

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4261299号
(P4261299)

(45) 発行日 平成21年4月30日(2009.4.30)

(24) 登録日 平成21年2月20日(2009.2.20)

(51) Int.Cl. F I
H03M 7/42 (2006.01) H03M 7/42

請求項の数 8 (全 25 頁)

(21) 出願番号	特願2003-328428 (P2003-328428)	(73) 特許権者	392026693
(22) 出願日	平成15年9月19日(2003.9.19)		株式会社エヌ・ティ・ティ・ドコモ
(65) 公開番号	特開2005-94652 (P2005-94652A)		東京都千代田区永田町二丁目11番1号
(43) 公開日	平成17年4月7日(2005.4.7)	(74) 代理人	100088155
審査請求日	平成18年4月11日(2006.4.11)		弁理士 長谷川 芳樹
		(74) 代理人	100092657
			弁理士 寺崎 史朗
		(74) 代理人	100114270
			弁理士 黒川 朋也
		(74) 代理人	100122507
			弁理士 柏岡 潤二
		(74) 代理人	100123995
			弁理士 野田 雅一

最終頁に続く

(54) 【発明の名称】 データ圧縮装置、データ復元装置およびデータ管理装置

(57) 【特許請求の範囲】

【請求項1】

型と値をそれぞれ有する複数の頂点と、該頂点間の参照情報とを有する入力データを、前記頂点間の参照情報を有する相互参照関係データと、前記型と値を有する複数の頂点からなる頂点群とに分離し、その分離された前記頂点群のデータを出力する分離手段と、

特定のパターンを有する前記頂点間の参照情報を、前記頂点間の参照情報が共有可能なテンプレートとして蓄積するテンプレート蓄積手段と、

前記分離手段により分離された前記相互参照関係データから、前記テンプレート蓄積手段に蓄積されているテンプレートと一致する箇所を検出するテンプレート一致箇所検出手段と、

前記分離手段により分離された相互参照関係データのうち、前記テンプレート一致箇所検出手段により検出された一致箇所を前記テンプレートで前記頂点間の参照情報を参照可能な状態に置換し、その置換された相互参照関係データを出力するテンプレート置換手段とを有することを特徴とするデータ圧縮装置。

【請求項2】

それぞれの値を有し、該各値が属性情報として型を有することが可能な複数の頂点と、該頂点間の参照情報とを有する入力データを、前記頂点間の参照情報を有する相互参照関係データと、前記値を有する複数の頂点からなる頂点群とに分離し、その分離された前記頂点群のデータを出力する分離手段と、

特定のパターンを有する前記頂点間の参照情報を、前記頂点間の参照情報が共有可能な

テンプレートとして蓄積するテンプレート蓄積手段と、

前記分離手段により分離された前記相互参照関係データから、前記テンプレート蓄積手段に蓄積されているテンプレートと一致する箇所を検出するテンプレート一致箇所検出手段と、

前記分離手段により分離された相互参照関係データのうち、前記テンプレート一致箇所検出手段により検出された一致箇所を前記テンプレートで前記頂点間の参照情報を参照可能な状態に置換し、その置換された相互参照関係データを出力するテンプレート置換手段とを有することを特徴とするデータ圧縮装置。

【請求項 3】

前記テンプレートが、前記頂点間の参照情報の一部または全部の参照方向を反転可能なことを特徴とする請求項 1 または 2 記載のデータ圧縮装置。

10

【請求項 4】

第 1 から第 N までの N 個の頂点を有し、前記第 1 の頂点と第 2 の頂点以外の連続番号を有する前記頂点は相互に参照し、前記第 1 の頂点が前記第 2 の頂点を参照し、かつ外部への参照を保持し、前記第 N の頂点が第 N - 1 の頂点を参照し、前記第 2 から第 N までの各頂点が、外部への参照を保持しないか、またはすべて同数の参照を保持する連続兄弟参照部を有する前記相互参照関係データに適用するための接続情報を有する連続兄弟参照用テンプレートが、前記テンプレート蓄積手段に蓄積されていることを特徴とする請求項 1 ~ 3 のいずれか一項記載のデータ圧縮装置。

【請求項 5】

20

前記頂点間の参照情報に、前記テンプレートを適用可能な親テンプレートが前記テンプレート蓄積手段に蓄積されていることを特徴とする請求項 1 ~ 4 のいずれか一項記載のデータ圧縮装置。

【請求項 6】

複数の前記入力データに共用可能な共用テンプレートが前記テンプレート蓄積手段に蓄積されていることを特徴とする請求項 1 ~ 5 のいずれか一項記載のデータ圧縮装置。

【請求項 7】

特定のパターンを有する複数の頂点間の参照情報を、前記頂点間の参照情報が共有可能なテンプレートとして蓄積するテンプレート蓄積手段と、

前記テンプレートにより置換され、圧縮された相互参照関係データを第 1 の入力データとして入力し、前記相互参照関係データから、前記テンプレートを用いて圧縮前の元の相互参照関係データを復元する展開手段と、

30

型と値をそれぞれ有する複数の前記頂点からなる頂点群のデータを第 2 の入力データとして入力し、前記頂点群のデータを前記展開手段により復元された前記相互参照関係データと合成したデータを出力する合成手段とを有することを特徴とするデータ復元装置。

【請求項 8】

圧縮可能なデータを蓄積する第 1 のデータ蓄積手段と、

請求項 1 ~ 6 のいずれか一項記載のデータ圧縮装置により圧縮されたデータを蓄積する第 2 のデータ蓄積手段と、

前記第 1 のデータ蓄積手段及び第 2 のデータ蓄積手段に蓄積されたそれぞれのデータの利用頻度を観測し、該観測された利用頻度に応じて移動要求を出力する利用頻度観測手段と、

40

該利用頻度観測手段からの移動要求にしたがい、前記利用頻度が高い高頻度データを請求項 7 記載のデータ復元装置により復元して前記第 1 のデータ蓄積手段へ格納し、前記利用頻度が前記高頻度データよりも低いデータを前記データ圧縮装置により圧縮して前記第 2 のデータ蓄積手段へ格納するようにして、データを移動させるための制御を行う制御手段と、

前記第 1 のデータ蓄積手段と第 2 のデータ蓄積手段のいずれかから、前記制御手段の指示に応じてデータを取得して出力する選択手段とを有することを特徴とするデータ管理装置。

50

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、電子データの圧縮装置、復元装置およびデータ管理装置に関する。

【背景技術】

【0002】

近年、WWW(World Wide Web)の普及により、HTML(Hyper Text Markup Language)やXML(EXTensible Markup Language)等、構造化文書を用いたデータ交換が増加している。特に、XMLはHTMLを補う次世代の言語として注目を集めており、今後インターネットにおける情報交換の場において最も普及していくことが予想される。

10

【0003】

XMLは要素の階層構造を表すデータ表現形式を伴った言語であり、XMLを用いた文書(XML文書)は、例えば図18のように記述される。図18は、XML文書10を示す図である。図18に示すとおり、XMLは大きくマークアップとテキスト情報に分けられる。図18に示すXML文書10では、マークアップは、要素開始記号(開始タグ)Ma、要素終了記号(終了タグ)Mb、空要素記号(空要素タグ)Mcからなっている。図18では、<book>、<title>、<authors>、<author>、<contents>および<chapter>が要素開始記号Maを表している。また、</book>、</title>、</authors>、</author>、</contents>および</chapter>が要素終了記号Mbを表し、<misc/>が空要素記号Mcを表している。これらの要素開始記号Maから対応する要素終了記号Mbまでの領域、または空要素記号Mcが要素(XMLの基本となる情報単位)を表している。

20

【0004】

要素開始記号Maと要素終了記号Mbの間には、別な要素記号の他、テキスト情報を記述することができる。例えば、図18に示すXML文書10では、要素<title>には、文字列“XMLの基礎”が、要素<authors>の中に現れる最初の要素<author>には、文字列“山田 太郎”がそれぞれテキスト情報として定義されている。

【0005】

要素やテキスト情報の間には、親子関係、兄弟関係が定義されている。図18に示すXML文書10の場合、要素開始記号Maが<book>で始まり、要素終了記号Mbが</book>で終了する要素(要素<book>)の中に、要素開始記号Maが<title>で始まり、要素終了記号Mbが</title>で終了する要素(要素<title>)が含まれている。このとき、要素<book>は要素<title>の親要素であるといい、要素<title>は要素<book>の子要素であるという。これが要素の親子関係である。

30

【0006】

また、要素<title>と要素<authors>とは、同一の親要素<book>を持ち、かつ連続している。このとき、要素<title>と要素<authors>とは兄弟であるといい、要素<title>は要素<authors>の前兄弟、要素<authors>は要素<title>の次兄弟であるという。これが要素の兄弟関係である。

【0007】

一般に、XMLは、コンピュータ間で通信を行う際や、ハードディスク装置やフラッシュメモリに蓄積する際には、図18に示すXML文書10のようにテキスト形式で表現されている。一方、コンピュータ内部で検索や修正用に利用するときは、解析されてコンピュータ内部に適したデータ構造に変換されている。

40

【0008】

図19は、図18に示すXML文書10を解析し、コンピュータの内部利用に適した形式に変換したデータ構造11を示す図である。図19では、各要素及びテキスト情報が型および値を有する頂点301~317として記述されている。型は各頂点301~317の左側に記述され、“E”であれば要素を表し、“T”であればテキスト情報を表している。例えば、頂点301では型301aは“E”である。また、値は頂点の右側に記述され、例えば、頂点301では値301bは“book”である。そして、頂点が要素を表

50

す場合は値に要素の名称（要素名）が記述され、テキスト情報を表す場合は文字列が記述される。例えば、頂点302では要素名<title>を表し、頂点306ではテキスト情報“XMLの基礎”を表している。

【0009】

また、各頂点301～317は、もとの（変換前の）XML文書10の親子関係および兄弟関係を表現するため、親参照、子参照、次兄弟参照および前兄弟参照の4つの参照を表す参照情報を有している。上述のXML文書10の場合、要素<title>は要素<book>の子要素であり、要素<book>は要素<title>の親要素であるから、図19に示すデータ構造11では、例えば、頂点301, 302については、それぞれ、<book>から<title>への子参照P1と、<title>から<book>への親参照P2を有し、それらが矢印によって表現されている。また、要素<book>は<title>の次の子要素として要素<authors>も有している。この場合、頂点302, 303については、要素<title>から要素<authors>への次兄弟参照P3、要素<authors>から要素<title>への前兄弟参照P4が保持されている。なお、兄弟関係にある要素では、先頭の子要素（例えば要素<title>）以外は親参照を直接に有しないものとされている。

