

(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2017-126185

(P2017-126185A)

(43) 公開日 平成29年7月20日(2017.7.20)

(51) Int.Cl.	F 1	テーマコード (参考)
G06F 17/22 (2006.01)	G06F 17/22 652	5B109
	G06F 17/22 617	
	G06F 17/22 647	

審査請求 未請求 請求項の数 11 O L (全 22 頁)

(21) 出願番号	特願2016-4797 (P2016-4797)	(71) 出願人	00005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番1号
(22) 出願日	平成28年1月13日 (2016.1.13)	(74) 代理人	110002147 特許業務法人酒井国際特許事務所
		(72) 発明者	出内 将夫 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	片岡 正弘 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	田尾 幸資 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		Fターム(参考)	5B109 NH20 QA06 SA08

(54) 【発明の名称】 符号化プログラム、符号化方法、符号化装置、復号化プログラム、復号化方法および復号化装置

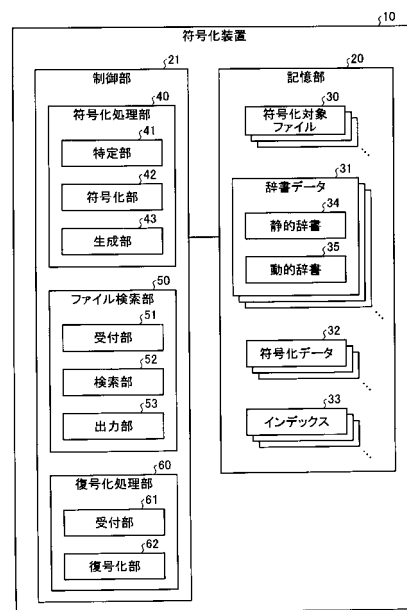
(57) 【要約】

【課題】 書構造に対応した符号化、復号化を行う符号化プログラム、符号化方法、符号化装置、復号化プログラム、復号化方法および復号化装置を提供する。

【解決手段】 特定部41は、構造化された第1の文書の文書構造を特定する。符号化部42は、文書構造を特定した第1の文書中の特定階層の文字列を、当該文書構造に対応した階層構造に応じた符号化方式により符号化する。

【選択図】 図3

符号化装置の構成の一例を示す図



- 【特許請求の範囲】
- 【請求項 1】
コンピュータに、
構造化された第 1 の文書の文書構造を特定し、
文書構造を特定した前記第 1 の文書中の特定階層の文字列を、当該文書構造に対応した階層構造に応じた符号化方式により符号化する
処理を実行させることを特徴とする符号化プログラム。
- 【請求項 2】
前記符号化する処理は、前記第 1 の文書中の文書構造を規定する文字列を、共通の符号化方式により符号化する
ことを特徴とする請求項 1 に記載の符号化プログラム。 10
- 【請求項 3】
前記符号化する処理は、データ属性が類似する階層の文字列を同じ符号化方式により符号化する
ことを特徴とする請求項 1 または 2 に記載の符号化プログラム。
- 【請求項 4】
前記符号化する処理は、前記特定階層の文字列を、当該特定階層に出現する文字列の特性に対応した符号化方式により符号化する
ことを特徴とする請求項 1 ~ 3 の何れか 1 つに記載の符号化プログラム。
- 【請求項 5】 20
前記符号化する処理は、1 またはデータ属性が類似する複数の階層ごとに、出現頻度の高いパターンを短い符号に変換する符号化方式により符号化する
ことを特徴とする請求項 1 ~ 4 の何れか 1 つに記載の符号化プログラム。
- 【請求項 6】
コンピュータに、
符号化方式ごとに、符号化した文字列に出現したパターンを示したインデックスを生成する
処理をさらに実行させることを特徴とする請求項 1 ~ 5 の何れか 1 つに記載の符号化プログラム。
- 【請求項 7】 30
コンピュータが、
構造化された第 1 の文書の文書構造を特定し、
文書構造を特定した前記第 1 の文書中の特定階層の文字列を、当該文書構造に対応した階層構造に応じた符号化方式により符号化する
処理を実行することを特徴とする符号化方法。
- 【請求項 8】
構造化された第 1 の文書の文書構造を特定する特定部と、
前記特定部により文書構造が特定された前記第 1 の文書中の特定階層の文字列を、当該文書構造に対応した階層構造に応じた符号化方式により符号化する符号化部と、
を有することを特徴とする符号化装置。 40
- 【請求項 9】
コンピュータに、
構造化された第 1 の文書が、当該第 1 の文書の文書構造に対応した階層構造に応じた符号化方式により符号化された符号化データの特定階層の文字列を、当該特定階層の符号化方式により復号化する
処理を実行させることを特徴とする復号化プログラム。
- 【請求項 10】
コンピュータが、
構造化された第 1 の文書が、当該第 1 の文書の文書構造に対応した階層構造に応じた符号化方式により符号化された符号化データの特定階層の文字列を、当該特定階層の符号化 50

方式により復号化する

処理を実行することを特徴とする復号化方法。

【請求項 1 1】

構造化された第 1 の文書が、当該第 1 の文書の文書構造に対応した階層構造に応じた符号化方式により符号化された符号化データの特定階層の文字列を、当該特定階層の符号化方式により復号化する復号化部

を有することを特徴とする復号化装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、符号化プログラム、符号化方法、符号化装置、復号化プログラム、復号化方法および復号化装置に関する。

