

(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2010-267092

(P2010-267092A)

(43) 公開日 平成22年11月25日(2010.11.25)

(51) Int.Cl.  
G06F 17/21 (2006.01)

F I  
G06F 17/21 501T

テーマコード(参考)  
5B109

審査請求 未請求 請求項の数 8 O L (全 13 頁)

(21) 出願番号 特願2009-118047 (P2009-118047)  
(22) 出願日 平成21年5月14日 (2009.5.14)

(71) 出願人 000001007  
キヤノン株式会社  
東京都大田区下丸子3丁目30番2号  
(74) 代理人 100076428  
弁理士 大塚 康徳  
(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

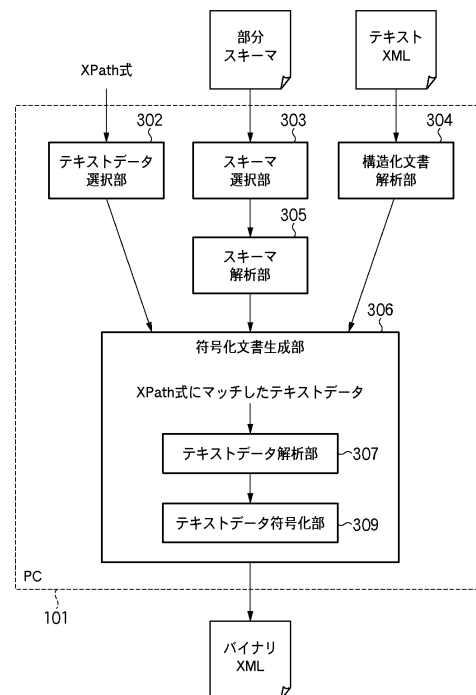
(54) 【発明の名称】 情報処理装置、情報処理方法

(57) 【要約】 (修正有)

【課題】 サイズの大きなテキストデータを属性値又は要素内容として含むテキストXMLからバイナリXMLへの変換効率を改善するための技術を提供する。

【解決手段】 テキスト形式の構造化文書を取得する手段と、前記構造化文書を解析し、テキストデータを検出する手段と、前記テキストデータを子要素に変換することで前記構造化文書を変換し、変換した構造化文書に基づいてバイナリ形式の構造化文書を生成する生成手段とを備えることを特徴とする情報処理装置。

【選択図】 図4



PC 101

**【特許請求の範囲】****【請求項 1】**

テキスト形式の構造化文書を取得する手段と、  
前記構造化文書を解析し、テキストデータを検出する手段と、  
前記テキストデータを子要素に変換することで前記構造化文書を変換し、変換した構造化文書に基づいてバイナリ形式の構造化文書を生成する生成手段と  
を備えることを特徴とする情報処理装置。

**【請求項 2】**

前記テキストデータは、前記構造化文書に含まれる属性値であることを特徴とする請求項 1 に記載の情報処理装置。

**【請求項 3】**

前記テキストデータは、前記構造化文書に含まれる要素内容であることを特徴とする請求項 1 に記載の情報処理装置。

**【請求項 4】**

前記生成手段は、  
前記テキストデータの構造を定義するスキーマを取得する手段と、  
前記スキーマを用いて前記テキストデータを子要素に変換することで前記構造化文書を変換する変換手段と  
を備えることを特徴とする請求項 1 乃至 3 の何れか 1 項に記載の情報処理装置。

**【請求項 5】**

前記変換手段は、  
前記テキストデータにおいて、前記スキーマのgroup要素以下で定義された列挙値に該当する変数については、当該列挙値をタグ名とする開始タグに変換する手段と、  
前記テキストデータにおいて、前記変数に後続する値については、前記開始タグに後続するデータとする手段と  
を備えることを特徴とする請求項 4 に記載の情報処理装置。

**【請求項 6】**

テキスト形式の構造化文書を取得する工程と、  
前記構造化文書を解析し、テキストデータを検出する工程と、  
前記テキストデータを子要素に変換することで前記構造化文書を変換し、変換した構造化文書に基づいてバイナリ形式の構造化文書を生成する生成工程と  
を備えることを特徴とする情報処理方法。

**【請求項 7】**

コンピュータを請求項 1 乃至 5 の何れか 1 項に記載の各手段として機能させるためのコンピュータプログラム。

**【請求項 8】**

請求項 7 に記載のコンピュータプログラムを格納した、コンピュータ読み取り可能な記憶媒体。

**【発明の詳細な説明】****【技術分野】****【0001】**

本発明は、構造化文書进行处理する技術に関するものである。

**【背景技術】****【0002】**

これまで、XMLデータのフォーマットは、一般的にテキスト形式であった。しかし、XMLデータは、データ構造を表現するために冗長なデータを必要とし、コンピュータが読み書きするのに時間がかかるという問題があるため、近年では、バイナリXML技術が注目されている。