【0010】

データ構造は、各頂点間の参照情報と、要素名やテキスト情報とを分離して管理することができ、例えば、そのそれぞれを図20(a)、図20(b)のように表現することができる。ここで、図20(a)は各頂点間の参照情報を有する相互参照関係データ400を示す図であり、図20(b)は要素とテキスト情報のいずれかに設定される型と値を有する複数の頂点の集合（頂点群ともいう）を示すテーブル450を示す図である。

【0011】

しかしながら、メモリ等の記憶装置の容量は有限であるため、データ構造を蓄積するときは、そのデータ構造を効率的に圧縮して蓄積することが求められる。この点に関し、非特許文献1には、図20(b)に示すような要素名やテキスト情報を圧縮する方法が開示されている。非特許文献1では、各頂点が保持する要素名やテキスト情報を辞書として別途蓄積し、各頂点には辞書のインデックスを持たせ、同じ文字列を複数蓄積しないようにすることで、圧縮する方法が開示されている。

【0012】

一方、非特許文献2には、XML文書中の部分的な構造を再利用することで、XML文書を圧縮する方法が開示されている。この方法は元のXML文書を構造、要素名情報、テキスト情報の3つに分離したのち、そのそれぞれをLZ77等の一般的な圧縮アルゴリズムで圧縮するというものである（LZ77について詳しくは、Jacob Ziv, Abraham Lempel: A Universal Algorithm for Sequential Data Compression. IEEE Transactions on Information Theory 23(3):337-343(1977)を参照）。

ここで、非特許文献2に開示されている圧縮方法について説明する。この圧縮方法ではまず、要素開始記号や空要素記号をそれぞれ「#1」、「#2」のような短い要素名で置換し、要素終了記号を「/」で置換する。また、テキスト情報は「C」で置換する。

以上の圧縮方法を分離したXML文書10に適用すると、分離後のデータ構造12、要素名情報13およびテキスト情報14はそれぞれ図21、図22、図23のように表現される。

【0013】

また、非特許文献2に記載の圧縮方法では、LZ77等に代表される圧縮アルゴリズムを用いてそれぞれを独立に圧縮するが、ここではその圧縮アルゴリズムの概要について説明する。LZ77等の圧縮アルゴリズムは元の入力情報に含まれる部分的なパターンを発見し、それをテンプレートとして繰り返し再利用することにより、圧縮を行う。例えば、図21に示すデータ構造12の圧縮について説明すると、テンプレートとして、テンプレートX, Y, Z, W, Vを用いるとし、それぞれのテンプレートの割り当てを、

X = “ # 1 # 2 C / # 3 ”, Y = “ # 4 C / ”,

Z = “ / # 5 ”, W = “ # 6 C / ”, V = “ / # 7 / / ” のように設定すると、

10

20

30

40

50

図 2 1 に示したデータ構造 1 2 は “ X Y Y Y Z W W V ” のように表せる。これは一部の文書構造をあらわすテンプレートとして、Y , W を複数回利用している。このように、テンプレートが繰り返し利用でき、元の文書を少ないテンプレートで表現することができれば、元の XML 文書を表す情報量が少なく済むから圧縮が可能になる。

【 0 0 1 4 】

【非特許文献 1】Mathias Neumuller and John N. Wilson: “ Compact In-Memory Representation of XML ” Internal Report of University of strathclyde

【非特許文献 2】Hartmut Liefke and Dan Suciu.: “ XMill: An Efficient Compressor for XML Data ” , In proceedings of ACM SIGMOD International Conference on Management of Data, 2000

10

【発明の開示】

【発明が解決しようとする課題】

【 0 0 1 5 】

しかし、上述した従来の技術では、テンプレートを用いて XML 文書のデータ構造を圧縮する際、テンプレート相互の情報が独立していたため、テンプレート数が増えると、その分だけ必要とされる記憶量が増えてしまい、メモリ等の記憶装置を圧迫するという問題があった。

【 0 0 1 6 】

そこで、本発明は上記課題を解決するためになされたもので、テンプレート数が増えても各テンプレートの格納に要するメモリ等の記憶装置を圧迫しないようにすることができる構成を備えたデータ圧縮装置、データ復元装置およびデータ管理装置を提供することを目的とする。

20

【課題を解決するための手段】

【 0 0 1 7 】

上記課題を解決するため、本発明は型と値をそれぞれ有する複数の頂点と、その頂点間の参照情報とを有する入力データを、頂点間の参照情報を有する相互参照関係データと、型と値を有する複数の頂点からなる頂点群とに分離し、その分離された頂点群のデータを出力する分離手段と、特定のパターンを有する頂点間の参照情報を、頂点間の参照情報が共有可能なテンプレートとして蓄積するテンプレート蓄積手段と、分離手段により分離された相互参照関係データから、テンプレート蓄積手段に蓄積されているテンプレートと一致する箇所を検出するテンプレート一致箇所検出手段と、分離手段により分離された相互参照関係データのうち、テンプレート一致箇所検出手段により検出された一致箇所をテンプレートで頂点間の参照情報を参照可能な状態に置換し、その置換された相互参照関係データを出力するテンプレート置換手段とを有するデータ圧縮装置を特徴とする。

30

このデータ圧縮装置は、テンプレート蓄積手段に蓄積されているテンプレートが頂点間の参照情報が共有可能な構成を有するので、テンプレート数を記憶するための記憶容量が少なく済むようになる。

【 0 0 1 8 】

また、本発明は、それぞれの値を有し、その各値が属性情報として型を有することが可能な複数の頂点と、その頂点間の参照情報とを有する入力データを、頂点間の参照情報を有する相互参照関係データと、値を有する複数の頂点からなる頂点群とに分離し、その分離された頂点群のデータを出力する分離手段と、特定のパターンを有する頂点間の参照情報を、頂点間の参照情報が共有可能なテンプレートとして蓄積するテンプレート蓄積手段と、分離手段により分離された相互参照関係データから、テンプレート蓄積手段に蓄積されているテンプレートと一致する箇所を検出するテンプレート一致箇所検出手段と、分離手段により分離された相互参照関係データのうち、テンプレート一致箇所検出手段により検出された一致箇所をテンプレートで頂点間の参照情報を参照可能な状態に置換し、その置換された相互参照関係データを出力するテンプレート置換手段とを有するデータ圧縮装置を提供する。

40

このデータ圧縮装置も、テンプレート蓄積手段に蓄積されているテンプレートが頂点間

50

の参照情報が共有可能な構成を有するので、テンプレート数を記憶するための記憶容量が少なく済むようになる。

【0019】

上記いずれのデータ圧縮装置も、テンプレートが、頂点間の参照情報の一部または全部の参照方向を反転可能なことが好ましい。

このようなテンプレートを有すると、1つのテンプレートを反転させて適用することができるから、テンプレート数を少なくすることができる。

【0020】

また、第1から第NまでのN個の頂点を有し、第1の頂点と第2の頂点以外の連続番号を有する頂点は相互に参照し、第1の頂点が第2の頂点を参照し、かつ外部への参照を保持し、第Nの頂点が第N-1の頂点を参照し、第2から第Nまでの各頂点が、外部への参照を保持しないか、またはすべて同数の参照を保持する連続兄弟参照部を有する相互参照関係データに適用するための接続情報を有する連続兄弟参照用テンプレートが、テンプレート蓄積手段に蓄積されているようにすることもできる。

このテンプレートは、各頂点間の参照情報を有しなくてもよいので、記憶容量が少なく済む。

【0021】

上記いずれのデータ圧縮装置も、頂点間の参照情報に、テンプレートを適用可能な親テンプレートがテンプレート蓄積手段に蓄積されているとよい。

このテンプレートは、テンプレートを定義するのに他のテンプレートの定義を利用できるため、テンプレートを表現するのに必要な記憶容量を削減することができる。

【0022】

さらに、複数の入力データに共有可能な共用テンプレートがテンプレート蓄積手段に蓄積されていることが好ましい。

共用テンプレートは、複数の入力データに共有されるから、テンプレート数を少なくすることができるようになる。

【0023】

そして、本発明は、特定のパターンを有する複数の頂点間の参照情報を、頂点間の参照情報が共有可能なテンプレートとして蓄積するテンプレート蓄積手段と、テンプレートにより置換され、圧縮された相互参照関係データを第1の入力データとして入力し、相互参照関係データから、テンプレートを用いて圧縮前の元の相互参照関係データを復元する展開手段と、型と値をそれぞれ有する複数の頂点からなる頂点群のデータを第2の入力データとして入力し、頂点群のデータを展開手段により復元された相互参照関係データと合成したデータを出力する合成手段とを有するデータ復元装置を提供する。

このようなデータ復元装置によれば、圧縮された相互参照関係データをテンプレート蓄積手段に蓄積されたテンプレートを用いて圧縮前の元の相互参照関係データを復元することができる。

【0024】

さらに、本発明は、圧縮可能なデータを蓄積する第1のデータ蓄積手段と、請求項1~6のいずれか一項記載のデータ圧縮装置により圧縮されたデータを蓄積する第2のデータ蓄積手段と、第1のデータ蓄積手段及び第2のデータ蓄積手段に蓄積されたそれぞれのデータの利用頻度を観測し、その観測された利用頻度に応じて移動要求を出力する利用頻度観測手段と、その利用頻度観測手段からの移動要求にしたがい、利用頻度が高い高頻度データを請求項7記載のデータ復元装置により復元して第1のデータ蓄積手段へ格納し、利用頻度が高頻度データよりも低いデータを上記データ圧縮装置により圧縮して第2のデータ蓄積手段へ格納するようにして、データを移動させるための制御を行う制御手段と、第1のデータ蓄積手段と第2のデータ蓄積手段のいずれかから、制御手段の指示に応じてデータを取得して出力する選択手段とを有するデータ管理装置を提供する。

このデータ管理装置は、利用頻度に応じて、第1のデータ蓄積手段と第2のデータ蓄積手段との間でデータを移動させることができ、データを圧縮済みまたは非圧縮の状態に適

10

20

30

40

50

宜変えて蓄積することができる。したがって、利用頻度の高いデータの利用速度が高いまま維持されるため、高速化が可能となる。

【発明の効果】

【0025】

本発明によれば、データ構造の圧縮に用いるテンプレート数が増えても、各テンプレートの格納に要するメモリ等の記憶装置を圧迫しないようにすることができる。

【発明を実施するための最良の形態】

【0026】

以下、図面を参照して本発明に係るデータ圧縮装置、データ復元装置およびデータ管理装置の実施の形態について、添付図面を用いて詳細に説明する。

10

データ圧縮装置の実施の形態

(第1の実施の形態)

図1は本実施の形態に係るデータ圧縮装置101の構成を示すブロック図である。図1に示すように、データ圧縮装置101はテンプレート蓄積手段102、分離手段103、テンプレート一致箇所検出手段104およびテンプレート置換手段105を有している。このデータ圧縮装置101は、入力データ106から、第1の出力データ107と第2の出力データ108を出力するようになっている。

