

(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2005-18672

(P2005-18672A)

(43) 公開日 平成17年1月20日(2005.1.20)

(51) Int. Cl. ⁷	F I	テーマコード (参考)
G06F 5/00	G06F 5/00 H	5B009
G06F 12/00	G06F 12/00 511A	5B082
G06F 17/21	G06F 17/21 501T	
	G06F 17/21 570G	

審査請求 未請求 請求項の数 19 O L (全 59 頁)

(21) 出願番号	特願2003-186103 (P2003-186103)	(71) 出願人	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成15年6月30日 (2003.6.30)	(74) 代理人	100075096 弁理士 作田 康夫
		(72) 発明者	室 啓朗 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内
		(72) 発明者	舟生 幸雄 神奈川県横浜市戸塚区戸塚町216番地 株式会社日立製作所ディフェンスシステム 事業部内
		Fターム(参考)	5B009 QA06 SA08 5B082 AA11 GA01

(54) 【発明の名称】 構造化文書の圧縮方法

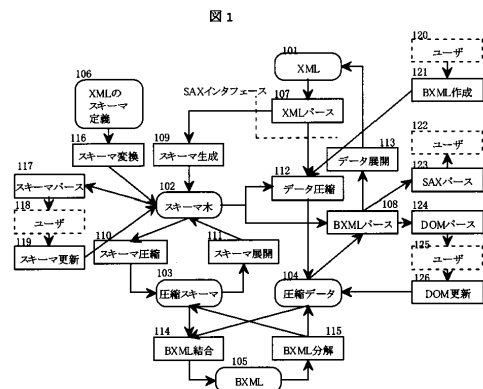
(57) 【要約】

【課題】 多量の図形情報が含まれる情報格納件数やデータサイズの巨大な構造化文書を、コンパクトに格納し、高速に解析・更新・作成するための方法を提供すること。

【解決手段】 構造化文書を圧縮スキーマ部と圧縮データ部に分離し、構造化文書の個別情報であるテキスト要素、属性値、重複数のみを圧縮データに、それ以外を圧縮スキーマ部に格納する。さらにテキスト要素の辞書化、可変長数値の適用、相対値管理を行うことでサイズを削減する。また標準インタフェースであるSAX、DOMを利用して構造化文書を解析・更新・作成する手段を提供すると共に、圧縮スキーマ更新により大量の圧縮データの更新なしで構造の一括変更手段を提供する。

【効果】 拡張性の高い構造化文書の作成、及び高速なデータ解析を可能とする。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

多量の図形情報が含まれる地図情報のように情報格納件数やデータサイズの巨大な構造化文書をコンパクトに格納し高速に通信するための構造化文書圧縮方法であり、構造化文書を圧縮スキーマ部と圧縮データ部に分離し、構造化文書のタグ名称、タグの性質を示すスキーマ型、タグが保有する子要素の数、タグが保有する属性の数を一つのタグに関する情報とし、構造化文書の属性キー名称、属性の性質を示すスキーマ型を一つの属性キーに関する情報とし、属性や子タグの情報をそれが所属する親タグの情報の直後に連続して再起的に記述する、という方法で構造化文書に出現する全種類のタグ情報・属性キー情報を圧縮スキーマ部に記述し、また構造化文書のタグや属性キーがその位置で何個連続して出現するかを示す重複数、テキスト要素、属性値を、圧縮スキーマで記述した順番通りに圧縮データ部に記述することにより構造化文書とその論理関係を保ったまま変形し、これにより構造化文書のサイズを削減し、また圧縮データ部の値を圧縮スキーマ部のタグ階層関係、タグと属性の関係、タグの重複関係、圧縮データ部のタグ重複数から類推することにより復元することを可能とする構造化文書圧縮方法。

【請求項 2】

請求項 1 の構造化文書圧縮方法において、圧縮スキーマ部のタグに関する情報として、タグが重複して存在するかを示す重複フラグを設定し、通常は重複フラグに「重複」と記述し、圧縮データ部にタグや属性キーがその位置で何個連続して出現するかを示す重複数を記載し、また必ず一回のみ出現するという事が保証されているタグや属性キーについては重複フラグに「単独」と記述し、圧縮データ部の重複数の記載を省略することにより、構造化文書のサイズを削減することを特徴とする構造化文書圧縮方法。

【請求項 3】

請求項 1 の構造化文書圧縮方法において、特定タグの下に複数の指定タグが順不同に複数回出現する可能性がある構造化文書を圧縮する場合において、圧縮スキーマ部における該特定タグのスキーマ型を重複選択型と記述し、圧縮データ部の対応する部分に子タグ総数を記述し、それぞれの子タグの情報として種類を示す識別子、子タグのテキスト要素をペアで記述し、あるいはまた特定タグの下に複数の指定タグのどれか一つのみが出現することが保証されている構造化文書を圧縮する場合において、圧縮スキーマ部における該特定タグのスキーマ型を選択型と記述し、圧縮データ部の対応する部分の子タグ総数の記述を省略し、それぞれの子タグの情報として種類を示す識別子、子タグのテキスト要素のみをペアで記述することにより、構造化文書のサイズを削減することを特徴とする構造化文書圧縮方法。

【請求項 4】

請求項 1 の構造化文書圧縮方法において、繰り返し出現する冗長なテキスト要素や属性値を識別子により代替することにより構造化文書のサイズを削減する方式であり、圧縮スキーマ部の特定タグの付加情報として、該特定タグに繰り返し出現するテキスト要素や属性キーの全部または断片の文字列を辞書として格納し、圧縮データ部の対応するテキスト要素・属性キーの格納位置に、文字列が識別子で代替表現されている事示すフラグと識別子のペア、文字列が識別子で代替表現されていない事示すフラグと文字列のペア、の順不同の組み合わせを記述する事により、構造化文書のサイズを削減することを特徴とする構造化文書圧縮方法。

【請求項 5】

請求項 1 の構造化文書圧縮方法において、テキスト要素や属性キーが汎用的な文字列ではなく、整数型や実数型、日付型など、特定の規則で表現される特殊な値であることが保証されている場合、圧縮スキーマ部においてスキーマ型を整数型や実数型、日付型など、該特定型と指定し、圧縮データにおいて該特定型に対応した記述を行うことにより、構造化文書のサイズを削減することを特徴とする構造化文書圧縮方法。

【請求項 6】

請求項 1 の構造化文書圧縮方法において、テキスト要素や属性キーが汎用的な文字列では

なく、整数型や実数型、日付型など、特定の規則で表現される数値であることが保証されており、かつ該数値が緯経度座標など巨大な数値であり、その変動値が該数値の絶対値と比較して微少である場合、圧縮スキーマ部においてスキーマ型を整数型や実数型、日付型など、該特定型と指定し、また圧縮スキーマ部において該数値の基準値を記載し、圧縮データにおいて該特定型に対応した該基準値からの相対値を可変長数値表現で記述し、圧縮データの数値記載部分のサイズを削減することにより、構造化文書のサイズを削減することを特徴とする構造化文書圧縮方法。

【請求項 7】

構造化文書を請求項 1 の構造化文書圧縮方法に従って効率よく圧縮するための方法であり、現在作成中のタグとその子要素に当たる次に到来すると期待されるタグを管理し、構造化文書を解析して開始タグ・終了タグ・テキスト要素・属性をその出現順に取得し、次に到来すると期待されるタグと該開始タグ、現在作成中のタグと該終了タグと比較することによりタグの親子関係・重複関係を決定する方式と、タグのテキスト要素と過去に登録したタグのスキーマ型、属性と過去に登録した属性のスキーマ型とを上書規則を用いて比較することによりスキーマ型を決定する方式からなるスキーマ生成ブロックと、要素単位でそのタグの重複数と該重複数の圧縮データ挿入位置を管理し、テキスト要素を圧縮データに書き込みつつ構造化文書を解析してその開始タグ、終了タグと比較することにより、重複数を勘定して、事後に圧縮データに書き込む方式からなるデータ圧縮ブロックから構成される構造化文書圧縮方法。

10

【請求項 8】

請求項 1 の構造化文書圧縮方法において圧縮された構造化文書を効率よく解析するための方法であり、現在参照しているタグに関する情報として、圧縮スキーマ部分の該タグに対応するスキーマ要素への参照、該タグと同名のタグが連続して何個出現しているかを示す重複数、該タグがそのうち何番目のタグであることを示す重複カウンタ、圧縮データにおける該タグに対応する部分へのポインタを管理し、現在参照しているタグが重複型であれば圧縮データから重複数を取得し、スキーマ要素のスキーマ型に従った型のデータを圧縮データから取得してテキスト要素とし、重複カウンタが重複数よりも少なければ次の要素も同じスキーマ型として繰り返し、スキーマ要素が子要素を保持する場合は現在参照しているタグを該子要素に移動して再起的に処理を行い、全ての子要素の解析を終了した時点で現在参照しているタグを親要素に移動して再起的に処理を行うことにより圧縮された構造化文書をパースし、また圧縮データにおいて重複数や各スキーマ型に従ったデータが出現した時に外部プログラムにその出現情報を伝えるメッセージを発行することにより処理を委譲できる構造化文書解析方法。

20

30

【請求項 9】

請求項 1 の構造化文書圧縮方法において圧縮された構造化文書を効率よく解析するための方法であり、現在参照しているタグを示すノードポインタとして、圧縮スキーマ部分の該タグに対応するスキーマ要素への参照、該タグと同名のタグが連続して何個出現しているかを示す重複数、該タグがそのうち何番目のタグであることを示す重複カウンタ、圧縮データにおける該タグに対応する部分へのポインタを保有し、子ノードへのポインタ移動を行う際には圧縮データの属性値部分をスキップし、該子要素へのスキーマ要素を返戻することにより実施し、弟ノードへのポインタ移動を行う際には圧縮データの自分自身並びに子要素部分を再起的にスキップし、該弟要素のスキーマ要素を返戻することにより実施し、また該ノードポインタの情報をその親ノード、祖父ノードと根ノードまで遡り、リストとして管理することにより、親ノードへのポインタ移動を可能とすることにより提供されるノード移動手段と、該圧縮データを該スキーマ要素のスキーマ型に従って復元することによるノードのテキスト要素参照手段を備えることにより、構造化文書の解析を高速に、かつ少ない記憶容量で行うことを特長とする構造化文書解析方法。

40

【請求項 10】

請求項 9 の構造化文書解析方法において、指定したタグ名の弟要素・子要素を高速に検索するための方法ならびに装置であり、圧縮スキーマ部分にスキーマ要素を一意に識別する

50

ための識別子である数値を記述し、該ノードポインタの情報としてさらに検索対象のタグ名を表す識別子である数値を管理し、弟ノードへのポインタ移動を行う際に請求項9の方法に加え、弟候補ノードの識別子と該検索対象の識別子と比較し不一致の場合は次の弟候補ノードを再起的に検索する比較検索方法を適用し、子ノードへのポインタ移動を行う際には請求項9の方法に加え、該比較検索方法を適用することにより、指定したタグ名の弟要素・子要素の高速検索を可能とすることを特徴とする構造化文書解析方法。

【請求項11】

請求項9の構造化文書解析方法において、タグを示すノードポインタを特定の分野に特化することにより利用者の利便性を向上させる方法であり、地理情報分野における構造化文書の解析手段として地物を示す地物ノード、図形を示す図形ノードを設け、これらのノードの内部情報として請求項9で示した汎用ノードを保有することにより汎用のノード移動手段・テキスト要素参照手段を提供し、更に固定個数の図形ノードのタグ名をあらかじめ記憶しておき、該タグ名と検索対象のタグ名を比較することにより、地物ノードからタグ名を指定しないで図形ノードへポインタ移動を行うことを可能とすることにより、利用者の利便性を向上させる事を特徴とする構造化文書解析方法。

10

【請求項12】

請求項1の構造化文書圧縮方法において圧縮された構造化文書を効率よく更新するための方法であり、要素群の挿入時には挿入対象要素群の圧縮データ挿入に先立ち挿入トークンと挿入要素のスキーマ要素を識別する識別子を挿入し、また要素群の削除時には削除対象要素群に先立ち削除トークンを挿入し、また更新された構造化文書の解析時には挿入トークン出現時には重複カウンタの増加を抑制しつつ挿入要素の解析を実行し、また削除トークン出現時には重複カウンタの増加のみを行い削除要素部分の解析をスキップすることにより、圧縮データの重複要素部分の修正を行わずに要素群の挿入・削除を実現することにより、高速な更新を行えることを特徴とし、さらに更新部分に挿入トークンと削除トークンを挿入することにより更新部分が判定できるため、更新の取り消しや更新部分のみの差分通信を可能とすることを特長とする構造化文書更新方法。

20

【請求項13】

請求項1の構造化文書圧縮方法において圧縮された構造化文書を効率よく更新するための方法であり、圧縮スキーマ部のタグに関する情報として、さらにタグを隠蔽するかどうかを識別する隠蔽フラグを設定し、該請求項8の構造化文書解析方法において隠蔽されたタグおよびその子要素については解析を中止し、スキップすることにより利用者に対して該タグを隠蔽する方法であり、オリジナルの構造化文書の圧縮データ部分を作り直すことなく利用者の要求する部分の抽出機能を実現することを特長とする構造化文書更新方法。

30

【請求項14】

請求項1の圧縮された構造化文書を作成するための方法であり、利用者へのインタフェースとして子要素作成メソッド、弟要素作成メソッド、親へ移動メソッドから構成されるノード作成メソッド群、属性の作成、テキスト要素の作成から構成されるデータ作成メソッド群を提供し、これに対して開始タグ、終了タグ、属性、テキスト要素を請求項7で示した構造化文書圧縮装置に提供することにより、利用者に簡易なインタフェースで、効率的に構造化文書を作成することを可能とする構造化文書作成方法。

40

【請求項15】

請求項14の構造化文章作成方法において、子要素作成メソッド、弟要素作成メソッド、親へ移動メソッドから構成されるノード作成メソッド群を一時的なバッファに記憶しておき、子要素作成メソッド発行後に属性の作成、テキスト要素の作成から構成されるデータ作成メソッド群が一度も発行されない状態で親へ移動メソッドが発行された場合に前回発行された子要素作成メソッドをキャンセルすることにより、0個の子要素を作成する際に不要な子要素を作成することをせず、これにより利用者の構造化文章作成手順を簡易化することを特長とする構造化文書作成方法。

【請求項16】

請求項1の圧縮された構造化文書同士を結合するための方法であり、結合対象の構造化文

50

書の圧縮スキーマが同一の場合は、該圧縮スキーマを持つ圧縮された構造化文書を請求項 14 の方法で新たに作成し、結合対象のノードを示す圧縮データをコピーし、また結合対象の構造化文書の圧縮スキーマが異なる場合は、双方の構造化文書の圧縮スキーマを含む圧縮スキーマを作成し、結合対象のノードを示す圧縮データを新たに作成した圧縮スキーマに適合する構造に変換したのちにコピーすることにより構造化文書同士を結合するための方法。

【請求項 17】

クライアントとサーバ間の低速な通信回線を通じて効率よく大量の構造化文書を交換するための方法であり、クライアントとサーバ間のクライアント装置側に請求項 7 で示した構造化文書圧縮装置と請求項 8 で示した構造化文書解析装置を持つ中継装置を設置し、クライアントの発行する構造化文書を圧縮してサーバに転送し、サーバの発行する構造化文書を解析して元の構造化文書に復元し、クライアントに転送することにより、本発明の構造化文書圧縮装置・解析装置を保有しない一般のクライアントに対し、効率的に大量の構造化文書を交換することを可能とする中継装置。

10

【請求項 18】

多量の構造化文書を格納し、クライアントの検索要求に従って構造化文書の部分集合を返戻するデータベース装置であり、表形式のデータベース装置を保有し、データベースの要素の一つとして請求項 1 で示した圧縮データを格納し、また請求項 16 で示した構造化文書結合装置を保有し、検索に該当した圧縮データを結合してクライアントに返戻することにより、高速な応答速度を実現することの可能なデータベース装置。

20

【請求項 19】

空間データを表現する多量の構造化文書を格納し、クライアントの検索要求に従って構造化文書の部分集合を返戻するデータベース装置であり、領域で分割された空間データを表現する請求項 1 で示した圧縮された構造化文書の集合を保有し、クライアントが要求する空間範囲に該当する圧縮された構造化文書の部分集合を収集し、請求項 13 で示すタグ隠蔽方法によりクライアントが要求していない種類の空間データを隠蔽し、また請求項 16 で示した構造化文書結合装置を保有し、検索に該当した圧縮データを結合してクライアントに返戻することにより、高速な応答速度を実現することの可能なデータベース装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、XML (eXtensible Markup Language) や SGML (Standard Generalized Markup Language) 等の構造化文書の圧縮・復元・解析・更新・作成に関わる技術に関する。

30

【0002】

【従来の技術】

近年、インターネットやイントラネット等、計算機間を接続するネットワークが普及するに従い、ネットワークを用いて計算機を相互接続する動きがある。

【0003】

複数の異種システム間でデータやプロトコルを統一して利用するための規格の一つとして、XML が 1998 年 2 月に W3C (World Wide Web Consortium) によって勧告された。XML は、同様の規格 SGML のサブセットとなっており、文書の中にタグを埋め込む形で、その文書の構造が記述される。XML や SGML により記述された文書は、一般に構造化文書と呼ばれる。

40

【0004】

XML では、「`<`」と「`>`」とで囲まれた領域をタグ、「`<` タグ名」を開始タグ、「`</` タグ名」を終了タグ、開始タグと終了タグで囲まれた領域をテキスト要素と呼ぶ。また開始タグの内部に「`キー = "値"`」の形式で属性を記述することができる。開始タグと終了タグでテキスト要素を囲むことにより、そのテキスト要素が何についての記述なのかという意味づけを行うことができる。また、開始タグと終了タグの間に再帰的に別の開始

50

タグ・終了タグを含むことにより階層化されたデータ構造を表現することができる。

【0005】

XMLの利用はWWW(World Wide Web)の分野を中心に広がっており、特にEDI(Electronic Data Interchange)、EC(Electronic Commerce)等で広く利用されつつある。XMLは階層化されたデータ構造を表現することのできる汎用的な仕組みであり、データの表現方式としても、またデータ要求・応答を実現する通信プロトコルとしても利用されている。

【0006】

XMLを解析・更新するためのインタフェースも標準化されている。たとえばXMLをイベント駆動形式で走査するSAX(Simple API for XML)、XMLを木状に展開して管理し、検索・更新を行うDOM(Document Object Model)が存在する。これらの実装としては、米国Microsoft社のMSXML、Apache Software FoundationのDOM/SAXパーサであるXercesなど、多数のメーカーや団体から有償・無償で公開されている。

10

【0007】

本発明では、XML文書を圧縮する際にSAXインタフェースを利用している。また本発明では、本発明により圧縮したXML文書に対する標準のインタフェースとして、SAXインタフェース、DOMインタフェースを提供している。そのため、以下にSAXインタフェースとDOMインタフェースの簡単な説明を行う。

【0008】

図2を用いてSAXパーサのアルゴリズムを示す。SAXパーサの目的は、XML文書を走査して、外部のプログラムに対しXML文書構造やテキスト要素の情報を提供することにある。SAXパーサは、XML文書を走査し、XML文書の構造に対応した複数のイベントをユーザプログラムに通知する。具体的には、XML文書を走査する直前に文書開始処理201を、XML文書の捜査が終了した時点で文書終了処理202を呼び出す。また、開始タグを発見した時には開始タグ処理203を、終了タグを発見した時には終了タグ処理204を、テキスト要素を発見した時にはテキスト要素処理205を、属性を発見した時には属性処理206を呼び出す。

20

ユーザは、上記201~206までのイベントに対応したイベント処理プログラムを実装することにより、XML文書に対する様々な操作を行うことができる。

30

【0009】

図3を用いて、DOMパーサについて説明する。DOMパーサの目的は、XML文書をメモリ空間上に木状に展開して、それに対する様々なアクセス方式をユーザに提供することによりXML文書の解析・更新機能を提供する。

【0010】

DOMパーサが作成する木を一般にはDOM木と呼ぶが、ここではXML文書のインスタンスをそのまま木状に表現することから、インスタンス木と呼ぶことにする。図3のXML文書301を木状に展開することにより、インスタンス木302を得る。インスタンス木302内部の円で囲まれる部分はノードと呼ばれ、それぞれXML文書のタグが囲む領域に相当する。XML文書301では、タグAの内部に2つのタグBと1つのタグCを保有する。また2番目のタグBの内部には、タグD、タグEが含まれる。タグBにはキー「id」で指定される属性が格納され、タグD、タグE、タグCにはテキスト要素が格納されている。インスタンス木302は、この階層構造を木状に表現しており、ノードAに2つのノードBと1つのノードCが接続し、また2番目のノードBに、ノードD、ノードEが接続するという構造を取る。

40

【0011】

ユーザにはインタフェースとして、インスタンス木の一つのノードを指すノードポインタが与えられる。ノードポインタはXML文書の階層構造を移動することができ、ノードに属するテキスト要素や属性を取得することができる。

【0012】

50

例えばノードポインタがノードB(304)を指しているとする。親へ移動(307)することにより、ノードポインタはノードA(303)に移動する。また弟へ移動(308)することにより、ノードポインタはノードC(305)に移動する。子へ移動(309)することにより、ノードポインタはノードD(306)に移動する。またノードB(304)で、キー「id」の属性を取得する(311)と、属性値「b2」が取得でき、またノードD(306)でテキスト要素を取得する(310)と、テキスト要素「d1」が取得できる。従来のDOMパーサは、メモリ空間上に展開されるインスタンス木の302のサイズが巨大であった。インスタンス木302全体が占めるメモリ空間は、元のXML文書301の数倍のサイズとなる場合がある。小さなサイズのXML文書では問題はないが、本発明が対象とする地理情報分野などの、ノード数が多くサイズも巨大なXML文書 10
 に対しては、大量のメモリ空間を占有し、かつメモリ空間への展開やXML文書への復元に時間がかかる。本発明では、XML文書よりサイズの小さな圧縮XML文書に対するDOMインタフェースを提供し、かつインスタンス木をメモリ空間に展開しない方式を用いることにより、メモリ空間の節約を図る。

【0013】

XML圧縮技術については何度も繰り返される長いタグ名をタグ辞書に格納し、その識別子をタグの代替として使用することにより、XML文書のサイズを削減するものが、「WAP Binary XML (WBXML) Encoding Specification」として1999年にエリクソン・IBM・モトローラ・Phone.comから構成されるWAPフォーラムによりW3Cに標準化案として提出されている。又、図5 20
 に開示する、XML文書の階層構造を示すタグ部分と、実質のデータを示すテキスト要素部分を分離するものがある(特許文献1)。

【特許文献1】特開2002-163248号公報

【発明が解決しようとする課題】

XMLのメリットは次の3点である。(1)データ構造が高い柔軟性を有するため、ベンダに中立なデータ構造となり、ベンダ間の容易な相互接続を可能とする。(2)データ構造が高い拡張性を有するため、ユーザの自由な拡張を可能とする。(3)テキストベースであるため、人間に対し高い視認性を有する。

又、XMLはテキストベースであることが原因となり、従来フォーマットに比べてデータサイズが大きくなる。又、テキストベースのため計算機が扱いにくい。例えばXML文書で記述された地図を画面上に描画するのは、従来フォーマットで記述された地図を描画するよりはるかに時間がかかる。特に地理情報処理分野におけるXML文書は(1)項目数・サイズ共に桁違いに巨大である、(2)データのほとんどの項目が図形形状を表現する座標数値である。よって、従来のXML文書を直接管理する方式では、容量・処理時間ともに実用に耐える時間では利用できず、本発明による圧縮・復元・解析方法が必要となる。

【0014】

本発明では、上記の課題を解消するため、XML等の構造化文書の意味的構造を保ったままバイナリ構造に圧縮することにより、データサイズを削減する。また圧縮XMLの圧縮した状態での解析・更新機能を提供することにより、計算機が扱いやすい、従ってXML 40
 文書と比較して高速な解析処理・加工処理を実現する機能を提供する。また圧縮XMLを直接作成する機能を提供することにより、高速な構造化文書作成機能を提供する。

【0015】

【課題を解決するための手段】

本願の開示する発明の概要を説明すると以下の通りである。

構造化文書を圧縮する方法であって、タグとデータを分離し、構造の型に分類し、該型に応じたデータ構造へ変換する。さらに、テキストデータ中の頻出単語を識別子で代替する。又、複数の数値表示型を採用し、座標をある基準点からの相対的な位置で表示する。さらに、通常 40
 の構造化文書を上記圧縮方法で作成される圧縮データ構造へと圧縮する方法、圧縮データ構造から通常 50
 の構造化文書への復元方法、通常 50
 の構造化文書の復元を行わず

に圧縮データ構造を直接解析する方法、直接更新する方法、圧縮データ構造を直接作成する方法について開示する。

構造化文書の圧縮は、構造化文書を解析してタグ間の関連やタグ内に格納されるデータの型を抽出するスキーマ生成工程と、生成したスキーマに従って構造化文書を適切な圧縮データ構造に変換するデータ圧縮行程で実現する。

構造化文書の復元は、スキーマ構造に従って圧縮データを解析し、重複・省略・条件選択されるタグを復元し、タグのデータ型に従ってデータを復元することにより実現する。

構造化文書の解析は、現在解析しているタグのスキーマ構造の位置・圧縮データの位置を管理するノードポインタを用意し、タグの親・子・弟への移手段、タグが保有するデータの取得手段を提供することにより実現する。

構造化文書の更新は、構造化文書の部分集合の削除と他構造化文書の挿入で行い、それぞれ削除位置に削除トークンを挿入し、挿入位置に挿入トークンを挿入し、トークンに付随するデータは全体の重複タグ数から無視するという方式で効率的に実現する。またスキーマ構造を修正することにより、圧縮データの修正なしに全体の構造を一括変換することを実現する。

構造化文書の作成は、子要素作成、弟要素作成、親へ移動、属性作成、テキスト要素作成というメソッドを用意し、これに対して開始タグ、終了タグ、属性、テキスト要素を発行することにより実現し、かつ一時バッファを用意してやり直しを許容する機能を実現する。

【0016】

【発明の実施の形態】

以下本願ではXMLを例にとって説明する。しかし、XMLに限らず、データが階層的に管理できる構造化言語であれば本願の技術を採用できる。

図6に、地理情報システムの分野で利用されるXML文書の例を示す。図6は全体をMapタグで囲み、Mapタグの子要素として複数のfeatureMemberタグが記述されているという構造を取る。Mapタグは複数の地物を含む地物集合を表し、具体的には地図を表現する。featureMemberタグの中にはRoadやHouseなど複数種類の地物を表すタグが記述され、これらが地図を構成する個々の道路や家屋情報を表現する。例えば図6のXML文書は、1/25、000縮尺の地図10km四方で、featureMemberタグが数千個から数万個になり、データサイズも10MB以上になる場合もある。個々の地物に着目すると、道路地物を表現するRoadタグの中には、道路名を表すNameタグ、車線数を示すLaneタグ、舗装状態を示すPavementタグ、道路中心線の形状情報を示すcenterLineOfタグが記述されている。centerLineOfタグ内には、折線形状を表すLineStringタグが記述され、さらにその中には座標系を示すSR属性と、折線の構成点を示す多数のcoordタグが記述されている。coordタグ内にはさらに構成点のX座標、Y座標を示すXタグ、Yタグが記述されている。例えば道路、河川、等高線など、複雑な形状の図形を表現する場合、coordタグが100点から1000点以上となる場合もある。本発明によるXMLの圧縮構造をBXML(Binary XML)と呼ぶことにする。初めにBXMLの圧縮方式とデータ構造の説明を行う。その後、(1)通常のXML文書をBXML文書に圧縮するための手段、(2)BXML文書を通常のXML文書に展開するための手段、(3)BXML文書を解析するための手段、(4)BXML文書を更新するための手段、(5)BXML文書を直接作成するための手段、について順番に説明する。