**【0003】**

例えば、W3Cの標準仕様としてEfficient XML Interchange (<http://www.w3.org/XML>

10

20

30

40

50

/EXI/)がある。また、ISOの標準仕様としては、Fast Infoset(<http://www.iso.org/iso/en/CatalogueDetailPage.CatalogueDetail?CSNUMBER=41327&scopeList=PROGRAMME>)がある。

【0004】

これらの技術は、XMLデータに含まれる要素名や属性名などの各ボキャブラリを、XMLデータ内での出現順に番号を振って符号化することによって、データサイズを小さくすることが可能であった。符号とボキャブラリとの対応を示す表は、符号化テーブルと呼ばれる。バイナリXML形式のデータは、従来のテキスト形式のXMLデータに比べて、高速に読み書きすることが可能であった。

【0005】

また、属性値や要素内容を従来のような文字列型のデータではなく、整数型、浮動小数点型や、独自形式の圧縮アルゴリズムを用いてエンコードすることにより、データサイズを更に小さくすることが可能であった。また、それらデータを解析するには、従来のようなアプリケーションによる文字列データから数値データへの変換処理の必要がなくなるので、解析にかかる処理時間を更に短くすることができた。

【0006】

また、XMLデータの文書構造と、含まれる属性値や要素内容の値のデータ型を、スキーマを用いることによって定義することができた。一般的には、それらスキーマは、XMLデータの文書構造と、値のデータ型が厳密にスキーマの定義内容と合っているかどうかを検証するために使用される。W3Cのスキーマの標準仕様としては、XML Schema(<http://www.w3.org/XML/Schema>)がある。また、ISOのスキーマの標準仕様としては、RELAX NG(<http://www.oasis-open.org/committees/relax-ng/>)がある。

【0007】

また、スキーマを利用した他の技術としては、データバインディングがある。これは、開発時に、スキーマからXMLデータを格納するクラスファイルを自動生成しておき、実行時には、アプリケーションがそのクラスを利用することによって、アプリケーションの開発コストを小さくするものである。スキーマから生成したクラスにXMLデータをバインディングするには、スキーマで定義した各データ型に基づいて、含まれるテキストデータを変換し、各メンバに格納する。よって、従来のテキストベースのXMLパーサとは異なり、アプリケーションはスキーマから生成したクラスを用いることによって、XMLデータに含まれるテキストデータを各データ型の値として扱うことが可能になる。一般的にはSun MicrosystemsのJAXB(<https://jaxb.dev.java.net/>)や、Relaxer(<http://www.relaxer.jp/>)が利用されている。

【0008】

また、KDDIのXEUSは、属性値や要素内容において、カンマなどで区切られた数値の配列に対しては、各数値をそれぞれ符号化することによって、XMLデータを効率的に圧縮することが可能であった。

【0009】

また、スキーマ情報を用いることによって、文書構造の異なる複数のXMLデータ間の変換用スタイルシートを自動生成するものがあった(特許文献1)。また、アプリケーションからXMLデータへのアクセス効率を改善するために、XMLデータに含まれる階層構造が浅くなるように、タグ名を変換するものがあった(特許文献2)。

【先行技術文献】

【特許文献】

【0010】

【特許文献1】特開2005-352945号公報

【特許文献2】特開2002-297569号公報

【発明の概要】

【発明が解決しようとする課題】

【0011】

10

20

30

40

50

従来のバイナリXMLでは、属性値や要素内容などに大きなサイズのテキストデータが含まれている場合、バイナリXML形式に符号化してもそれらの圧縮効率が良くないので、データサイズを小さくすることが難しかった。

【0012】

ここで、W3C標準のSVGを用いて図形オブジェクト(図1(b))を記述したテキスト形式のXMLデータ(図1(a))について説明する。このテキスト形式のXMLデータをバイナリXML形式に符号化しても、path要素のd属性の値はそのままテキストデータとして符号化され、バイナリXMLデータ内の文字列テーブルに格納される。よって、XMLデータ内に繰り返し出現する可能性の低い、サイズの大きなテキストデータは、従来のバイナリXML技術を用いても効率的に圧縮することが難しかった。特にSVGにおいては、path要素のd属性の値に長いテキストデータを用いて図形オブジェクトを記述する機会が多いので、バイナリXMLによる圧縮効果を得ることが難しかった。

10

【0013】

本発明は以上の問題に鑑みてなされたものであり、サイズの大きなテキストデータを属性値又は要素内容として含むテキストXMLからバイナリXMLへの変換効率を改善するための技術を提供することを目的とする。