【0027】

図1におけるデータ圧縮装置101はテンプレート蓄積手段102、分離手段103、テンプレート一致箇所検出手段104、テンプレート置換手段105までが一体化された単一の装置として構成されているが、必ずしも単一の装置として実現される必要はなく、各手段を図示しない通信手段により接続して実現することも可能である。例えば、データ圧縮装置101からテンプレート蓄積手段102を取り除き、テンプレート蓄積手段102を別体の装置として実現し、データ圧縮装置101とテンプレート蓄積手段102とを通信手段によって接続する構成にしてもよい。そうすると、テンプレート蓄積手段102を含まない複数のデータ圧縮装置を複数設け、それらの間で1つのテンプレート蓄積手段102を共有するといったことも可能になる。

20

【0028】

なお、本実施の形態においては、図2に示すXML文書20を圧縮する手順をもって発明の詳細を説明するが、本実施の形態におけるデータ圧縮装置101による圧縮の対象はXML文書20のようなXML文書に限定されるものではなく、型と値を有する頂点、または値を有し、その各値が属性情報として型を有する頂点と、その頂点間の参照情報とを有するような一般的なグラフ構造によって表現されるデータにも適用することができる。値を有し、その各値が属性情報として型を有する頂点とは、例えば値として“1”、属性情報として“整数”といった情報を有するような頂点であり、属性情報から値の型が判定可能であることから、型と値を有する頂点を有するデータと同様に、本実施の形態におけるデータ圧縮装置101により、圧縮を行うことが可能である。

30

また、一般的なグラフ構造によって表現されるデータとは、XML文書のような根付順序木に限定されず、例えば閉路を有するような、より一般的な無向・有向グラフのことをいう。

40

【0029】

図2(a)はXML文書20のテキスト表現の一例を示す図である。XML文書20は既知の手法(例えば<http://xml.apache.org/xerces2-j/>において示されるXercesなど)により、図2(b)に示すようなコンピュータ内部に適したデータ構造21に変換可能である。そこで、以下では、XML文書20を図2(b)に示すデータ構造21に変換した後の圧縮工程について説明する。このデータ構造21は、型と値をそれぞれ有する複数の頂点と、その頂点間の参照情報とを有している。

【0030】

一方、図1に示した分離手段103は、図2(b)に示すデータ構造21を入力データ106として入力し、入力データ106を各頂点間の参照情報を有する相互参照関係デー

50

タと、型と値を有する複数の頂点からなる頂点群とに分離する。すなわち、分離手段103は、各頂点を一意に識別可能なID（頂点ID）を各頂点901～923に順に割り振ったものを相互参照関係データ900とし、割り振った頂点IDと、対応する頂点901～923がもともと有していた型および値との組を列記してテーブル950を生成し、そのテーブル950を型と値を持つ頂点901～923の集合（頂点群）とすることによって、データ構造21を分離している。各頂点901～923への頂点IDの割り振り方には、幅優先探索や、深さ優先探索等があるが、ここでは幅優先探索を用いている。また、分離手段103は、分離して得られる頂点群のデータを第2の出力データ108として出力する。

【0031】

そして、XML文書20から分離された相互参照関係データ900、およびテーブル950はそれぞれ、図3(a)、図3(b)に示す通りである。ここで、図3(b)は、頂点ID950a、型950bおよび値950cを一行とする形式で表現されている。

テンプレート蓄積手段102は、圧縮に先立ちあらかじめテンプレートとテンプレート実体を蓄積している。このとき、テンプレート蓄積手段102は、テンプレートやテンプレート実体として、例えば、あらかじめ高い頻度で適用されることが分かっている高頻度のものを蓄積している。このようなテンプレートとテンプレート実体としては、それぞれ例えば、図4(a)、図4(b)にそれぞれ示すテンプレート1003とテンプレート実体1020とがある。

【0032】

テンプレート1003は、テンプレートID1001、接続情報1002およびパターン情報1004を有している。テンプレートID1001はテンプレート蓄積手段102に複数のテンプレートが蓄積された際に、その各テンプレートを一意に識別するために用いられる。パターン情報1004はテンプレートによって表現される頂点間の参照情報のパターンを表し、複数の頂点とそれら相互の参照情報とを有している。パターン情報1004に含まれる参照情報には、親参照、子参照、次兄弟参照、前兄弟参照の4つの参照が設けられている。なお、接続先の頂点がない参照のうち、テンプレートを適用し、後述するようにして、相互参照関係データ900を圧縮する際に利用されない参照については、その旨がパターン情報1004に記述されている。これは、例えば、無効な頂点を定義しておき、その頂点への参照とすることで実現可能である。接続情報1002には、テンプレート1003を適用して、相互参照関係データ900を圧縮する際に他のテンプレートや頂点との接続を示す接続参照が列挙されている。

【0033】

テンプレート1003は、接続参照を有する接続情報1002と、参照情報を有するパターン情報1004とを区別して構成しているから、異なるテンプレート1003同士でパターン情報1004を共有することができる。つまり、接続情報1002を異ならせることにより、接続され得る頂点や他のテンプレートを異ならせ、パターン情報1004が同じでも、別テンプレートのようにして利用することができる。すると、テンプレート内に含まれる頂点間の参照情報が省略可能となり、テンプレート蓄積手段102のメモリ使用量（記憶領域）を効率よく利用することが可能となる。

【0034】

図4(c)は、テンプレート1003の具体的な一例となる第1のテンプレート1011を示す図である。図4(c)に示す第1のテンプレート1011におけるテンプレートID1012には、“1”が設定されている。第1のテンプレート1011におけるパターン情報1019は、5つの頂点1014～1018と、それらの間の参照とにより構成され、参照は矢印で記述されている。なお、参照の種類は矢印に対し、親参照はp、子参照はc、次兄弟参照はns、前兄弟参照はpsとして記述されている。例えば、頂点1014の子参照cは頂点1016を指定しており、頂点1016の次兄弟参照は頂点1017を指定している。

【0035】

10

20

30

40

50

また、テンプレートを適用し相互参照関係データ900を圧縮する際に利用されないことを示す参照は端点を「x」で記述し、テンプレートを適用し相互参照関係データ900を圧縮する際に他のテンプレートや頂点と接続されることを示す参照は端点を「」で記述している。後者に該当する4つの参照、すなわち、頂点1014の親参照、頂点1016、頂点1017、頂点1018の子参照については、接続情報1013に頂点IDと参照の種類が列挙されている。

【0036】

テンプレート実体1020は、入力データ106に対し、相互参照関係データ900を圧縮する際に、テンプレートを適用したこと(テンプレート適用済み)を表すために用いられる。このテンプレート実体1020は、テンプレート実体ID1005と、反転フラグ1006と、テンプレート独自情報1007とを有している。また、テンプレート独自情報1007は適用するテンプレートを表す利用テンプレートID1008と、実体接続情報1009および実体情報1010とを有している。圧縮後の相互参照関係データにおいて、テンプレート実体1020を参照することにより、テンプレートの適用内容が把握できるようになる。

10

【0037】

テンプレート実体ID1005は、テンプレートを適用して相互参照関係データ900を圧縮した際、そのテンプレートの適用箇所を一意に特定するために用いられる。反転フラグ1006は、テンプレートを適用して相互参照関係データ900を圧縮した際に兄弟関係の方向を反転して利用するか否かを示す。本実施の形態では、反転フラグ1006が「偽」のときに次兄弟参照と、前兄弟参照は文字通りの意味を持ち、反転フラグ1006が「真」のときに次兄弟参照が前兄弟参照の意味を、前兄弟参照が次兄弟参照の意味を持つものとしている。

20

【0038】

テンプレート独自情報1007に含まれる実体情報1010は、テンプレートを適用して相互参照関係データ900を圧縮した際に接続する先の頂点を列挙するために設けられている。この実体接続情報1010には、テンプレートを適用して相互参照関係データ900を圧縮した際にテンプレート内に内包される頂点のIDが蓄積されている。実体接続情報1009については後述する。なお、反転フラグ1006は、同じテンプレートを反転させて利用する場合には必要であるが、そのようなテンプレートの適用を行わない場合は設けなくてもよい。

30

【0039】

次に、テンプレート一致箇所検出手段104は、入力データ106から分離手段103により分離された相互参照関係データ900から、テンプレート蓄積手段102に蓄積されているテンプレートに一致する箇所を検出する。テンプレート蓄積手段102には、複数のテンプレートが蓄積されることが予想されるため、テンプレート一致箇所検出手段104による検出結果は複数通り存在すると考えられる。ただし、例えば後述する図8に示す一致箇所検出手順によれば、検出結果は一意に定まる。

【0040】

本実施の形態では、テンプレート蓄積手段102に第1のテンプレート1011のみが蓄積されているときに、相互参照関係データ900に対して、図8に示す手順により求めた一致箇所を示す。そのテンプレートの一致箇所1501は、例えば図9(a)に示すように、利用テンプレートID1502、反転フラグ1503およびテンプレートの頂点から元の相互参照関係データ900の頂点への割り当てを表す頂点对応情報1504によって表すことができる。

40

【0041】

そして、テンプレート一致箇所検出手段104により、相互参照関係データ900からテンプレート蓄積手段102に蓄積されているテンプレートの一致箇所を検出した結果を表すテンプレート一致箇所情報1505は図9(b)に示す通りである。このテンプレート一致箇所情報1505は、相互参照関係データ900から、テンプレート蓄積手段10

50

2に第1のテンプレート1011のみが蓄積されているとき、図8に示す手順にて検出した結果であり、ここでは、第1、第2、第3の一致箇所1506, 1507, 1508を有し、一致箇所が3箇所あったことを示している。

【0042】

第1の一致箇所1506及び第2の一致箇所1507は第1のテンプレート1011を反転せず、テンプレートの頂点と元の相互参照関係データ900の頂点をそれぞれ頂点对応情報1509、頂点对応情報1510に示すように対応させることで一致することを表している。また、第3の一致箇所1508は第1のテンプレート1011を反転し、テンプレートの頂点と元の相互参照関係データ900の頂点を頂点对応情報1511に示すように対応させることで一致することを表している。

10

そして、テンプレート一致箇所検出手段104は、このようなテンプレート一致箇所情報1505をテンプレート置換手段105に伝達(入力)する。

【0043】

テンプレート置換手段105は、テンプレート一致箇所検出手段104からテンプレート一致箇所情報1505を入力し、そのテンプレート一致箇所情報1505を用いて、元の相互参照関係データ900に対してテンプレートを適用し、テンプレート実体1020を用いて、頂点間の参照情報を参照可能な状態のまま置換し、置換した結果を第1の出力データ107として出力する。