【0017】

初めに、BXMLの圧縮方式とデータ構造について説明する。

BXMLは、図7に示す次の5種類の方式を取ることによりXML文書を圧縮する。すなわち、「タグとデータの分離」圧縮701、「データ型の適用」圧縮702、「文字列の辞書化」圧縮703、「可変長数値の適用」圧縮704、「相対値の適用」圧縮705、である。これらのうち、「タグとデータの分離」701が必須の圧縮方式であり、残りの方式はオプションである。オプションの各圧縮方式は、単独あるいは図7に示す依存関係

10

20

30

40

50

に従った組合せで適用することができる。矢印の矢尻に記された方式を適用するためには、矢印の先頭に記された方式を適用する必要がある。どの組合せを適用するかは、圧縮対象のXML文書の特徴や、圧縮速度と圧縮サイズのどちらを優先するかなどのユーザ目的、XMLを圧縮するかBXMLを直接作成するかなどの使用条件に応じ、BXML作成ユーザが決定することができる。

「タグとデータの分離」圧縮701は、BXML文書を圧縮スキーマと圧縮データの二つの部分で構成し、XML文書で何度も繰り返し記述される冗長なタグ名やタグ同士の階層関係を圧縮スキーマにまとめて記述し、個別に必要なテキスト要素や属性値、タグの重複回数などを圧縮データに記述するという方法である。

「データ型の適用」圧縮702は、テキスト要素や属性値を単なる文字列として格納するのではなく、整数型・実数型・日付型などの多数の型に分類し、それぞれの型に応じて最適な格納方式を適用することにより、圧縮データのサイズを削減し、また検索速度を向上させる方式である。

「文字列の辞書化」圧縮703は、テキスト要素を単語に分解し、頻出単語を圧縮スキーマ部分で辞書化し、圧縮データ部分で頻出単語を識別子で代替する方式である。頻出する長い単語が、短い識別子で代替されることにより、圧縮データのサイズを削減することができる。

「可変長数値の適用」圧縮704は、上記の整数型・実数型のデータを、例えばJava言語（Javaは登録商標である）のshort型整数（2byte）、long型整数（4byte）、double型浮動小数点数（8byte）のように固定長数値を利用するのではなく、0から100までは1byte、101から10000までは2byteのように、数値の大きさにより格納サイズが可変となる数値格納方式を適用することにより、圧縮データのサイズを削減する方式である。

「相対値の適用」圧縮705は、特に多数の地図図形の座標をXML文書で管理する場合に適用し、座標をある起点からの相対値で表現する方式である。前述の「可変長数値の適用」704と併用し、相対値表現により小さくなった数値を可変長数値により少ないサイズで記述することにより、圧縮データのサイズを削減することができる。

【0018】

図6を用いて1番目の圧縮方式である「タグとデータの分離」圧縮701について説明する。

【0019】

図6のXML文書は、<Road>などの開始タグ、</Road>などの終了タグなどで構成されるタグ部分と、「Route 6」など、タグとタグの間に挟まれたデータ部分に分けられる。ここで、<Road>タグが多数回繰り返される場合を考えると、データ部分はそれぞれ異なるが、タグ部分は同じ構造が繰り返されるという冗長な構造となっていることが分かる。一般にXML文書は、このような冗長なタグ部分が全体サイズの半分近くを占めている。「タグとデータの分離」圧縮701では、タグ部分とデータ部分を分離し、最低限必要なデータ部分は個別に記述するが、冗長なタグ部分は1回しか記述しない、という方式を取ることで、XML文書のサイズを削減する。

【0020】

図8を用いて、「タグとデータの分離」圧縮701について詳細に説明する。

【0021】

BXML文書803において、分離したタグ部分を表現する構造を圧縮スキーマ804、データ部分を表現する構造を圧縮データ810と呼ぶことにする。XML文書を構成する情報のうち、圧縮スキーマに記述する情報と、圧縮データに記述する情報は次のように分類される。

【0022】

圧縮スキーマに記述される情報は次の7つである。(1)タグ名称805、(2)重複フラグ806、(3)スキーマ型807、(4)子要素保有数808、(5)保有子要素一覧、(6)属性保有数809、(7)保有属性一覧。

10

20

30

40

50

(1) タグ名称 805 は、XML 文書に出現するタグの名称である。これを圧縮スキーマに記述することにより、XML 文書の半分近くを占める開始タグ・終了タグを削減することができ、圧縮サイズは約半分になる。

(2) 重複フラグ 806 は、そのタグが重複するか、すなわち連続した兄弟の位置に同じタグが複数並ぶ可能性があるかどうかを示すフラグである。これは「重複」と「単独」の2つの値のいずれかを取り、「単独」の場合はそのタグが必ず1回のみ出現する、「重複」はそれ以外の場合であり、そのタグが0回以上複数回出現する可能性があることを示す。なお、タグが省略される可能性がある場合も「重複」となる。一方、個々のタグが何個出現したかという重複数の情報はデータに依存するため、圧縮データ部分に記述し、圧縮スキーマには記述しない。なお、重複フラグ 806 が「単独」と指定されたタグは、重複数は必ず1回であることが保証されているため、圧縮データ部分に記述する必要はない。

(3) スキーマ型 807 は、そのタグがどのような性質を持つのかを記述する。例えばそのタグが子要素を持つか、あるいは子要素を持たず、テキスト要素を持つかを指定する。さらに子要素を持つ場合、指定した子要素群を順番通りに持つか (SEQUENCE 型)、順不同に持つか (MCHOICE 型)、いずれか一つのみを持つか (CHOICE 型)、あるいはテキスト要素を持つ場合、任意文字列が入るか、整数や実数のみしか入らないか、などを細かく指定することができる。これは XML 文書に明示的に記述されていない情報であるが、本発明により効率のいい圧縮を行うため、本発明の圧縮工程で付与される情報である。なお、XML Schema 等の XML 文書定義言語で圧縮対象 XML 文書の構造が提供される場合は、ここからスキーマ型 807 を取得することも可能である。

(4) 子要素保有数 808 は、そのタグに出現する可能性のある全子要素の種類数である。これは種類の数であり、実際の子要素数ではないため、個々のタグで子要素が省略されたり子要素が重複して出現したりして、実際の子要素数と異なっても構わない。個々のタグにおいてどの子要素が省略されたか、またどの子要素が何個重複して出現したかという情報はデータに依存するため、圧縮データ部分に記述し、圧縮スキーマには記述しない。

(5) 保有子要素一覧は、子要素保有数 808 で指定した数の、子要素情報の集合である。それぞれの子要素情報も再帰的に、上述した (1) タグ名称、(2) 重複フラグ、(3) スキーマ型、(4) 子要素保有数、(5) 保有子要素一覧、(6) 属性保有数、(7) 保有属性一覧を保有する。

(6) 属性保有数 809 は、そのタグに出現する可能性のある全属性の種類数である。これは種類の数であり、実際の属性数ではないため、個々のタグで属性が省略されてこの属性保有数よりも少ない属性を持つ場合があっても構わない。個々のタグにおいてどの属性が省略されたかという情報はデータに依存するため、圧縮データ部分に記述し、圧縮スキーマには記述しない。

(7) 保有属性一覧は、属性保有数 809 で指定した数の、属性情報の集合である。属性情報とは、属性キーと属性値であるが、属性キーをタグ名称、属性値をテキスト要素と置き換えることにより、特殊な子要素と考えることができる。従って属性情報は、上述した子要素情報と同様に (1) キー名称、(2) 重複フラグ、(3) スキーマ型、を持つ。属性は階層構造を持たないため子要素保有数は必ず0、属性は属性を持たないため、属性保有数も必ず0である。そのため、子要素保有数 808、保有子要素一覧、属性保有数 809、保有属性一覧は省略可能である。

【0023】

圧縮データ 810 には、次の4つの情報が記述される。(1) テキスト要素、(2) 属性値、(3) タグ重複数、(4) 選択タグ識別子。

(1) テキスト要素は、XML 文書においてタグとタグの間に囲まれた部分である。XML 文書のテキスト要素がそのまま記述されたり、図7に示す他の圧縮方式 702 ~ 705 により圧縮されたテキスト要素が記述されたりする。

(2) 属性値は XML 文書の保有する属性の値である。これもテキスト要素と同様、XML 文書の属性値がそのまま記述されたり、図7に示す他の圧縮方式 702 ~ 705 により圧縮された属性値が記述されたりする。

10

20

30

40

50

(3) タグ重複数は、XML 文書の兄弟の位置に存在する、同名のタグがいくつ連続で続くかを示す数値である。これは XML 文書により変化し、また一つの XML 文書内でも場所により変化するため、圧縮スキーマではなく圧縮データに直接記述する。

(4) 選択タグ識別子は、子要素として複数種類のタグのどれか一つが出現する可能性がある場合、どのタグが出現したかを示すための識別子である。圧縮スキーマ 804 の子要素一覧として、あるタグの子要素として出現する可能性のあるタグ一覧を記述するが、その記述順の通し番号を識別子とすることができる。属性の場合も同様に、複数種類のキーのどれか一つが出現する可能性がある場合、どのキーが出現したかを選択タグ識別子として記述しても構わない。

【0024】

XML 文書は、図 8 に示す 7 個のパターンの組み合わせで表現できる。(1) 子要素パターン 811、(2) 兄弟要素パターン 812、(3) 重複要素パターン 813、(4) 省略要素パターン 814、(5) 選択要素パターン 815、(6) 重複選択要素パターン 816、(7) 属性パターン 817、である。以下、上記 7 個のパターンを例に取り、XML 文書がどのような構造の圧縮スキーマ、圧縮データで表現されるかを説明する。

(1) 子要素パターン 811 は、タグの中に子要素が存在するパターンである。XML サンプル 802 では、タグ A の子要素としてタグ B が、タグ B の子要素としてタグ C が格納されている。タグ C は子要素を持たず、その代わりにテキスト要素として文字列 text を持つ。

【0025】

圧縮スキーマでは、最初にルートタグ A の情報を記述する。タグ A は 1 回のみ出現するため重複タグ 806 は「単独」となる。また、タグ B のみを子要素として持つため、スキーマ型 807 は「SEQUENCE」、子要素保有数 808 は 1 個となる。タグ B、タグ C の情報はタグ A の情報に連続して記述する。タグ B も 1 回のみ出現するため重複タグ 806 は「単独」、また、タグ C のみを子要素として持つため、スキーマ型 807 は「SEQUENCE」、子要素保有数 808 は 1 個となる。タグ C も 1 回のみ出現するため重複タグ 806 は「単独」、また、タグ C は子要素を持たず、テキスト要素のみを持つため、スキーマ型 807 は「文字列」となる。子要素パターンの圧縮データは、タグ C のテキスト要素である text 1 のみを記述する。これは、圧縮スキーマでの定義で、タグ A、タグ B がテキスト要素を持たないことを示しているため、最初に出現するテキスト要素がタグ C のものであるということが自明だからである。

(2) 兄弟要素パターン 812 は、一つのタグの中に複数の異なるタグが順番に格納されるパターンである。XML サンプル 802 では、タグ A の中にタグ B とタグ C が順番に格納されている。兄弟要素パターンを示す圧縮スキーマでは、その兄弟タグの親タグのスキーマ文に続けて、兄弟タグをその出現順に記述する。タグ A は 1 回のみ出現するため重複タグ 806 は「単独」となる。また、タグ B、タグ C を子要素として持つため、スキーマ型 807 は「SEQUENCE」、子要素保有数 808 は 2 個となる。タグ B、タグ C の情報はタグ A の情報に連続して記述する。タグ B も 1 回のみ出現するため重複タグ 806 は「単独」、子要素を持たず、テキスト要素のみを持つため、スキーマ型 807 は「文字列」となる。タグ C も 1 回のみ出現するため重複タグ 806 は「単独」、子要素を持たず、テキスト要素のみを持つため、スキーマ型 807 は「文字列」となる。

【0026】

兄弟要素パターンを示す圧縮データでは、テキスト要素が出現順に記述される。圧縮スキーマによれば各テキスト要素がどのスキーマに対応するか自明だからである。

(3) 重複要素パターン 813 は、一つのタグの中に複数の同じタグが格納されるパターンである。XML サンプル 802 では、タグ A の中に 2 個のタグ B が重複要素として格納され、その後タグ C が兄弟要素として格納されている。

【0027】

重複要素パターンを示す圧縮スキーマでは、重複する要素の重複フラグに「重複」を設定する。タグ A は 1 回のみ出現するため重複タグ 806 は「単独」となる。また、タグ B、

10

20

30

40

50

タグCを子要素として持つため、スキーマ型807は「SEQUENCE」、子要素保有数808は2個となる。タグBは2回出現するため重複タグ806は「重複」、子要素を持たず、テキスト要素のみを持つため、スキーマ型807は「文字列」となる。タグCは1回のみ出現するため重複タグ806は「単独」、子要素を持たず、テキスト要素のみを持つため、スキーマ型807は「文字列」となる。

【0028】

重複要素パターンを示す圧縮データでは、重複するタグの位置に重複数を格納する。本例の場合、タグBが2回重複して出現しているため、重複数として2を格納する。テキスト要素は出現順に記述する。スキーマの解釈によれば、タグBの重複フラグは「重複」であるため、圧縮データ先頭の整数「2」はタグBの重複数であることが分かる。それに続く文字列text1、text2は、タグBのスキーマ型が「文字列」であり、重複数が2であることから、それぞれ重複する2つのタグBのテキスト要素であることが分かる。さらに続く文字列text3は、タグBの次に定義されているタグCのテキスト要素であることが分かる。

10

(4) 省略要素パターン814は、あるタグの子要素が、ある場所では存在するが、別の場所では省略される場合があるパターンである。これは、省略されたタグは重複フラグが「重複」のタグであり、重複数が0であるとみなすことにより、重複要素パターンと同一パターンとみなすことができる。XMLサンプル802では、最初のタグAの子要素としてタグB、タグCが存在するが、最初のタグAではタグCが省略されており、次のタグAではタグCが省略されていない。

20

【0029】

省略要素パターンを示す圧縮スキーマでは、省略する要素の重複フラグに「重複」を設定する。XMLサンプルでは、タグCが省略される場合があるため、タグCのスキーマ文に「重複」を設定する。

タグAは2回のみ出現するため重複タグ806は「重複」となる。また、タグB、タグCを子要素として持つため、スキーマ型807は「SEQUENCE」、子要素保有数808は2個となる。タグBは必ず1回のみ出現するため重複タグ806は「単独」、子要素を持たず、テキスト要素のみを持つため、スキーマ型507は「文字列」となる。タグCは省略される場合もあるため、重複タグ806は「重複」、子要素を持たず、テキスト要素のみを持つため、スキーマ型807は「文字列」となる。

30

【0030】

省略要素パターンを示す圧縮データでは、最初にタグAの重複数である2を格納する。次に最初のタグAの子要素であるタグBのテキスト要素「text1」を記述する。次にタグCの重複数を記述する。この場合タグCは省略されているため、重複数は0である。次に2番目のタグAの子要素であるタグBのテキスト要素「text2」を記述する。次にタグCの重複数を記述する。この場合タグCは省略されていないため、重複数は1である。引き続きタグCのテキスト要素「text3」を記述する。スキーマの解釈によれば、タグAの重複フラグは「重複」であるため、最初の整数「2」はタグAの重複数であることが分かる。タグAの最初の子要素であるタグBの重複フラグは「単独」であるため、次の文字列text1はタグBの保有である。タグAの次の子要素であるタグCの重複フラグは「重複」であるため、次の整数「0」はタグCの重複数である。次のテキスト要素text2は二番目のタグAの先頭子タグBのテキスト要素である。次の整数「1」は、タグCの重複フラグであり、値が1であるため続く文字列text3はタグCのテキスト要素であることが分かる。

40

(5) 選択要素パターン815は、あるタグの子要素が複数の候補のうちの一つであるというパターンである。XMLサンプル802では、最初のタグAの子要素としてタグBが来ており、次のタグAの子要素としてタグCが来ている。

【0031】

選択要素パターンを示す圧縮スキーマでは、子要素が選択要素となるタグに対し、スキーマ型を「CHOICE」と設定する。そして選択される可能性のある子要素情報を連続し

50

て記述する。なお、選択される子要素の重複数は必ず「単独」とする。

タグ A は 2 回のみ出現するため重複タグ 806 は「重複」となる。また、タグ B、タグ C のいずれか一つを子要素として持つ可能性があるため、スキーマ型 807 は「CHOICE」、子要素保有数 808 は 2 個となる。タグ B は選択される子要素のため重複タグ 806 は「単独」、子要素を持たず、テキスト要素のみを持つため、スキーマ型 507 は「文字列」となる。タグ C も同様に、重複タグ 806 は「単独」、子要素を持たず、テキスト要素のみを持つため、スキーマ型 807 は「文字列」となる。

【0032】

選択要素パターンを示す圧縮データでは、タグ A の重複数の後に、実際に到来したタグの識別子を記述し、その後到来したタグのデータを記述する。ここでいうタグの識別子は、同じ階層のタグ群から選択できればよいため、圧縮スキーマ部で定義したタグの出現順序を 0 で始まる通し番号として記述するので構わない。例えば XML サンプルでは、A の重複数の後、タグ B が出現したことを示すタグ B の識別子「0」を格納し、タグ B のテキスト要素を続いて格納する。その後、次にタグ C が出現したことを示す識別子「1」を格納し、続いてタグ C のデータを記述する。

(6) 重複選択要素パターン 816 は、重複要素と選択要素が同時に生じるパターンである。XML サンプル 802 では、タグ A の子要素として、タグ B とタグ C が順不同で複数回出現している。子要素が重複選択要素となるタグに対し、スキーマ型を「MCHOICE」と設定する。そして選択される可能性のある子要素情報を連続して記述する。なお、選択される子要素の重複数は必ず「単独」とする。

タグ A は 1 回のみ出現するため重複タグ 806 は「単独」となる。また、タグ B、タグ C が順不同で複数回出現する可能性があるため、スキーマ型 807 は「MCHOICE」、子要素保有数 808 は 2 個となる。タグ B は選択される子要素のため重複タグ 806 は「単独」、子要素を持たず、テキスト要素のみを持つため、スキーマ型 507 は「文字列」となる。タグ C も同様に、重複タグ 806 は「単独」、子要素を持たず、テキスト要素のみを持つため、スキーマ型 807 は「文字列」となる。

【0033】

重複選択要素パターンを示す圧縮データでは、タグ A の保有する子タグ群の総数を記述したあとに、選択要素パターンと同様、実際に到来したタグの識別子と、タグのデータのペアを、子タグ群の総数と同じ回数、連続して記述する。例えば XML サンプルでは、タグ A の子要素であるタグ B、タグ C を合わせた重複数「4」を記述し、最初に到来したタグ B の識別子「0」、タグ B のテキスト要素「text1」、2 番目に到来したタグ C の識別子「1」、タグ C のテキスト要素 text2、3 番目に到来したタグ B の識別子「0」、タグ B のテキスト要素「text3」、4 番目に到来したタグ C の識別子「1」、タグ C のテキスト要素 text4、を順に記述する。

(7) 属性パターン 817 は、タグに属性が格納されるパターンである。XML サンプル 802 では、タグ A の属性としてキー a、値 text1 という属性が格納されている。

【0034】

属性パターンを示す圧縮スキーマでは、属性を含むタグのスキーマ文に、属性保有数 809 として格納される属性の数を記述し、さらに保有する属性の情報を連続して記述する。属性を示すスキーマ文は、子要素を示すスキーマ文とほとんど同じ構造を持つ。タグ名として、属性のキーを持つ。属性は省略可能だが同じキーの属性が複数来ることはない、すなわち重複数は必ず 0 か 1 であるため、重複フラグ 506 には特別に「属性」という値を設定しても構わない。また、属性は子要素や属性を持たないため、子要素保有数 508、属性保有数 509 は省略して構わない。なお、属性を示すスキーマ文は、子要素を示すスキーマ文に先行して記述する。

【0035】

タグ A は 1 回のみ出現するため重複タグ 806 は「単独」となる。タグ B のみを子要素として持つため、スキーマ型 807 は「SEQUENCE」、子要素保有数 808 は 1 個となる。また、属性 a を持つため、属性保有数 809 は 1 個となる。属性 a の重複タグ 80

6は「属性」、属性のスキーマ型807は「文字列」となる。タグBは1回のみ出現するため重複タグ806は「単独」、また、タグBは子要素を持たず、テキスト要素のみを持つため、スキーマ型807は「文字列」となる。

【0036】

属性パターンを示す圧縮データでは、属性の値とテキスト要素を、その出現順に記述する。属性が出現する可能性があるかどうかは圧縮スキーマに記述され、属性が省略されたかどうかは圧縮データに記述され、また属性は子要素に先行して記述するというルールから、属性の値を一意に取得することができる。

図8の圧縮データ810の構造について、図9のサンプルXML文書901～905を用いて詳細に説明する。

【0037】

XML文書901では、複数のタグBが存在し、さらにタグBの中に複数のタグCが存在する。これに対する圧縮データは906の構造を取る。最初にタグBの重複数Nb（この場合は重複数2）を記述し、その後タグBの子要素の圧縮データ913、914を記述する。圧縮データ913は重複する二つのタグCを表現するため、再起的に最初のタグCの重複数Nc1（この場合は重複数2）を記述し、その後最初のタグCのテキスト要素であるc1、次のタグCのテキスト要素であるc2を記述する。c1、c2は単なる文字列が記述される場合もあるし、図7の圧縮方式で説明したデータ型の利用702、文字列の辞書化703、可変長数値の適用704、相対値の適用705を適用し、さらに圧縮された構造が格納される場合がある。c1、c2がどの形式の構造であるかはテキスト要素が含まれるタグの圧縮スキーマ804のスキーマ型807を参照することにより判断できる。また圧縮データのある位置の情報がNbなどの重複数であるか、c1、c2などのテキスト要素であるかの判断も、圧縮スキーマ804の重複フラグ806を参照することにより判断できる。

【0038】

XML文書902は、XML文書901に対し、タグCの兄弟要素としてタグDが追加された構造を持つ。これに対する圧縮データは907の構造を取る。最初にタグBの重複数Nb（この場合は重複数2）を記述し、その後タグBの子要素の圧縮データ915、916を記述する。圧縮データ915はさらにタグCとタグDから構成されているため、出現順番通りタグCの圧縮データ917を記述した後、タグDの圧縮データ918を記述する。圧縮データ917は、タグCが重複要素であるため、最初にタグCの重複数Nc1を記述し、2つのタグCのテキスト要素c1、c2を記述する。同様に圧縮データ918は、タグDの重複数Nd1を記述した後、2つのタグDのテキスト要素d1、d2を記述する。圧縮データ916も同様である。

【0039】

XML文書903は、XML文書902の特別な構造であり、タグC、タグDが重複せず、単独で出現する。これに対する圧縮データは、圧縮データ907と同様な考え方で908の構造を取ってもよいが、タグC、タグDは必ず1個しか出現しないことが明らかになっている場合、重複数Nc1、Nd1などは全て1個となり、冗長である。そのため909の構造に示すように、重複数を取り除いた構造を取ることができる。あるタグが重複数を取るかどうかの判断は、図8におけるそのタグに対応する圧縮スキーマ804において、重複フラグ806が「単独」の場合重複数は省略され、「重複」の場合重複数が記述されるという方法で判定できる。圧縮データ909の構造は、最初のタグBの重複数Nbを除き、全てがXML文書のテキスト要素のみであり、XML文書のタグ階層構造を示す情報がほとんど除去された構造となっている。多くのXML文書の特徴として、単独要素の数が重複要素の数よりも多いため、重複フラグ806が「単独」の場合重複数を省略する、という圧縮方式は、効率のいい圧縮が期待できる。

【0040】

XML文書904では、タグBの中にタグC、タグD、タグEのいずれか一つが格納される。これに対する圧縮データは、圧縮データ907と同様な考えかたで910の構造を取

10

20

30

40

50

ってもよいが、 $N_{c i}$ 、 $N_{d i}$ 、 $N_{e i}$ ($i = 1, 2, 3, \dots$) の全ての組において、 $N_{c i}$ 、 $N_{d i}$ 、 $N_{e i}$ のいずれか一つを除いて全て 0 個であることが保証されるため、冗長である。そのため 9 1 1 の構造に示すように、どの種類のタグが出現するかという識別子 K_c 、 K_d 、 K_e の何れかを記述し、その後その種類のタグに関するデータのみを記述する、という方法を取ることににより、圧縮データのサイズを低減することができる。ここで識別子は、圧縮スキーマにおいて定義されている子要素 C、D、E の出現順の通し番号を付ければよい。あるタグ B の子要素を表現する圧縮データが 9 1 1 の構造になっていることを判定するためには、タグ B に対応する圧縮スキーマにおいて、スキーマ型 8 0 7 が S E Q U E N C E 型の場合 9 1 0 の構造、C H O I C E 型の場合 9 1 1 の構造である、という形で判定できる。

10

【0041】

X M L 文書 9 0 5 では、タグ B の中にタグ C、タグ D、タグ E が順不同に 0 個以上格納される。これに対する圧縮データは、タグ B の子要素の圧縮データとして、最初に順不同に出現するタグ C、タグ D、タグ E の総数 $N_{c d e}$ (この場合 6) を記述し、その後圧縮データ 9 1 1 と同様な方式で出現するタグの識別子、テキスト要素の組み合わせを記述する、という方式で記述できる。圧縮データ 9 1 2 の構造を取ることを判定するためには、タグ B に対応する圧縮スキーマにおいて、スキーマ型 8 0 7 が M C H O I C E 型の場合 9 1 2 の構造である、という形で判定できる。

【0042】

次に、図 10 を用いて、X M L 圧縮構造の 2 番目の方法であり、最初のオプションである「文字列の辞書化」方式について説明する。本圧縮方式は、圧縮データ領域を対象とする。前述したように、圧縮データ領域は主に X M L 文書のテキスト要素を記述する部分である。ただし図 8 の 7 番目のパターン「属性」で示したとおり、本発明の圧縮方式では属性部分をテキスト要素を持つ子要素と同一の方式で管理している。そのため本発明におけるテキスト要素の圧縮方式は、属性の値部分の圧縮にも容易に適用することができる。

20

【0043】

X M L 文書のデータ部にも何度も繰り返し出現する冗長な部分が存在する。例えば図 10 の X M L 文書 1 0 0 1 の例では、 $\langle T y p e \rangle$ タグの中に道路区分として、「国道」「県道」の 2 種類のいずれかが記述されている。一般に、同じ名称のタグが多量に存在し、そのタグの持つテキスト要素として取りうる値の種類がタグの総数よりも十分小さければ、そのテキスト要素を表す識別子を用意し、その識別子でテキスト要素を置き換えることにより冗長性が減少し、圧縮効率が向上する。

30

【0044】

B X M L 文書 1 0 0 2 では、圧縮スキーマ 1 0 0 3 において、指定したタグのテキスト要素として含まれる文字列の候補を辞書として持っておき、圧縮データ 1 0 0 4 において文字列を直接記述する代わりに、圧縮スキーマで保持している辞書データの識別子を持たせることができる。