【課題を解決するための手段】

【0014】

本発明の目的を達成するために、例えば、本発明の情報処理装置は以下の構成を備える。即ち、テキスト形式の構造化文書を取得する手段と、前記構造化文書を解析し、テキストデータを検出する手段と、前記テキストデータを子要素に変換することで前記構造化文書を変換し、変換した構造化文書に基づいてバイナリ形式の構造化文書を生成する生成手段とを備えることを特徴とする。

20

【発明の効果】

【0015】

本発明の構成により、サイズの大きなテキストデータを属性値又は要素内容として含むテキストXMLからバイナリXMLへの変換効率を改善することができる。

【図面の簡単な説明】

【0016】

【図1】テキスト形式のXMLデータとその描画結果を示す図。

30

【図2】実施形態に係るシステムの構成例を示す図。

【図3】PC101のハードウェア構成例を示すブロック図。

【図4】PC101の処理機能ブロック図。

【図5】PC101が行う処理のフローチャート。

【図6】ステップS113の詳細を示すフローチャート。

【図7】部分スキーマの一例を示す図。

【図8】(a)はXPath式の例を示す図、(b)はバイナリ形式のXMLデータを示す図

。

【図9】各イベントの符号を示す図。

【図10】名前テーブルを示す図。

40

【図11】(a)は値テーブルを示す図、(b)は符号のテーブルを示す図。

【発明を実施するための形態】

【0017】

以下、添付図面を参照し、本発明の好適な実施形態について説明する。なお、以下説明する実施形態は、本発明を具体的に実施した場合の一例を示すもので、特許請求の範囲に記載の構成の具体的な実施例の1つである。

【0018】

本実施形態では、比較的大きなサイズのテキストデータを属性値として有する属性を含む要素が記された構造化文書としてのテキストXMLを、圧縮効率良くバイナリXMLに変換する為の技術について説明する。

50

## 【 0 0 1 9 】

先ず、本実施形態に係るシステムについて、図 2 を用いて説明する。図 2 に示す如く、本実施形態に係るシステムは、本実施形態に係る情報処理装置としての P C ( パーソナルコンピュータ ) 1 0 1、ネットワークの一例である L A N 1 0 2、デジタルカメラ 1 0 3、複合機 1 0 4、ファイルサーバ 1 0 5 により構成されている。

## 【 0 0 2 0 】

デジタルカメラ 1 0 3 は、P C 1 0 1 で生成した X M L データと共に用いる画像を撮像により取得するためのもので、デジタルカメラ 1 0 3 により撮像された画像のデータは、L A N 1 0 2 を介して P C 1 0 1 やファイルサーバ 1 0 5 に転送される。もちろん、デジタルカメラ 1 0 3 の使用方法是これに限定するものではない。

10

## 【 0 0 2 1 】

複合機 1 0 4 は、P C 1 0 1 で生成した X M L データと共に用いる画像を複写により取得するためのもので、複合機 1 0 4 により複写された画像のデータは、L A N 1 0 2 を介して P C 1 0 1 やファイルサーバ 1 0 5 に転送される。もちろん、複合機 1 0 4 の使用方法是これに限定するものではない。

## 【 0 0 2 2 】

ファイルサーバ 1 0 5 は、デジタルカメラ 1 0 3 や複合機 1 0 4 から転送された画像のデータを保持すると共に、P C 1 0 1 から転送された X M L データも保持する。ファイルサーバ 1 0 5 が保持する画像のデータや X M L データに対しては、L A N 1 0 2 を介して適宜アクセスすることができる。

20

## 【 0 0 2 3 】

P C 1 0 1 は、本実施形態に係る情報処理装置として用いるもので、テキスト形式の X M L データをバイナリ形式の X M L データに変換する後述の処理を実行する。次に、図 3 を用いて、P C 1 0 1 のハードウェア構成例について説明する。

## 【 0 0 2 4 】

C P U 2 0 1 は、R O M 2 0 2 や R A M 2 0 3 に格納されているコンピュータプログラムやデータを用いて P C 1 0 1 全体の制御を行うと共に、P C 1 0 1 が行うものとして後述する各処理を実行する。

## 【 0 0 2 5 】

R O M 2 0 2 は、コンピュータ読み取り可能な記憶媒体の一例であり、P C 1 0 1 の設定データやブートプログラムなどが格納されている。R A M 2 0 3 は、コンピュータ読み取り可能な記憶媒体の一例であり、記憶部 2 0 4 からロードされたコンピュータプログラムやデータ、L A N I / F 2 0 7、U S B I / F 2 0 9 を介して外部から受信したデータ等を一時的に記憶する為のエリアを有する。また、R A M 2 0 3 は、C P U 2 0 1 が各種の処理を実行する際に用いるワークエリアも有する。即ち、R A M 2 0 3 は、各種のエリアを適宜提供することができる。