テンプレートを用いて元の相互参照関係データ900を置換する置換手順は、図10に示すとおりで、また、置換した結果は図5に示す相互参照関係データ1100のようになる。この相互参照関係データ1100は圧縮済みの相互参照関係データ(以下「圧縮済み参照データ」ともいう)である。

20

【0044】

図10は、置換手順を示す図である。置換手順は、処理開始後、ステップ1に進み*i*に0をセットして、ステップ2に進み、テンプレート一致箇所情報(テンプレート一致箇所情報1505)に含まれるすべての一致箇所(上述の場合は第1、第2、第3の一致箇所1506, 1507, 1508)について、それぞれ1つずつステップ3以下の処理を繰り返す。

選択した一致箇所は*M_i*とする。

ステップ3では、テンプレート実体を1つ作成し、実体ID=*i*とする。このテンプレート実体を*O_i*とし、以下の処理を行う。

30

利用テンプレートID、反転フラグ*M_i*の利用、テンプレートID、反転フラグよりそれぞれ複製する。

実体情報を*M_i*の頂点对応情報より複製する。

実体接続情報は実体情報に記述された対応関係より、元の参照をそのまま代入する。

次に、ステップ4に進んで*i = i + 1*を計算する、

続くステップ5では、作成済のテンプレート実体を1つずつ選択し、以下の処理を繰り返す。選択したテンプレート実体を*O_i*とする。

次に、ステップ6に進み、実体接続情報に記述された参照の接続先頂点が他のテンプレート実体に含まれる場合はテンプレート実体IDとテンプレートの頂点の組に置換する。

40

【0045】

図5において、テンプレート実体として、3つのテンプレート実体1124、テンプレート実体1133、テンプレート実体1140が存在している。各テンプレート実体1124、テンプレート実体1133、テンプレート実体1140はそれぞれ実体ID1125、実体ID1131、実体ID1138を有し、そのそれぞれが、“1”、“2”、“3”の値を持っていることで識別可能になっている。各テンプレート実体は、すべて利用テンプレートID1128, 1135, 1142に“1”が設定され、反転フラグ1126, 1132, 1139は、前2者が偽、後者が真の値を有している。これにより、テンプレート実体1124, 1133は図4(c)に示す第1のテンプレート1011がそのまま適用され、テンプレート実体1140は第1のテンプレート1011が反転して適用

50

されたことを表す。

【 0 0 4 6 】

図 5 における各テンプレート実体 1 1 2 4 , 1 1 3 3 , 1 1 4 0 の実体情報は、テンプレートが内包する頂点とテンプレート適用前の相互参照関係データの頂点との対応を示す情報が設定されている。そのため、置換された相互参照関係データ 1 1 0 0 では、頂点間の参照情報が残り、これらが参照可能になっている。例えば、テンプレート実体 1 1 2 4 の場合、実体情報 1 1 3 0 には、第 1 のテンプレート 1 0 1 1 の頂点 1 0 1 4 , 1 0 1 5 , 1 0 1 6 , 1 0 1 7 および 1 0 1 8 がそれぞれ図 3 (a) に示す圧縮前の相互参照データ 9 0 0 の頂点 9 0 2 , 9 0 3 , 9 0 4 , 9 0 5 , 9 0 6 にそれぞれ一致することを表す情報が設定されている。

10

また、各テンプレート実体の実体接続情報には、他のテンプレート実体や頂点との接続関係を示す情報が設定されている。各テンプレート実体が適用しているテンプレートは第 1 のテンプレート 1 0 1 1 であるが、第 1 のテンプレート 1 0 1 1 は外部と接続できる参照を 4 つ保持していることがその接続情報 1 0 1 3 に記述されている。

【 0 0 4 7 】

そこで、各テンプレート実体の実体接続情報には、これらの参照先がどの頂点となるのかを記述する。例えばテンプレート実体 1 1 2 4 の場合、実体接続情報 1 1 2 9 には、頂点 1 0 1 4 の親参照は頂点 9 0 7 へ、頂点 1 0 1 6 、頂点 1 0 1 7 の子参照はそれぞれテンプレート実体 ID が “ 2 ” の頂点 1 0 1 4 、テンプレート実体 ID が “ 3 ” の頂点 1 0 1 4 へ、頂点 1 0 1 8 の子参照はどの頂点にも接続しないことを示す情報が設定されている。

20

【 0 0 4 8 】

図 3 に示すように、テンプレート適用前の相互参照関係データ 9 0 0 は、各頂点 9 0 2 ~ 9 1 7 がそれぞれ 4 つの参照を持っていたが、図 5 に示すテンプレート適用後の相互参照関係データ 1 1 0 0 は、各頂点間の参照情報を持たないテンプレート実体により置換されている。このような置換を行うことにより、テンプレートの圧縮が可能となっている。また、適用するテンプレートも、従来技術とは以下のような相違がある。つまり、テンプレート一致箇所 1 5 0 8 は、第 1 のテンプレート 1 0 1 1 と一致しなかったため、従来技術では、テンプレート一致箇所 1 5 0 8 のための別なテンプレートが必要であったが、本実施の形態では、上述したように、第 1 のテンプレート 1 0 1 1 を反転させてテンプレート一致箇所 1 5 0 8 に一致させることができるから、テンプレート一致箇所 1 5 0 8 のための別なテンプレートを設ける必要がない。そのため、テンプレート蓄積手段 1 0 2 のメモリ使用量（記憶領域）を効率よく利用することが可能となる。

30

【 0 0 4 9 】

なお、本実施の形態では、図 3 (b) に示す入力データから分離された型と値を持つ頂点の集合である第 2 の出力データ 1 0 8 については、分離手段 1 0 3 により分離された後に出力されるだけで圧縮については何ら触れられていない。第 2 の出力データ 1 0 8 は、例えば、非特許文献 1 に示される方法等と組み合わせることにより圧縮することが可能である。

また、本実施の形態では、反転フラグは各テンプレート実体に 1 つずつ用意しているが相互参照関係データ全体で 1 つとしてもよいし、また両方を設定してもよい。

40

【 0 0 5 0 】

一方、上述した本実施の形態で示すデータ圧縮装置 1 0 1 を複数種類の入力データに適用する場合、テンプレートについてはその複数種類の入力データ間で共用する共用テンプレートとすることができる。その共用テンプレートは、複数種類の入力データに適用可能であるから、それぞれの入力データに対応してテンプレートを設けることを要しない。したがって、テンプレート蓄積手段 1 0 2 のメモリ使用量の効率化が可能である。

【 0 0 5 1 】

例えば、図 1 1 (a) に示す XML 文書 3 0 の場合、図 8 に示す一致箇所検出手順により、図 2 に示す XML 文書 2 0 と同様、図 4 (c) に示す第 1 のテンプレート 1 0 1 1 を

50

適用することができる。そこで、テンプレート蓄積手段102は双方のXML文書20, 30に適用したテンプレートを区別せずに同一の共用テンプレートとして蓄積することにより、テンプレート蓄積手段102におけるメモリ使用量の効率化(メモリ利用効率の向上)を図ることができる。

【0052】

なお、図8に示す一致箇所検出手順は以下のとおりである。

処理開始後ステップ11で、パターン蓄積手段に蓄積されたパターンから、頂点の数が多い順に1つずつ選択し、以下の処理を繰り返す。

選択したパターンをP_jとする。

次にステップ12に進んで、反転フラグの値を偽、真のそれぞれに対し、以下を繰り返す。

続くステップ13では、相互参照関係データに含まれる頂点から、選択したパターンPの頂点の数と一致する頂点を選択する組み合わせをX₁, X₂, X_mとし、その中から1つずつ選択して、以下を繰り返す。

選択した組み合わせをX_kとする。

次に、ステップ14に進み、X_kに含まれる頂点はすべて置換済みマークが無いが否かを判断し、無ければステップ15に進み、そうでなければ処理を終了する。ステップ15に進むと、X_kがP_jと同型か否かを判断し、同型であればステップ16に進み、そうでなければ処理を終了する。ステップ16に進むと、X_kを一致箇所として登録し、X_kに含まれる頂点は置換済みとしてマークする。

【0053】

(第2の実施の形態)

図6(a)に示すような連続する複数の兄弟参照を有する相互参照関係データ1204を圧縮する場合について説明する。データ圧縮装置101によれば、この相互参照関係データ1204は、図6(a)に示す連続兄弟参照部1200を図6(b)に示す第2のテンプレート1201を用いて圧縮する。図6(b)に示す第2のテンプレート1201はテンプレートID1202と、接続情報1203を有するが、第1のテンプレート1101とは異なり、パターン情報を有していない。この第2のテンプレート1201は、相互参照関係データ1204のような連続兄弟参照部を有する相互参照関係データを圧縮するために設けた連続兄弟参照用テンプレートである。

【0054】

ここで、相互参照関係データ1204では、図6(a)に示すように、連続兄弟参照部1200が、第1、第2、第3から第NまでのN個の頂点2a, 2b, 2c, …, 2nを有し、第1の頂点2aと第2の頂点2b以外の連続番号を有する各頂点は必ず相互に参照し、第1の頂点2aが第2の頂点2bを参照し、かつ連続兄弟参照部1200の外部にある頂点2pへの参照を保持している。また、第Nの頂点2nは、図示しない第N-1の頂点を参照し、さらに、第2から第Nまでの各頂点が、連続兄弟参照部1200の外部への参照をまったく保持しないようになっている(または、例えば、図3における頂点912, 913のように同数の参照を保持するようになっていてもよい)。

このような相互参照関係データ1204を圧縮するには、少なくとも、連続兄弟参照部1200を構成する頂点の個数と、第1の頂点の外部への参照がわかればよいので、第2のテンプレート1201における接続情報1203には、連続兄弟参照部を構成する頂点の個数Nと、テンプレートの親参照pが設定されている。なお、テンプレートID1202は“2”を有している。

【0055】

図7は、図6(b)に示す第2のテンプレート1201を用いて、相互参照関係データ1204を圧縮した後の相互参照関係データ1300を示す図である。この場合の圧縮では、テンプレートにおける一致箇所の検出やテンプレートの適用は第1の実施の形態と同様にすることで可能である。

図7において、相互参照関係データ1300には、テンプレート実体として、テンプレ

ート実体 1 3 1 8 , 1 3 2 3 , 1 3 2 9 の 3 つが存在し、実体 ID 1 3 1 7 , 実体 ID 1 3 2 2 , 実体 ID 1 3 2 8 は、それぞれ “ 1 ” , “ 2 ” , “ 3 ” の値を有している。各テンプレート実体は、すべて利用テンプレート ID 1 3 2 6 , 1 3 3 4 , 1 3 3 3 を有していて、いずれも “ 2 ” が設定されている（これは、第 2 のテンプレート 1 2 0 1 を用いて圧縮したことを意味している）。なお、本実施の形態では、反転フラグを用いないため各テンプレート実体 1 3 1 8 , 1 3 2 3 , 1 3 2 9 には反転フラグが設けられていない。