【0045】

テキスト要素を列挙させる場合において、例外的なテキスト要素が存在する場合がある。たとえば 1 0 0 件のテキスト要素が存在し、その中の 9 0 件は 5 種類の値のいずれかしか取らないが、残りの 1 0 件はまったく異なる値を取る、という場合がある。指定したタグの全てのテキスト要素を辞書に登録し、その識別子でしか参照しない場合、このような例外的なテキスト要素を圧縮スキーマに登録するのは無駄である。通常のテキスト要素と辞書への識別子参照を混在することにより、この問題を解決することができる。

40

【0046】

テキスト要素の種類がそのタグの総数に比べ十分に小さくない場合でも、そのテキスト要素の部分を切り出せば十分まとめることが可能な場合がある。この例の代表例としては、X M L の識別子 (I D) が挙げられる。I D が世界中で完全にユニークであるための仕掛けとして、W 3 C では U R I (U n i f o r m R e s o u r c e I d e n t i f i e r) という形式を提唱している。これはある組織が I D を決定する場合において、I D を

50

その組織の保有するURL (Uniform Resource Locator) にその組織内で固有のIDを組み合わせるという方式である。例えば、「http://crl.hitachi.co.jp/gis/」というURLを持つ組織が「f00001」という固有IDを持つ時、URIは「http://crl.hitachi.co.jp/gis/f00001」となる。この方式の欠点として、IDが非常に長いものになるという点が挙げられる。例えば地図上の道路や家屋など、数千個の地物全てにURI形式でIDを付加すると、IDの大部分を占めるURL部分が同じという冗長な構造でデータ量が肥大する、という問題がある。ここでIDをURLの部分と組織内固有IDの部分に分割し、URLの部分を辞書に登録することによりデータ量の圧縮が実現できる。また、XML内部のデータをIDをキーとして高速に検索する場合、一般に文字列比較よりも数値比較の方が高速なため、IDは文字列ではなく、数値であるほうが都合がよい。このような場合もURLの部分を辞書化しておき、組織固有IDを数値で管理しておけば、IDの高速検索が実現できる。

10

【0047】

本発明による「文字列の辞書化」では、テキスト要素を辞書化しないタグ、テキスト要素を全て辞書化するタグ、部分的に辞書化するタグが混在するという構造を取る。図11に辞書化のパターンを示す。

【0048】

パターン1「インライン文字列」では、辞書を用いず、入力テキスト要素をそのまま記載する。圧縮データには、辞書に登録されたIDか文字列かを判定するため、「次のデータは文字列であり、終端する」ことを示すトークン「T_WORD」を記述し、その後文字列を記載する。

20

【0049】

パターン2「辞書ID」では、テキスト要素を全て辞書に格納する例である。辞書に識別子と登録文字列のペアを格納しておき、圧縮データには、「次のデータは辞書に登録された文字列の識別子であり、終端する」ことを示すトークン「T_ID」を記述し、その後識別子である数値「1」を記載する。

【0050】

パターン3「辞書ID + インライン文字列」では、最初の文字列「国道」を辞書に格納し、それ以降の文字列「20号線」を圧縮データに記載する。それを表現するため、圧縮データには、「次のデータは辞書に登録された文字列の識別子であり、テキスト要素は継続する」ことを示すトークン「C_ID」、辞書の識別子である数値「1」、「次のデータは文字列であり、終端する」ことを示すトークン「T_WORD」を記述し、その後文字列を記載する。

30

【0051】

パターン4「辞書ID + 数値」では、「T_WORD」の代わりに「次のデータは数値であり、終端する」ことを示すトークン「T_NUM」を記述し、その後数値を記載する。

【0052】

パターン5「インライン文字列 + 辞書ID」では、「次のデータは文字列であり、テキスト要素は継続する」ことを示すトークン「C_WORD」、文字列、「次のデータは辞書に登録された文字列の識別子であり、終端する」ことを示すトークン「T_ID」を記述し、その後識別子である数値「1」を記載する。

40

3番目の圧縮方式である、「データ型の適用」方式について説明する。本圧縮方式は、圧縮データ領域を対象とする。図4に示すように、「タグとデータの分離」方式では、テキスト要素しか持たないタグのスキーマ型を「文字列型」と設定していた。本圧縮方式では、従来の「文字列型」に加え、さらに7種類に拡張する。すなわち、「整数型」、「実数型」、「2次元整数型」、「2次元実数型」、「2次元整数配列型」、「2次元実数配列型」、「日付型」である。

整数型や実数型は、テキスト要素に整数や実数のみしか格納されていない場合に適用する

50

。「文字列型」に対応する圧縮データが、ASCIIやSHIFT-JIS、UNICODE等の文字コードで記載されているのに対し、たとえば計算機のメモリ上での整数格納形式や、不動小数点数格納形式(IEEE754)などの、より計算機に利用しやすい形で格納する。これは圧縮データのサイズ削減に貢献する。

例えば、整数32767を文字列として格納する場合、「3」「2」「7」「6」「7」「¥0」の6つの文字となる。最後の¥0は文字列の終端を表すヌル文字である。ASCII文字セットでは文字は1BYTEで表現されるため、整数32767は6BYTEの領域を占める。UNICODE文字セットでは文字は2BYTEで表現されるため、整数32767は12BYTEの領域を占める。一方、Java言語のshort型では、数値32767は0x7F、0xFFの2BYTEで表現される。同様に実数3.1415926は、「3」「.」「1」「4」「1」「5」「9」「2」「6」「¥0」となり、ASCII文字セットで10BYTE、UNICODE文字セットで20BYTEの領域を占めるが、Java言語のfloat型では4BYTEで表現される。

10

【0053】

また本方式は、XMLの解析処理速度を向上させるという効果もある。XML文書を利用するXMLアプリケーションは、例えば数値の大小判定を行う、四則演算を行うなど、数値を利用する場合、数値を計算機のメモリ上で理解できる数値型で保持しておく必要がある。XMLでは数値は文字列として格納されているため、内部処理として文字列を数値型に変換する処理が必要となり、処理時間が必要となる。一方、本方式で圧縮したBXMLをそのまま解析することにより、文字列・数値変換が不要となるため、処理速度が向上する。

20

例えば地図をあらわしたXML文書は、その大部分が図形を示す点・折れ線・ポリゴンの座標となる。データの大部分を占める数値情報を文字列ではなく数値として管理することにより、大きなデータ量削減効果が得られる。また用途として、地図を示すXML文書全体を解析して全ての座標を取得して画面上に描画する、XML文書全体を解析して、指定した矩形領域内に存在する地物を抽出する、などの処理を行うにあたり、数値情報を文字列ではなく数値として管理することにより、大きな処理速度向上効果が得られる。

タグのスキーマ型として、整数型、実数型について説明したが、残る「2次元整数型」、「2次元実数型」、「2次元整数配列型」、「2次元実数配列型」、「日付型」について説明する。

30

「2次元整数型」は二つの整数の組、「2次元実数型」とは二つの実数の組を表現する構造である。これは、圧縮元となるXML文書において、例えば<座標>x、y</座標>のように、数値がテキスト要素として分離していない場合に利用する。上記のような表記のXML文書は、「整数型」「実数型」では表現できない。これを文字列型として表現するとデータサイズ、解析速度が悪化する。このような構造は座標を管理する地図用XMLでは頻出する可能性があるため、特別な型として定義する。これは圧縮スキーマにおいて「2次元整数型」あるいは「2次元実数型」と指定し、付加情報として二つの数値を分離するセパレータ文字(空白やカンマ文字)を定義し、圧縮データにおいて数値を連続して並べることにより実現できる。2次元整数、2次元実数は、他にも地物の生成時間と消滅時間のような時間的な期間を表すなど、様々なデータに利用できる。同様に、立体図形を表現するための座標(x、y、z)、あるいは空間上の移動量・移動方向を表すベクトル(x、y、z)を表現する「3次元整数型」、「3次元実数型」への拡張も可能である。

40

「2次元整数配列型」は2次元整数型の0個以上の組、「2次元実数配列型」は2次元実数型の0個以上の組を表す。これは例えば<座標列>x1、y1 x2、y2 x3、y3</座標列>のように、数値がテキスト要素として分離していない場合に利用する。このような構造は座標を管理する地図用XMLでは頻出する可能性があるため、特別な型として定義する。これは圧縮スキーマにおいて「2次元整数配列型」あるいは「2次元実数配列型」と指定し、付加情報とし座標x、yを分離するセパレータ文字、座標と座標を分離するセパレータ文字(上記の例ではカンマと空白)を定義し、圧縮データにおいて最初に配列の個数を、その後数値を連続して並べることにより実現できる。これは同様に「

50

3次元整数配列型」、「3次元実数配列型」への拡張も可能であり、<数値列> a 1、a 2、a 3、a 4 </数値列>のように整数や実数の配列である「1次元整数配列型」、「1次元実数配列型」も可能である。

【0054】

「日付型」は日付を指定するために利用する。W3CではW3C-DTFとして、またISO8601で日付の書式が定義されており、例えば2002年12月1日0時0分0秒を、「2002-12-01T00:00:00」という文字列で表現する。一方、計算機内部で広く使われる日付の管理方式として「1970年1月1日午前0時(UTC)を起算点とする通算の秒数」で日付を管理する場合がある。この方式は日付のデータに4byteしか利用せず、また日付同士の比較、日付と期間の加算が容易であるというという特徴を持つ。では、圧縮スキームにおいて「日付型」と指定し、圧縮データにおいてこの通算秒などよりコンパクトな日付の格納を行うことにより実現できる。

4番目の圧縮方式である、「可変長数値の適用」方式について説明する。本圧縮方式は、圧縮データ領域を対象とする。本圧縮方式では、整数と実数を、本発明で独自に定義した、MBI(Multi Byte Integer:可変長整数)、MBF(Multi Byte Float:可変長実数)で表現する。

従来の計算機内で利用されている整数・実数の管理方式は、char型整数は1byte、short型整数は2byte、long型整数は4byte、float型実数は4byte、double型実数は8byteと固定されている、固定長数値である。本方式の可変長数値では、例えば数値10は1byte、100は2byte、10、000は3byte、1、000、000は4byteと可変長で管理できる。

【0055】

固定長数値の欠点は、表現できる桁数が限られるという点にある。例えばchar型整数は127まで、short型整数は32、767まで、long型整数は2、147、483、646までしか表現できない。float型実数は仮数部の精度7桁、double型実数は精度15桁までしか表現できない。XML文書を「データ型の適用」を用いて圧縮する時、固定長数値を利用する場合、どの数値型を選択する必要があるかを判定する必要がある。全ての整数をlong型の4byteで管理するとデータ圧縮上効率が悪い。利用者が出現する整数の取りうる範囲を考慮して適切な数値型を割り当てるのは手間がかかり、また汎用的ではない。また例外的に極端に大きな数値が出現した時に破綻する可能性がある。本発明の可変長数値を利用することにより、小さな整数や精度の低い実数は必要最低限の領域で格納でき、またメモリ領域やファイルシステム領域の制限の許す限り大きな整数、精度の高い実数を格納することができる。

【0056】

図12に本発明の可変長整数であるMBIの構造を示す。MBIは、W3Cの標準規格の一つであるWAP Binary XMLで定義されるMBI(以降WAP-MBIと呼ぶ)を、負数が表現できるように拡張したものである。

【0057】

図12(a)は、WAP-MBIにおける、整数を表現する複数のバイト列の、一つのバイトに着目した構造である。最上位ビットは継続フラグであり、1の場合は後続バイトにデータが続くことを示し、0の場合はそのバイトでデータが終端することを示す。残りの7bitに数値の情報が格納される。数値の情報を7bitずつ区切り、格納する。

【0058】

図12(b)は、本発明によるMBIであり、WAP-MBIを、負数が表現できるように拡張したものである。

最上位ビットは継続フラグというのはWAP-MBIと同様である。先頭バイトは例外的にデータ領域が5bitであり、継続フラグとデータ領域の間に余分の2bitが存在する。最上位ビットの次のビットは常に0が格納され、MBIであることを識別するために存在する。次のビットは符号フラグであり、0または正数の場合は0、負数の場合は1が格納される。データ領域は正数の絶対値が2進数表現で格納される。

10

20

30

40

50

【0059】

本発明の可変長実数であるMBFは、実数を仮数部と指数部の二つの整数に分割し、それぞれ二つのMBIとして連続して格納することにより表現する。MBFと二つの連続したMBIをそれぞれで区別する方法はない。圧縮スキーマ部でスキーマ型が「整数型」の場合、圧縮データから1つのMBIを取得して、これを整数として認識する。スキーマ型が「実数型」の場合、圧縮データから連続する2つのMBIを取得して、これをMBFとして認識する、という方式で区別を行う。「2次元整数型」、「2次元実数型」、「2次元整数配列型」、「2次元実数配列型」等の型も同様な処理を行う。

【0060】

5番目の圧縮方式である、「相対値管理」について説明する。本圧縮方式は、圧縮データ領域を対象とする。 10

特に地図図形の座標を多数XML文書で管理する場合に適用され、座標をある起点からの相対値で表現する方式である。

地図上の地物の位置や形状を一意に表現するためには、緯度経度などの絶対座標が用いられる。しかし地図上の多数の座標を全て絶対座標で表現すると、データ量が増大する。ある領域内に限定した地図をみると、その上位桁はほとんど同じで、下位桁のみ変化しているという冗長性のある構造を取っている。

従来は地図管理フォーマットでは、図面という概念を導入して、座標を図面原点からの相対座標で管理するという方式を取っている。この目的は、(1)座標格納領域のデータ量を削減する、(2)座標をshort型など比較的小さいサイズの固定長数値で表現できるようにする、という二点である。図面管理方式の欠点として、座標を固定長数値で表現するため、図面のサイズがある程度固定されてしまう点にある。図面として規定されたサイズを超える図形や、図面をまたがる図形は、複数の図形断片に分断されてしまうことになる。 20

文字列ベースであるXMLでは固定長数値の制約はないため、自由に大きな桁数の座標が記述できる。そこでXMLをベースにした空間データフォーマットは、座標を緯経度などの絶対座標で格納し、図面という制約を仕様から外しているものが多い。ただしこれがデータサイズ増大の原因となっている。

【0061】

本発明の圧縮方法の一つである「相対値管理」は、このような絶対座標で記述された空間データをBXML内部で図面と同様の相対座標管理を行うことができる。相対座標管理はBXMLの内部構造で行っており、空間データフォーマット上では相対座標であることが隠蔽されている。これにより、従来は図面管理方式と同様に、座標情報のデータ量削減が実現できる。また、4番目の圧縮方法の「可変長数値の適用」で示したとおり、本発明のデータ構造は固定長数値の制約が存在しないため、相対座標の原点から極端に外れている座標も問題なく管理できる。 30

【0062】

相対値管理は、圧縮スキーマ部において、特定のタグに起点となる数値を記述し、圧縮データ部においてその起点数値からの相対座標を記述することにより実現できる。

【0063】

次に、図13を用いて、本発明によるBXML文書に関連し、本発明が提供する機能について説明する。本発明が提供する機能は、圧縮処理1301、展開処理1302、解析処理1303、更新処理1304、作成処理1305に分類される。 40

【0064】

圧縮処理1301は、XML文書1306を読み込み、これをバイナリ構造に圧縮したBXML文書1307に変換する。展開処理1302は、BXML文書1307をXML文書1306に展開する。

【0065】

解析処理1303は、BXML文書1309のタグ構造の移動と、BXML文書の任意の位置での情報取得機能を、ユーザの作成するプログラム1308に提供する。解析処理1 50

303はXML文書の解析のための標準インタフェースであるSAXとDOMに準拠したインタフェースをユーザプログラム1308に提供する。

【0066】

更新処理1304は、BXML文書の任意の位置でのデータの更新機能を、ユーザプログラム1308に提供する。更新処理1304はXML文書の更新のための標準インタフェースであるDOMに準拠したインタフェースをユーザプログラム1308に提供する。

【0067】

作成処理1305は、BXML文書1311を直接作成する機能をユーザプログラム1310に提供する。ユーザプログラム1310は、XML文書を作成することなく、直接BXML文書1311を作成することができる。

図1に、図13の機能ブロックを実現するためのシステムブロック図を示す。ここで、矩形はそれぞれの機能を実現する各処理ブロックを示す。また破線で構成される矩形は、ユーザが作成するプログラムを示す。角の取れた矩形はデータを示す。また矢印の方向は、矢印の始点から終点方向にデータが参照され、あるいはデータが変換されることを示す。最初に、図1に記載されているデータ構造群について説明する。XML文書101は本発明が圧縮対象とするXML文書であり、BXML文書105は本発明によりバイナリ構造に圧縮されたXML文書である。

XMLのスキーマ定義106は、XMLの構造を定義した文書である。一般に、XML文書の構造を定義する標準的な形式として、DTD(Document Type Definition)、XML Schema、RELAX NGなどが規定されているが、スキーマ定義106はこれらの形式の文書に相当する。これは圧縮処理1302や生成処理1305に利用される場合がある。

圧縮スキーマ103および圧縮データ104は、本発明の処理で一時的に利用されるデータであり、この二つを組み合わせることによりBXML文書105が構成される。圧縮スキーマ103にはXML文書のタグ名やタグ間の階層関係など、XML文書全体で共通に利用される情報が、圧縮データ104にはテキスト要素や属性値、重複数など、データ固有の情報が格納される。

スキーマ木102は圧縮スキーマ103を計算機のメモリ空間に計算機に処理しやすいように、木状に展開したものである。

図14にスキーマ木の構造を示す。スキーマ木は、XML文書のルート要素を根とする木構造を取り、木の要素としてXMLのタグを表すスキーマ要素1401を持つ。重複するタグは同じスキーマ要素となる。スキーマ要素に含まれる変数は、図8に示す圧縮スキーマ804に記述される情報とほぼ同じであり、スキーマ要素識別子1404、タグ名1405、重複フラグ1406、スキーマ型1407、属性保有数1408、子要素保有数1409を持つ。また、各子要素を示すスキーマ要素には、0個以上の属性を示すスキーマ要素1402が接続し、また0個以上の子要素を示すスキーマ要素1403が接続する。これを実現するため、各スキーマ要素には、属性を表すスキーマ要素1402への0個以上のポインタ、子要素を表すスキーマ要素1403への0個以上のポインタ、親要素を表すスキーマ要素への1個のポインタを保持する。

【0068】

図8の圧縮スキーマ804に存在しない新たな情報としてスキーマ要素識別子1404がある。これは、スキーマ要素を識別するための情報であり、例えば数値で記載される。この識別子の決定方法は様々な方法が考えられるが、例えばスキーマ木で出現する全てのタグ、属性を示すスキーマ要素群の出現順の通し番号としてもよい。

属性を示すスキーマ要素1402の構造は、タグを表すスキーマ要素1401と同じ構造を持つ。属性のキー名1411はタグ名1405で表現する。スキーマ要素識別子1410、重複フラグ1412、スキーマ型1413はスキーマ要素1401と同じである。属性に付随する属性は存在せず、また属性は階層構造を持たないため、属性を示すスキーマ要素1402には属性保有数1408、子要素保有数1409は存在しない。同様に、属性を示すスキーマ要素には、属性を表すスキーマ要素へのポインタ、子要素を表すスキーマ要素へのポインタを保持する。

10

20

30

40

50

マ要素へポインタは持たず、親要素を表すスキーマ要素への1個のポインタのみを保持する。

図15を用いてスキーマ木について、より詳細に説明する。図3において、従来のDOMで使用される木構造であるインスタンス木302について説明したが、スキーマ木は、インスタンス木と違い、重複したタグそれぞれについてのノードを持たない。重複したタグのスキーマ構造は全て同じであり、これを持つことは冗長であるためである。例えば図15の最初の例1501では、タグRの中にタグAが3個、タグBが2個格納されている。インスタンス木ではそれぞれのタグについて全て木の構成要素を持つが、スキーマ木ではタグA、タグBの種類につき一つだけの構成要素しか持たない。図15の2番目の例1502では、タグRの中にタグAが3個存在し、最初のタグAにはタグBが、次のタグAにはタグCが、3番目のタグAにはタグDが格納されている。インスタンス木ではそれをそのまま格納しているが、スキーマ木ではタグAの下にタグB、タグC、タグDが格納される、という構造を取る。これによりタグAのスキーマ構造を重複せず持つことができる。図15の3番目の例1503では、タグCがタグA、タグBの両方の下に出現している。これを表現するスキーマ木としては、スキーマ木1504、スキーマ木1505の2通りの構造が取れる。スキーマ木1504の方がタグCのスキーマ構造を重複せず持つことができるので効率的である。ただしスキーマ木構造にループができてしまうため、圧縮方式が複雑になり、圧縮処理やスキーマ構造の変更処理における計算速度が低下する可能性がある。そのため、スキーマ木を表現する圧縮スキーマのサイズが圧縮データのサイズと比較して無視できるほど十分に小さい場合、冗長であるが簡単なスキーマ木1505の構造を取っても構わない。

【0069】

次に、図1のシステムブロック図を用いて、図13の機能ブロック図の各機能を実現する手段について説明する。図13の機能ブロック図に示す各機能は、図1のシステムブロック図に示す各処理ブロックの組合せにより実現される。見易さのため、図1のシステムブロック図を、図13の機能ブロック図に示す5つの機能に合わせて抜粋した図を、図16、図17、図18、図19、図20に示す。

【0070】

図16は、図13の圧縮処理1301を実現する処理ブロック群を、図1から抜粋したブロック抜粋図である。図16を用いてXMLからBXMLへの圧縮の処理フローについて説明する。XMLのBXMLへの圧縮は、3種類の方法がある。すなわち、(1)XML文書のみを利用した圧縮、(2)スキーマ定義を利用したXML文書の圧縮、(3)BXML文書を利用したXML文書の圧縮である。

(1)「XML文書のみを利用した圧縮」では、XML文書をBXML文書に圧縮する際、圧縮対象のXML文書以外の情報を利用しないで行う圧縮方法である。最初にXML文書101全体をXMLパースブロック107によりパースし、スキーマ生成ブロック109により、スキーマ木102を作成する。そして再度XML文書101全体をXMLパースブロック107によりパースし、データ圧縮ブロック112により既に作成したスキーマ木102を利用しながら圧縮データ104を作成する。次にスキーマ圧縮ブロック110によりスキーマ木102を圧縮スキーマ103に圧縮し、BXML結合ブロック114により圧縮スキーマ103と圧縮データ104を結合し、最終成果であるBXML文書105を作成する。

(2)「スキーマ定義を利用した圧縮」では、処理時間のかかるスキーマ生成ブロック109を省略することにより圧縮時間を短縮する方式である。XML SchemaなどのXMLスキーマ定義文書106から直接スキーマ木102を作成することによりスキーマ生成ブロック109を省略することができる。スキーマ変換ブロック116により、XMLのスキーマ定義106からスキーマ木102を作成する。後は「XML文書のみを利用した圧縮」で説明したのと同様な流れでBXML文書105を作成する。

(3)「BXMLを利用した圧縮」も上記と同様に、処理時間のかかるスキーマ生成ブロック109を省略することにより圧縮時間を短縮する方式である。同じ構造のXML文書

10

20

30

40

50

を繰り返し圧縮する場合に、スキーマ木を作成する方法として、以前に作成した同じスキーマ構造であることが分かっている B X M L からスキーマ木を取得する。

図 1 において、あらかじめ作成してある B X M L 文書 1 0 5 を用意し、B X M L 分解ブロック 1 1 5 により B X M L 文書 1 0 5 から圧縮スキーマ 1 0 3 を取り出す。圧縮スキーマ 1 0 3 をスキーマ展開ブロック 1 1 1 により変換することでスキーマ木 1 0 2 を得る。あとは「X M L 文書のみを利用した圧縮」で説明したのと同様な流れで B X M L 文書 1 0 5 を作成する。上記説明において、B X M L 文書 1 0 5 ではなく圧縮スキーマ 1 0 3 をあらかじめ保持しておくこともできる。

図 1 7 は、図 1 3 の展開処理 1 3 0 2 を実現する処理ブロック群を、図 1 から抜粋したブロック抜粋図である。

10

【 0 0 7 1 】

B X M L から X M L への展開は、2 種類の方法がある。すなわち、(1) B X M L 文書の展開、(2) 既存圧縮スキーマを利用した圧縮データの展開、である。

(1) 「B X M L 文書の展開」では、B X M L 文書 1 0 5 から圧縮スキーマ 1 0 3 と圧縮データ 1 0 4 を取り出し、それを用いて X M L 文書に展開する。B X M L 分解ブロック 1 1 5 により B X M L 文書 1 0 5 から圧縮スキーマ 1 0 3、圧縮データ 1 0 4 を取得し、スキーマ展開ブロック 1 1 1 により圧縮スキーマ 1 0 3 からスキーマ木 1 0 2 を取得する。そして B X M L パースブロック 1 0 8 でスキーマ木 1 0 2 を利用して圧縮データ 1 0 4 をパースし、データ展開ブロック 1 1 3 により X M L 文書 1 0 1 を取得する。

(2) 「既存圧縮スキーマを利用した圧縮データの展開」では、同じ構造の B X M L 文書を繰り返し展開する時に利用する事ができる方式である。同じ構造の B X M L 文書 1 0 5 が多数存在する場合、その圧縮スキーマ 1 0 3 は全て共通であり、従って圧縮スキーマ 1 0 3 の部分が冗長である。そこで、複数の B X M L 文書 1 0 5 を管理するのではなく、単独の圧縮スキーマ 1 0 3 と複数の圧縮データ 1 0 4 を管理する方が、データサイズの削減となり、その結果、データ管理領域の削減やデータ転送速度の向上が実現できる。例えば、クライアント・サーバシステムにおいて、あらかじめクライアントに共通の圧縮スキーマ 1 0 3 を送信しておく。そしてクライアントの要求に応じて、圧縮データ 1 0 4 を複数回クライアントに送信する。クライアントではあらかじめ、取得した圧縮スキーマ 1 0 3 をスキーマ展開ブロック 1 1 1 により圧縮スキーマ 1 0 3 をスキーマ木 1 0 2 に展開しておき、圧縮データ 1 0 4 が繰り返し送信されてきた時に、あらかじめ取得しておいたスキーマ木 1 0 2 を用いて繰り返し送付される圧縮データ 1 0 4 を B X M L パース 1 0 8 でパースし、データ展開ブロック 1 1 3 により、X M L 文書 1 0 1 を取得する。

20

30

図 1 8 は、図 1 3 の解析処理 1 3 0 3 を実現する処理ブロック群を、図 1 から抜粋したブロック抜粋図である。

【 0 0 7 2 】

一般に、X M L 文書の解析を行うための標準インタフェースとして、S A X、D O M が規定されている。解析処理 1 3 0 3 では、B X M L 文書を解析するインタフェースとして標準インタフェースである S A X、D O M を提供する。これにより、ユーザは X M L 文書を解析すると同様の方式で B X M L 文書を解析することができる。

【 0 0 7 3 】

B X M L 文書 1 0 5 の解析に先立ち、まず B X M L 分解ブロック 1 1 5 により B X M L 文書 1 0 5 から圧縮スキーマ 1 0 3、圧縮データ 1 0 4 を取得し、スキーマ展開ブロック 1 1 1 により圧縮スキーマ 1 0 3 からスキーマ木 1 0 2 を取得しておく。そして B X M L パースブロック 1 0 8 でスキーマ木 1 0 2 を利用して圧縮データ 1 0 4 をパースする。