30

## 【 0 0 2 6 】

記憶部 2 0 4 は、コンピュータ読み取り可能な記憶媒体の一例であり、ハードディスクドライブ装置に代表される、大容量情報記憶装置である。記憶部 2 0 4 には、O S ( オペレーティングシステム ) や、P C 1 0 1 が行うものとして後述する各処理を C P U 2 0 1 に実行させるためのコンピュータプログラムやデータが保存されている。このコンピュータプログラムやデータは、C P U 2 0 1 による制御に従って適宜 R A M 2 0 3 にロードされ、C P U 2 0 1 による処理対象となる。

40

## 【 0 0 2 7 】

操作部 2 0 5 は、キーボードやマウスなどにより構成されており、P C 1 0 1 の操作者が操作することで、各種の指示を C P U 2 0 1 に対して入力することができる。表示部 2 0 6 は、C R T や液晶画面などにより構成されており、C P U 2 0 1 による処理結果を画像や文字などでもって表示することができる。

## 【 0 0 2 8 】

L A N I / F 2 0 7 は、P C 1 0 1 を L A N 1 0 2 に接続する為のもので、P C 1 0

50

1はこのLAN I/F 207を介してLAN 102に接続されている機器との通信を行う。USB I/F 209は、PC 101をUSB回線 210に接続する為のもので、PC 101はこのUSB I/F 209を介してUSB回線 210に接続されている機器との通信を行う。上記各部は何れも共通のバスに接続されており、互いに通信を行うことができる。

#### 【0029】

なお、以上で説明したシステムの構成やPC 101の構成は一例であって、後述する処理を実現することができ、且つ後述する処理によって得られるバイナリ形式のXMLデータを用いた応用が可能であれば、他の構成を適用しても良い。

#### 【0030】

次に、テキスト形式のXMLデータからバイナリ形式のXMLデータへの効率の良い圧縮符号化処理について、図5を用いて説明する。なお、図5に示したフローチャートに従った処理をCPU 201に実行させるためのコンピュータプログラムやデータは記憶部 204に保存されている。係るコンピュータプログラムやデータは、CPU 201による制御に従って適宜RAM 203にロードされ、CPU 201による処理対象となる。CPU 201がこのロードされたコンピュータプログラムやデータを用いて処理を実行することで、PC 101は、以下に説明する各処理(図5に示した各ステップにおける処理)を実行することになる。

#### 【0031】

また、以下では、説明を簡単にするために、圧縮符号化対象のテキスト形式のXMLデータは、図1に示したXMLデータであるものとする。もちろん、以下に説明する処理の本質は、他のテキスト形式のXMLデータであっても同じである。

#### 【0032】

先ずステップS 102では、CPU 201は、バイナリ形式のXMLデータへの圧縮符号化対象となるテキスト形式のXMLデータを選択する。係る選択は、PC 101上で動作しているアプリケーションによって行っても良いし、操作部 205を介して操作者が入力した選択指示に基づいて行っても良い。上述の通り、ここでは、図1に示したXMLデータが選択されたものとする。係るXMLデータには、効率的に圧縮することの難しいテキストデータ(X path式にマッチするテキストデータ)を属性値として有する属性(図1ではd属性)を含む要素(図1ではpath要素)が含まれている。

#### 【0033】

次にステップS 103では、効率的に圧縮することの難しいテキストデータに関するデータ型を記述した部分スキーマと、そのデータ型の情報のルートとなる定義名と、を選択し、この部分スキーマを解析し、RAM 203内にDOMツリーを作成する。

#### 【0034】

なお、部分スキーマと、そのデータ型の情報のルートとなる定義名は、表示部 206にこれらの選択を行うためのユーザインタフェース画面を表示し、ユーザがこの画面を見ながら操作部 205を用いて選択指示を入力するようにしても良い。また、部分スキーマと、そのデータ型の情報のルートとなる定義名を、あらかじめアプリケーションロジックに組み込むことで、これらを選択するようにしても良い。

#### 【0035】

上記の通り、本実施形態では、効率的に圧縮することの難しいデータは、図1におけるpath要素のd属性の属性値である。また、RELAX NGを用いて記述した、図1におけるpath要素のd属性のデータ型を定義した部分スキーマの一例について、図7を用いて説明する。

#### 【0036】

図7に例示した部分スキーマでは、define要素を用いて複数のデータ型を記述しているが、d属性の値のデータ型のルートとなる定義名は、"SVG.PathData.datatype"である。然るに、図7の場合、ステップS 103で選択する定義名は、"SVG.PathData.datatype"となる。