【 0 0 5 6 】

各テンプレート実体 1 3 1 8 , 1 3 2 3 , 1 3 2 9 の実体接続情報 1 3 2 0 , 1 3 2 5 , 1 3 3 1 は連続兄弟参照部を構成する頂点の個数 N （それぞれ、 $N = 4 , 3 , 2$ ）と、各テンプレート実体の親参照 p を記録している。例えばテンプレート実体 1 3 1 8 の場合では、 $N = 4$ 、 p は 4 0 1 になっている。各テンプレート実体の実体情報 1 3 2 1 , 1 3 2 7 , 1 3 3 2 には、テンプレートにより内包される頂点を示す情報が示されている。例えば、テンプレート実体 1 3 1 8 の場合は、実体接続情報 1 3 2 0 より示されている 4 つの頂点、すなわち、頂点 4 0 2、4 0 3、4 0 4、4 0 5 を示している。

以上のように相互参照関係データ 1 2 0 4 は第 2 のテンプレート 1 2 0 1 を用いて圧縮可能である。その圧縮に用いる第 2 のテンプレート 1 2 0 1 はパターン情報を有していないため、第 2 のテンプレート 1 2 0 1 を記憶するのに必要な記憶容量が少なく済む。そのため、各テンプレート蓄積手段 1 0 2 のメモリ使用量を削減することが可能である。

【 0 0 5 7 】

データ復元装置の実施の形態

次に、データ復元装置 1 8 0 1 について、図 1 2 を用いて説明する。図 1 2 は本実施の形態に係るデータ復元装置 1 8 0 1 の構成を示すブロック図である。このデータ復元装置 1 8 0 1 は、図 5 に示すような圧縮後の相互参照関係データ（圧縮済み参照データ） 1 1 0 0 と、図 3（b）に示す型と値を有する複数の頂点からなる頂点群のデータとから、それぞれが型と値を有する複数の頂点と、頂点間の参照情報とを有する元の入力データを復元する。データ復元装置 1 8 0 1 は、テンプレート蓄積手段 1 8 0 2 と、合成手段 1 8 0 3 と、テンプレート展開手段 1 8 0 4 とを有している。

【 0 0 5 8 】

なお、図 1 2 において、データ復元装置 1 8 0 1 は各手段が一体化された単一の装置として構成されているが、必ずしも単一の装置として実現される必要はなく、複数の装置を図示しない通信手段により接続して実現することも可能である。例えば、データ復元装置 1 8 0 1 からテンプレート蓄積手段 1 8 0 2 を分離した上で、テンプレート蓄積手段 1 8 0 2 を別な単一装置として実現し、両装置間を図示しない通信手段により接続する構成にしてもよい。そうすると、テンプレート蓄積手段 1 8 0 2 を有しない複数のデータ圧縮装置間でテンプレート蓄積手段 1 8 0 2 を共有するといったことも可能になる。

テンプレート展開手段 1 8 0 4 は第 1 の入力データ 1 8 0 6 として与えられた圧縮後の相互参照関係データをテンプレート蓄積手段 1 8 0 2 に蓄積されたテンプレートを用いて展開する。その展開は、例えば上述した図 1 3 に示した復元手順で行うことができる。テンプレート展開手段 1 8 0 4 により復元された相互参照関係データは、図 3 に示す相互参照関係データ 9 0 0 のようになる。

【 0 0 5 9 】

合成手段 1 8 0 3 は展開された相互参照関係データと、図 3（b）に示す第 2 の入力データとして与えられた複数の型と値を有する頂点群のデータとを合成し、合成されたデータを出力データ 1 8 0 5 として出力する。その合成は、図 3（b）に示すテーブル 9 5 0 における型と値を持つ頂点群において、各頂点に頂点 ID が割り振ってあるため、頂点 ID が一致する相互参照関係データの頂点に、型と値をあてはめていくことによって行う。

以上のような手順により、それぞれが型と値を有する元の複数の頂点と、頂点間の参照情報からなる入力データを復元することが可能である。

復元手順は以下のとおりである。

図 1 3 において、開始後のステップ 2 1 で、圧縮済みの相互参照関係データに含まれる

10

20

30

40

50

すべてのテンプレート実体を X_1 , X_2 , X_n とし、すべてについて以下を行う。

選択したテンプレート実体を X_i とする。

次にステップ 22 に進み、テンプレート実体 X_i が利用するテンプレートが持つ頂点間の参照情報を複製し、テンプレート実体 X_i の実体情報に記述される頂点の ID を割り振る。

次いでステップ 23 に進み、テンプレート実体 X_i の実体接続情報に記述された頂点が他のテンプレート実体 X_m に含まれる頂点の場合、テンプレート実体 X_m に記述される頂点 ID で置換する。

【 0 0 6 0 】

データ管理装置の実施の形態

10

本発明によるデータ圧縮装置により、相互参照関係データを圧縮すると、その圧縮後のデータへのアクセス速度の若干の低下が見込まれる。そのため、データの利用頻度を観測しておいて、その時々で利用頻度の高いものは非圧縮とし、いったん圧縮した相互参照関係データについても、利用頻度が高くなれば非圧縮の状態に戻し、逆に利用頻度が低くなれば再度圧縮する、といった方法でデータ管理を行うことが好ましい。このようなデータ管理を行うデータ管理装置を設ければ、装置全体の高速化と省メモリ化を両立させることも可能である。

【 0 0 6 1 】

図 14 は、このようなデータ管理を行えるデータ管理装置 2000 の構成を示すブロック図である。データ管理装置 2000 は、第 1 のデータ蓄積手段 2001 と、第 2 のデータ蓄積手段 2004 と、データ圧縮装置 2002 と、データ復元装置 2003 とを有している。また、データ管理装置 2000 は、利用頻度観測手段 2006 と、制御手段 2005 と、選択手段 2008 とを有している。

20

【 0 0 6 2 】

なお、本実施の形態におけるデータ管理装置 2000 は、各装置が一体化された単一の装置とされているが、本発明によるデータ管理装置は、必ずしも単一の装置として実現される必要はなく、各装置を図示しない通信手段により接続して実現することもできる。例えば、データ管理装置 2000 より、第 1 のデータ蓄積手段 2001 を分離し、データ圧縮装置 2002 から、後述のテンプレート蓄積手段 102 を取り除いた上で、第 1 のデータ蓄積手段 2001 を別な単一装置として実現し、両装置間を通信手段により接続する構成をとることができる。そうすると、第 1 のデータ蓄積手段 2001 を有しない複数のデータ管理装置間で第 1 のデータ蓄積手段 2001 を共有するといったことも可能になる。また、その他の構成手段についても同様である。

30

【 0 0 6 3 】

第 1 のデータ蓄積手段 2001 は、圧縮可能なデータとして、圧縮前のコンピュータに適した形式のデータ（例えば図 20 に示した相互参照関係データ 400 等）を蓄積している。第 2 のデータ蓄積手段 2004 は、圧縮されたデータ（例えば、図 5 に示す相互参照関係データ 1100 等）を蓄積している。ここで、データ圧縮装置 2002 は上述した本発明によるデータ圧縮装置 101 と同じ構成を有し、データ復元装置 2003 は上述したデータ復元装置 1801 と同じ構成を有している。

40

【 0 0 6 4 】

制御手段 2005 は第 1 のデータ蓄積手段 2001、第 2 のデータ蓄積手段 2004、データ圧縮装置 2002、データ復元装置 2003、選択手段 2008 をシステム外部から入力されるデータ指定 2009 に基づいて制御する。この制御手段 2005 は、利用頻度観測手段 2006 からの移動要求にしたがい、データ指定 2009 の指定に対応するデータを移動させるための制御を行う。選択手段 2008 は、制御手段 2005 の指示にしたがい、第 1 のデータ蓄積手段 2001 と第 2 のデータ蓄積手段 2003 のいずれかから蓄積されているデータを取得して出力する。

【 0 0 6 5 】

利用頻度観測手段 2006 は、第 1 のデータ蓄積手段 2001 または第 2 のデータ蓄積

50

手段 2004 に蓄積されているデータ（相互参照関係データ）の利用頻度を観測し、観測した利用頻度に応じて後述する移動要求を出力する。この利用頻度観測手段 2006 は、利用履歴リスト 2007 を内部に保持している。この利用履歴リスト 2007 には、例えば、利用要求のあったデータの識別 ID を利用要求のあった順に複数個（N 個）線形リスト（図示せず）として保存している。

本実施の形態におけるデータ管理装置 2000 は、内部に蓄積するデータを一意に識別するための識別 ID を各データに割り振っている。データ指定 2009 は、そのための識別 ID を外部から入力する手段である。

【0066】

以降、実際の動作内容について説明する。

データ管理装置 2000 の場合、データは第 1、第 2 いずれかのデータ蓄積手段 2001、2004 に保存されているが、初期状態では、第 1 データ蓄積手段 2001、第 2 のデータ蓄積手段 2004 のいずれにデータを蓄積しておいてもよい。以下の説明では、すべて第 2 のデータ蓄積手段 2004 に蓄積しておくことを想定している。

制御手段 2005 は、データ指定 2009 により、外部から識別 ID が入力されると、それを受けて利用頻度観測手段 2006 に指示を入力する。利用頻度観測手段 2006 は制御手段 2005 の指示を受けて、該当するデータが第 1 のデータ蓄積手段 2001 と第 2 のデータ蓄積手段 2004 のいずれに蓄積されているか、および、両手段の間でのデータの移動があるか否かを通知する情報を制御手段 2005 に入力する。

【0067】

ここで、利用頻度観測手段 2006 は、第 1、第 2 のデータ蓄積手段 2001、2004 の指定について、利用履歴リスト 2007 を参照し、データ指定 2009 により指定されるデータがその利用履歴リスト 2007 に有るか否かを判断する。そして、例えばそのデータが有れば第 1 のデータ蓄積手段 2001 に蓄積されているとし、無ければ第 2 のデータ蓄積手段 2004 から蓄積されている、というようにして返答する。

【0068】

さらに、利用頻度観測手段 2006 は、第 1、第 2 のデータ蓄積手段 2001、2004 の間でデータ移動の有無については、次のようにして制御手段 2005 に返答する。例えば、利用履歴リスト 2007 が更新された時、その利用履歴リスト 2007 に新規に載った（記録された）データは第 2 のデータ蓄積手段 2004 から第 1 のデータ蓄積手段 2001 へ移動したとし、利用履歴リスト 2007 から外れたデータは第 1 のデータ蓄積手段 2001 から第 2 のデータ蓄積手段 2004 へ移動したとして返答する。