40

【 0 0 7 4 】

ユーザプログラム 1 2 2 が標準インタフェースである S A X で B X M L 文書 1 0 5 を解析する場合、B X M L パースブロック 1 0 8 は S A X パースブロック 1 2 3 を呼び出し、S A X パースブロック 1 2 3 がユーザプログラム 1 2 2 に対し X M L を構成する部品、すなわち開始タグ、終了タグ、テキスト要素、属性の情報などを提供する。

【 0 0 7 5 】

50

ユーザプログラム 125 が標準インタフェースである DOM で BXML 文書 105 を解析する場合、DOM パースブロック 124 は BXML パースブロック 108 を利用して XML 文書のノード間の移動を行い、また任意のノードにおけるタグ情報、テキスト要素、属性情報をユーザプログラム 125 に提供する。

図 19 は、図 13 の更新処理 1304 を実現する処理ブロック群を、図 1 から抜粋したブロック抜粋図である。

【0076】

BXML 文書 105 の更新に先立ち、まず BXML 分解ブロック 115 により BXML 文書 105 から圧縮スキーマ 103、圧縮データ 104 を取得し、スキーマ展開ブロック 111 により圧縮スキーマ 103 からスキーマ木 102 を取得しておく。そして BXML パースブロック 108 でスキーマ木 102 を利用して圧縮データ 104 をパースする。

10

【0077】

更新処理 904 では、BXML 文書 105 の更新機能として 3 つの機能を提供する。すなわち、(1) インスタンス変更機能、(2) 部分 BXML 文書抽出機能、(3) スキーマ変更機能である。

【0078】

(1) インスタンス変更機能では、BXML 文書 105 の任意の位置に存在するテキスト要素や属性値の追加・削除・置換、任意のノードの挿入・削除・置換を行う。ユーザがテキスト要素や属性、ノードを指定する方法として、前に説明した解析機能 1303 の DOM インタフェースを用いる。DOM パースブロック 124 は BXML パースブロック 108 を利用して XML 文書のノード間の移動を行う。DOM 更新ブロックはユーザのテキスト要素、属性情報の追加・削除・置換、ノードの挿入・削除・置換処理を受け付け、圧縮データ 104 に書き込む。そして更新した圧縮データ 104 を BXML 結合ブロック 114 により圧縮スキーマ 103 と結合し、BXML 文書 105 を出力する。

20

【0079】

(2) 部分 BXML 文書抽出機能では、BXML 文書 105 のユーザの指定した任意のノードにおいて、そのノード以下の部分 BXML 文書を抽出し出力する。DOM パースブロック 124 は BXML パースブロック 108 を利用してユーザが抽出したいノードへの移動を行う。DOM 更新ブロックは圧縮データの現在参照している位置と、現在参照しているスキーマ木の位置とを利用して、BXML 結合ブロック 114 を実行することにより、部分 BXML 文書 105 を取得する。

30

【0080】

(3) スキーマ変更機能は、スキーマ木を変更することにより、BXML 文書 105 のユーザが指定したタグ名称を一括変更する機能、ユーザが指定したノード以下の部分 XML 文書を一括削除する機能などを提供する。スキーマパースブロック 117 は、スキーマ木に対してユーザが変更したいスキーマ要素へ移動する機能を提供する。スキーマ更新ブロック 119 は、ユーザのタグ名一括変更、ノード一括削除要求を受け付け、スキーマ木の変更を行う。スキーマ圧縮ブロック 110 により変更されたスキーマ木 102 を圧縮スキーマ 103 に変換し、BXML 結合ブロック 114 により圧縮データと結合し、BXML 文書 105 を出力する。

40

【0081】

図 20 は、図 13 に示す BXML 文書作成処理 1305 を実現する処理ブロック群を、図 1 から抜粋したブロック抜粋図である。生成処理 1305 では BXML 文書を新規に作成する機能を提供する。

【0082】

XML 文書を新規に作成する従来方法としては、(1) 専用のプログラムを作成して、テキストの XML 文書を直接記述する、(2) 標準の DOM で XML 文書を作成し、完成したインスタンス木から XML 文書を生成する、という 2 通りの方法が存在する。しかし最初の方法では、ユーザプログラムのミスで開始タグと終了タグの対応付けが正しく記述されないというエラーが、XML 文書が完成して整合性の確認を行うまで発見できない。ま

50

たテキストベースのXML文書を直接作成するため、巨大なXML文書を作成するには大量の処理時間とメモリ空間を必要とする。2番目の方法は、従来のDOMパーサでは、作成するインスタンス木は目的のXML文書より数倍のサイズのメモリ空間を必要とするため、巨大なXML文書を作成するには大量のメモリ空間を必要とする。本発明の作成処理1305は、XML文書より小さなサイズのBXML文書を直接作成すること、従来のDOMパーサの作成するインスタンス木のような大量のメモリ空間を必要とする中間生成物を必要としないことにより、従来の上記2方法と比較し占有メモリ空間が少なく、かつ高速にBXML文書を生成することができる。

【0083】

BXML文書を新規に作成する手順は、(1)スキーマ木102の作成、(2)圧縮データ104の作成、の2手順で行う。この二つのデータが完成すると、圧縮処理1301で説明したとおり、スキーマ圧縮ブロック110、BXML結合ブロック114を経由して最終結果であるBXML文書105を作成することができる。

(1)スキーマ木102の作成方法は、次に示す4つの方法が存在する。(1)作成対象のBXML文書と同じ構成のXML文書やテンプレート101から、スキーマ生成ブロック109を経てスキーマ木102を作成する。(2)XML Schemaなどのスキーマ定義文書106からスキーマ変換ブロック116を経てスキーマ木102を作成する。(3)あらかじめ用意しておいたBXML文書105や圧縮スキーマ103からスキーマ展開ブロック111を経てスキーマ木を作成する。(4)スキーマパーサブロック117、スキーマ更新ブロック119を用いてスキーマ木を最初から構築する。また、これら4つの方法のいずれかで作成したスキーマ木は、スキーマパーサブロック117、スキーマ更新ブロック119を用いたユーザプログラム118でスキーマ要素の追加・置換・削除を行ったり、別のBXML文書が持つスキーマ木の一部を挿入したりすることにより更新することもできる。

(2)圧縮データ104の作成は、圧縮行程1301で説明したデータ圧縮ブロック113をそのまま利用する。圧縮処理1301で説明した通り、データ圧縮ブロック113はXMLパーサブロック107からSAXインタフェースで呼び出される。すなわち、XMLパーサブロック107はXML文書を走査し、「開始タグ」、「終了タグ」、「属性」、「テキスト要素」などのイベントをデータ圧縮ブロック113に発行し、データ圧縮ブロック113はそれに対応して圧縮データ112を作成していく。

図20のBXML作成ブロック121はXMLパーサブロック107と同様に、SAXインタフェースを用いてデータ圧縮ブロック113を呼び出す。ユーザプログラム122が、「開始タグ」、「終了タグ」、「属性」、「テキスト要素」などのSAXイベントを直接データ圧縮ブロック113に発行すると圧縮データ104を作成することができるが、これはユーザにとって使いにくい。そのためBXML作成ブロックは、ユーザプログラム122に対し、「親に移動」「子に移動」「弟に移動」「テキスト要素の挿入」などのインタフェースを提供し、これをSAXイベントに変換してデータ圧縮ブロック113を呼び出す、という方式をとる。また、「テキスト要素の挿入」のインタフェースの派生として、「整数の挿入」「実数の挿入」など、図7の「データ型の適用」圧縮702で説明した複数のデータ型を直接挿入するインタフェースを提供することにより、ユーザプログラム122が持つデータ型をXMLの文字列に変換し、BXML文書の内部構造へ変換する、という手間が省け、その結果処理速度を向上させることができる。

【0084】

BXML作成ブロック121のもう一つの特長として、他のBXML文書の指定したノード以下の部分集合を一括挿入することができる、という点がある。挿入元のBXML文書の圧縮データを一括して挿入先の圧縮データにコピーすることにより、タグ単位で挿入する方式より高速な挿入が実現できる。またこの機能を利用して、複数のBXML文書の高速な結合処理が実現できる。

本発明の実施形態として、図1に示す本発明を構成するシステムブロック図の各処理ブロックについて詳細に説明する。その後、本発明の実施例として、本発明を利用したアプ

10

20

30

40

50

リケーションについて説明する。

最初に、図1に示す本発明を構成するシステムブロック図の各処理ブロックについて、(1)XMLパースブロック107、(2)スキーマ生成ブロック109、(3)スキーマ変換ブロック116、(4)データ圧縮ブロック112、(5)スキーマ圧縮ブロック114、(6)スキーマ展開ブロック111、(7)BXML結合ブロック114、(8)BXML分解ブロック115、(9)BXMLパースブロック108、(10)データ展開ブロック113、(11)SAXパースブロック123、(12)DOMパースブロック124、(13)DOM更新ブロック126、(14)BXML作成ブロック121、(15)スキーマパースブロック117、(16)スキーマ更新ブロック117、の順番で説明を行う。

10

最初に、図1のXMLパースブロック107について詳細に説明する。

XMLパースブロック107の目的は、入力で与えられたXML文書101の初めから最後までを走査し、外部のプログラムに対しXML文書構造やテキスト要素の情報を提供することにある。XMLパースブロック107は、図2で説明した、標準のXML解析インタフェースであるSAXを利用することにより実現する。

【0085】

次に、図1のスキーマ生成ブロック109について詳細に説明する。

【0086】

図1のスキーマ生成ブロック109の目的は、XML文書101からスキーマ木102を作成することである。前述の6つの文書圧縮方式、すなわち(1)タグとデータの分離、(2)文字列の辞書化、(3)データ型の適用、(4)可変長数値の適用、(5)相対値の適用、に必要な情報を同時に作成し、スキーマ木101に登録する。

20

「タグとデータの分離」を行うため、図21に示すアルゴリズムを用いて処理を行う。

[STEP 1] 最初にXML文書をメモリ空間に読み込む。またスキーマ木のルートスキーマ要素を設定する。ここでルートスキーマ要素とは、XML文書のルート要素ではなく、その親の位置にある仮想的なものである。図14に示すように、スキーマ要素には属性を示すスキーマ要素が0個以上、子タグを示すスキーマ要素が0個以上接続される構造を取るが、初期状態ではいずれも0個とする。また各スキーマ要素には、処理のための一時的な変数として、「期待タグ」として、次に来ると期待される子タグのスキーマ要素を指すポインタ、「期待キー」として、次に来ると期待される属性キーのスキーマ要素を指すポインタを持つ。またスキーマ生成部が保持する一時変数「参照要素」として、スキーマ木の現在作成中のスキーマ要素を記憶しておく。参照要素の初期状態は、ルートスキーマ要素を指している。

30

[STEP 2] 文書終端に達するまで、[STEP 3]以降を繰り返す。

[STEP 3] XML文書の部品の一つを取得する。ここでXML文書の部品とは、開始タグ、終了タグ、テキスト要素、属性である。これは標準のSAXでの動作と同様である。

[STEP 4] 取得した部品である開始タグ、終了タグ、テキスト要素、属性に対し、それぞれ処理を分岐させる。ここで分岐する処理は開始タグ、終了タグ、属性である。実際はテキスト要素に対しても処理を行うが、これは「タグとデータの分離」の圧縮のためではなく、その他のオプションの圧縮方針を実現するために利用するため、ここでは割愛する。

40

[STEP 5] 部品が開始タグである場合、本処理を実行する。以降、参照要素の期待タグが指すスキーマ要素のことを期待要素と命名する。期待要素のタグ名称と、今回取得した開始タグのタグ名称とを比較し、同一であればSTEP 6を実行する。それ以外の場合、参照要素の期待タグを、子タグ一覧の次に指すタグにして、STEP 5の先頭に戻る。子タグ一覧が空である場合、または子タグ一覧に継続する要素が存在しない場合、STEP 8を実行する。

[STEP 6] 取得した開始タグが期待要素のタグ名称と一致している場合、本処理を実行する。期待要素の重複数を1増加させる。重複数は初期設定時点では0となってい

50

る。この時点で重複数が2以上になっている場合、その期待要素は図4の重複要素パターンであることが確定されるため、重複フラグを「重複」に設定する。

[STEP 7] 参照要素を現在の参照要素の期待要素に移動する。

[STEP 8] 取得した開始タグが期待要素のタグ名称と一致していない場合、あるいは期待要素が存在しない場合、本処理を実行する。新しいスキーマ要素を作成し、参照要素の持つ子タグ一覧の末尾に挿入する。参照要素の期待タグはここで作成したスキーマ要素を指すようにする。

[STEP 9] 参照要素には一時変数として処理済フラグというフラグを設け、一度処理を行ったスキーマ要素にはTRUEが、処理を行っていない要素にはFALSEが設定されるようにしておく。そして処理済フラグがTRUEの場合、STEP 10を実行する。これはSTEP 8で到来したタグが、以前は省略されていたことを示す。具体例を図15の例1502を用いて説明する。最初のタグAに存在するタグBが到来した時、タグAの処理済フラグはFALSEであり、タグBの重複フラグは「単独」となる。2番目のタグAの中に存在するタグCが到来した時、タグAの処理済フラグはTRUEであるため、タグCの重複度は「重複」となる。これは、タグAが2回目以降の処理であり、前回の処理でタグCが省略された、すなわちタグCは省略可能なタグである、と判断したことに基づく。

[STEP 10] STEP 8で新規に挿入したスキーマ要素の重複フラグを「重複」に設定する。

[STEP 11] 参照要素の要素型を「SEQUENCE」に設定する。これは参照要素に子要素が存在しているためである。

[STEP 12] 部品が終了タグである場合、本処理を実行する。参照タグの全ての子要素について、STEP 13を実行する。

[STEP 13] 重複数が1ではない子要素について、子要素の重複フラグを「重複」に設定する。重複数が0の子要素は、参照タグ内においてその子要素が省略されたことを示し、重複数が2以上の子要素は、参照タグ内においてその子要素が2回以上出現したことを示す。いずれも「子要素が必ず1回のみ出現する」という条件から外れるため、重複フラグは「重複」となる。

[STEP 15] 参照要素の一時変数である処理済フラグをTRUEにする。

[STEP 16] 部品が属性キーである場合、本処理を実行する。属性キーが参照要素の属性として登録済みの場合、何も処理は行わない。属性キーが登録されていない場合、新たな属性を生成し、属性一覧の末尾に挿入する。属性の重複フラグは必ず「属性」とし、これは属性が省略されるか、1回出現するかのいずれかを示す。XMLの定義によると属性が2回以上重複することはないため、子要素のような重複フラグ判定処理は省略可能である。

【0087】

なおここでは簡単のため、図8に示すパターンのうち兄弟要素パターン812でのスキーマ木生成方式について示したが、選択要素パターン815、重複選択要素816の生成も可能である。具体的には、STEP 12の終了タグ処理において、図8パターン816のXMLサンプル802のタグBのように、登録された子要素のタグのうち、同名のタグが連続せずに存在する場合、重複選択要素816のパターンであることが判定でき、その場合参照要素Aの要素型を「SEQUENCE」から「MCHOICE」に変更する。また、図8パターン815のXMLサンプル802のタグB、タグCのように、子要素数は必ず1個であり、そのタグ名は参照要素Aによって異なる場合、選択要素815のパターンが判定でき、その場合、参照要素Aの要素型を「SEQUENCE」から「CHOICE」に変更する。

残りの圧縮方式、すなわち、(2)文字列の辞書化、(3)データ型の適用、(4)可変長数値の適用、(5)相対値の適用、を行うため、図22に示すアルゴリズムを用いて処理を行う。これは図21で示した「タグとデータの分離」と同時に行われるが、説明を簡単にするため、図を分けている。

10

20

30

40

50

[STEP 1] ~ [STEP 4]までは、前述した、図21の説明と同じである。ここでは、前回割愛した、テキスト要素についての説明のみを抜き出して示す。

[STEP 5] 部品がテキスト要素である場合、本処理を実行する。ここでは、テキスト要素が前述した、「整数型」、「実数型」、「2次元整数型」、「2次元実数型」、「2次元整数配列型」、「2次元実数配列型」、「日付型」、「文字列型」のいずれかであるかを判定する。

【0088】

整数型の判定規則としては、その記述形式（先頭のプラスマイナス記号と0から9までの文字が1個以上でのみ構成されるなど）、記述範囲（long値を超える極端に大きな数値は整数ではなく実数や文字列で管理するなど）の規則を適用できる。また「002」のような特殊な記述がある場合、「2」として整数で記述するか、文字列で記述するか、また他の拡張形式（例えば固定長3桁の整数など）で記述するかをユーザ指定で方針を切り替えても構わない。

10

【0089】

実数型の判定規則としては、その記述形式（小数点1個を含んでもよい、E+02などの指数部表現を許可するなど）、記述範囲（指数部の桁数や仮数部の精度桁数がdouble値を超えるものは文字列で管理するなど）の規則を適用できる。

2次元整数型、2次元実数型、ならびにそれを拡張した整数配列型、実数配列型は、テキスト要素が整数や実数と、それを区切る区切り文字でできていることを判定する。区切り文字はあらかじめカンマ、空白、タブなどのいずれかで判定してもよく、テキスト要素を解析することで区切り文字を導出してもいい。2次元整数配列型、2次元実数配列型のように区切り文字が二種類存在する場合も同様である。

20

【0090】

日付型は、その記述形式（ISO8601に従うことなど）、記述範囲（日付を1970年1月1日午前0時（UTC）を起算点とする通算の秒数で管理する場合、その起算点からlong値の限界まで）の規則を適用できる。

【0091】

これらの判定規則には優先順位を設ける。たとえば「2」という文字列は整数型でも実数型でも文字列型でもあるが、整数型とした方が圧縮効率がいい。そのため「整数型でも実数型でもある場合は整数型にする」、という優先順位を設ける。優先順位は、「整数型」、「実数型」、「2次元整数型」、「2次元実数型」、「2次元整数配列型」、「2次元実数配列型」、「日付型」、「文字列型」の順となる。

30

[STEP 6] 現在の参照要素の期待タグが示すスキーマ要素（以降は期待要素と呼ぶ）に対し、スキーマ型としてSTEP 5で判定した型を代入する。ただし上書規則を設定する必要がある。例えば2個の重複要素が存在し、最初のタグのテキスト要素に整数、次のタグのテキスト要素に実数が入っている場合、そのタグのスキーマ要素の型は「実数」となる。最初のタグのテキスト要素に実数、次のタグのテキスト要素に整数が入っている場合も、そのタグのスキーマ要素の型は「実数」となる。

【0092】

図20に上書規則を示す。あるスキーマ要素のスキーマ型が、既に図20の列に示す型として判定されている状態で、新たに左端に示すスキーマ型で上書しようとする場合の規則を示す。ここで、「OK」は上書可能、「NG」は上書不可能、「文字列」は元の型や上書する型にかかわらず文字列型に置換されてしまうことを示す。

40

[STEP 7] 文字列辞書化圧縮方式のオプションが選択されており、かつ期待要素のスキーマ型が文字列型と判定された場合、本処理を実行する。本処理では、頻出する文字列を抜き出し、スキーマ木のスキーマ要素毎に存在する辞書に、文字列と、それに対応する識別子を登録することにより、データサイズの削減を図る。

【0093】

文字列辞書化のアルゴリズムは多数考えられる。例えば、（1）単純辞書化、（2）頻出単語登録、（3）事前辞書準備、（4）数値分離、の方式が考えられる。

50

(1) 単純辞書化は、全てのテキスト要素を分割せずにそのまま辞書登録を行う。辞書登録量の爆発を防止するため、辞書の登録件数が閾値以上(例えば100件以上)の場合、そのスキーマ要素は辞書化が困難であると判定して辞書登録を停止する、という戦略が可能である。

(2) 頻出単語登録は、テキスト要素を単語に分割して、その出現数と共に一時バッファに登録する。単語の分割は既存の形態素解析等を利用することで行う。そして単語の出現数が閾値以上(例えば上位10件)のみを登録することにより、辞書の肥大化を防止することができる。

(3) 事前辞書準備は、あらかじめ外部から辞書を提供する方式である。一般に、XML文書を作成するユーザはそのXML文書内に頻出する用語を理解している場合があり、この場合には本方式が有効となる。 10

(4) 数値分離は、テキスト要素を共通の文字列と固有の数値に分離できると仮定して、テキスト要素を文字データの前半部分と数値の後半部分に分離し、前半部分を辞書化する方法である。例えば前述のURIのように、識別子の前半がURLのような文字データ、後半が組織固有の数値というように分離できる場合に有効である。

[STEP 8] 相対座標管理方式のオプションが選択されており、かつ期待要素のスキーマ型が整数型、実数型などの数値型であると判定された場合、本処理を実行する。本処理では、数値の相対値管理を行うための、起点となる数値を取得するために実行する。

【0094】

相対値管理のアルゴリズムは多数考えられる。例えば次の3つが考えられる。すなわち、 20

(1) 差分管理、(2) 起点計算、(3) 起点入力、である。

(1) 差分管理は、テキスト要素として、同じスキーマ要素の前出したテキスト要素からの相対値管理を行う方式である。例えば3つのテキスト要素、10000、10001、10002がある場合、格納される相対値としては10000、1、1となる。

(2) 起点計算は、同じスキーマ要素に対し、起点を計算する方式である。同じスキーマ要素の全てのタグに対して最小となる数値を計算し、その数値を起点とする。あるいは同じスキーマ要素の全てのタグに対して平均となる数値を計算し、その数値を起点とする。後者の方が、極端に大きいあるいは小さい数値を含む場合に全体として安定している。例えば3つのテキスト要素、10000、10001、10002がある場合、起点を平均値10001とすると、圧縮データに格納される相対値は、-1、0、1となる。起点情報 30

は図29の2932のように、圧縮スキーマに記述する。

(3) 起点入力は、あらかじめ外部から起点を入力する方式である。例えばユーザの設定により起点を入力する。また、ある種のXML文書には数値範囲を記録している場合がある。たとえば地理情報分野に置ける多数の地物を表現するXML文書の中には、地物の位置を示す座標群の他に、全ての地物を外接する矩形を格納している場合がある。この矩形の端点、あるいはこの矩形の中心点を相対値管理の起点として採用することもできる。

【0095】

次に、図1のスキーマ変換ブロック116について詳細に説明する。スキーマ変換ブロック116の目的は、スキーマ生成ブロック809の代替ブロックとして、XML SchemaなどのXMLスキーマ定義文書106からスキーマ木102を作成することである 40

【0096】

たとえばXML Schemaでは、あるタグの中にどのタグが含まれるか、タグの最小出現回数と最大出現回数、テキスト要素の型を定義することができるため、スキーマ木への変換を行うことができる。また頻出単語や相対値など、XML Schemaで表現できない項目については、その部分だけ別の形式で管理することもでき、またXML Schemaの拡張を行い、上記の情報を格納するようにすることもできる。

【0097】

逆に、XML Schemaでは表現可能であるが、BXMLの圧縮スキーマでは表現できない構造も存在する。本発明の圧縮XMLの目的は、XMLのサイズを削減し、高速な 50

処理を実現することであり、XML Schemaで示される構造のチェック (validation)を行う必要は必ずしもない。そのため、以下のようなXML構造の簡略化を行う。

(1) XML Schemaでは、「正の整数型」「単精度32bit浮動小数点型」など、あらかじめ44種類のデータ型が定められている。本発明では、「正の整数型」を「整数型」に変換するなど、型変換規則を適用することにより、型の違いを吸収する。

(2) XML Schemaでは、あるタグの子要素として出現する組み合わせを自由に定義できる。例えば、タグAの子要素としてタグB、タグC、タグD、タグEが存在する可能性がある場合、タグB、タグC、タグD、タグEが連続して1回だけ来るのか、タグB、タグC、タグD、タグEの何れかが1回だけ来るかなど、XML Schemaでは

10

様々な組み合わせが定義できる。一方、本発明では、図8の兄弟要素パターン812、選択要素パターン815、重複選択要素パターン816の三種類のみを定義している。本発明では、兄弟要素パターン812、選択要素パターン815に当てはまらないパターンを全て重複選択要素パターン816とすることにより組み合わせの違いを吸収する。

【0098】

次に、図1のデータ圧縮ブロック112について詳細に説明する。データ圧縮ブロック113の目的は、スキーマ木102を利用してXML文書101をパースし、圧縮データ104を作成することである。

20

【0099】

初めに、「タグとデータの分離」圧縮701に対する圧縮データの作成方法について説明する。その後、「文字列辞書の適用」圧縮703、「相対値の適用」圧縮705に対する圧縮データの作成方法について説明する。

【0100】

「タグとデータの分離」圧縮701に対する圧縮データの作成方法について説明する。本方式では、スキーマ木を参照しながら、XML文書を解析し、目的の圧縮データを作成する。

【0101】

データ圧縮ブロック113は、SAXパーサであるXMLパースブロック107を利用して実現する。データ圧縮ブロック113では、各SAXのイベントに対し、XML文書で現在参照している位置がスキーマ木のどの部分であるかを把握しておく必要がある。そのため、SAXイベントに対し、スキーマ木を移動していく手段が必要となる。

30

【0102】

初めに、図25を用いて、XML文書をパースする際、現在参照中のタグやテキスト要素が、スキーマ木のどのスキーマ要素に対応するかを判定する方法について説明する。

【0103】

XML文書2501と、そのスキーマ構造を表すスキーマ木2502が存在する場合を考える。スキーマ木2502の矩形A、B、CなどはXML文書2501の同名のタグが囲む領域を表すスキーマ要素である。一時変数として、スキーマ木外部に、現在参照中のスキーマ要素を表す参照要素2504を1つ設定する。また、全てのスキーマ要素に対し、現在参照している子要素を示す期待タグ(2503、2505)と、現在参照している属性を示す期待キーを持つ。以降の説明では簡単のため、期待タグの指すスキーマ要素を期待要素と呼ぶことにする。

40

【0104】

スキーマ木移動のアルゴリズムを図26に示す。これを図25の具体例、すなわちサンプルXML文書2501と、そのスキーマ構造を示すスキーマ木2502を用いて説明する。

[STEP 1] XML文書2501の終端になるまで、STEP 2以降を繰り返す

50

。

[STEP 2] XML文書を構成する部品、すなわち開始タグ、終了タグ、テキスト要素、属性を一つずつ取り出す。

[STEP 3] 部品の種類により、STEP 4、STEP 6、STEP 7、STEP 8のいずれかに分岐する。

[STEP 4] 本ステップは開始タグが到来した時に実行される。本ステップでは、到来したタグが、現在の参照要素の指すスキーマ要素が持つ期待要素と一致するかを判定し、一致しない場合期待要素を弟要素に進めていく。

本ステップは、重複したり省略されたりするXML文書のタグを正しく処理するために存在する。タグの出現パターンは、(1)単独、(2)重複、(3)省略、(4)エラータグの4つに分類される。XML文書2501において、次の4つのパターンで説明する。

(1) 開始タグD(2506)が到来した時、参照要素2104はスキーマ要素B(2510)を指しており、スキーマ要素B(2510)の期待要素は最初の子要素であるスキーマ要素D(2511)を指している。到来したタグと期待要素のタグ名が一致するため、そのままSTEP 5に移行する。

(2) 開始タグE(2507)が到来した時、スキーマ要素B(2510)の期待要素は最初の子要素であるスキーマ要素D(2511)を指しており、到来したタグのタグ名と一致しない。そこで期待要素を一つ進めてスキーマ要素E(2512)に移動する。この結果タグ名が一致したため、STEP 5に移行する。これは単独タグD 2506を正しく処理したことを示す。