10

20

30

40

50

## 【 0 0 3 7 】

図 5 に戻って、次に、ステップ S 1 0 4 では、RELAX NG仕様書の第 4 章単純化に基づいて、選択した部分スキーマの文書構造を単純化する。そして、ステップ S 1 0 5 では、XPath式を用いることにより、定義した部分スキーマに対応するテキストデータを選択する。テキスト形式のXMLデータに出現する全てのpath要素のd属性を表現したXPath式の例は、図 8 ( a ) に示すようになっている。なお、XPath式を用いたテキストデータの選択については、表示部 2 0 6 に選択を行うためのユーザインタフェース画面を表示し、ユーザがこの画面を見ながら操作部 2 0 5 を用いて選択指示を入力するようにしても良い。また、XPath式をあらかじめアプリケーションロジックに組み込むことで、係る選択を行うようにしても良い。

10

## 【 0 0 3 8 】

次に、ステップ S 1 0 6 では、テキスト形式のXMLデータを解析する。係る解析は、テキスト形式のXMLデータの先頭から順次行うものとする。次に、ステップ S 1 0 7 では、テキスト形式のXMLデータからバイナリ形式のXMLデータへの符号化が完了したか否かを判断する。係る判断の結果、完了している場合には本処理は終了し、完了していない場合には処理はステップ S 1 0 8 に進む。

## 【 0 0 3 9 】

ステップ S 1 0 8 では、テキスト形式のXMLデータ中に、XPath式にマッチするテキストデータが含まれているか否かを判断する。係る判断の結果、含まれている場合(テキスト形式のXMLデータからXPath式にマッチするテキストデータを検出可能である場合)には処理をステップ S 1 1 0 に進め、含まれていない場合には処理をステップ S 1 0 9 に進める。ステップ S 1 0 9 では、テキスト形式のXMLデータをそのまま従来通りにバイナリ形式のXMLデータに符号化する。

20

## 【 0 0 4 0 】

ステップ S 1 1 0 では、テキスト形式のXMLデータにおいて、XPath式にマッチしないデータを符号化する。図 1 のテキスト形式のXMLデータの場合、最初の開始タグsvgはXPath式にマッチするテキストデータを含んでいないので、そのまま開始タグとして符号化する。そして、次のpath要素は、XPath式にマッチするd属性を含んでいるので、開始タグpathと、fill、stroke、stroke-width属性を符号化する。なお、XPath式にマッチしたd属性とその属性値は、符号化しないでRAM 2 0 3 に保存しておく。

30

## 【 0 0 4 1 】

次に、ステップ S 1 1 1 では、XPath式にマッチしたテキストデータの所属先が、属性値か否かを判断する。係る判断の結果、XPath式にマッチしたテキストデータの所属先が属性値である場合には、処理をステップ S 1 1 2 に進め、XPath式にマッチしたテキストデータの所属先が属性値ではない場合には、処理をステップ S 1 1 5 に進める。

## 【 0 0 4 2 】

ステップ S 1 1 2 では、XPath式にマッチしたテキストデータの所属先としての属性値の属性名をタグ名とする開始タグを符号化する。図 1 のテキスト形式のXMLデータの例では、XPath式にマッチしたテキストデータは属性値に所属し、その属性名が " d " であるので、 " d " という名前を持つ開始タグ < d > を符号化する。

40

## 【 0 0 4 3 】

次に、ステップ S 1 1 3 では、XPath式にマッチしたテキストデータを先頭から逐次解析し、選択した部分スキーマのデータ型の情報を用いることによって、そのテキストデータを子要素に変換する。

## 【 0 0 4 4 】

ここで、ステップ S 1 1 3 における処理の詳細について、図 6 を用いて説明する。まずステップ S 2 0 2 では、部分スキーマのDOMツリー内の内部状態を初期化する。次に、ステップ S 2 0 3 では、テキストデータの解析途中の状態を記憶しておくために、部分スキーマ内の任意のノードを指し示すポインタ p を、NULL や最初のノードへのポインタ値等に初期化する。

50

## 【 0 0 4 5 】

次に、ステップ S 2 0 4 では、XPath式にマッチしたテキストデータについての解析が完了したか否かを判断する。係る判断の結果、完了した場合にはステップ S 1 1 3 の処理は終了し、ステップ S 1 1 4 に処理を進める。一方、係る判断の結果、完了していない場合には、処理をステップ S 2 0 5 に進める。なお、テキストデータの解析処理は、テキストデータに含まれる空白文字を区切りの識別子とし、この識別子でテキストデータを分割し、それぞれの分割テキストデータ毎に解析する。

## 【 0 0 4 6 】

ステップ S 2 0 5 では、分割テキストデータが、部分スキーマのgroup要素以下で定義された列挙値に該当する変数であるか否かを判断する。係る判断の結果、該当する場合には、処理をステップ S 2 0 6 に進め、該当しない場合には、処理をステップ S 2 0 7 に進める。