そして、制御手段 2005 は利用頻度観測手段 2006 からの上述した返答に基づき、第 1 のデータ蓄積手段 2001 または第 2 のデータ蓄積手段 2004 を制御して、記憶しているデータを出力させ、選択手段 2008 により、いずれかから得たデータをデータ管理装置 2000 の外部に出力する。

【0069】

また、制御手段 2005 は利用頻度観測手段 2006 からデータの移動要求があったときに、データ蓄積手段 2001、2004 の間でデータを移動させるための制御を行う。例えば、第 1 のデータ蓄積手段 2001 から第 2 のデータ蓄積手段 2004 に移動するような移動要求があったときは、そのデータの利用頻度が低いため、第 1 のデータ蓄積手段 2001 からデータを取り出し、そのデータをデータ圧縮装置 2002 を用いて圧縮した上で、第 2 のデータ蓄積手段 2004 に格納するように、データ移動の制御を行う。第 1 のデータ蓄積装置 2001 からは取り出したデータを削除するように制御する。

【0070】

逆に、第 2 のデータ蓄積手段 2004 から第 1 のデータ蓄積手段 2001 に移動するように、移動要求があったときは、そのデータの利用頻度が高いため、第 2 のデータ蓄積手段 2004 からデータを取り出し、その圧縮されているデータをデータ復元装置 2003 を用いて復元し、その復元されたデータを第 1 のデータ蓄積手段 2001 に格納するように制御する。また、第 2 のデータ圧縮装置 2004 からは取り出したデータを削除するよ

10

20

30

40

50

うに制御する。

【0071】

以上のように、データ管理装置2000によると、利用頻度観測手段2006からの移動要求に応じて制御手段2005がデータの移動を制御することにより、第1のデータ蓄積手段2001と第2のデータ蓄積手段2004との間でデータを移動させることができるから、利用頻度に応じて、相互参照関係データを圧縮済みか、非圧縮の状態に適宜変えて蓄積することができる。すると、利用頻度の高いデータに関しては、非圧縮の状態で蓄積することにより、利用速度が高いまま維持されるため、動作速度を高速にしつつ全体としてのメモリ使用量を抑えることができる。

【0072】

(その他の実施の形態)

本発明によれば、図15に示すようなテンプレートをテンプレート蓄積手段に格納することもできる。図15は第3のテンプレート2100を示す図である。図15に示す第3のテンプレート2100のパターン情報2103において、エリア2104, 2105における頂点間の参照情報は、図4(c)に示す第1のテンプレート1011のパターン情報に一致することが分かる。したがって、テンプレートやテンプレート実体を拡張し、内部に適用されたテンプレートの数や、テンプレート内部の各テンプレートを一意に識別できるID等を記述する情報を付加することにより、テンプレート内部のパターン情報に対しても、前述までのテンプレートを適用することが可能になる。

【0073】

ここでは説明のためにテンプレート内部のパターン情報に前述までのテンプレートを適用するとき、元のテンプレートを親テンプレート、テンプレート内部のパターン情報に適用されるテンプレートを子テンプレートと呼ぶことにする。

テンプレート内部のパターン情報にも、テンプレートの適用を可能とするための親テンプレート2300と親テンプレート実体2320の構成例をそれぞれ図17(a), (b)に示す。親テンプレート2300と親テンプレート実体2320は、図4に示す第1の実施の形態で用いたテンプレート1003及びテンプレート実体1020に対し、前者に内部テンプレート情報2301を追加し、後者に内部テンプレート実体接続情報2302及び内部テンプレート実体情報2303を追加している。

【0074】

内部テンプレート情報2301には、例えば、子テンプレートの数や、子テンプレートのIDを記述する。内部テンプレート実体情報2303には、子テンプレートに内包される頂点が親テンプレートを実際に適用するとき、親テンプレート適用前の頂点とどのように対応するののかの対応関係を各頂点毎に記述する。内部テンプレート実体接続情報2302には、親テンプレートを実際に適用した際、子テンプレートが外部の親テンプレートと接続するときの接続情報を記述する。

【0075】

そして、図16に、図17(a)に示すテンプレートの具体的な親テンプレート2200を示す。この親テンプレート2200は、内部テンプレート情報2216に2種類の子テンプレートを2箇所に適用していることから、ID=1、ID=2としている。パターン情報2203は、第1の実施の形態に示すように、相互参照関係データを圧縮する要領で圧縮されており、第1のテンプレート1011が適用されている。

このように構成した親テンプレート2200を相互参照関係データに適用すれば、テンプレートを定義するのに他のテンプレートの定義を利用できるため、図15に示すような第3のテンプレート2100に比べて、テンプレートの情報量を削減することが可能である。

【図面の簡単な説明】

【0076】

【図1】本発明の実施の形態に係るデータ圧縮装置の構成を示すブロック図である。

【図2】(a)はXML文書の一例を示す図、(b)は(a)のXML文書のデータ構造

10

20

30

40

50

を示す図である。

【図3】(a)は図2のXML文書から分離された相互参照関係データを示す図、(b)は頂点の集合のテーブルを示す図である。

【図4】(a)はテンプレートの構成を示すブロック図、(b)はテンプレート実体の構成を示すブロック図、(c)は第1のテンプレートの構成を示すブロック図である。

【図5】圧縮後の相互参照関係データを示す図である。

【図6】(a)は別の相互参照関係データを示す図、(b)は第1のテンプレートの構成を示すブロック図である。

【図7】図6における圧縮後の相互参照関係データを示す図である。

【図8】一致箇所の検出手順の一例を示す図である。

10

【図9】(a)はテンプレート的一致箇所を示すブロック図、(b)はテンプレート一致箇所情報を示すブロック図である。

【図10】置換手順の一例を示す図である。

【図11】(a)は別のXML文書を示す図、(b)は(a)のXML文書のデータ構造を示す図である。

【図12】本実施の形態に係るデータ復元装置の構成を示すブロック図である。

【図13】テンプレートから元の頂点の参照情報を復元する手順の一例を示す図である。

【図14】データ管理装置の構成を示すブロック図である。

【図15】第1のテンプレートの構成を示すブロック図である。

【図16】内部にテンプレートを有する親テンプレートの一例を示すブロック図である。

20

【図17】(a)は親テンプレートを示すブロック図、(b)は親テンプレート実体を示すブロック図である。

【図18】XML文書の別の一例を示す図である。

【図19】図18のXML文書のデータ構造を示す図である。

【図20】(a)は図19のXML文書から分離された相互参照関係データを示す図、(b)は頂点の集合のテーブルを示す図である。

【図21】図18のXML文書から分離されたデータ構造を示す図である。

【図22】図18のXML文書から分離された要素名情報を示す図である。

【図23】図18のXML文書から分離されたテキスト情報を示す図である。

【符号の説明】

30

【0077】

20, 30 ... XML文書

21, 31 ... データ構造

101 ... データ圧縮装置

102, 1802 ... テンプレート蓄積手段

103 ... 分離手段

104 ... テンプレート一致箇所検出手段

105 ... テンプレート置換手段

106 ... 入力データ、107 ... 第1の出力データ

108 ... 第2の出力データ

40

900, 1100 ... 相互参照関係データ

1204, 1300 ... 相互参照関係データ

901 ... 頂点、950 ... テーブル

1003 ... テンプレート

1020, 1124, 1133 ... テンプレート実体

1140, 1323, 1329 ... テンプレート実体

1011 ... 第1のテンプレート

1201 ... 第2のテンプレート

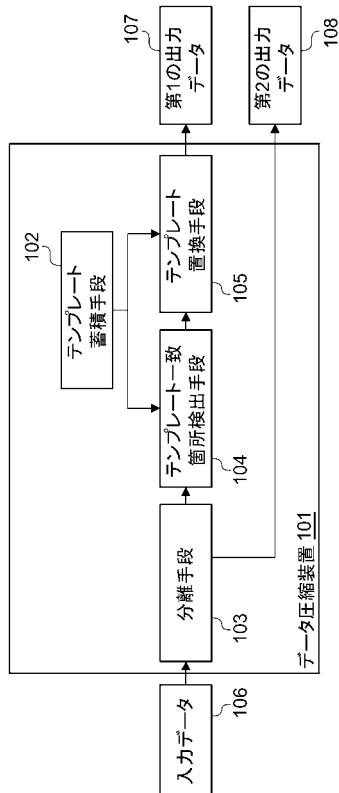
1501 ... 一致箇所

1505 ... テンプレート一致箇所情報

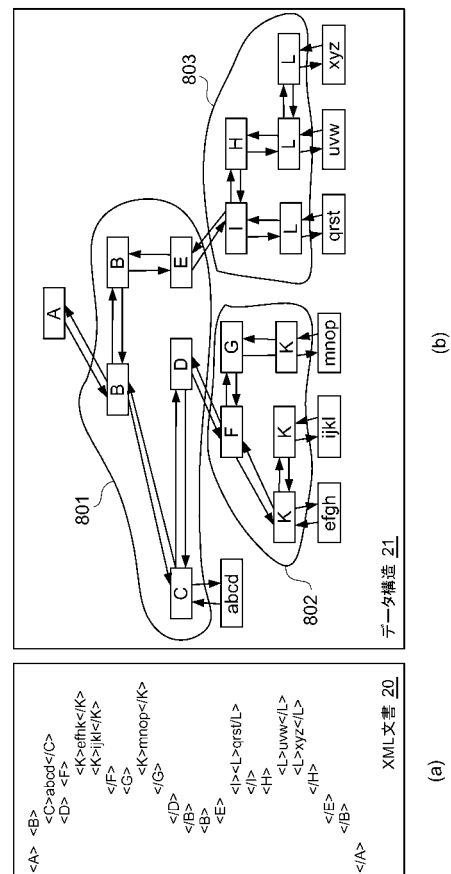
50

- 1 8 0 1 , 2 0 0 3 ... データ復元装置
- 1 8 0 3 ... 合成手段
- 1 8 0 4 ... テンプレート展開手段
- 2 0 0 0 ... データ管理装置
- 2 0 0 1 ... 第1のデータ蓄積手段
- 2 0 0 2 ... データ圧縮装置
- 2 0 0 4 ... 第2のデータ蓄積手段
- 2 0 0 5 ... 制御手段、2 0 0 6 ... 利用頻度観測手段
- 2 0 0 8 ... 選択手段、2 1 0 0 ... 第3のテンプレート
- 2 2 0 0、2 3 0 0 ... 親テンプレート
- 2 3 2 0 ... 親テンプレート実体

