依存 XPath	XPath	依存レベル	列長
1	X10	X8	/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[1]/tbody[1]/tr[1]/td[1]/table[2]	3	15
2	X11	X8	/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/tbody[2]/tbody[1]/tr[1]/td[1]/table[1]/tbody[1]/tr[1]	3	13
3	X8	X4	/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[2]/table[1]	1	11
4	X16	X4	/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[1]/table[2]	1	11
5	X13	X4	/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/tbody[2]/tbody[1]/tr[1]/td[1]/table[1]	2	11
6	X12	X4	/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/table[2]/tbody[1]/tr[1]/td[3]	1	10
7	X14		/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/tbody[2]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]	1	7
8	X17	X4	/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/tbody[4]	2	7
9	X15	X4	/html[1]/tbody[1]/table[2]/tbody[1]/tr[1]/td[3]/tbody[1]/tbody[1]/tbody[1]/tbody[1]	1	7
10	X6	X4	/html[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]	1	7
11	X5		/html[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]	1	7
12	X4	X1	/html[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]	0	6
	X3	X1	/html[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]	0	5
	X18	X18	/html[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]/tbody[1]	0	11
			/html[1]/tbody[1]/tbody[1]	0	13

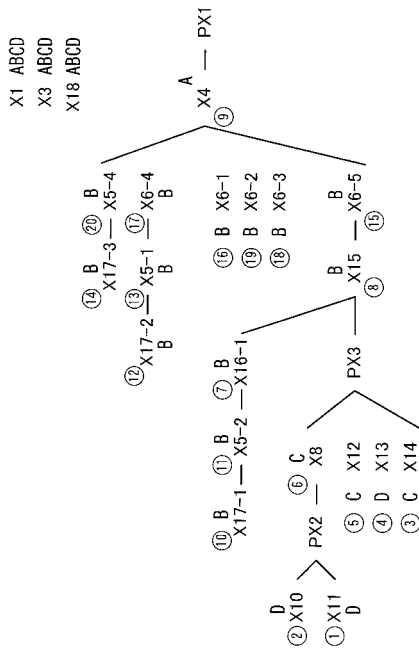
【 図 2 5 】

XPath ID	依存XPath	XPath
PX1		/body/table[2]/tbody[1]/tr[1]
PX2	X8	[X8]/tbody[1]
PX3		[X15-1]/tbody[1]/tr[1]
X10	PX2	{PX2}/tr[1]/td[1]/table[2]
X1-1		/body/table[1]/tbody[1]/tr[1]
X11-1	PX2	{PX2}/tr[2]
X12	PX3	{PX3}/td[1]/table[2]
X13	PX3	{PX3}/td[2]/table[1]
X14	PX3	{PX3}/td[3]
X15-1	X6-5	[X4]/table[2]
X16	X15-1	[X4]/table[4]
X17-1	X6-2	[X4]/table[7]
X17-2	X5-1	[X4]/font[11]
X17-3	X4	[X4]/a[2]
X18		/body/table[3]
X3		{PX1}/td[1]/table[1]/tbody[1]/tr[1]/td[1]/table[1]
X4		{PX1}/td[3]
X5-1	X6-4	[X4]/font[6]
X5-2	X16	[X4]/table[6]
X5-4	X4	[X4]/a[1]
X6-1	X4	[X4]/script[1]
X6-2	X4	[X4]/b[1]
X6-3	X4	[X4]/br[1]
X6-4	X4	[X4]/font[1]
X6-5	X4	[X4]/table[1]
X8	PX3	{PX3}/td[1]/table[1]

【 図 2 6 】



【 図 2 7 】



【 図 2 8 】

優先順位	XPath ID	依存XPath	XPath	依存レベル
1	X11	PX2	{PX2}/tr[2]	7
2	X10	PX2	{PX2}/tr[1]/td[1]/table[2]	7
3	X14	PX3	{PX3}/td[3]	5
4	X13	PX3	{PX3}/td[2]/table[1]	5
5	X12	PX3	{PX3}/td[1]/table[2]	5
6	X8	PX3	{PX3}/td[1]/table[1]	5
7	X16	X15	[X4]/table[4]	4
8	X15	X6-5	[X4]/table[2]	4
9	X4		{PX1}/td[3]	1
10	X17-1	X5-2	[X4]/table[7]	6
11	X5-2	X16	[X4]/table[6]	5
12	X17-2	X5-1	[X4]/font[11]	4
13	X5-1	X4	[X4]/font[6]	4
14	X17-3	X4	[X4]/a[2]	3
15	X6-5	X4	[X4]/table[1]	2
16	X6-1	X4	[X4]/script[1]	2
17	X6-4	X4	[X4]/font[1]	2
18	X6-3	X4	[X4]/br[1]	2
19	X6-2	X4	[X4]/b[1]	2
20	X5-4	X4	[X4]/a[1]	2
分割 XPath	PX3	X8	[X15-1]/tbody[1]/tr[1]	
	PX2	X8	/body/table[2]/tbody[1]/tr[1]	
必須 XPath	X1		/body/table[1]/tbody[1]/tr[1]	
	X18		/body/table[3]	
	X3		{PX1}/td[1]/table[1]/tbody[1]/tr[1]/td[1]/table[1]	

【 図 29 】

/E1 [C1]/E2 [C2]/E3 [C3].../En [Cn] [predicates]/any XPath expression

バート 1

バート 2

フロントページの続き

- (72)発明者 中村 宏明
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内
- (72)発明者 百合山 まどか
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内
- (72)発明者 高木 啓伸
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内

審査官 高瀬 勤

- (56)参考文献 特開2000-293191(JP,A)
Aaron Skonnard, XPathを使ってInfosetにアクセスする, Microsoft Developer Network Magazine日本語版, 日本, 株式会社アスキー, 2000年 8月18日, 第5号, 第149~156頁
村上隆一, Mastering XML 第15回, Java WORLD, 日本, (株)IDGジャパン, 2001年 5月1日, 第5巻, 第5号, 第179-183頁

(58)調査した分野(Int.Cl., DB名)

G06F 12/00
G06F 17/21-25
G06F 17/30