10

## 【 0 0 4 7 】

ステップ S 2 0 6 では、部分スキーマのgroup要素以下で定義された列挙値に該当する変数としての分割テキストデータを、その列挙値をタグ名とする開始タグに符号化する。図 7 に示した部分スキーマを用いる場合、列挙値は、" M "、" m "、" L "、" l " であるので、最初に解析する分割テキストデータ " M " は、" moveType " を名前として持つ define要素以下に対応する。然るにこのような場合には、最初に解析する分割テキストデータ " M " を、" M " をタグ名とする開始タグに符号化する。そしてその後、ステップ S 2 0 8 では、解析途中の状態を記憶しておくために、例えば図 7 の場合には、部分スキーマ内の名前が " moveType " の define要素以下の choice要素を指し示すように、ポインタ p を更新する。

20

## 【 0 0 4 8 】

一方、ステップ S 2 0 7 では、ポインタ p が指し示す部分スキーマ内の参照先のノードを解析する。そして、分割テキストデータを、そのテキストデータのデータ型に基づいて、1つのテキストイベントとして符号化する。図 7 に示した部分スキーマを用いる場合、最初の分割テキストデータ " M " を開始タグとして符号化した後、" M " に後続する分割テキストデータ " 1 0 0 " は、部分スキーマで定義されている float型で1つのテキストイベントとして符号化する。

30

## 【 0 0 4 9 】

そしてその後ステップ S 2 0 8 では、現在ポインタ p が指し示していたノードの兄弟ノードを参照するように、ポインタ p を更新する。従って、分割テキストデータ " 1 0 0 " を符号化した後、ポインタ p は、" moveType " を名前に持つ define要素以下の < data type = " float " > ノードを参照することになる。このようにして、1つのテキストデータを構成する全ての分割テキストデータについて順次、符号化する。

## 【 0 0 5 0 】

図 5 に戻って、次に、ステップ S 1 1 4 では、ステップ S 1 1 2 で符号化した開始タグに対応する終了タグを符号化する。本実施形態では、" d " というタグ名を有する終了タグ < /d > を符号化する。一方、ステップ S 1 1 5 では、ステップ S 1 1 2 , S 1 1 4 で符号化した開始タグと終了タグを出力しないで、ステップ S 1 1 3 と同様の処理を行う。

40

## 【 0 0 5 1 】

以上説明した符号化処理を図 1 に示したテキスト形式の XML データに対して行うことで、図 8 ( b ) に示すバイナリ形式の XML データを作成することができる。図 8 ( b ) に示したバイナリは、それぞれの単位に対応するバイナリである。

## 【 0 0 5 2 】

図 1 のテキスト形式の XML データの path要素に含まれる d属性のデータは、path要素の子要素として含まれるように変換される。また、d属性の値のテキストデータは、" M " や " L " といったコマンドをグループ ( テキストデータの構造単位 ) の単位として子要素として変換される。その要素内容には、" M " や " L " のコマンドに対して指定する座標情報が、1つ1つのテキストイベントとして格納される。アプリケーションが、" M " や " L "

50



”のコマンドに対して指定した座標情報をデコードした際には、各座標情報を別々のテキストイベントとして取得することができる。

【0053】

図8(b)に示したバイナリ形式のXMLデータに関する各イベントの符号を図9に示す。アプリケーションは事前に各イベントの符号を知っているものとする。また、図8(b)で参照される名前テーブルを図10に示す。この名前テーブルは、要素や属性の名前と、各名前に割り当てる符号と、各名前のデータサイズを格納したテーブルである。

【0054】

また、図8(b)に示したバイナリ形式のXMLデータで参照される値テーブルを図11(a)に示す。この値テーブルは、要素内容や属性値の値と、各値に割り当てる符号と、各値のデータ型、各値のデータサイズを格納したテーブルである。図10, 11(a)に示したテーブルは、変換元となるテキスト形式のXMLデータに依存するものであるため、テキスト形式のXMLデータの内容が変われば、図10, 11(a)に示したテーブルの内容もそれに依存して変わる。図11(b)は、図11(a)に示したテーブルの各値のデータ型として参照する符号のテーブルであり、アプリケーションは事前に知っているものとする。

10

【0055】

次に、PC101の処理機能ブロックについて、図4を用いて説明する。構造化文書解析部304は、バイナリ形式のXMLデータへの圧縮符号化対象となるテキスト形式のXMLデータを取得し、解析する。

20

【0056】

スキーマ選択部303は、効率的に圧縮することの難しいテキストデータに関するデータ型を記述した部分スキーマと、そのデータ型の情報のルートとなる定義名と、を選択する。