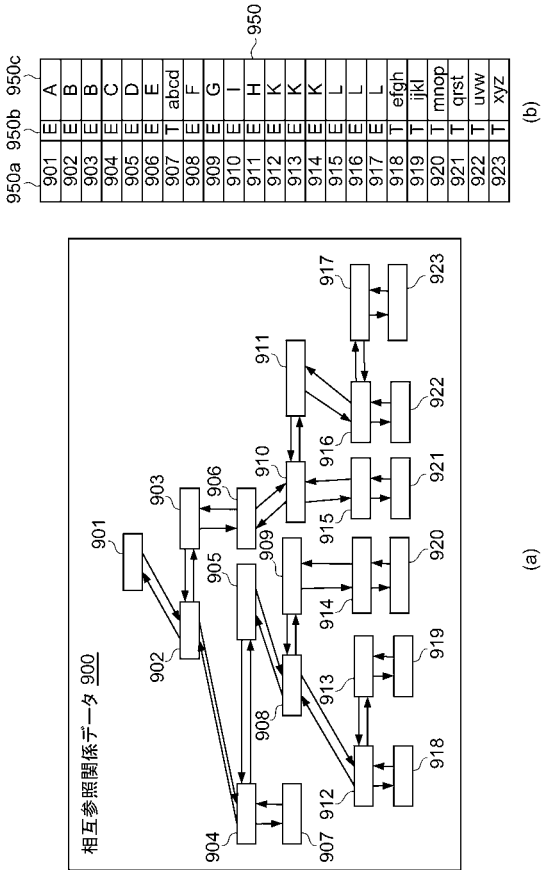
【 図 1 】



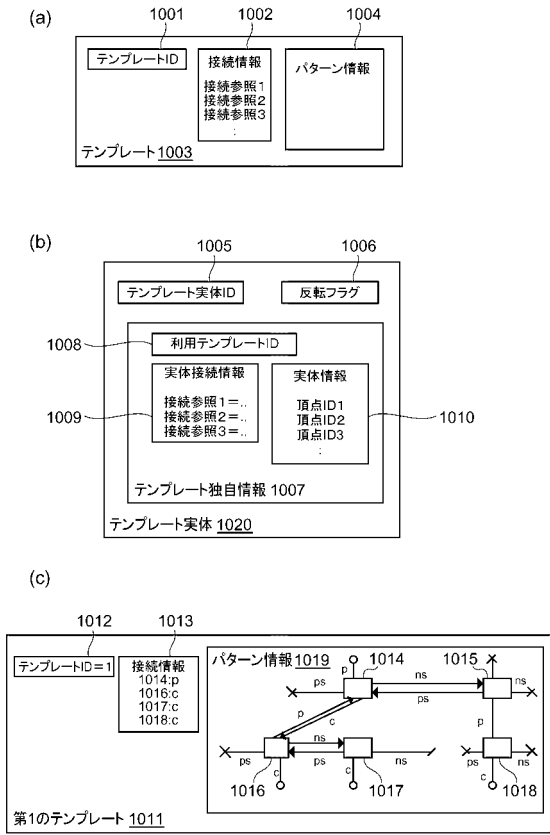
【 図 2 】



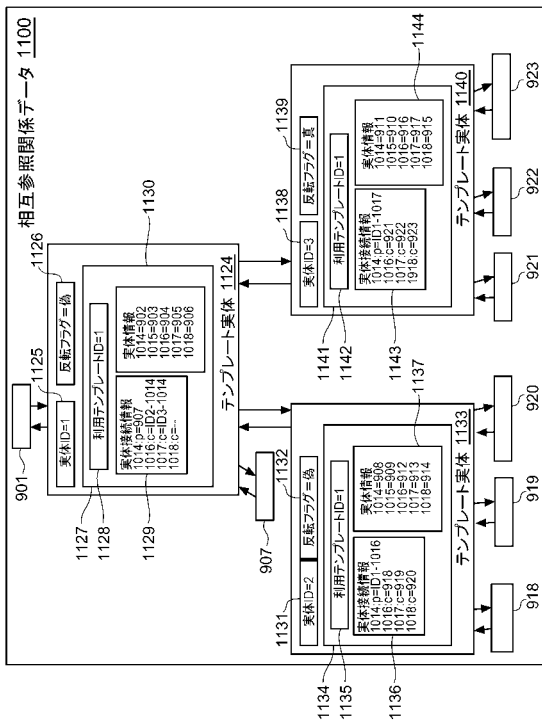
【 図 3 】



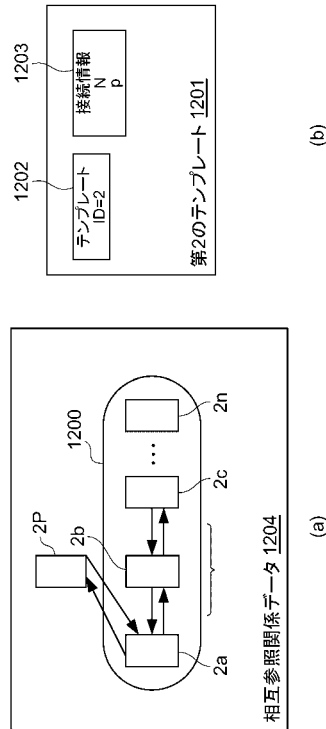
【 図 4 】



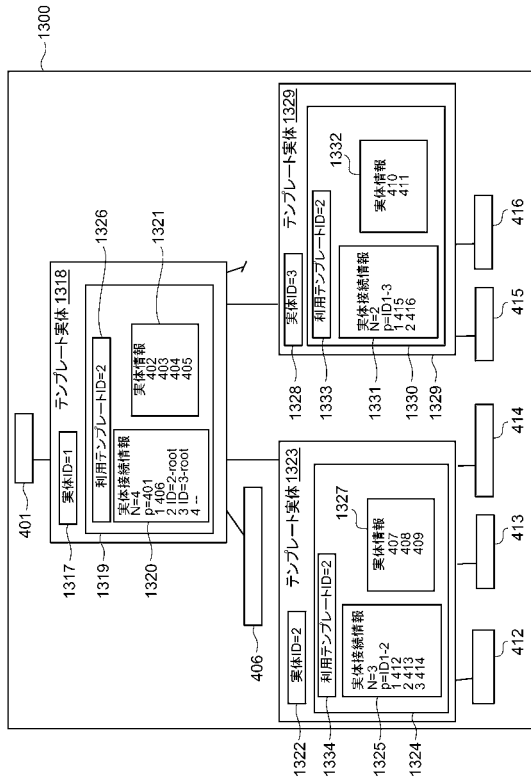
【 図 5 】



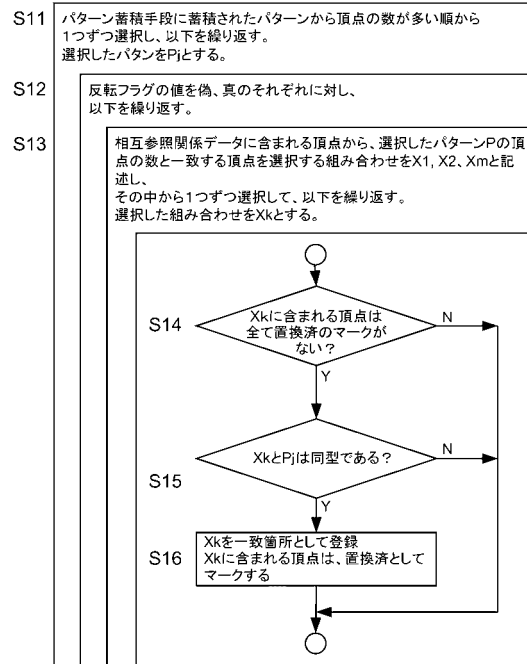
【 図 6 】



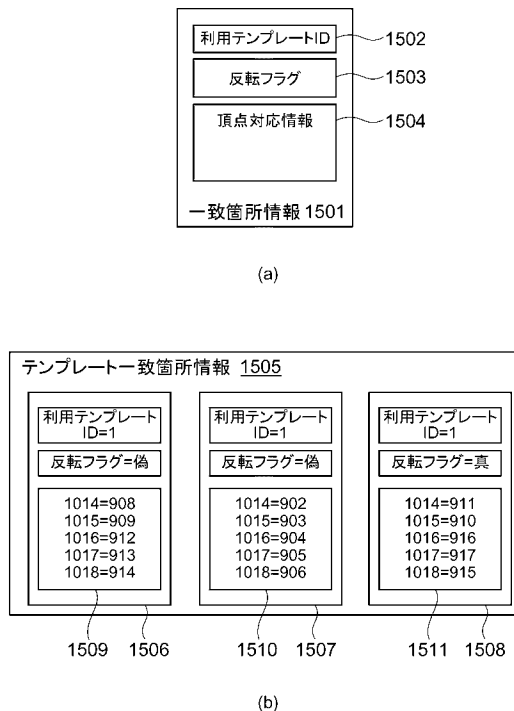
【図7】



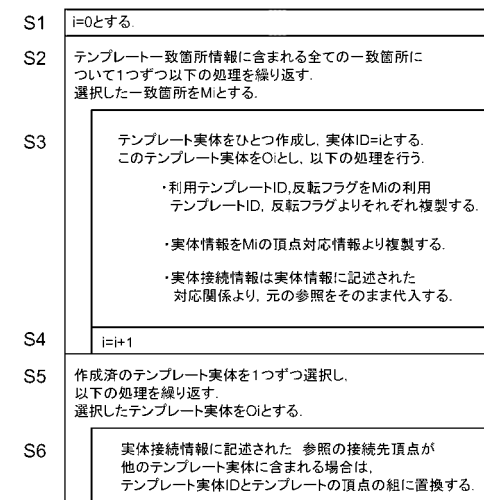
【図8】



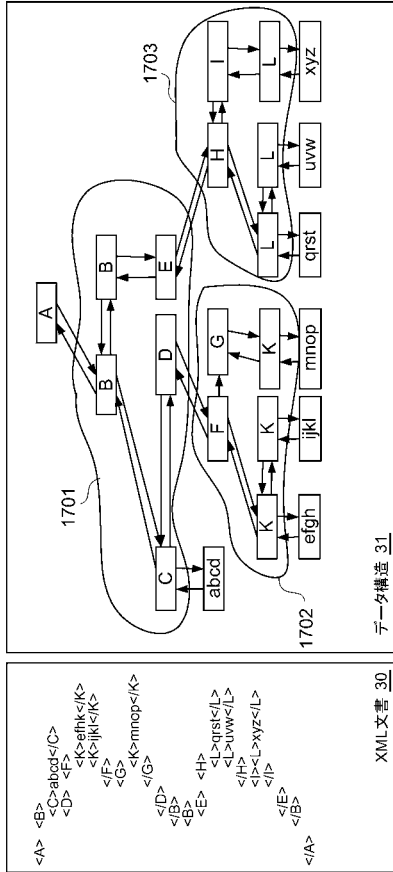
【図9】



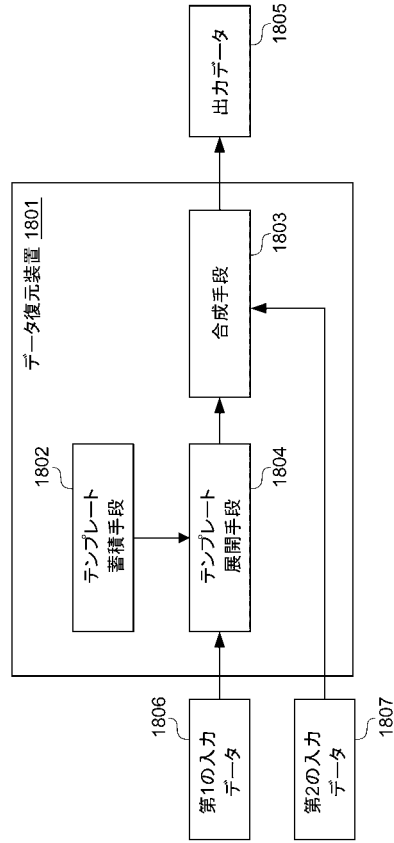
【図10】



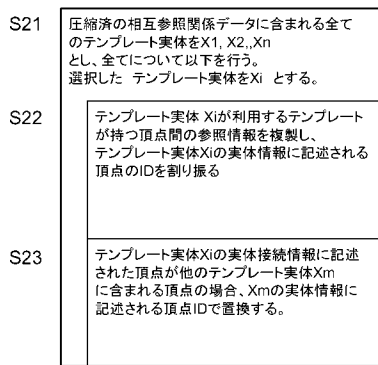
【図 1 1】



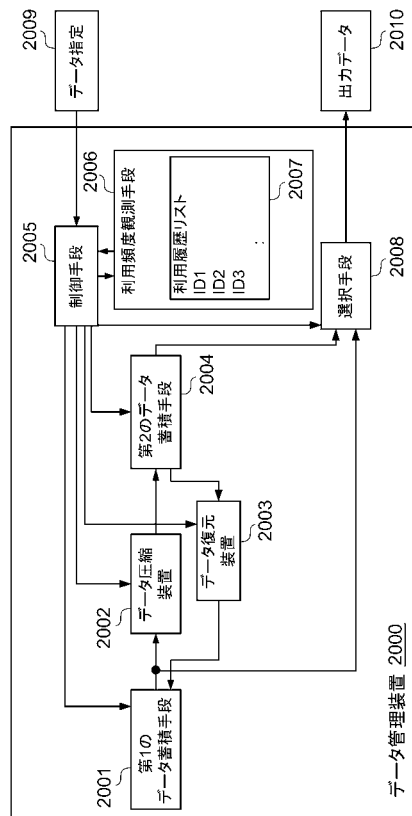
【図 1 2】



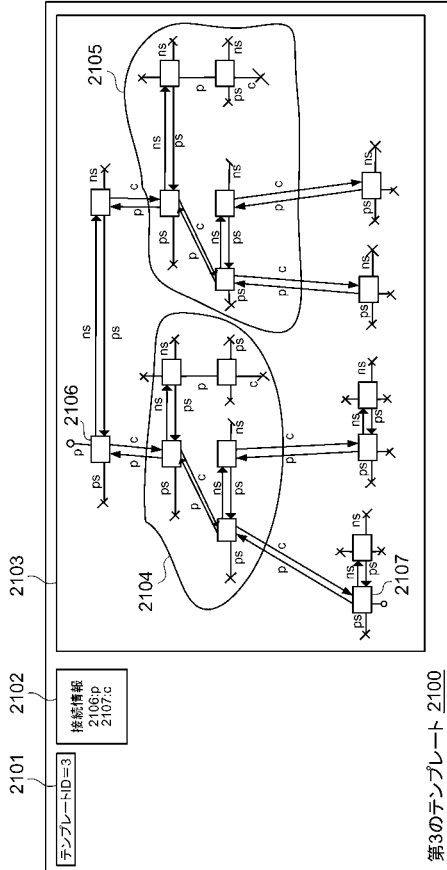
【図 1 3】



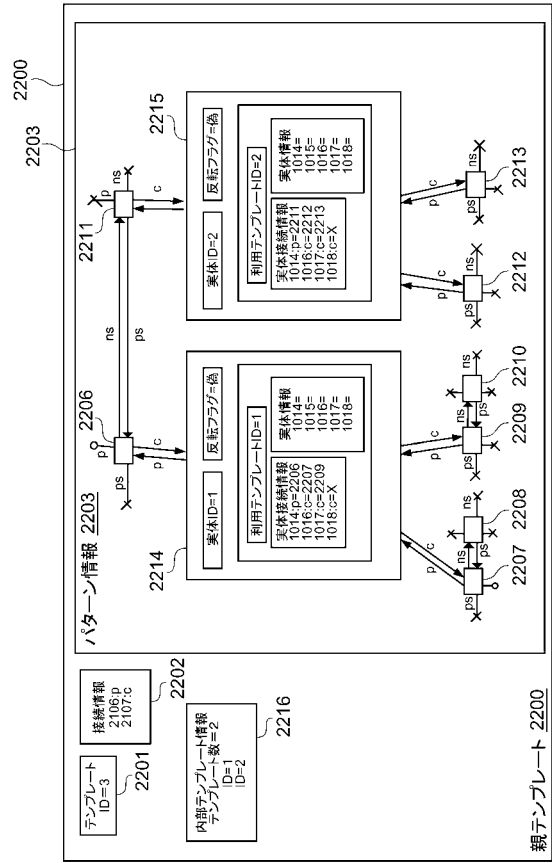
【図 1 4】



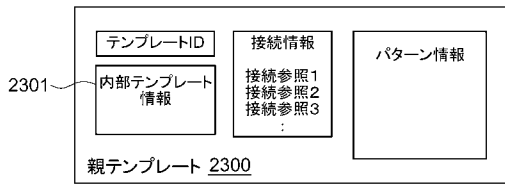
【図15】



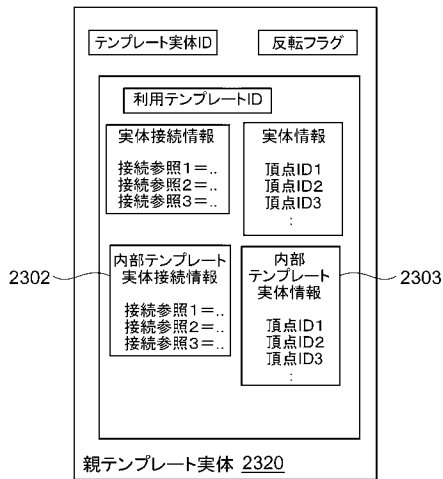
【図16】



【図17】



(a)

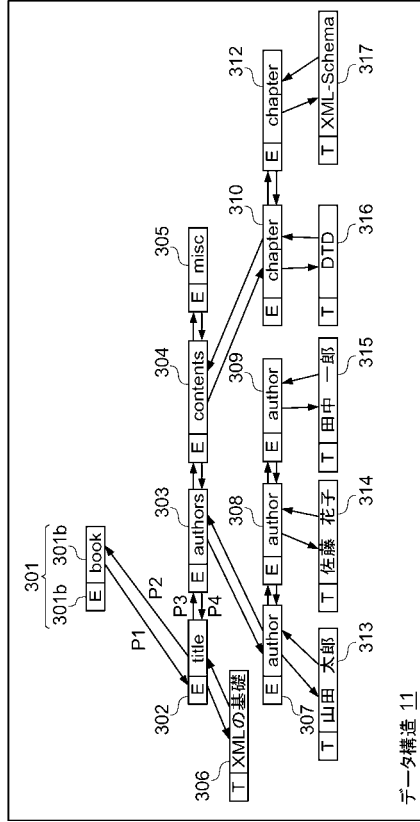


(b)

【図18】



【図 19】