(3) 2番目の開始タグE(2508)が到来した時、スキーマ要素B(2510)の期待要素はスキーマ要素E(2512)を指しており、到来したタグと期待要素のタグ名が一致するため、そのままSTEP 5に移行する。これは重複タグE(2507、2508)を正しく処理したことを示す。

(4) XML文書2501において仮に2つのタグE(2507、2508)が省略され、終了タグD(2514)の直後に開始タグF(2509)が来た場合を考える。スキーマ要素B(2510)の期待要素は最初の子要素であるスキーマ要素D(2511)を指しており、到来した開始タグF(2509)と一致しない。期待要素を次のスキーマ要素E(2511)にするが、まだ開始タグF(2509)と一致しない。期待要素をさらに次のスキーマ要素F(2513)にすると、開始タグF(2509)と一致するため、STEP 5に移行する。これはタグE(2507、2508)が省略された状況を正しく処理したことを示す。

(5) スキーマ木に存在しないタグXが来た時、期待タグは全ての子要素を探索するが、到来タグと一致する要素を発見することができない。この場合本ステップは例外を発行し、XML文書2501が正しい構文ではないことを報告する。

[STEP 5] 参照要素を、現在の参照要素の期待要素に移動する。例えば上記の例のパターン(1)で開始タグD(2506)が到来した時、スキーマ要素B(2510)を指していた参照要素2104をスキーマ要素D(2511)に移動する。

[STEP 6] 本ステップは終了タグが到来した時に実行される。本ステップでは、現在の参照要素をその参照要素の親要素に移動する。例えば、終了タグD(2514)が到来した時、スキーマ要素D(2511)を指していた参照要素2504をスキーマ要素B(2510)に移動する。

【0105】

上記説明したとおり、STEP 4、STEP 5とSTEP 6をXML文書の開始タグ・終了タグに順次適用することにより、重複タグ、省略タグを含む任意のXML文書に対し、参照要素2104をスキーマ木の適切な位置に移動することができる。

[STEP 7] 本ステップはテキスト要素が到来した時に実行される。スキーマ木移動アルゴリズムでは、本ステップでは何も行わない。後述するデータ圧縮ブロック113では、本ステップでテキスト要素圧縮処理を行う。本ステップにおいて、参照要素2504はそのテキスト要素を含むタグのスキーマ要素を指している。例えばテキスト要素 "d

1" (2515) が到来した時の参照要素 2504 はスキーマ要素 D (2511) を指している。

[STEP 8] 本ステップは属性が到来した時に実行される。図 21 では簡単のため省略したが、図 14 に示すように、スキーマ要素 1401 は、子要素を示す 0 個以上のスキーマ要素 1403 への接続を持つと同様に、属性を示す 0 個以上のスキーマ要素 1402 への接続を持つ。属性のスキーマ木移動も上述した子要素のスキーマ木と同様にできる。本ステップでは、到来した属性キーが、現在の参照要素の指すスキーマ要素が持つ期待キーと一致するかを判定し、一致しない場合期待キーを弟要素に進めていく。

【0106】

図 27、図 28、図 29 を用いて「タグとデータの分離」による圧縮について説明する。

10

【0107】

図 27 は、図 8 で示した XML 文書のパターンのうち、子要素パターン 811、兄弟要素パターン 812 のみを用いた重複・省略なしの単純な XML 文書に対するデータ圧縮ブロック 113 のアルゴリズムである。

【0108】

XML 文書 2701 から前に説明したスキーマ生成ブロック 110 を用いて、スキーマ木 2702 を得る。タグ A はタグ B、タグ C を保有し、タグ B、タグ C はテキスト要素を保有する。タグ A、タグ B、タグ C はいずれも、必ず 1 回のみ出現するため、重複フラグは「単独」である。圧縮データ 2703 はテキスト要素 b1、b2 が連続する構造を取る。

【0109】

図 26 に示すスキーマ木移動のアルゴリズムをベースに、データ圧縮を実現する。具体的には、図 26 の STEP 7 のテキスト要素処理においてテキスト文字列を取得し、これを圧縮データに追加していく。STEP 7 で説明したとおり、STEP 7 を実行する時点でそのテキスト要素が含まれるタグのスキーマ要素は、参照要素 2504 からアクセスすることができる。そのスキーマ要素のスキーマ型に対し、以下のルールを適用することにより圧縮データを作成する。

20

[RULE 1] スキーマ要素が複合型 (SEQUENCE 型、CHOICE 型、MCHOICE 型) の場合、テキスト要素は無視する。これは具体的には、XML 文書 2501 のタグ <A> と 、 と <C>、</C> と の間のテキスト要素を無視するという意味を意味する。上に上げたテキスト要素は全てスキーマ要素 A (2704) のテキスト要素であり、スキーマ要素 A のスキーマ型は SEQUENCE 型であるため無視される。

30

[RULE 2] スキーマ型が文字列型の場合、圧縮データに文字データを格納する。

【0110】

BXML 圧縮方式として、いくつかのオプションである圧縮方式を前述している。これらの圧縮方式により、以下の拡張ルールを適用する。

「文字列の辞書化」による圧縮を行う場合、次のルールを適用する。

[RULE 3] スキーマ型が文字列型で語彙辞書が存在する場合、テキスト要素をこの辞書と比較して、図 11 に示したようなテキスト要素と識別子の組合せを圧縮データに格納する。

40

「データ型の適用」「可変長数値の適用」圧縮を行う場合は次のルールを適用する。

[RULE 4] スキーマ型が整数の場合、テキスト要素を図 12 で示した MBI 形式整数に変換して圧縮データに追加する。

[RULE 5] スキーマ型が実数の場合、テキスト要素を図 12 で示した MBF 形式実数に変換して圧縮データに追加する。

[RULE 6] スキーマ型が 2 次元整数の場合、テキスト要素を連続する 2 つの MBI 形式整数に変換して圧縮データに追加する。

[RULE 7] スキーマ型が 2 次元実数の場合、テキスト要素を連続する 2 つの MBF 形式実数に変換して圧縮データに追加する。

【0111】

50

同様に他のデータ型に対しても、それに特有の格納方法で格納することができる。

【0112】

また、「相対値の適用」圧縮を行う場合、スキーマ木の各スキーマ要素に相対値の起点が格納されている。そして次のルールを適用する。

[RULE 8] スキーマ型が数値(整数・実数等)であり、相対値の起点が格納されている場合、テキスト要素の数値から起点までの差分を圧縮データに追加する。

【0113】

図28は、図27に対し、図8で示したXML文書のパターンのうち、さらに重複要素813が追加されたXML文書に対するデータ圧縮ブロック113のアルゴリズムである。

【0114】

XML文書2801から前に説明したスキーマ生成ブロック110を用いて、スキーマ木2802を得る。タグAはタグB、タグCを保有し、タグB、タグCはテキスト要素を保有する。タグA、タグB、タグCはいずれも、複数回出現するため、スキーマ要素2805、2806の重複フラグは「重複」である。圧縮データ2803は、タグBの重複数 N_b 、タグBのテキスト要素 b_1 、 b_2 、タグCの重複数 N_c 、タグCのテキスト要素 c_1 、 c_2 が連続する構造を取る。

【0115】

図27と異なる点は、圧縮データに重複要素の重複数を記述するという点である。圧縮データ2403において、タグBの重複数 N_b はタグBのテキスト要素 b_1 、 b_2 より先に記述する。重複数 N_b はテキスト要素 b_1 、 b_2 の処理が終わった時点で判明するため、先に仮の重複数を記述しておき、テキスト要素の処理が終了した後に正しい重複数を記述する、という方式を取る。

【0116】

図26に示すスキーマ木移動のアルゴリズムをベースに、データ圧縮を実現する。初めに、スキーマ木の各スキーマ要素が持つ一時変数として、重複カウンタを設置する。重複カウンタはそのスキーマ要素が同じ階層で何回出現したかを計測するカウンタであり、開始タグが出現する度に1ずつ増加し、その親タグの終了タグが出現した時に上述した重複数を設定するために用いられる。

具体的には、図26で説明したステップに、図27で説明した処理に加え、以下の処理を加える。

[STEP 4] 開始タグ処理の拡張：到来した開始タグが期待要素と一致する場合、期待要素の重複カウンタを1つ増加させる。期待要素の重複フラグが「重複」であり、重複カウンタが1(初めて出現したタグ)の場合、圧縮データに重複数を追加する。ただしこの段階では重複数の値は不明であるため、仮の値を追加しておき、後で変更を行うため、スキーマ要素の一時変数として「重複数登録位置」を格納しておく。

[STEP 6] 終了タグ処理の拡張：現在の参照要素の持つ全ての子要素について、その重複カウンタを取得し、これを最終的な重複数として圧縮データに記述する。記述する場所は、それぞれのスキーマ要素の一時変数「重複数登録位置」に記録されている値を利用する。

【0117】

図29は、図27に対し、図8で示したXML文書のパターンのうち、省略要素パターン814が追加されたXML文書に対するデータ圧縮ブロック113のアルゴリズムである。

【0118】

XML文書2901から前に説明したスキーマ生成ブロック110を用いて、スキーマ木2902を得る。タグRは3個のタグAを保有する。タグA以下のスキーマ構造は、図28と同様である。最初のタグAは子要素としてタグB、タグCを持つが、2番目のタグAは子要素のタグBが省略されている。3番目のタグAは子要素のタグCが省略されている。

圧縮データ2803は、タグAの重複数 N_a 、最初のタグAにおけるタグBの重複数 N_b

10

20

30

40

50

、タグBのテキスト要素b1、タグCの重複数Nc、タグCのテキスト要素c1、そして2番目のタグAにおけるタグBの重複数Nb2、ここで2番目のタグAにおけるタグBは省略されているのでタグBのテキスト要素はなく、次にタグCの重複数Nc2、タグCのテキスト要素c2、そして3番目のタグAにおけるタグBの重複数Nb3、タグBのテキスト要素b3、タグCの重複数Nc3、3番目のタグAにおけるタグCは省略されているのでタグCのテキスト要素はなし、というデータが連続する構造を取る。

【0119】

図27と異なる点は、省略により開始タグが出現しない場合がある、という点である。STEP6の終了タグ処理で、「重複数設定位置」に子要素の重複数を記述するが、省略されたタグについてはSTEP4の開始タグ処理を通らないため、「重複数設定位置」が設定されていない。 10

これに対応するため、図26で説明したステップに、図27、図28で説明した処理に加え、以下の処理を加える。

[STEP4] 開始タグ処理の拡張：到来した開始タグが期待要素と一致しない場合、期待要素を弟要素に進めていくが、その前に一致しなかった期待要素の重複数0を圧縮データに追加する。

[STEP6] 終了タグ処理の拡張：終了タグの子要素の内、期待要素より後に並ぶ要素群は、その終了タグの子要素として出現しなかったことになる。出現しなかった子要素の重複数0を圧縮データに追加する。

【0120】

次に、図1のスキーマ圧縮ブロック110について詳細に説明する。スキーマ圧縮ブロック110は、計算機のメモリ空間上に展開されているスキーマ木102を、ファイルとして保管・通信できるような形である圧縮スキーマ103に変換する。 20

【0121】

図30の圧縮スキーマ構造を用いてスキーマ圧縮ブロック110について説明する。

【0122】

圧縮スキーマ103は、「タグとデータの分離」圧縮701を行う場合に使用する、スキーマ木の構造を表す構造スキーマ3001、「文字列の辞書化」圧縮703を行う場合に使用する辞書スキーマ3002、「相対値の適用」圧縮705を行う場合に使用する相対値スキーマ3003から構成される。各スキーマはスキーマ開始トークン3005、3006、3007より開始し、圧縮データ3004の開始を表すデータトークン3008と区別する。またそれぞれのスキーマの種類を判定するため、構造スキーマ3001には構造トークン3009、辞書スキーマ3002には辞書トークン3010、相対値スキーマには相対値トークン3011を格納する。 30

【0123】

構造スキーマ3001の説明を行う。図14のスキーマ木を構成するスキーマ要素1401は、タグ名1405、重複フラグ1406、スキーマ型1407、属性保有数1408、子要素保有数1409、で構成され、属性保有数で示される個数(0個以上)の、属性を示すスキーマ要素1402へのリンクポインタ、子要素保有数で示される個数(0個以上)の、子要素タグを示すスキーマ要素1403へのリンクポインタから構成されている。属性を示すスキーマ要素1402も同じ形式で表現でき、タグ名をキー名1411として使用し、属性保有数や子要素保有数が0個、属性や子要素へのリンクポインタが存在しない構造で表現できる。 40

スキーマ木のタグを示すスキーマ要素1401は次のような手順で圧縮スキーマ3019に変換することができる。まずタグ名3012は文字列であるため、0で終端する文字列として格納する。重複フラグ3013は「重複」「単独」のいずれかを取るため、それぞれの列挙に識別子を割り振ることにより登録する。スキーマ型3015は前述の「SEQUENCE」や「CHOICE」、「MCHOICE」、「整数型」、「実数型」などと列挙される項目のいずれかを取るため、それぞれの列挙に識別子を割り振ることにより登録する。子要素保有数3014や属性保有数3016は整数であるため、前述のMBI形 50

式の整数で格納する。また、スキーマ型が子要素を含まないことが自明である場合、具体的にはSEQUENCE型、CHOICE型、MCHOICE型以外、すなわち文字列型、整数型、実数型、ならびに「データ型の適用」702において述べた各種スキーマ型である場合、子要素保有数3014は必ず0個であり、省略できる。

属性を示すスキーマ要素1402は、次のような手順で圧縮スキーマ3020に変換することができる。キー名3017、スキーマ型3018を、タグ名3012、スキーマ型3015と同様な方法で格納する。属性保有数、子要素保有数は必ず0個であり、省略できる。また重複フラグ1412も必ず「属性」であるため、省略できる。

【0124】

スキーマ木の、スキーマ要素間のリンク関係は次のような手順で圧縮スキーマに変換することができる。まずルート要素のスキーマ要素3019を圧縮スキーマに記述する。属性保有数3016を1個以上指定した場合、属性を示す1個以上のスキーマ要素3020を次に連続して記述する。次に子要素保有数を1個以上指定した場合、子要素タグを示す1個以上のスキーマ要素3021を次に連続して記述する。子要素タグを示すスキーマ要素がそれぞれ更に属性や子要素を持つ場合、そのスキーマ要素3021の弟を記述する前に再帰的に属性、子要素を記述する。この方式では、全てのスキーマ要素に対し、自分自身の情報として属性保有数3016、子要素保有数3014を持ち、その後に属性・子要素が順番に記述される、という構造を持つため、リンク関係は完全に保存される。

【0125】

次に、辞書スキーマ3002の説明を行う。「文字列の辞書化」の圧縮を行う場合、圧縮スキーマには、辞書情報を格納する必要がある。辞書情報としては、スキーマトークン3006、辞書トークン3010の後に辞書の総数3022を格納する。その後、辞書情報3023を、辞書総数3022と同数繰り返して記述する。辞書情報3023は、スキーマ要素識別子3024、登録単語数3025、を格納する。その後、単語情報3026を、登録単語数3025と同数繰り返して記述する。単語情報3026は、単語を識別する識別子3027、単語の文字列3028を格納する。

【0126】

スキーマ要素識別子3024は、辞書が対象とするスキーマ要素を示す識別子を記述する。この識別子は、図14のスキーマ要素識別子1404を記載する。図30の構造スキーマ3001にスキーマ要素識別子1404を明示的に記述してもよいが、「スキーマ要素識別子はスキーマ木の出現順の通し番号とする」などの決定方法が明確になっていれば、構造スキーマ3001への記載は省略することが出来る。

【0127】

次に、相対値スキーマ3003の説明を行う。「相対値の適用」の圧縮を行う場合、圧縮スキーマには相対値の起点となる数値を格納する必要がある場合がある。起点情報としては、スキーマトークン3007、辞書トークン3011の後に相対値格納数3029を格納する。その後、相対値情報3030を、相対値格納数3029と同数繰り返して記述する。相対値情報3030は、スキーマ要素識別子3031、起点数値3032、を格納する。

【0128】

スキーマ要素識別子3031は、相対値管理を行うスキーマ要素を示す識別子であり、辞書スキーマにおけるスキーマ要素識別子3024と同様、図14のスキーマ要素識別子1404を記載する。また起点数値3032は、整数、実数、2次元整数、すなわち座標(x、y)を示す二つの数値、2次元実数のいずれかを格納する。起点数値がどの形式で格納されるかは、そのスキーマ要素識別子3031が示すスキーマ要素のスキーマ型3015により決定することができる。

【0129】

次に、図1のスキーマ展開ブロック111について詳細に説明する。スキーマ展開ブロック111は、スキーマ圧縮ブロック110の逆の機能ブロックであり、ファイルや通信で取得した圧縮スキーマ103をスキーマ木102として計算機のメモリ空間上に展開する

。これは、スキーマ圧縮ブロック 110 で説明した手順の逆操作を行うことにより実現できる。

【0130】

スキーマ展開ブロックでは、初めにスキーマトークン 3005、3006、3007 を識別子、それに続く構造トークン 3009、辞書トークン 3010、相対値トークン 3011 を判別し、それぞれに従った展開を行う。本発明の圧縮構造を拡張してあらたなスキーマ構造を追加する場合、スキーマトークンの後に拡張したスキーマを識別するトークンを記述することにより容易に拡張することが出来る。そしてデータトークンが出現した位置で圧縮データが開始されると判定することが出来る。

【0131】

次に、図 1 の B X M L 結合ブロック 114 について詳細に説明する。B X M L 結合ブロック 114 では、圧縮スキーマ 103 および圧縮データ 104 を結合することにより、B X M L 文書 105 を取得する。B X M L 結合ブロック 114 は、単純に圧縮スキーマ 103 の後に圧縮データ 104 を連結させることにより実現できる。

【0132】

次に、図 1 の B X M L 分解ブロック 115 について詳細に説明する。B X M L 分解ブロック 115 では、B X M 文書 105 を圧縮スキーマ 103 と圧縮データ 104 に分解する。B X M L 分解ブロック 115 は、B X M L 文書 105 において、圧縮スキーマ 103 の開始位置と終了位置、圧縮データ 104 の開始位置と終了位置を調査し、開始位置から終了位置までを別データとして保存することにより実現できる。実際の利用として、文字通り分解しなくても、圧縮データの開始位置を取得するだけで十分な場合がある。なお、圧縮スキーマの開始位置は B X M L 文書の開始位置と同じである。圧縮データの開始位置は、スキーマ展開ブロック 111 で説明した通り、データトークン 3008 を判定することにより確定できる。

【0133】

次に、図 1 の B X M L パースブロック 108 について詳細に説明する。X M L パースブロック 107 が X M L 文書をパースし、外部ハンドラルーチンに対し X M L 文書の構造を表すイベントを発行すると同様、B X M L パースブロック 108 は、B X M L 文書をパースし、外部ハンドラルーチンにたいし B X M L 文書の構造を表すイベントを発行する。外部ハンドラルーチンとしてデータ展開ブロック 113 を利用することにより B X M L 文書を X M L 文書に展開することができる。また S A X パースブロック 123 を利用することにより B X M L 文書の標準 S A X インタフェースを実現することができる。

【0134】

図 3 1 に、B X M L パースブロックのアルゴリズムを示す。一時変数として、図 1 4 に示すスキーマ木において現在参照しているスキーマ要素 1401 を指す「参照要素」、圧縮データにおいて現在参照している位置を指す「参照ポインタ」を用意する。「参照要素」の初期位置は X M L 文書のルート要素を指すスキーマ要素を、「参照ポインタ」の初期位置は圧縮データの先頭を指しておく。

[STEP 1] 参照要素の重複フラグ 1406 を取得し、「単独」の場合は STEP 3 以降のノード処理を実行する。重複フラグが「重複」の場合は STEP 2 に進む。

[STEP 2] 圧縮データから子要素の重複数を取得し、重複数と同じ回数、STEP 3 以降のノード処理を実行する。

[STEP 3] 開始タグ名を引数として、ユーザの定義する開始タグ処理を実行する。開始タグ名は参照要素のタグ名 1405 から取得できる。

[STEP 4] 参照要素から属性保有数 1408 を取得し、属性保有数と同じ回数、STEP 5 以降を繰り返す。

[STEP 5] 圧縮データから属性の重複数を取得し、重複数が 1 の場合のみ、STEP 6 以降を実行する。

[STEP 6] 圧縮データから属性値として文字列を取得する。「文字列の辞書化」圧縮 703 を行っている場合は、図 1 1 の方法を用いて、単語識別子 3027 から文字列

10

20

30

40

50

3028への復元を行う。

[STEP 7] 属性キーと属性値を引数として、ユーザの定義する属性処理を実行する。属性キーは参照要素のキー名1411から取得できる。

[STEP 8] 参照要素のスキーマ型1407により、処理を分岐する。簡単のため、スキーマ型がSEQUENCE型の場合と、スキーマ型が文字列型の場合について説明するが、CHOICE型、MCHOICE型、整数型、実数型等、他のスキーマ型でも同様である。

[STEP 9] 本ステップは、スキーマ型が「SEQUENCE型」の場合に実行する。参照要素から子要素保有数1409を取得し、子要素保有数と同じ回数、STEP 1以降を再帰的に繰り返す。ここで繰り返されるSTEP 1以降では、参照要素は対象とする子要素に移動している。

10

[STEP 10] 本ステップは、スキーマ型が「文字列型」の場合に実行する。圧縮データからテキスト要素として文字列を取得する。「文字列の辞書化」圧縮を行っている場合は、図11の方法を用いて、単語識別子3027から文字列3028への復元を行う。

[STEP 11] 文字列型のテキスト要素を引数に、スキーマ型が「文字列型」のテキスト要素に対するユーザ定義処理を実行する。

[STEP 12] 本ステップは、スキーマ型が「整数型」の場合に実行する。圧縮データからテキスト要素として整数を取得する。「相対値の適用」圧縮705を行っている場合は、ここで相対値から元の数値への復元を行う。

20

[STEP 13] 整数型のテキスト要素を引数に、スキーマ型が「整数型」のテキスト要素に対するユーザ定義処理を実行する。

[STEP 14] 終了タグ名を引数として、ユーザの定義する終了タグ処理を実行する。終了タグ名は参照要素のタグ名1405から取得できる。

【0135】

次に、図1のデータ展開ブロック113について詳細に説明する。データ展開ブロック113の目的は、スキーマ木102を利用して圧縮データ104をXML文書101に展開することである。データ展開ブロック113は、図31に示したBXMLパーズブロック108を利用することにより実装できる。具体的には、図31の、各ユーザ定義ステップで以下の処理を実現する。

30

[STEP 3] 開始タグ処理：前回の開始タグの閉じ括弧が出力されていない場合、閉じ括弧「>」を出力する。また、開始タグの括弧「<」を出力し、開始タグ名を取得し、XML文書の開始タグ名として出力する。

[STEP 7] 属性処理：属性キーと属性値を取得し、XML文書の属性部分を出力する。

[STEP 11] テキスト処理（文字列）：前回の開始タグの閉じ括弧が出力されていない場合、閉じ括弧「>」を出力する。また、文字列を取得し、XML文書のテキスト要素として出力する。

[STEP 13] テキスト処理（整数）：前回の開始タグの閉じ括弧が出力されていない場合、閉じ括弧「>」を出力する。また、整数を取得し、文字列に変換して、XML文書のテキスト要素として出力する。

40

[STEP 14] 終了タグ処理：前回の開始タグの閉じ括弧が出力されていない場合、閉じ括弧「>」を出力する。また、終了タグ名を取得し、XML文書の終了タグ部分を出力する。

【0136】

次に、図1のSAXパーズブロック123について詳細に説明する。SAXパーズブロック123の目的は、BXML文書に対して標準のXML文書解析インタフェースの一つであるSAXに準拠したインタフェースをユーザに提供することである。

SAXパーズブロック123は、図31に示したBXMLパーズブロック108を利用することにより実装できる。具体的には、図31の、各ユーザ定義ステップで以下の処理を

50

実現する。

[STEP 3] 開始タグ処理：前回の開始タグが未処理の場合、ユーザにSAX開始タグイベントを発行する。また開始タグ名を記憶する。

[STEP 7] 属性処理：属性キーと属性値を記憶する。

[STEP 11] テキスト処理（文字列）：前回の開始タグが未処理の場合、SAX開始タグイベントを発行する。また、文字列を取得し、SAXテキスト要素イベントを発行する。

[STEP 13] テキスト処理（整数）：前回の開始タグが未処理の場合、SAX開始タグイベントを発行する。また、整数を取得し、文字列に変換して、SAXテキスト要素イベントを発行する。

[STEP 14] 終了タグ処理：前回の開始タグが未処理の場合、SAX開始タグイベントを発行する。また、終了タグ名を取得し、SAX終了タグイベントを発行する。

【0137】

次に、図1のDOMパースブロック124について詳細に説明する。DOMパースブロック124の目的は、XML文書の圧縮構造であるBXML文書に対して、図3で説明した標準XML文書解析インタフェースの一つであるDOMに準拠したインタフェースを提供することにより、BXML文書の解析・更新手段をユーザに提供することである。本発明のDOMパースブロックの特徴は、従来方式と比較し、DOMパースに必要なメモリ空間が小さいこと、実行速度が速いことが挙げられる。メモリ空間を小さくする手段としては、(1)解析する対象がXML文書ではなく、本発明により圧縮されたBXML文書であること、(2)従来のDOMパーサのようにインスタンス木をメモリ上に展開しない、という点が挙げられる。実行速度を高速化するための手段としては、文字列をできるだけ使用しない方式を提供することが挙げられる。たとえばBXML文書において整数や実数は文字列で格納されておらず、整数や実数を表す適切な型で格納されている。本方式では文字列と整数・実数変換を介することなく取得し、また更新する手段を提供する。またXML文書のタグの判定や移動も、文字列ではなく適切な識別子で判定できる手段を提供する。

【0138】

図32を用いて、DOMパースブロック124の実装方式について説明する。XML文書3201は、本発明ではスキーマ木3202と圧縮データ3203で表現される圧縮構造を取る。本発明では、図3で示したインスタンス木302をメモリ空間に展開しない。インスタンス木のノード304に相当する本発明のノードとして、ノード3207のようなBXMLノード構造を定義する。例えば、ノードB1(3207)はXML文書3201の最初のタグBを、ノードB2(3208)は、2番目のタグBを、ノードC1(3209)は最初のタグCを表すノードである。

【0139】

BXMLノードは、内部の変数として、参照スキーマ要素3211、重複数3212、重複カウンタ3213、圧縮データ位置3214を持つ。

【0140】

参照スキーマ要素3211は、そのノードが対応するスキーマ木3202のスキーマ要素を指す。ノードB1(3207)の参照スキーマ要素3211は、スキーマ木3202の対応するスキーマ要素3205を指している。

【0141】

重複数3212は、ノード3207と連続する兄弟の位置に存在する、同名のタグの個数である。XML文書3201にはタグBが2つ連続して存在するため、重複数3212は2である。これは、圧縮スキーマ3203のNbから取得される値である。

【0142】

重複カウンタ3213は、ノード3207が同名の連続するタグの、0から始まる何番目のタグかを示すカウンタである。ノード3207は2つあるタグBの最初のタグであるため、重複カウンタは0である。一方、ノード3208の重複カウンタは、2番目のタグで