【0057】

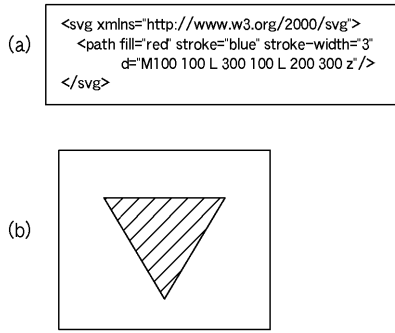
スキーマ解析部305は、この部分スキーマを解析し、RAM203内にDOMツリーを作成する。そして更に、スキーマ解析部305は、RELAX NG仕様書の第4章単純化に基づいて、選択した部分スキーマの文書構造を単純化する。

【0058】

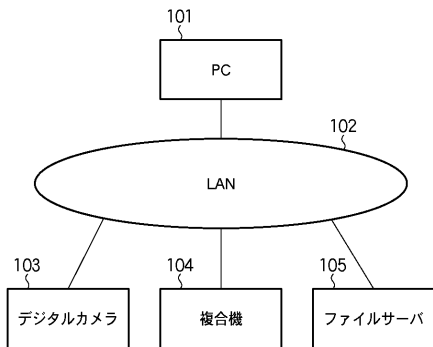
テキストデータ選択部302は、XPath式を取得する。符号化文書生成部306は、XPath式を用いることにより、定義した部分スキーマに対応するテキストデータを選択する。符号化文書生成部306は、テキストデータ解析部307、テキストデータ符号化部309を有しており、各部は以下のように動作する。テキストデータ解析部307は、XPath式にマッチしたテキストデータの解析を行い、テキストデータ符号化部309は、上記各符号化処理を行う。