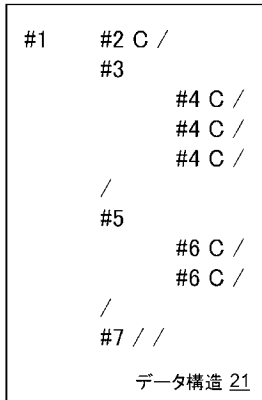


【図 20】

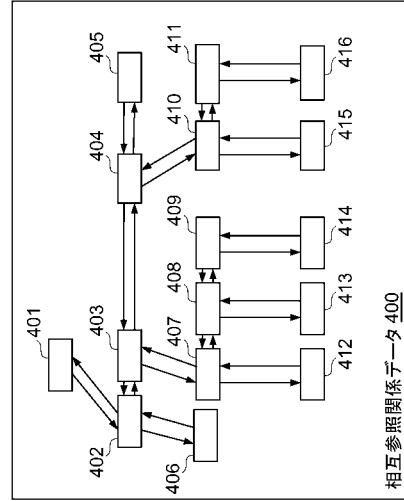
401	E	book
402	E	title
403	E	authors
404	E	contents
405	E	misc
406	T	XMLの基礎
407	E	author
408	E	author
409	E	author
410	E	chapter
411	E	chapter
412	T	山田 太郎
413	T	佐藤 花子
414	T	田中 一郎
415	T	DTD
416	T	XML-Schema

(b)

【図 21】



【図 22】



(a)

- #1: <book>
- #2: <title>
- #3: <authors>
- #4: <author>
- #5: <contents>
- #6: <chapter>
- #7: <misc>

要素名情報 13

【 図 2 3 】

XMLの基礎
山田 太郎
佐藤 花子
田中 一郎
DTD
XML-Schema
テキスト情報 14

フロントページの続き

- (72)発明者 行友 英記
東京都千代田区永田町二丁目11番1号 株式会社エヌ・ティ・ティ・ドコモ内
- (72)発明者 中山 雄大
東京都千代田区永田町二丁目11番1号 株式会社エヌ・ティ・ティ・ドコモ内
- (72)発明者 金野 晃
東京都千代田区永田町二丁目11番1号 株式会社エヌ・ティ・ティ・ドコモ内
- (72)発明者 竹下 敦
東京都千代田区永田町二丁目11番1号 株式会社エヌ・ティ・ティ・ドコモ内

審査官 渡辺 未央子

- (56)参考文献 特開2001-282516(JP,A)
特開平09-130616(JP,A)
特開2001-251617(JP,A)
特開2002-163248(JP,A)
特開2003-044459(JP,A)
特開2004-032774(JP,A)

(58)調査した分野(Int.Cl., DB名)
H03M 7/30