10

20

30

40

50

あるため、重複カウンタは1である。

【0143】

圧縮データ位置3214は、そのノードの情報が圧縮データ3203のどの位置に存在するかを示すポインタである。ノードB1(3207)の圧縮データ位置3214は、圧縮データ3203のb1の位置、ノードB2(3208)の圧縮データ位置3214は、圧縮データ3203のb2の位置を指している。

【0144】

図32の構成で、図3に示す標準DOMインタフェースの5つのメソッド、すなわち(1)属性値の取得、(2)テキスト要素の取得、(3)子へ移動、(4)弟へ移動、(5)親へ移動、の実装アルゴリズムを以下に示す。

図33を用いて、DOMパースブロック124における「属性値取得」のアルゴリズムについて説明する。「属性値取得」は、引数として属性のキーを入力し、そのノードにおけるキーに対応する属性値を返す処理である。

[STEP 1] ノードの参照スキーマ要素3211から、そのノードに対応したスキーマ要素を取得し、それに接続する全属性についてSTEP 2以降を実行する。

[STEP 2] 圧縮データから重複数を取得する。

[STEP 3] 重複数が0の場合はその属性は省略されていることになる。STEP 1に戻り、次の属性を探索する。全ての属性を探索し尽くせば、全ての属性が省略されていることになり、「属性値取得」は失敗する。重複数が1の場合は、属性が発見されたことになり。現在の属性のキー名と引数キーとを比較し、一致する場合、圧縮データから文字列を取得し、終了する。

図34を用いて、DOMパースブロック124における「テキスト要素取得」のアルゴリズムについて説明する。

[STEP 1] ノードの変数である圧縮データ位置3214は、ノードを示す圧縮データの先頭位置を指している。データ先頭にはそのノードの属性情報が記述されているため、これを読み飛ばす。属性スキップ処理は、図33を利用することにより実現できる。

[STEP 2] ノードの参照スキーマ要素3211から、そのノードに対応したスキーマ要素を取得し、そのスキーマ型に応じた方法で圧縮データからテキスト要素を取得する。

図35を用いて、DOMパースブロックにおける「子へ移動」のアルゴリズムについて説明する。

[STEP 1] ノードの変数である圧縮データ位置3214は、ノードを示す圧縮データの先頭位置を指している。ノードの先頭にはそのノードの属性情報が記述されているため、これを読み飛ばす。

[STEP 2] ノードの参照スキーマ要素3211から、そのノードに対応したスキーマ要素を取得し、それに接続する全子要素についてSTEP 3以降を実行する。子要素が存在しない場合、「子への移動」は失敗する。

[STEP 3] 子要素の重複数を取得する。子要素の重複フラグが「単独」の場合、重複数は1である。子要素の重複フラグが「重複」の場合、圧縮データから子要素の重複数を取得する。

[STEP 4] 重複数が0の場合、その子要素は省略されているため、STEP 2に戻り、次の子要素を探索する。全ての子要素を探索し尽くせば、全ての子要素が省略されていることになり、「子への移動」は失敗する。一方、重複数が0より大きい場合、子要素が発見されたことになり。パラメータを設定し、「子への移動」は成功終了する。

図36を用いて、DOMパースブロック124における「弟へ移動」のアルゴリズムについて説明する。

[STEP 1] ノードの変数である圧縮データ位置3214を、自ノードの最後までパースした状態に移動する。自ノードが子要素を含む場合、子要素の情報もスキップする。これは図1のXMLパースブロック108を利用することにより実現できる。

[STEP 2] ノードの変数である重複カウンタ3213が、同じノードの変数であ

10

20

30

40

50

る重複数 3 2 1 2 より小さい場合、同名の弟タグが残っていることになる。パラメータを設定し、「弟への移動」は成功終了する。

[STEP 3] ノードの参照スキーマ要素 3 2 1 1 から、そのノードに対応したスキーマ要素を取得し、そのスキーマ要素の後に存在する全弟要素についてSTEP 4以降を実行する。弟要素が存在しない場合、「弟への移動」は失敗する。

[STEP 3] 弟要素の重複数を取得する。弟要素の重複フラグが「単独」の場合、重複数は1である。弟要素の重複フラグが「重複」の場合、圧縮データから弟要素の重複数を取得する。

[STEP 4] 重複数が0の場合、その弟要素は省略されているため、STEP 3に戻り、次の弟要素を探索する。全ての弟要素を探索し尽くせば、全ての弟要素が省略されていることになり、「弟への移動」は失敗する。一方、重複数が0より大きい場合、弟要素が発見されたことになる。パラメータを設定し、「弟への移動」は成功終了する。

図32を用いて、DOMパースブロック124における「親に移動」のアルゴリズムについて説明する。

【0145】

「親へ移動」の機能を実装するためには、たとえば図32のノード3207において、その親ノードに対する状態変数、すなわち参照スキーマ要素3211、重複数3212、重複カウンタ3213、圧縮データ位置3214を取得する必要がある。そこで、全てのノードにおいて、ルート要素からそのノードに到るまでの全ての状態変数を可変長配列として記憶する。その結果、配列の先頭にはルート要素の状態変数が、配列の最後尾には自分のノードの状態変数が格納されるという構造になる。「親へ移動」の機能は、状態変数配列の最後尾の状態変数を取り除き、最後から2番目の状態変数を自分のノードの状態変数とすることで実現できる。また前述した「子への移動」機能では、子への移動の結果の状態変数を状態変数配列の最後尾に追加する。「弟への移動」機能では、配列の最後尾の状態変数を弟への移動の結果の状態変数に置き換える。これにより「親へ移動」機能が実現できる。

【0146】

このような状態変数配列をノードに持たせると速度低下の原因となる。そこで、「親へ移動」の機能が不要な場合は、状態変数配列を持たないバージョンのノードを利用することにより高速化を図ることができる。

【0147】

また、「子に移動」、「弟に移動」のアルゴリズムを利用して、「指定タグ名を持つ最初の子に移動」、「指定タグ名を持つ最初の弟に移動」、「指定タグ名を持つ次の弟に移動」という高機能ノード移動機能を実現できる。

この実現手段としては、ノード3207の保有する情報として、さらに検索対象タグ名を付加する。「指定タグ名を持つ最初の子に移動」を実現するには、「子に移動」を実行し、検索対象タグ名と移動したノードのタグ名を比較し、それが一致するまで繰り返し「弟に移動」を実行すればよい。「指定タグ名を持つ最初の弟に移動」を実現するには、「弟に移動」を実行し、検索対象タグ名と移動したノードのタグ名を比較し、それが一致するまで繰り返し「弟に移動」を実行すればよい。「指定タグ名を持つ次の弟に移動」を実現するには、前回の検索で設定した検索対象タグ名を保持しておき、検索対象タグ名とノードのタグ名が一致するまで繰り返し「弟に移動」を実行すればよい。

【0148】

同様に、「指定タグ名を持つ最初の孫に移動」、「指定タグ名を親に持つ最初の孫に移動」、「指定タグ名を親に、別の指定タグ名を持つ最初の孫に移動」、「指定条件の次の孫に移動」という二階層のノード移動機能も実現できる。この実現手段としては、ノード3207の保有する情報として、検索対象タグ名に加え、検索対象親タグ名を付加し、判定対象ノードのタグ名・親タグ名と検索対象タグ名・検索対象親タグ名を比較すればよい。

【0149】

次に、図37を用いて、XML文書のタグ名の判定を文字列ではなく識別子で判定する手

10

20

30

40

50

段について説明する。例えばXML文書3701で、タグ<A>が多数存在する状況において、タグ<A>の下にタグを検索し、そのテキスト要素を取得し、何らかのユーザ処理を行う場合を考える。従来の検索方式では、タグ名文字列による判定3703を行っていた。すなわち、ブロック3705において、ユーザの指定したタグ名の文字列と、XML文書3701に記述された文字列との文字列比較を行う。本発明が提供する識別子による判定方式3604では、あらかじめ検索したいタグ名のスキーマ要素識別子を取得しておき(3606)、ループの中のブロック3607において、ユーザの取得したスキーマ要素識別子と、BXML文書3601の参照スキーマに記述されるスキーマ要素識別子とを比較する。スキーマ要素識別子は、図14のスキーマ要素識別子1404の説明で示したとおり数値で表現できるため、文字列比較より高速な比較が可能となる。

10

【0150】

XMLの名前空間に対応すると、速度差はより顕著となる。図3602に示す名前空間に対応したXML文書において、タグを判定する場合を考える。タグ<X:B> 3711、タグ<Y:B> 3712、タグ<Z:B> 3713の三つのタグにおいて、タグ<X:B> 3711とタグ<Y:B> 3712は異なるタグ、タグ<X:B> 3711、タグ<Z:B> 3713は同じタグであると判定しなければならない。従来の判定手段は、タグ<X:B> 3711の修飾子「X」の名前空間定義3708とタグ<Y:B> 3712の修飾子「Y」の名前空間定義3709が異なるためタグ<X:B> 3711とタグ<Y:B> 3712は異なるタグ、タグ<X:B> 3711、タグ<Z:B> 3713は修飾子「X」「Z」が異なるが、その名前空間定義3708と3710は同じであるため同じタグ、と判定していた。このような処理をブロック3705で多数回行うと時間がかかる。

20

一方、本方式では、タグ<X:B> 3711とタグ<Z:B> 3713のスキーマ要素識別子は同じ値、<Y:B> 3712のみ異なる値となる。ループの外の処理ブロック3706で1回のみ名前空間の判定を行い、多数回繰り返されるループ内の処理ブロック3707でスキーマ要素識別子の比較だけを行うことにより、高速なタグ名判定が実現できる。

【0151】

また、DOMパースブロック124の応用例として、図3で説明した標準XML文書解析インタフェースであるDOMインタフェースの他に、特定の分野に特化した高機能なインタフェースを提供することも考えられる。高機能なインタフェースの実現方式として、図3のノード304のように、XML文書の一つのタグに対応したノードではなく、意味を持ったタグの集合を一つのノードとして管理し、ノード移動機能を提供することが考えられる。

30

【0152】

図38を用いて、地理情報分野に特化したノードの説明を行う。図38は、図6と同様、地理情報分野のXML文書の例であり、一枚の地図を表す。地図である地物集合を表すMapタグの孫の位置に、個々の地物を表すRoadタグ3801、3807、Houseタグ3808が存在する。Roadタグ3801に着目すると、孫の位置に道路中心線(centerLineOf)を表す折線図形(LineString)タグ3802、交差点(junction)を表す点集合(MultiPoint)タグ3803が存在する。さらに点集合タグ3803の孫の位置には、点集合を構成する個々の点(Point)タグ3804、3805、3806が存在する。

40

【0153】

ここで、地物を表す地物ノード、図形を表す図形ノードを定義する事が出来る。ノード3801、3807、3808は地物ノードである。また、ノード3802、3803、3804、3805、3806は図形ノードである。地物ノードの機能として、「最初の地物ノードに移動」、「最初の道路地物ノードに移動」、「次の地物ノードに移動」、「地物に含まれる最初の道路中心線を表す折線を示す図形ノードを取得」、「地物に含まれる最初の折線図形ノードを取得」、「次の折線図形ノードに移動」などのインタフェースが

50

実現できる。また図形 3803 のような複数の図形を含む図形に対し、「図形に含まれる最初の図形ノードを取得」などが定義できる。これらのインタフェースは、前述の「子に移動」「孫に移動」等の基本インタフェースを組み合わせることにより実現できる。

さらに地物ノード、図形ノードのように、地理情報分野に特化したノードでは、汎用の構造化文書が持たない特徴を利用したインタフェースが構築できる。例えば地物ノード 3801 は、道路・家屋・河川など利用者が自由に定義できるが、図形ノード 3802 は、点・折線・ポリゴンなど限られた固定個数の種類に限られる。この性質を利用して、「地物に含まれる最初の図形ノードに移動」というインタフェースが実現できる。汎用の構造化文書に対するノードではタグ名を指定する必要があったが、「地物に含まれる最初の図形ノードに移動」というインタフェースは、あらかじめ固定個数の図形タグを記憶しておき、判定時にその全ての図形タグと比較を行うことにより実現できる。また、「地物の属性を取得する」、「図形の座標系を取得する」、「図形の重心を取得する」等の地理情報分野に特有なインタフェースを追加することにより、地理情報分野の XML 文書を効率的に処理することの出来るインタフェースを提供することが出来る。

10

【0154】

次に、図 1 の DOM 更新ブロック 126 について詳細に説明する。DOM 更新ブロック 126 の目的は、既存の XML 文書の更新を行うことである。

【0155】

XML 文書の更新は、図 32 で説明した DOM パースブロック 124 のインタフェースを用いて更新対象ノードに移動し、DOM 更新ブロック 126 が提供するインタフェースである、(1) 属性値変更、(2) テキスト要素変更、(3) 属性の挿入、(4) 属性の削除、(5) 子要素の挿入、(6) 子要素の削除、(7) 子要素の置換を実行することにより行う。基本的に、テキスト要素変更、子要素の置換は、子要素の削除・子要素の挿入で表現することが出来る。また属性値変更、属性の挿入、属性の削除は、本発明では属性は子要素の特別な形としているため、同様に子要素の削除・子要素の挿入で表現することが出来る。そこで、子要素の挿入、子要素の削除方式について説明する。

20

【0156】

更新ブロックの実装は様々な方式が考えられるが、更新機能に必要な条件として、次の 3 点が考えられる。(1) 多数のノードを高速に更新できること、(2) 更新のやりなおし (UNDO) ができること、(3) 更新情報を差分として抽出できること。特に三番目の更新差分は、地理情報サーバが地図を端末に配信し、端末側で地図を編集・更新し、それを地理情報サーバに反映するという利用を考えると、端末が更新差分のみをサーバに送信することにより通信コストが大きく低減するため、重要である。以下では、これらの条件を満たす更新方式について説明する。

30

【0157】

図 39 を用いて DOM 更新ブロック 126 の説明を行う。初めに、「子要素の挿入」の例として、XML 文書 3901 を XML 文書 3902 に更新する場合を考える。これは二つのタグ C の間に新たなタグ C が挿入された状態を示す。

【0158】

「子要素の挿入」機能の実装例を 2 つ示す。第 1 の実装例を 3907 に示す。挿入されたタグ C のテキスト要素 3922 が挿入され、タグ C の重複数 3921 が元のデータ 3906 と比べ 1 増加している。しかしこの方式では、更新時に更新データ実体 3922 の挿入と重複数 3921 の更新と二カ所の変更が必要となり、複数のノードを同時更新する場合に処理が煩雑となり、処理速度が落ちる。また更新後にやりなおしや差分抽出が出来ない。

40

【0159】

「子要素の挿入」機能の第 2 の実装例を 3908 に示す。3908 は更新された XML 文書 3902 を表現する圧縮データである。挿入されたタグ C の部分に挿入データ 3915 を挿入する。挿入データ 3915 は、挿入トークン 3923、挿入タグのスキーマ要素識別子 3924、挿入データ実体 3925 から構成される。挿入トークン 3923 は、デー

50

タグが挿入されたことを示すトークンである。スキーマ要素識別子 3924 は、挿入されたデータがタグ C であることを示す識別子であり、図 14 のスキーマ要素識別子 1404 と同じものである。挿入データ実体 3925 は挿入されたノードの圧縮データであり、3922 と同じものである。

【0160】

スキーマ要素識別子 3924 を設ける理由は、挿入されたノードが子要素群の最後に存在する場合への対処である。XML 文書 3903 はタグ群 C の最後に新しいタグ C を挿入している。これを表す圧縮データは、圧縮データ 3909 における挿入データ 3916 である。また、XML 文書 3904 は、タグ C の親タグ B の最後に新しいタグ B を挿入している。これを表す圧縮データは、圧縮データ 3910 における挿入データ 3917 である。挿入データ 3916 と挿入データ 3917 は、スキーマ要素識別子がないと挿入位置の区別が付かない。これを区別するために挿入タグ識別子 3924 を設けている。

10

【0161】

次に、「子要素の削除」の例として、XML 文書 3901 を XML 文書 3905 に更新する場合を考える。これは、二番目のタグが削除された状態を示す。

【0162】

「子要素の削除」機能の実装例を 4 つ示す。第 1 の実装例を 3911 に示す。元のデータ 3906 と比較し、重複数 Nc が 1 減少しており、また 2 番目のタグの情報 c2 が削除されている。これも挿入の場合と同様な問題を持つため、別の方式が必要となる。

【0163】

「子要素の削除」機能の第 2 の実装例を 3912 に示す。これは削除ノードの情報全体を削除トークン 3926 で置き換えるという方式である。この方式は、データがコンパクトになり、また重複数 Nc を更新する必要がないので高速な削除処理を行うことができる。

20

【0164】

「子要素の削除」機能の第 3 の実装例を 3912 に示す。これは削除ノードの前方に削除トークン 3927 を挿入するという方式である。削除対象のデータ c2 は、データ 3926 のように実際は削除されていない。この方式は、削除データ情報全体を削除トークン 3926 で置き換えるという方式である。この方式の利点は、削除されるノードの情報が残されるため、「更新のやりなおし」や更新差分抽出が行えることである。

【0165】

「子要素の削除」機能の第 4 の実装例を 3914 に示す。これは削除ノードの位置に削除トークン 3928 と、削除ノードのサイズ 3929 を記述するという方式である。削除ノードの残りの情報はごみ 3930 として残るが、利用者は削除トークン 3928 が出現したら 3929 で示されるサイズだけデータの読み込みをスキップすることでごみ 3930 を無視することが出来る。この方式の利点は、更新処理が高速となるという点である。実装例 3911、3912、3913 はいずれも、「子要素の削除」処理により削除データ部分のサイズ変更が発生し、削除データより後方の圧縮データの位置をシフトするためにコピーが発生する。第 4 の実装例では、削除トークン 3928 と削除サイズ 3929 の合計が元のデータよりも大きくなるというわずかな例外を除いて、圧縮データの位置シフトが発生しないため、全体的な更新パフォーマンスの向上につながる。

30

40

【0166】

上記に示した手法で更新を行う場合、圧縮データに挿入トークン・削除トークン等の余分なデータが蓄積されていく。そのため定期的にガベージコレクションを実施することにより挿入トークン・削除トークンを削除する。

【0167】

図 40 を用いて、挿入・削除されたノードに対応した BXML パースブロック 108 のアルゴリズムを説明する。これは図 31 で説明した BXML パースブロック 108 を拡張したものである。

【0168】

STEP 1 ~ STEP 14 までの処理は図 31 の BXML パースブロック 108 と同

50

一であるため、説明を割愛する。ここでは、新たに挿入されたSTEP 15 ~ STEP 18について説明を行う。

【STEP 15】本ステップでは、図39のXML文書3902のように、新しいノードが重複する子要素群の最後尾を除く任意の位置に挿入された場合の処理を行う。STEP 2で子要素の重複数だけノード処理を実行するが、挿入された子要素は重複数には含まれないため、重複カウンタを進めてはならない。圧縮データからデータを取得し、それが挿入トークン3923である場合、重複カウンタを進めずに繰り返しノード処理を実行する。

【STEP 16】図39のXML文書3903のように、新しいノードが重複する子要素群の最後尾に挿入された場合の処理を行う。圧縮データからデータを取得し、それが挿入トークン3923であり、スキーマ識別子3924が現在の参照要素と同一である限り、繰り返しノード処理を実行する。

【STEP 17】本ステップは、子要素が削除された場合に対応する。圧縮データからデータを取得し、それが削除トークン3926の場合、圧縮データの参照ポイントを削除トークンの次に合わせ、ノード処理を終了する。

【STEP 18】本ステップは、属性が削除された場合に対応する。圧縮データからデータを取得し、それが削除トークン3926でない場合のみ、STEP 6以降を実行する。

【0169】

次に、図1のBXML作成ブロック121について詳細に説明する。BXML作成ブロック121の目的は、BXML文書を直接作成することである。BXML作成ブロック121は、図20で説明したとおり、データ圧縮ブロック112に対し図2に示したSAXイベントを発行することにより、圧縮データを直接作成する。BXML作成ブロック121の機能は、ユーザプログラム120にBXML作成のためのインタフェースを提供し、これをデータ圧縮ブロック112が受け付けるSAXイベントに変更することである。

【0170】

図41に、BXML作成ブロック121がユーザプログラム120に提供するインタフェースを示す。BXML作成ブロック121のインタフェースは、以下の12個のメソッド、すなわち、文書作成系の(1)文書開始メソッド、(2)文書終了メソッド、ノード作成系の(3)子要素作成メソッド、(4)弟要素作成メソッド、(5)親へ移動メソッド、(6)先祖へ移動メソッド、データ作成系の(7)属性の作成メソッド、(8)文字列の作成メソッド、(9)整数の作成メソッド、(10)実数の作成メソッド、(11)空文字の作成メソッド、(12)ノードのコピー、から構成される。

(1)文書開始メソッド4108は、BXML文書を作成開始する時に実行する。文書開始メソッドは、BXML文書作成のための初期処理を行う。これはデータ圧縮ブロック112に対し、図2の文書開始処理201を呼び出すことで実現する。

(2)文書終了メソッド4109は、BXML文書の作成を完了する時に実行する。文書終了メソッドは、BXML文書の終了タグのないタグ全てに終了タグを付加し、BXML文書作成の終了処理を行う。これは、スキーマ木の参照要素から親要素をさかのぼり、ルート要素に到るまで繰り返しデータ圧縮ブロック112に対し終了タグ処理204を呼び出し、その後文書終了処理202を呼び出すことで実現する。

(3)子要素作成メソッド4110は、作成途中のBXML文書における現在の位置に、新たな子要素を追加する時に実行する。子要素作成メソッドは、引数で指定したタグ名の子要素をBXML文書に追加する。これは、データ圧縮ブロック112に対し、引数で指定されるタグ名を用いて開始タグ処理203を呼び出すことにより実現できる。ここで、作成途中のBXML文書は、図41の4104と4105に示す2つの状態がある。これをそれぞれ開状態、閉状態と呼ぶことにする。開状態4104は、最後に開始タグが記述されている状態であり、閉状態4105は最後に終了タグが記述されている状態である。子要素作成メソッドは開状態4104で実行されることを前提としており、閉状態4105で実行するとエラーを返す。また別の実装方式として、閉状態4105で実行された時

10

20

30

40

50

に、現在のノード B と同じ弟ノード B を作成し、その弟ノード B に対して子ノード C を追加しても構わない。この場合、スキーマ木の参照要素のタグ名を用いて開始タグ処理 203 を呼び出し、その後引数で指定されるタグ名 C を用いて開始タグ処理 203 を呼び出すことにより実現できる。

(4) 弟要素作成メソッド 4111 は、作成途中の BXML 文書における現在の位置に、新たな弟要素を追加する時に実行する。弟要素作成メソッドは、引数で指定したタグ名の弟要素を BXML 文書に追加する。これは、閉状態 4107 において、データ圧縮ブロック 112 に対し、引数で指定されるタグ名を用いて開始タグ処理 203 を呼び出すことにより実現できる。閉状態 4106 で実行するとエラーを返す。別の実装として、閉状態 4106 で実行された時に、現在のノードを空タグとして閉じ、引数で指定したタグ名の弟要素を BXML 文書に追加しても構わない。この場合、スキーマ木の参照要素のタグ名を用いて終了タグ処理 204 を呼び出し、その後引数で指定されるタグ名 C を用いて開始タグ処理 203 を呼び出すことにより実現できる。

(5) 親へ移動メソッド 4112 は、現在のノードの書き込みを終了し、現在のノードから見て親の弟、あるいは親の親の弟など、上位階層のノードの書き込みを開始する直前に実行する。親へ移動メソッドは引数を持たない。閉状態 4109 では親のノードを閉じる。これはスキーマ木の参照要素における親のタグ名を用いて終了タグ処理 204 を呼び出すことにより実現できる。また閉状態 4108 では、現在のノードを空タグとして閉じ、その後親のノードを閉じる。これはスキーマ木の参照要素のタグ名を用いて終了タグ処理 204 を呼び出し、その後親のタグ名を用いて終了タグ処理 204 を呼び出すことにより実現できる。

(6) 先祖へ移動メソッド 4113 は、指定したノードまで親へ移動メソッドを繰り返し発行する。これはスキーマ木の参照要素から指定したタグ名のノードまで繰り返し遡って検索し、遡った回数を利用して親へ移動メソッドを繰り返し実行することにより実現できる。

【0171】

また、図 41 の BXML 作成ブロック 121 が提供するインタフェースの異なる実装として、子要素作成メソッド 4110、弟要素作成メソッド 4111、親へ移動メソッド 4112、先祖へ移動メソッド 4113 で示されるノード作成メソッド群に対して直接データ圧縮ブロック 112 に対して処理を実行するのではなく、一時的なバッファにメソッドを記憶しておき、属性の作成メソッド 4114、文字列の作成メソッド 4115 等のデータ作成メソッド群が発行された時点で一度に記憶しておいたノード作成メソッド群を発行するという方式が考えられる。例えば 0 個以上の子要素を作成したい利用者は、最初に子要素作成メソッド 4110 を発行し、文字列の作成メソッド 4115 などのデータ作成メソッドと弟要素作成メソッド 4111 を作成したい子要素数と同数回発行し、最後に親へ移動メソッド 4112 を発行する。この際、作成すべき子要素が 0 個であった場合、本発明の最初の実装では、子要素作成メソッド 4110 が必ず発行されるため、不要な子要素が出現してしまう。ここでノード作成メソッド群を一時的なバッファに記憶しておき、子要素作成メソッド 4110 を発行後、データ作成メソッド群が一度も発行されることなく親へ移動メソッド 4112 が発行された場合、バッファに記憶された最後の子要素作成メソッド 4110 をキャンセルする動作を実現すれば、上記のような不要な子要素の出現を抑制することができる。

(7) 属性の作成メソッド 4114 は、作成途中の BXML 文書における現在のノードに、指定したキーと値を持つ属性を挿入する。これは指定したキーと値を用いてデータ圧縮ブロック 112 に対し、属性処理 206 を呼び出すことにより実現できる。属性の作成メソッドは閉状態 4112 でのみ許可され、閉状態ではエラーとなる。

(8) 文字列の作成メソッド 4115 は、作成途中の BXML 文書における現在のノードに、テキスト要素として指定した文字列を挿入する時に使用する。これは指定した文字列を用いてデータ圧縮ブロック 112 に対し、テキスト要素処理 205 を呼び出し、その後スキーマ木の参照要素のタグ名を用いて終了タグ処理 204 を呼び出すことにより実

10

20

30

40

50

現できる。文字列の作成メソッドは開状態 4 1 1 2 でのみ許可され、閉状態ではエラーとなる。

(9) 整数の作成メソッド 4 1 1 6 は、作成途中の B X M L 文書における現在のノードに、テキスト要素として指定した整数を挿入する。これは文字列の作成メソッドにおけるテキスト要素処理 2 0 5 の呼び出しの代わりに、指定した整数を用いて整数要素処理を呼び出すことにより実現できる。整数要素処理とテキスト要素処理 2 0 5 の内部処理は全く同じであり、引数として整数を用いることにより余分な文字列変換を実行しない、高速化のためのメソッドである。

(10) 実数の作成メソッド 4 1 1 7 は作成途中の B X M L 文書における現在のノードに、テキスト要素として指定した整数を挿入する。実現手段は整数の作成メソッドと同様である。またここでは整数と実数の場合のみを示したが、図 4 の「データ型の適用」圧縮 4 0 2 で説明した多数のデータ型についても同様なメソッドが定義できる。