30

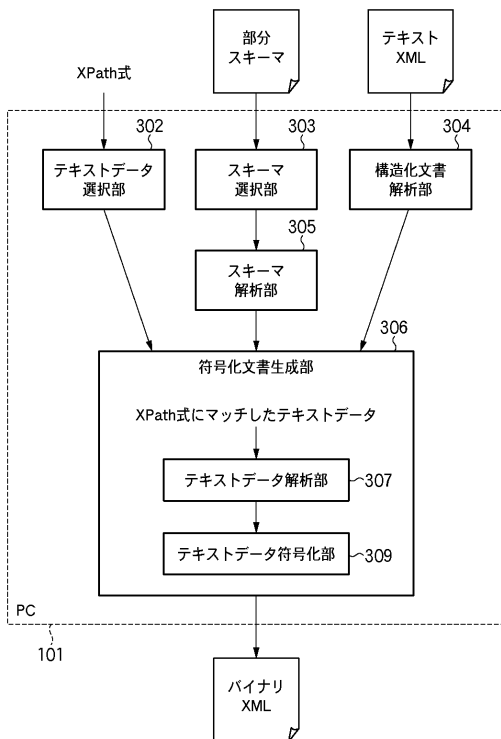
【 図 1 】



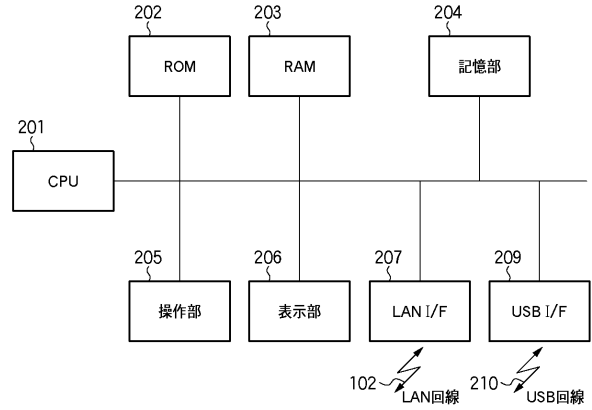
【 図 2 】



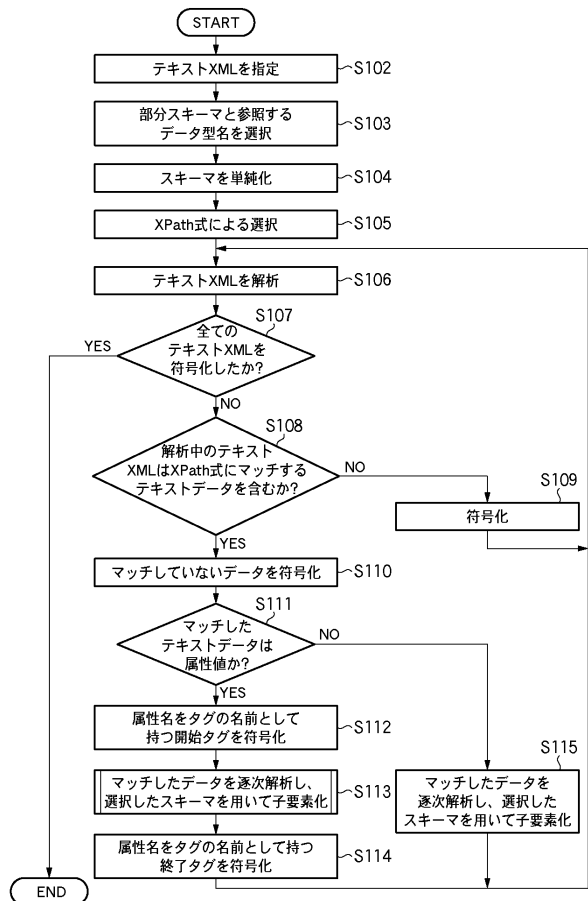
【 図 4 】



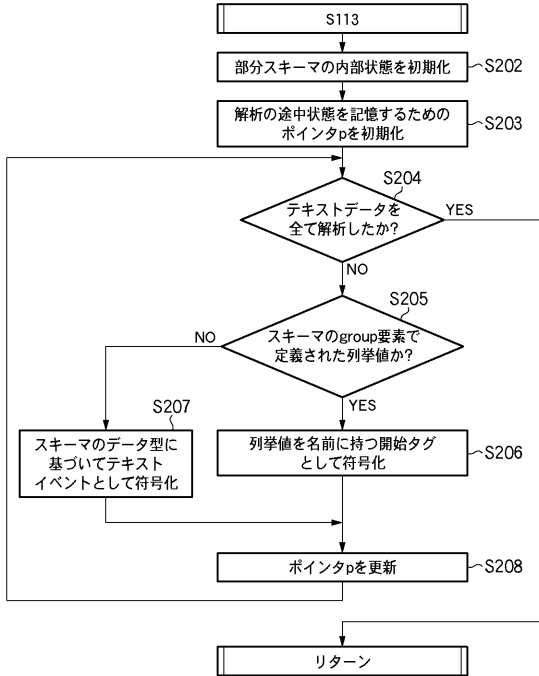
【 図 3 】



【 図 5 】



【 図 6 】



【 図 7 】

```

<grammar xmlns="http://www.w3.org/2000/svg">
  <define name="SVG.PathData.datatype">
    <list>
      <oneOrMore>
        <choice>
          <ref name="moveType"/>
          <ref name="lineType"/>
        </choice>
      </oneOrMore>
      <ref name="closepathType"/>
    </list>
  </define>
  <define name="moveType">
    <group>
      <choice>
        <value>M</value>
        <value>m</value>
      </choice>
      <data type="float"/>
      <data type="float"/>
    </group>
  </define>
  <define name="lineType">
    <group>
      <choice>
        <value>L</value>
        <value>l</value>
      </choice>
      <data type="float"/>
      <data type="float"/>
    </group>
  </define>
  <define name="closepathType">
    <group>
      <optional>
        <value>z</value>
      </optional>
    </group>
  </define>
</grammar>
  
```

【 図 8 】

(a) //path@d

```

(b) <svg xmlns="http://www.w3.org/2000/svg">
0x00 0x10 0x02 0x11 0x20
  <path fill="red" stroke="blue" stroke-width="3">
0x00 0x12 0x02 0x13 0x21 0x02 0x14 0x22 0x02 0x15 0x23
    <d>
0x00 0x16
      <M> 100 100 </>
0x00 0x17 0x03 0x24 0x03 0x24 0x01
      <L> 300 100 </>
0x00 0x18 0x03 0x25 0x03 0x24 0x01
      <L> 200 300 </>
0x00 0x18 0x03 0x26 0x03 0x25 0x01
      <z> </>
0x00 0x19 0x01
    </>
  </>
0x01
</>
0x01
  
```

【 図 9 】

符号	イベントの種類
0x00	開始タグ
0x01	終了タグ
0x02	属性
0x03	要素内容

【 図 10 】

符号	サイズ	名前
0x10	3	svg
0x11	4	xmlns
0x12	4	path
0x13	4	fill
0x14	6	stroke
0x15	12	stroke-width
0x16	1	d
0x17	1	M
0x18	1	L
0x19	1	z

【 図 1 1 】

(a)

符号	データ型	サイズ	値
0x20	0x30	26	<a href="http://www.w3.org/2000/svg">http://www.w3.org/2000/svg</a>
0x21	0x30	3	red
0x22	0x30	4	blue
0x23	0x30	1	3
0x24	0x35	4	100
0x25	0x35	4	300
0x26	0x35	4	200

(b)

符号	データ型
0x30	string
0x31	base64
0x32	short
0x33	int
0x34	long
0x35	float
0x36	double

フロントページの続き

(72)発明者 内田 均

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

Fターム(参考) 5B109 NH08