(11) 空タグの作成メソッド 4 1 1 8 は、現在のタグを閉じる時に使用する。これは例えば属性のみを持ち、子要素やテキスト要素を持たないタグを記述する場合に利用する。これはスキーマ木の参照要素のタグ名を用いて終了タグ処理 2 0 4 を呼び出すことにより実現できる。

(12) ノードのコピーメソッド 4 1 1 9 は、現在のタグ位置に別の B X M L 文書のノード群を挿入する時に利用する。

【0 1 7 2】

図 4 2 を用いてノードのコピーメソッド 4 1 1 9 を詳細に説明する。作成中文書 4 2 0 1 に対し、挿入元文書 4 2 0 3 の中の、2 つのタグ B で構成されるノード群 4 2 0 4 を挿入し、文書 4 2 0 2 にする例を考える。図 2 8 で説明したとおり、データ圧縮ブロック 1 1 2 ではスキーマ木 2 8 0 2 を利用して圧縮データ 2 8 0 3 を作成していく。そのためコピーメソッド 4 1 1 9 の実装は、スキーマ木の挿入対象タグ B のスキーマ要素 4 2 0 5、作成中の圧縮データ 4 2 0 8 に対し、挿入元圧縮データ 4 2 1 0 の 2 つのタグ B で構成されるノード群を示す圧縮データ 4 2 1 4 を挿入し、スキーマ要素 4 2 0 6、圧縮データ 4 2 0 9 にする操作ということになる。これは以下のステップで実行する。

[STEP 1] 挿入するノード数(本例の場合 2)を、挿入先スキーマ要素 4 2 0 5 の一時変数である重複カウンタに加える。

[STEP 2] 圧縮データ 4 2 1 4 を作成中の圧縮データにコピーする。

【0 1 7 3】

本例では、挿入元圧縮データ 4 2 1 4 であるタグ B は単純な文字列型としたが、タグ B がさらに子要素を持つ複雑な構造を持っていても同様にコピーすることが出来る。また、挿入元文書 4 2 0 3 のタグ B のスキーマ木が挿入先文書のタグ B のスキーマ木と異なる場合、挿入元圧縮データ 4 2 1 4 を単純にコピーするのではなく、挿入先文書のタグ B のスキーマ木に一致する構造に変換を行った後にコピーを行うことになる。このような構造変換は、図 1 の B X M L パースブロック 1 0 8 を用い、挿入元文書のスキーマ木 1 0 2 と挿入先文書のスキーマ木型変換ブロックを実装することにより実現できる。

【0 1 7 4】

またノードのコピーメソッド 4 1 1 9 を利用して、複数の B X M L 文書を結合する処理を実現することができる。結合対象の複数の B X M L 文書の圧縮スキーマが等しい場合、図 2 0 で示した生成処理 1 3 0 5 を用いて同じ圧縮スキーマを持つ新たな B X M L 文書を作成する。スキーマ木 1 0 2 はあらかじめ用意しておいた B X M L 文書 1 0 5 からスキーマ展開ブロック 1 1 1 を経て作成する。そして結合したいノード群に対し、ノードのコピーメソッド 4 1 1 9 を利用して圧縮データ 4 2 1 4 を作成中の圧縮データにコピーすればよい。結合対象の複数の B X M L 文書の圧縮スキーマが異なる場合、図 2 0 で示したスキーマパース処理 1 1 7、スキーマ更新処理 1 1 9 を利用して結合する圧縮スキーマの双方の構造を含んだ最小公倍数的な圧縮スキーマを作成し、さらに結合したいノード群に対し、ノードのコピーメソッド 4 1 1 9 を利用して、挿入先文書のスキーマ木に一致する構造に変換を行った後にコピーを行うことにより実現できる。

10

20

30

40

50

【0175】

次に、図1のスキーマパースブロック117について詳細に説明する。スキーマパースブロック117の目的は、スキーマ木の移動インタフェースをユーザに提供することである。

【0176】

スキーマパースブロック117は、ユーザプログラム118に対し、以下の3つの移動系メソッドを提供する。(1)親へ移動、(2)長男に移動、(3)次の弟に移動。これは図25に示す参照要素2504を用いて実装する。例えば「親へ移動」メソッドは、図14で説明した「親スキーマ要素へのポインタ」を利用し、参照要素2504を親に移動させる。「長男に移動」メソッドは、「子スキーマ要素へのポインタ」を利用し、参照要素2504を最初の子に移動させる。「次の弟に移動」メソッドは、親要素の持つ期待要素2505を一つ後ろの要素に移動させ、これを参照要素とする。また、スキーマパースブロック117は、図14に示すスキーマ木の各変数を参照する機能を提供する。これにより、「指定したタグ名の子に移動」「指定したタグ名の弟に移動」メソッドは、スキーマ要素のタグ名1405が一致するまで「次の弟へ移動」メソッドを実行することにより実現できる。

10

【0177】

次に、図1のスキーマ更新ブロック119について詳細に説明する。スキーマ更新ブロック119の目的は、スキーマの変更を行うことにより、BXML文書の一括変換を実現することである。

20

【0178】

スキーマ更新ブロック119は、上述したスキーマパースブロック117を用いて、更新の必要なスキーマ要素に移動し、図14に示すスキーマ要素の状態変数を修正する機能を提供する。

【0179】

スキーマ要素1401の、タグ名1405を変更することにより、XML文書全体の指定した開始タグ・終了タグを一括変更することができる。XML文書内に同じ名称のタグが多数存在していても、本機能によりスキーマ木の1箇所を修正するだけで全てのタグを変更することができるため、高速な変換を実現することができる。

【0180】

これは例えば名前空間の変換に利用できる。図37のXML文書3702において、タグ<X:B> 3711に対し、名前空間修飾子のない形のタグに変換する場合、また異なる修飾子<Z:B>に変換する場合、本機能によりスキーマ木の1箇所を変更することで全てのタグ<X:B>の変更ができる。

30

【0181】

また、特定のスキーマ要素に隠蔽の設定を行うことにより、XML文書の一部を隠蔽することができる。図14のスキーマ要素1401において、新たに隠蔽フラグを格納する。スキーマ更新ブロック119により隠蔽フラグをTRUEにしておき、BXMLパースブロック108において隠蔽されたスキーマ要素を持つタグをスキップすることにより隠蔽機能は実装できる。それにより例えば、図43のXML文書4301において、Roadタグ以外を隠蔽し、XML文書4302に変更することができる。またXML文書4303において、Roadタグの子要素であるLaneタグ、Pavementタグを隠蔽し、XML文書4304に変更することが出来る。

40

【0182】

隠蔽機能により、XML文書の高速な加工処理が行える。例えば地理情報サーバがXML文書4301を保有している時、地理情報クライアントが地理情報サーバに対し、道路を示すRoadのみを要求する場合を考える。地理情報サーバがXML文書4301を加工し、XML文書4301を作成するには時間がかかるが、隠蔽機能を用いることにより、XML文書内に同じ名称のタグが多数存在していても、スキーマ木の1箇所を修正するだけで全ての関連タグを隠蔽することができるため、高速な加工を実現することができる。

50

【0183】

図44を用いて、隠蔽されたノードに対応したBXMLパースブロック108のアルゴリズムを説明する。これは図31で説明したBXMLパースブロック108を拡張したものである。

【0184】

STEP 1、STEP 2、並びにノード処理は図31と同一であるため、説明を割愛する。ここでは、新たに挿入されたSTEP 3では、現在参照しているスキーマ要素の隠蔽フラグを確認し、隠蔽状態であればノードのスキップを行う。

【0185】

以上により、図1の全処理ブロックの詳細な説明を終了する。次に、本発明の実施例として、圧縮XMLを利用した地理情報システムについて説明する。 10

【0186】

図45を用いて、本発明の実施形態について説明する。図45は一般的なインターネットのWebサービスの構成を示す三階層モデルで構成された地理情報システムである。三階層モデルは、ユーザが利用する端末であるクライアント層4501、データを加工するアプリケーション層4502、データを格納するデータベース層4503、から構成される。

【0187】

クライアント層4501は、多数の一般ユーザが利用する汎用クライアント4504と、少数の特定ユーザが利用する専用クライアント4506に分類できる。汎用クライアント 20
では主に、地図の参照や利用を行う。汎用クライアントは多数の一般ユーザが利用するため、複数のベンダ製品での実装が前提となる。そのため、汎用クライアントとサーバ間のデータ形式や通信プロトコルは標準化がされている必要がある。従来技術で説明したとおり、標準化されたデータ形式として、JIS規格のG-XMLやOGCのGML、標準化された通信プロトコルとして、OGCのWFS、WMSが存在する。一方専用クライアント4506では主に、地図の保守や更新を行う。専用クライアントのユーザはデータベース層4503やアプリケーション層4502の管理者や、それに近い組織のユーザとなるため、専用クライアントはデータベース層4503やアプリケーション層4502と同じベンダ製品での実装も可能であり、必ずしも標準化がされている必要はない。

通信環境としては、一般に、アプリケーション層4502とクライアント層4501の間 30
は電話回線や無線回線などの低速回線4509で、アプリケーション層4502とデータベース層4503はLAN(Local Area Network)などの高速回線4512で接続される。ただし例えば県庁向け地理情報システムにおいて、地図データが各市町村のデータベースに分散され、県庁がゲートウェイとなり、各市町村のデータベース層を統合して検索する場合など、データベース層4503とアプリケーション層4502が異なる組織に分散し、組織間が低速回線4513で接続されている場合もある。

【0188】

本発明によるXML文書の圧縮構造であるBXML文書をデータベース層、アプリケーション層、クライアント層に適用する方法について説明する。

【0189】

本発明は、アプリケーション層4502に適用できる。一般にWebサービスでは、データベース層・アプリケーション層で大量のデータのやり取りを行い、大量の複雑な計算を行い、その結果として少量のデータをクライアント層に返戻する、というパターンを取る。データベース4515からBXML文書を取得し、アプリケーション4510において本発明の解析処理1303を用いることにより、大量のデータを用いた大量の複雑な計算処理を高速に実現することが出来る。例えば地図検索サービスでは、ユーザが指定する領域の地図をデータベース4515から検索し、アプリケーション4510において「道路は赤色で描画する」などのユーザが指定する描画規則で地図を描画し、地図画像として汎用クライアント4504に返戻する方式が考えられる。

多数の一般ユーザが利用する汎用クライアント4504と異なり、データベース層とアプ 50

リケーション層は同一のベンダでシステム構築を行うことが比較的容易である。そのため標準化されたXML文書ではなく、独自仕様である本発明のBXML文書を用いてもかまわない。またBXML文書はXML文書を圧縮した構造であるため、XML文書のメリットである自由度や拡張性を全て継承し、また可逆圧縮であるため変換により情報が劣化することはない。上述した県庁ゲートウェイと市町村データベースの場合のようにデータベース層4503とアプリケーション層4502が低速回線の場合も、コンパクトなBXML文書を利用することにより通信時間を削減することが出来る。

データベース層4503とアプリケーション層4502が低速回線4513で接続されており、かつデータベース4515やアプリケーション4510が標準のXML文書の入出力しか出来ない場合、ブリッジ4511、4514を利用することにより、途中の低速通信回線4513をコンパクトなBXML文書で通信を行うことにより通信時間を削減することが出来る。ブリッジは、たとえばHTTP(Hyper Text Transfer Protocol)通信におけるPROXYサーバとして実装される。データベース4515が標準のXML文書を出力する時、ブリッジ4514は圧縮処理1301を実行することにより入力文書をBXML文書に圧縮し、低速回線4513において高速に通信を行う。BXML文書は、ブリッジ4511で展開処理1302を実行することにより元のXML文書に展開され、XML文書のみを受け付けるアプリケーション4510に入力される。

【0190】

また本発明は、専用クライアント4506とアプリケーション4510との間の通信に利用できる。クライアントとアプリケーション層の間が電話回線や無線通信など低速回線の場合、サイズのコンパクトな本発明のBXML文書を通信に利用することにより、通信時間の短縮を実現することが出来る。また専用クライアント4506は保守・更新など、大量のデータを利用する機会が多く、本発明の解析処理1303、更新処理1304により高速にデータ解析・更新を実現することが出来る。多数の一般ユーザが利用する汎用クライアント4504と異なり、専用クライアント4506はデータベース層やアプリケーション層と同一のベンダでシステム構築を行うことが比較的容易である。そのため標準化されたXML文書ではなく、独自仕様である本発明のBXML文書を用いることもできる。また現在Webブラウザなどにおいて、JavaやMicrosoft社のActiveXなど、クライアント側で実行するソフトウェアをサーバ側から動的にダウンロードして使用する方式が利用されている。このようなクライアントでは、標準化されたXML文書ではなく、独自仕様である本発明のBXML文書を用いてもよい。

【0191】

また本発明は、クライアント層4501がモバイル端末の場合に利用できる。モバイル端末とアプリケーション層との間の通信は無線通信となり、低速回線となる。また通常のデスクトップ端末に対し、モバイル端末自身のメモリ容量は小さく、処理能力も低い。サイズのコンパクトな本発明のBXML文書を通信に利用することにより、通信時間の短縮を実現することが出来る。またモバイル端末に本発明の解析処理1303を組み込むことにより、小メモリ容量・低処理能力のモバイル端末でも地図表示・利用が可能となる。

【0192】

また本発明を、汎用クライアント4505に適用する方法について説明する。第1の方法は、本発明の装置4507を汎用クライアント4505に組み込み、標準板フェースを用いて接続する方法である。本発明は図18で説明したとおり、解析処理1303として標準のXML文書解析用インタフェースであるSAX、DOMを提供する。そのためSAX、DOMを利用した汎用クライアントでは、SAX、DOMのパース部品を本発明の解析処理1303と入れ替えることにより本発明のBXML文書を解析する事が可能となる。

【0193】

本発明を汎用クライアント4505に適用する第2の方法は、本発明の装置をクライアント側のデコーダ4508として組み込む方法である。デコーダは、たとえばHTTP通信におけるクライアント側のPROXYサーバとして実装される。アプリケーション451

10

20

30

40

50

0 が B X M L 文書を出力する時、デコーダ 4 5 0 8 は展開処理 1 3 0 2 を実行することにより元の X M L 文書に展開され、X M L 文書のみを受け付ける汎用クライアント 4 5 0 6 に入力される。

【 0 1 9 4 】

次に、図 4 5 のデータベース層 4 5 1 5 に対する本発明の実施例について説明する。

【 0 1 9 5 】

データベース層 4 5 1 5 の実装方式は、R D B (R e l a t i o n a l D a t a b a s e) を利用した方式と、図面を利用した方式の二種類に分類される。図 4 6 を用いて、R D B を利用した方式に対し、本発明を適用した実施例について説明する。

データベース層 4 5 1 5 の目的は、X M L 文書 4 6 0 2 のような構造を持つデータを多数、R D B に格納し、検索条件に合致するデータを抽出して返戻することである。 10

X M L 文書のデータ管理方式は木構造であり、R D B のデータ管理方式は表構造である。そのため、木構造の X M L 文書を表構造の R D B に格納する場合何らかの対応を行う必要がある。例えば X M L 文書 4 6 0 2 のタグ B、タグ C、タグ D のように、一回しか出現しないことが保証されているタグは、テーブル 4 6 0 3 の列 B、C、D のように列に展開することが出来る。しかしタグ E のように複数回重複して出現するタグや、さらにタグ F、タグ G のように、重複タグの中に子要素が存在する場合、単一のテーブル 4 6 0 3 では表現できない。そのためテーブル 4 6 0 3 の子供である別テーブル 4 6 0 4 を用意し、重複するタグを別テーブル 4 6 0 4 の行として管理する。テーブル 4 6 0 3 とテーブル 4 6 0 4 の関連付けは、テーブル 4 6 0 3 とテーブル 4 6 0 4 において、列 i d で示される共通のキーを持ち、テーブル 4 6 0 4 ではさらに、i d 2 で示される、重複数だけ存在するキーを設けることで実現できる。 20

【 0 1 9 6 】

しかしこの方式は検索時間がかかる。X M L 文書が単一のテーブル 4 6 0 3 のみで管理される場合、1 回の S Q L 検索要求のみで条件に該当する全ての X M L 文書 4 6 0 2 を作成するのに必要な情報が収集できる。しかし重複タグ E を含む X M L 文書 4 6 0 2 のような例では、初めにテーブル 4 6 0 3 の検索を行い、さらに検索結果の回数分テーブル 4 6 0 4 の検索を行う必要がある。X M L 文書 4 6 0 2 の構造が複雑化し、テーブル 4 6 0 4 にさらに子供である別テーブルが関連するなどの結果、検索時間は加速度的に悪化する。 30

【 0 1 9 7 】

R D B 4 6 0 5 の構成は、R D B 4 6 0 1 を改良した構成であり、本発明を利用して検索時間の高速化を行っている。子テーブル 4 6 0 8 はテーブル 4 6 0 4 と同一である。主テーブル 4 6 0 7 は主テーブル 4 6 0 3 に対し、新たな列として圧縮データ列 4 6 0 9 が付加される構造を取る。圧縮データ列 4 6 0 9 にはテーブル 4 6 0 7、4 6 0 8 で管理される情報をあらかじめ B X M L の圧縮データ 1 0 4 としたものを格納する。情報の二重管理となるため R D B 4 6 0 6 のデータサイズは増加するが、検索を高速化することが出来る。また圧縮データ列 4 6 0 9 に格納するものは本発明によるコンパクトな構造を取っているため、X M L 文書をそのまま格納する場合と比較して全体のデータ量を削減することが出来る。 40

【 0 1 9 8 】

ユーザから X M L 文書 4 6 0 6 の取得要求があった場合、テーブル 4 6 0 7 の圧縮データ列 4 6 0 9 を取得し、検索結果の取得された複数の圧縮データ 4 6 0 9 を、図 4 2 のノードコピー方式で結合し、一つの B X M L 文書として返戻することができる。なおその時に必要となる圧縮スキーマは別管理しておく。圧縮スキーマの数は、テーブル 4 6 0 7 の数と同じ数だけ存在する。 40

【 0 1 9 9 】

次に、図 4 7 を用いて、図面を利用したデータベース層 4 5 0 3 の実装方式に対し、本発明を適用した実施例について説明する。

【 0 2 0 0 】

クライアント層 4 5 0 1 が地理情報を要求するパターンは、地物検索と地図検索の 2 つに 50

分類できる。

地物検索パターンは、例えば「指定した二点を結ぶ最短経路を取得せよ」「指定した道路の、騒音の恐れのある周辺50mの範囲の領域に存在する全ての病院や幼稚園を抽出せよ」というように、ユーザが空間データの解析をデータベース層に依頼する場合に生じる。この検索パターンは検索条件が複雑であり、かつ返戻結果が比較的少ない、という特徴を持つ。また、比較的長い検索時間も許容される。これは図46で説明したRDBベースのデータベース層4515が得意とするパターンである。

一方、地図検索パターンは、「指定した矩形領域に含まれる地図を取得せよ」というように、データベース層に複雑な解析を依頼しない場合に生じる。取得した地図はクライアント層で地図として描画されたり、クライアント層側の機能で解析されたりする。この検索パターンは検索条件が比較的単純であり、かつ返戻結果が膨大となる、という特徴を持つ。またユーザがストレスなく地図の様々な領域をスクロールしたり拡大縮小したりするために、高速な検索時間が要求される。これは図面ベースのデータベース層4515が得意とするパターンである。

【0201】

図面ベースのデータベース層で地図を管理するためには、地図を図面に分割する必要がある。従来方式では、地図を図面に分割する際、図面端の図形を分断していた。例えば図47において、地図4701を図面4704と図面4705に分割する際、図面端の家屋図形4702は家屋4706と家屋4707に分断され、道路4703は道路4708と道路4709に分断されていた。これらの図面をクライアント層に取り込み、例えば家屋数を数える場合や道路ネットワーク網の解析を行う場合、分断図形の再結合という複雑な処理を必要とした。

【0202】

従来方式で図面を分割していた理由の一つとして、図形構成点の座標の容量を低減するため、相対値管理を行いたいという点がある。座標を緯度経度などの絶対座標で管理すると、表現したい精度に応じて例えば64bit整数などの大サイズの固定長整数で管理する必要があり、データ量が増大する。図面内の座標をローカル座標で管理することにより例えば16bit整数などの小サイズの固定長整数で管理することが出来る。この副作用として、例えば図面4710の道路4713のように、図面枠を大きくはずれる座標は小サイズの固定長整数でオーバーフローする可能性があるため、図面枠を超える図形を分断する必要があった。

【0203】

本発明では、図12に示したとおり、座標値を可変長数値で管理することが出来るため、図面枠を大きくはずれる座標のオーバーフローの心配はない。また図7の相対値の適用705で示したとおり、圧縮方式として相対値管理を行っており、図面レベルで明示的に相対値管理を行う必要はない。

そこで本発明を適用した図面ベースデータベース層4515では、図面分割の際、図形を分断しない。地図4701を図面4710と図面4711に分割する際、図面端の家屋図形4702は家屋4712、道路4703は道路4713としていずれかの図面に格納する。また座標は絶対座標をそのまま管理する。各図面は外接矩形4714を属性として持ち、図面内の全ての図形が外接矩形4714に含まれることを保証する。

【0204】

【効果】

本発明の効果は以下の通りである。

(1) 本発明の構造化文書圧縮処理ならびに展開処理を利用することにより、巨大なサイズの構造化文書をその意味的構造を完全に保ったままサイズを約1/10に削減できる。これにより、多量かつ巨大なサイズの構造化文書をデータベースシステムやファイルシステムに格納する場合、格納容量を低減することができる。またモバイル環境などの低容量通信路において巨大なサイズの構造化文書を伝送する場合、通信負荷を軽減することができる。また利用者の許容できる応答時間でより大量な情

10

20

30

40

50

報を伝達することができる。またモバイル機器など、記憶装置や二次記憶装置のサイズが小さいクライアントに対し、クライアントの記憶装置の負荷を軽減できる。

(2) 本発明の構造化文書解析処理を利用することにより、巨大なサイズの構造化文書を従来の非圧縮構造化文書解析処理の約10倍の速度で解析できる。また本発明の構造化文書解析処理は従来のDOM構造のように構造化文書を不要にメモリに展開しないため、巨大なサイズの構造化文書をより少ない記憶容量で解析できる。

(3) 本発明の構造化文書更新処理を利用することにより、巨大なサイズの構造化文書を従来の非圧縮構造化文書更新処理より高速に更新できる。また本発明の更新処理では更新部分に挿入トークン、削除トークンが挿入されるため、更新後に更新位置を特定することができる。これにより、更新の取り消しを実現することができる。また、クライアント・サーバシステムにおいて、サーバから取得した巨大な構造化文書をクライアントにおいて更新した場合、更新部分のみを差分情報としてサーバに転送することにより、クライアント・サーバ間の通信負荷を軽減できる。

(4) 本発明の構造化文書作成処理を利用することにより、巨大なサイズの構造化文書を非圧縮の形を経由せずに直接作成できる。そのため巨大なサイズの構造化文書を従来の非圧縮構造化文書の作成処理より高速かつ記憶装置の負荷をかけずに更新できる。

(5) 昨今の地理情報の標準化により、空間データを構造化文書で表現する方法が提案されてきている。しかし地理情報分野では街の地図など多量の空間データを扱う場合が多く、従来の非圧縮の構造化文書ではサイズが巨大であり解析処理も遅く、実用的に使用することはできなかった。そのため構造化文書による標準方式は、処理時間を問わないオフラインでのデータベース間のインポートに利用されるのみで、ネットワークを介したクライアント・サーバ間では利用されなかった。本発明による構造化文書圧縮処理・解析処理を利用することにより、利用者が許容できる時間で構造化文書の伝送・解析が可能となる。従来の地理情報システムでは、クライアント単体の機能としては、地図の表示と交差判定や包含判定などの空間演算しか行えなかった。本発明により、クライアントで構造化文書が利用可能となることから、「道路周辺の、騒音が到達する20m範囲での、防音設備のない病院や幼稚園を検索せよ」などの属性と空間データを交えた検索や、シミュレーション等さまざまなアプリケーションが開発可能となり、地理情報分野の応用範囲が大きく広がる。

【図面の簡単な説明】

【図1】本発明を構成するシステムブロック図である。

【図2】SAXパーサのアルゴリズムである。

【図3】DOMパースブロックの基本インタフェースである。

【図4】WAP Binary XMLの圧縮方式である。

【図5】特開2002-163248の圧縮方式である。

【図6】地物集合を表すXML文書の例である。

【図7】BXMLの圧縮方式である。

【図8】XML構造のパターンと圧縮方式である。

【図9】圧縮データのサンプルである。

【図10】「文字列の辞書化」の概念的な構造である。

【図11】辞書化のパターンである。

【図12】可変長整数の構造である。

【図13】本発明の機能ブロック図である。

【図14】スキーマ木の構造である。

【図15】スキーマ木のサンプルである。

【図16】圧縮処理1301を実現するブロック抜粋図である。

【図17】展開処理1302を実現するブロック抜粋図である。

【図18】解析処理1303を実現するブロック抜粋図である。

【図19】更新処理1304を実現するブロック抜粋図である。

【図20】生成処理1305を実現するブロック抜粋図である。

10

20

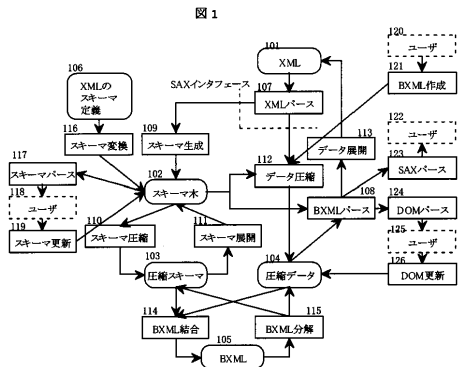
30

40

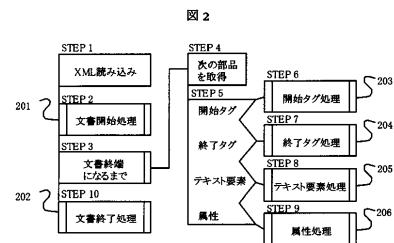
50

- 【図 2 1】スキーマ生成ブロックのアルゴリズム（タグとデータの分離）である。
- 【図 2 2】スキーマ生成ブロックのアルゴリズム（要素型、文字列辞書、相対値）である。
- 【図 2 3】スキーマ型の上書規則である。
- 【図 2 4】XML Schema の構造定義例と、その変換規則である。
- 【図 2 5】スキーマ木移動を説明する図である。
- 【図 2 6】スキーマ木移動のアルゴリズムである。
- 【図 2 7】データ圧縮ブロックのアルゴリズム（単独型）である。
- 【図 2 8】データ圧縮ブロックのアルゴリズム（重複型）である。
- 【図 2 9】データ圧縮ブロックのアルゴリズム（省略型）である。 10
- 【図 3 0】圧縮スキーマの構造である。
- 【図 3 1】XML パースブロック 1 0 8 のアルゴリズムである。
- 【図 3 2】DOM パースブロックにおけるノード情報である。
- 【図 3 3】DOM パースブロック 1 2 4 における「属性値取得」アルゴリズムである。
- 【図 3 4】DOM パースブロック 1 2 4 における「テキスト要素取得」アルゴリズムである。
- 【図 3 5】DOM パースブロック 1 2 4 における「子に移動」アルゴリズムである。
- 【図 3 6】DOM パースブロック 1 2 4 における「弟に移動」アルゴリズムである。
- 【図 3 7】スキーマ要素識別子を用いた高速タグ判定方式である。
- 【図 3 8】地理情報分野に特化したノード例である。 20
- 【図 3 9】更新された圧縮データのデータ構造である。
- 【図 4 0】更新された XML 文書に対する XML パースブロック 1 0 8 のアルゴリズムである。
- 【図 4 1】XML 作成ブロックが提供するインタフェースである。
- 【図 4 2】XML 部分ノードのコピー方式である。
- 【図 4 3】スキーマの隠蔽機能を説明する図である。
- 【図 4 4】隠蔽された XML 文書に対する XML パースブロック 1 0 8 のアルゴリズムである。
- 【図 4 5】三階層モデルと XML 圧縮を説明する図である。
- 【図 4 6】RDB をベースとしたデータベース層 4 5 1 5 の実装例である。 30
- 【図 4 7】図面をベースとしたデータベース層 4 5 1 5 の実装例である。

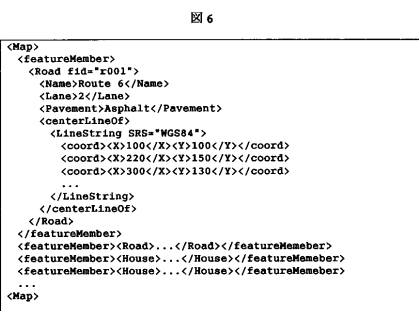
【 図 1 】



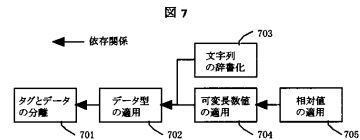
【 図 2 】



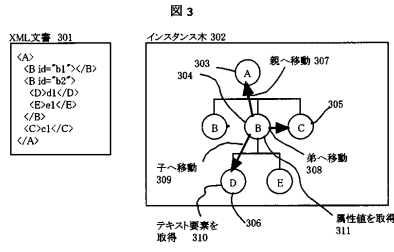
【 図 6 】



【 図 7 】



【 図 3 】



【 図 4 】

XML 文書	圧縮 XML 文書
<pre> <RoadGroup> <Road> <Name>国道 6 号線</Name> <Type>国道</Type> </Road> </RoadGroup> </pre>	<pre> (タグ略す) <RoadGroup> ::= 1 <Road> ::= 2 <Name> ::= 3 <Type> ::= 4 (圧縮 XML) <1> <2> <3>国道 6 号線</> <4>国道</> </> <2> <3>国道 20 号線</> <4>国道</> </> </pre>

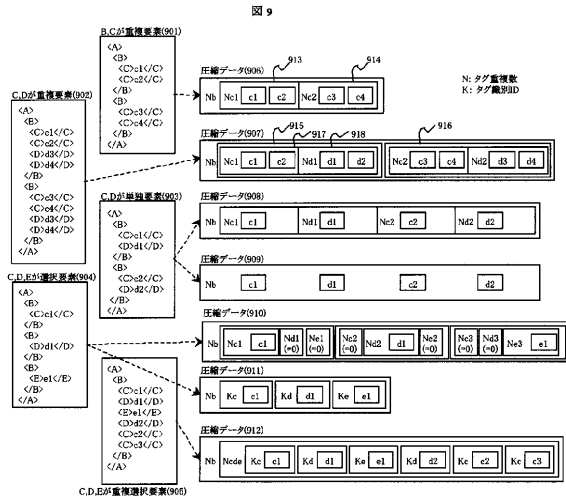
【 図 5 】

XML 文書	圧縮 XML 文書
<pre> <RoadGroup> <Road> <Name>国道 6 号線</Name> <Type>国道</Type> </Road> </RoadGroup> </pre>	<pre> (圧縮スキーマ) <RoadGroup> <Road> <Name></Name> <Type></Type> </Road> </RoadGroup> (圧縮データ) "..国道 6 号線".. "国道".. "国道 20 号線".. "国道".. </pre>

【 図 8 】

#	パターン	XML サンプル	BXML 文書 803						
			タグ名 805	重複フラグ 806	圧縮スキーマ 804	子要素保存数 808	属性保存数 809	圧縮データ 810	
1	子要素 811	<pre> <A> <C>text1</C> </pre>	A: 単独, B: 単独, C: 単独		SEQUENCE	1	0	0	text1
2	兄弟要素 812	<pre> <A> text1 <C>text2</C> </pre>	A: 単独, B: 単独, C: 単独		SEQUENCE	2	0	0	text1, text2
3	重複要素 813	<pre> <A> text1 text2 <C>text3</C> </pre>	A: 単独, B: 重複, C: 単独		SEQUENCE	2	0	0	B の重複数=2 text1 text2 text3
4	省略要素 814	<pre> <A> text1 <A> text2 <C>text3</C> </pre>	A: 単独, B: 重複, C: 重複		SEQUENCE	1	0	0	A の重複数=2 text1 C の重複数=0 text2 C の重複数=1 text3
5	選択要素 815	<pre> <A> text1 <A> <C>text2</C> </pre>	A: 単独, B: 重複, C: 単独		SEQUENCE	1	0	0	A の重複数=2 B が出現 text1 C が出現 text2
6	重複選択要素 816	<pre> <A> text1 <C>text2</C> text3 <C>text4</C> </pre>	A: 単独, B: 重複, C: 単独		CHOICE	2	0	0	B または C の重複数=4 B が出現 text1 C が出現 text2 B が出現 text3 C が出現 text4
7	属性 817	<pre> text2 </pre>	A: 単独, B: 属性, C: 単独		SEQUENCE	1	1	1	a の重複数=1 text1 text2

【 図 9 】

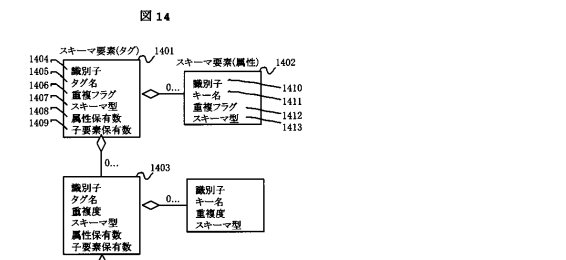


【 図 10 】

図 10

XML 文書 1001	BXML 文書 1002	圧縮データ 1004
<pre> <RoadGroup> <Road<Type>国道</Type></Road> <Road<Type>県道</Type></Road> <Road<Type>県道</Type></Road> <Road<Type>県道</Type></Road> </RoadGroup> </pre>	<pre> 圧縮スキーマ 1003 <Type>に含まれる文字列 ::= 1: 国道, 2: 県道 <RoadGroup> <Road<Type>1</Type></Road> <Road<Type>2</Type></Road> <Road<Type>2</Type></Road> <Road<Type>2</Type></Road> </RoadGroup> </pre>	<pre> <RoadGroup> <Road<Type>1</Type></Road> <Road<Type>2</Type></Road> <Road<Type>2</Type></Road> <Road<Type>2</Type></Road> </RoadGroup> </pre>

【 図 14 】



【 図 15 】

図 15

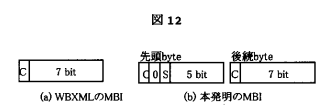
XML 文書	インスタンス木	スキーマ木
<pre> 1501 <R> <A>text1 <A>text2 <A>text3 text4 text5 </R> </pre>	<pre> R ├── A │ ├── A │ ├── A │ └── B └── B ├── B └── B </pre>	<pre> R ├── A └── B </pre>
<pre> 1502 <R> <A>text1 <A> <C>text3</C> <C>text3</C> <A> <D>text4</D> <D>text5</D> </R> </pre>	<pre> R ├── A │ ├── A │ ├── A │ └── A ├── B │ ├── B │ ├── C │ └── D └── D ├── D └── D </pre>	<pre> R ├── A ├── B ├── C └── D </pre>
<pre> 1503 <R> <A> <C>text1</C> <C>text2</C> <A> <C>text3</C> <C>text4</C> </R> </pre>	<pre> R ├── A │ ├── A │ └── B ├── C │ ├── C │ ├── C │ └── C └── C ├── C └── C </pre>	<pre> R ├── A ├── B ├── C └── C </pre>

【 図 11 】

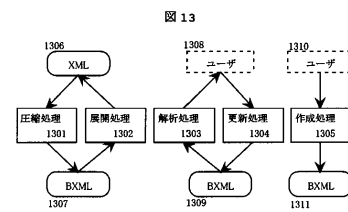
図 11

#	辞書化パターン	入力例	辞書	圧縮例
1	インライン文字列	* 国道*	なし	* 国道*
2	辞書 ID	* 国道*	1=* 国道*	T_ID. 1
3	辞書 ID + インライン文字列	* 国道 20 号線*	1=* 国道*	C_ID. 1. T_WORD. *20 号線*
4	辞書 ID + 数値	http://a.com/f2256 *	1=* http://a.com/f *	C_ID. 1. T_NUM. 2256
5	インライン文字列 + 辞書 ID	*user@ri.hitachi.co.jp*	1=*@ri.hitachi.co.jp*	C_WORD. *user*. T_ID. 1

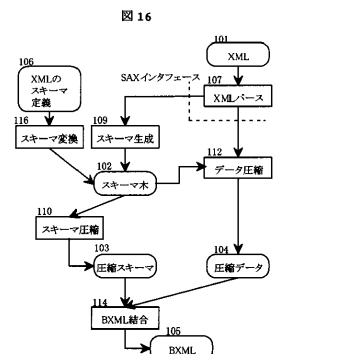
【 図 12 】



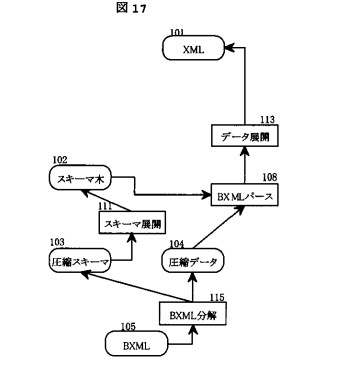
【 図 13 】



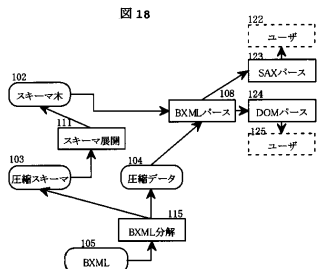
【 図 16 】



【 図 17 】

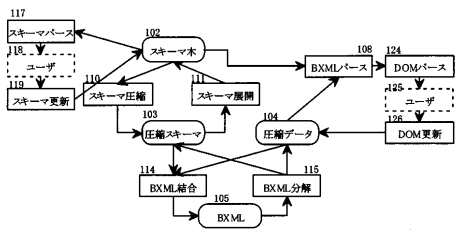


【 図 18 】

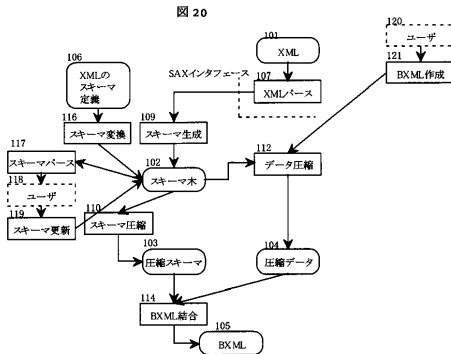


【 図 19 】

図 19

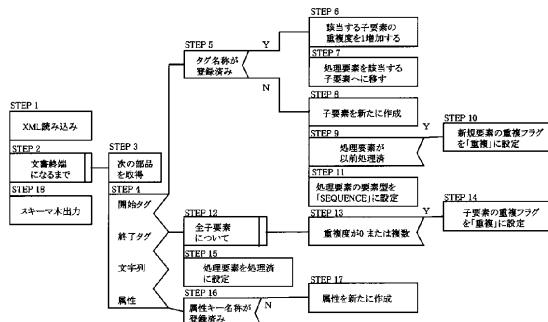


【 図 20 】



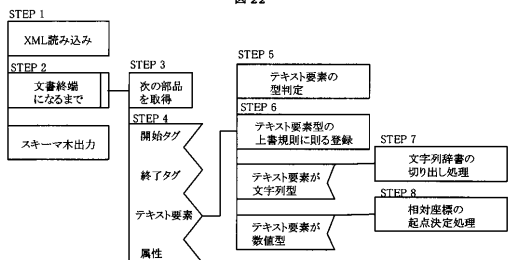
【 図 21 】

図 21



【 図 22 】

図 22



【 図 23 】

図 23

元	未定義	SEQUENCE CHOICE MCHOICE	文字列	実数	整数	2次元実数	2次元整数
上書き							
未定義	OK	NG	NG	NG	NG	NG	NG
SEQUENCE CHOICE MCHOICE	OK	OK	OK	OK	OK	OK	OK
文字列	OK	NG	OK	OK	OK	OK	OK
実数	OK	NG	NG	OK	OK	文字列	文字列
整数	OK	NG	NG	OK	OK	文字列	文字列
2次元実数	OK	NG	NG	文字列	文字列	OK	NG
2次元整数	OK	NG	NG	文字列	文字列	OK	OK

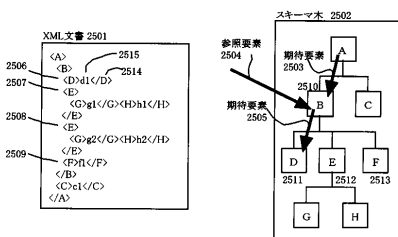
【 図 24 】

図 24

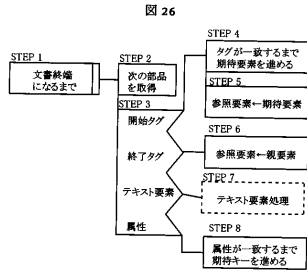
#	XML Schema	サンプル 2402	圧縮スキーマ	要素型 2404
1	BCDE	BCDE 固定	BCDE	SEQUENCE
2	B C D E	BD or CE or CD or CB	B C D E	CHOICE
3	*(BCDE)	BCDEBCDE	*(B C D E)	MCHOICE
4	*(B C D E)	BCDE 任意組合せ		
5	*(B C D E)	BCDB		
6	(BC) (DE)	BCDE	BCDE	SEQUENCE
7	(BC) (DE)	BC or DE	*(B C D E)	MCHOICE
8	(B C) (D E)	BD or BE or CD or CE		
9	(B C) (D E)	B or C or D or E		
10	*(BC)*(DE)	BCBCDEDE		
11	*(*(BC)) (*(DE))	BCBC or DEDE		
12	*(B C)*(D E)	BBBDDD or BBBBEE or CCCCDD or CCCCEE		
13	*(*(B C)) (*(D E))	BC 任意組合せ or DE 任意組合せ		
14	(BC) (D E)	BCD BCE		
15	(BC) (D E)	BC D E		
16	*(*(BC)) (*(D E))	BCBC D E		
17	(BC) (D E)	BC DEDEDE		
18	*(*(BC)) (*(D E))	BCBC DEDE		
19	*(*(BC))*(DE)	BCDEBCDE		
20	*(*(BC))*(DE)	BCBCDEBCBCDEDE		
21	*(BC) (D E)	BCBCD BCBC		
22	BC*(D E)	BCCCD BCCE		
23	*(BC)*(D E)	BCBCDDD BCBCBEE		
...				

【 図 25 】

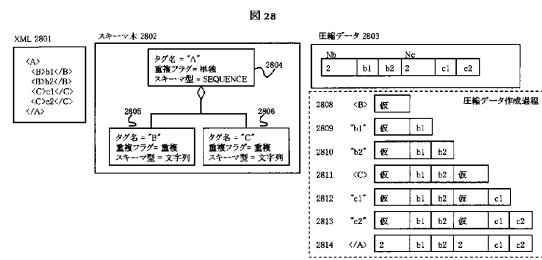
図 25



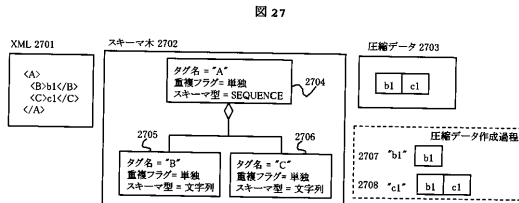
【 図 2 6 】



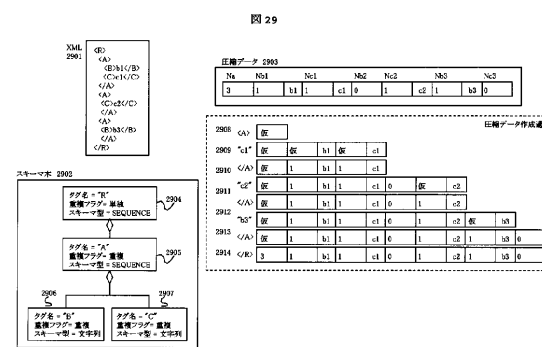
【 図 2 8 】



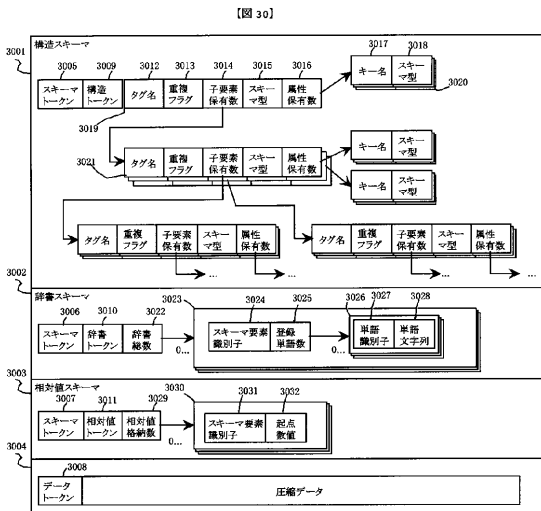
【 図 2 7 】



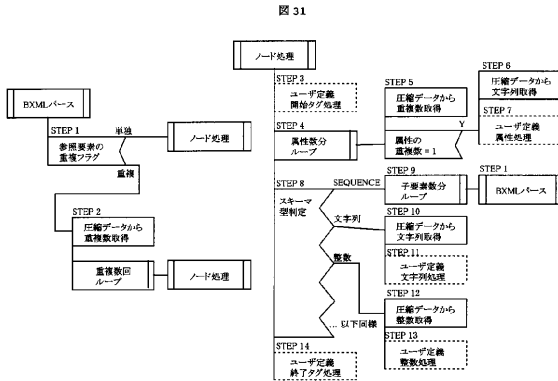
【 図 2 9 】



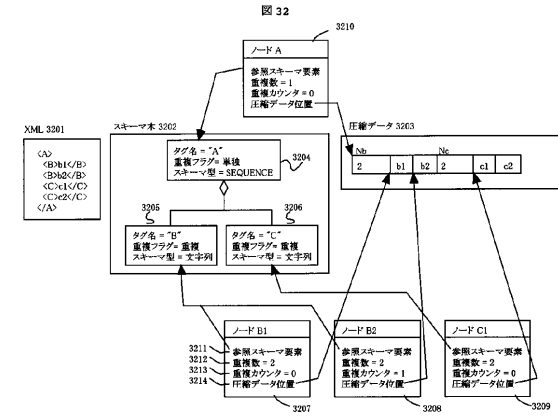
【 図 3 0 】



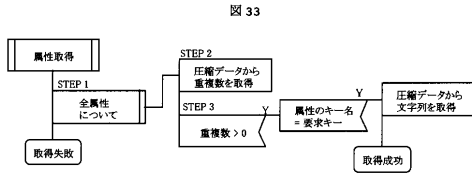
【 図 3 1 】



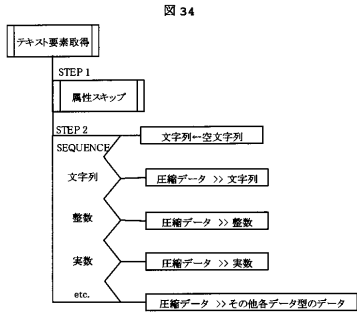
【 図 3 2 】



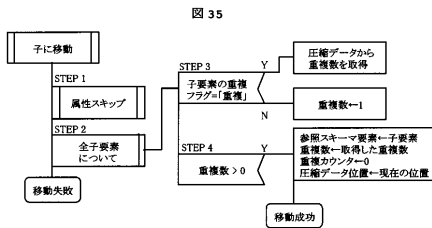
【 図 3 3 】



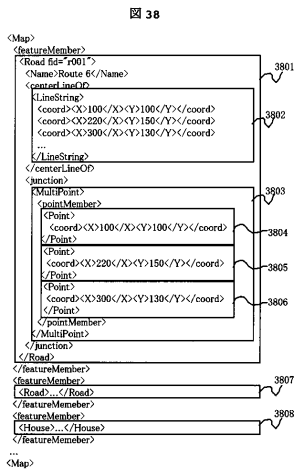
【 図 3 4 】



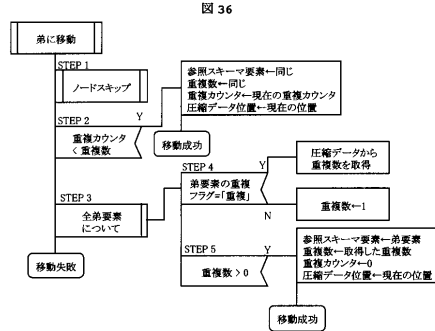
【 図 3 5 】



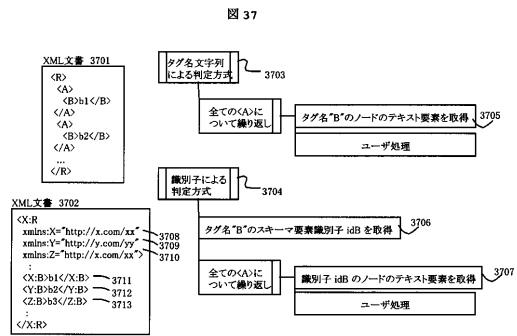
【 図 3 8 】



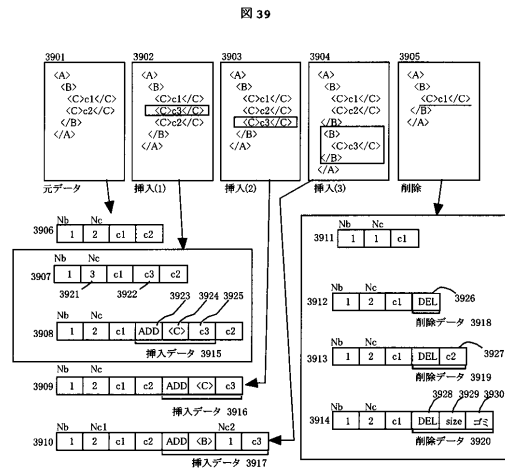
【 図 3 6 】



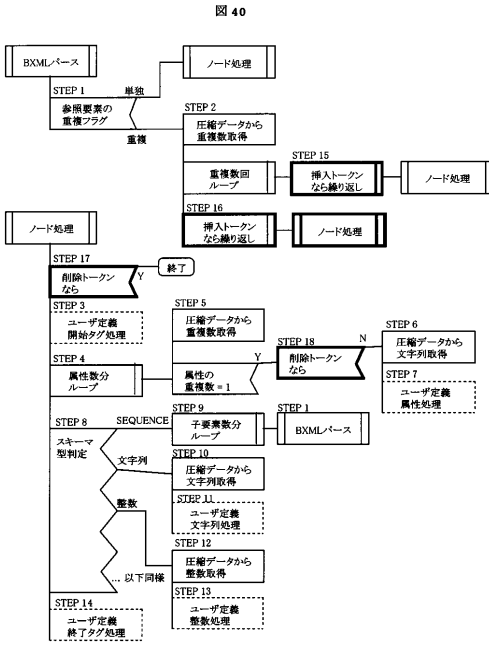
【 図 3 7 】



【 図 3 9 】



【 図 4 0 】

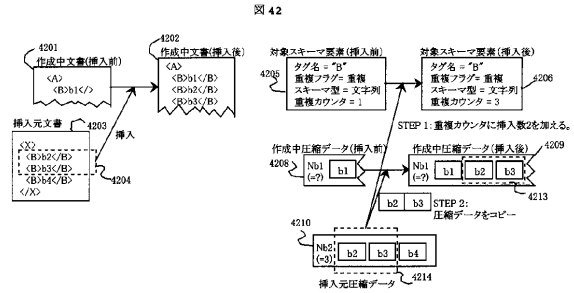


【 図 4 1 】

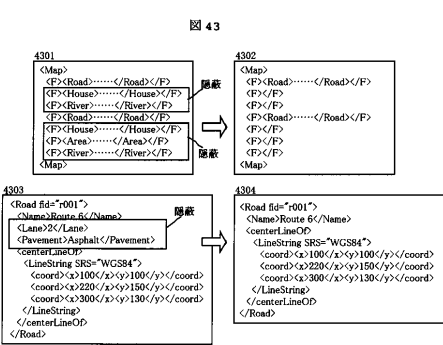
図 41

#	分野	メソッド	書式	開始状態	終了状態
1	文書作成	文書開始 (4108)	startDocument();	開 (4101)	初期状態
2	文書作成	文書終了 (4109)	endDocument();	開 (4102)	閉 (4103)
3	ノード作成	子要素作成 (4110)	moveToChild("C");	開 (4104)	閉 (4105)
4	ノード作成	弟要素作成 (4111)	moveToSibling("C");	開 (4106)	閉 (4107)
5	ノード作成	親へ移動 (4112)	moveToParent();	開 (4108)	閉 (4109)
6	ノード作成	先祖へ移動 (4113)	moveToAncestor("B");	開 (4110)	閉 (4111)
7	データ作成	属性の作成 (4114)	setAttribute("id", "1");	開 (4112)	閉 (4113)
8	データ作成	文字列の作成 (4115)	setText("text");	開 (4114)	閉 (4115)
9	データ作成	整数の作成 (4116)	setInteger(8);	開 (4116)	閉 (4117)
10	データ作成	実数の作成 (4117)	setDouble(3.14);	開 (4117)	閉 (4118)
11	データ作成	空タダの作成 (4118)	setNull();	開 (4118)	閉 (4119)
12	データ作成	ノードのコピー (4119)	setNode(node);	開 (4119)	閉 (4120)

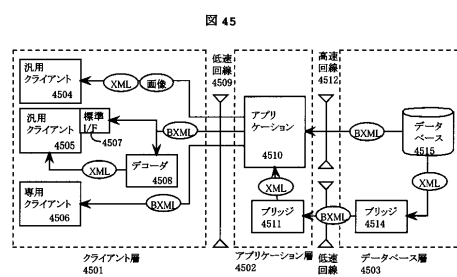
【 図 4 2 】



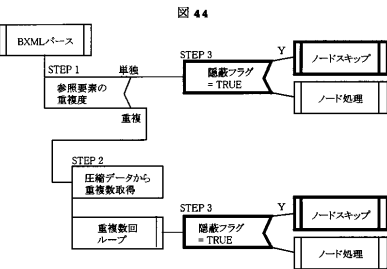
【 図 4 3 】



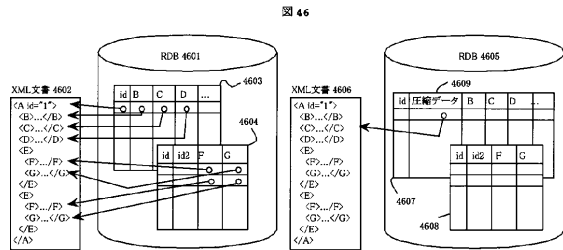
【 図 4 5 】



【 図 4 4 】



【 図 4 6 】



【 図 4 7 】