【背景技術】

【0002】

従来から、例えば、XML (Extensible Markup Language) などの構造化された文書データが利用されている。例えば、異なるシステム間でデータを交換する共通フォーマットとして、XML が広く普及している。また、XML など構造化された形式で各種の文書データが公開される。この際、保存や通信の際のデータ量を減らすため、構造化された文書データは、例えば、zip などの圧縮形式で全体が圧縮されて保存される。圧縮された文書データを活用する場合は、圧縮された文書データ全体の伸長を行って文書データを復元し、復元した文書データに対して各種の解析が行われる。例えば、文書データの特定の階層に特定の単語を含むかを検索する場合は、復元した文書データに対して字句解析や構造解析が行われる。

【先行技術文献】

【特許文献】

【0003】

【特許文献 1】特開 2005 - 215951 号公報

【特許文献 2】特開 2002 - 297568 号公報

【特許文献 3】特開 2005 - 18672 号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、zip などの圧縮形式で全体が圧縮された文書データを活用する場合、圧縮された文書データ全体を伸長した後、各種の解析を行うため、処理量が多い。圧縮された文書データは、携帯端末など処理能力の低い端末で活用される場合もあり、活用する際の処理量が多いと、処理に時間がかかる。また、構造化された文書において、文書構造の一部のみを利用する場合でも、zip などの圧縮形式で全体が圧縮されている場合は、文書全体を伸長する。

【0005】

一つの側面では、文書構造に対応した符号化、復号化を行う符号化プログラム、符号化方法、符号化装置、復号化プログラム、復号化方法および復号化装置を提供することを目的とする。

【課題を解決するための手段】

【0006】

第 1 の案では、符号化プログラムは、コンピュータに、構造化された第 1 の文書の文書構造を特定する処理を実行させる。符号化プログラムは、コンピュータに、文書構造を特定した第 1 の文書中の特定階層の文字列を、当該文書構造に対応した階層構造に応じた符号化方式により符号化する処理を実行させる。

【発明の効果】

【0007】

10

20

30

40

50

本発明の1実施態様によれば、文書構造に対応した符号化を行うことができるという効果を奏する。

【図面の簡単な説明】

【0008】

【図1】図1は、符号化処理の流れを概略的に示した図である。

【図2A】図2Aは、検索処理の流れを概略的に示した図である。

【図2B】図2Bは、インデックスが生成されない場合の検索処理の流れを概略的に示した図である。

【図3】図3は、符号化装置の構成の一例を示す図である。

【図4】図4は、符号の割当ての一例を示す図である。

10

【図5】図5は、符号の割当ての一例を示す図である。

【図6】図6は、スキーマの概略的構成を示した図である。

【図7A】図7Aは、タグにより文書構造を示した文書の一例を示す図である。

【図7B】図7Bは、文書の一部にタグによりメタデータを付与した文書の一例を示す図である。

【図8A】図8Aは、符号化の一例を示す図である。

【図8B】図8Bは、符号化の一例を示す図である。

【図9】図9は、符号化の流れを概略的に示した図である。

【図10A】図10Aは、検索の一例を示した図である。

【図10B】図10Bは、検索の一例を示した図である。

20

【図11】図11は、符号化処理の手順の一例を示すフローチャートである。

【図12】図12は、検索処理の手順の一例を示すフローチャートである。

【図13】図13は、検索処理の手順の一例を示すフローチャートである。

【図14】図14は、復号化処理の手順の一例を示すフローチャートである。

【図15】図15は、符号の割当ての一例を示す図である。

【図16】図16は、符号化プログラムを実行するコンピュータの一例を示す図である。

【図17】図17は、検索プログラムを実行するコンピュータを示す図である。

【図18】図18は、復号化プログラムを実行するコンピュータの一例を示す図である。

【発明を実施するための形態】

【0009】

30

以下に、本願の開示する符号化プログラム、符号化方法、符号化装置、復号化プログラム、復号化方法および復号化装置の実施例を図面に基づいて詳細に説明する。なお、この実施例によりこの権利範囲が限定されるものではない。各実施例は、処理内容を矛盾させない範囲で適宜組み合わせることが可能である。

【実施例1】

【0010】

[符号化処理]

最初に、図1を用いて符号化処理の概要について説明する。図1は、符号化処理の流れを概略的に示した図である。以下では、構造化された文書が記憶された符号化対象ファイル30を符号化する場合を例に説明する。

40

【0011】

符号化対象ファイル30には、例えば、XMLにより、構造化された文書が記憶されている。XMLでは、テキストにより文書が記述され、タグにより要素が区切れ、構造化されている。図1の例では、病院で患者のカルテのデータをXMLにより、構造化された文書とした場合を示している。図1の例では、要素名「体温」のタグにより、体温が36.0と記録されている。また、要素名「概要」のタグにより、患者に関する概要「XXX・・・」が記録されている。なお、符号化対象ファイル30は、構造化された文書であれば、何れであってもよい。

【0012】

符号化対象ファイル30の符号化を行う場合、符号化装置10の符号化処理部40は、

50

符号化対象ファイル 30 に記憶された文書を読み出し、文書の文書構造を特定する（図 1（1））。符号化処理部 40 は、例えば、符号化対象ファイル 30 に対応して XML のスキーマ（Schema）が定義されている場合、符号化対象ファイル 30 に対応するスキーマから文書構造を特定してもよく、文書を解析して文書構造を特定してもよい。

【0013】

符号化処理部 40 は、読み出した文書を、文書構造に対応した階層構造に応じた符号化方式により符号化し、符号化したデータを符号化データ 32 として格納する（図 1（2））。

【0014】

例えば、符号化処理部 40 は、文書構造を規定する文字列を、共通の符号化方式により符号化する。図 1 の例では、文書構造を示すタグ「<体温>」を符号 A1 に符号化し、タグ「</体温>」を符号 A2 に符号化し、タグ「<概要>」を符号 A3 に符号化し、タグ「</概要>」を符号 A4 に符号化している。終了タグの符号は、開始タグの符号と別な符号としてもよく、タグの終了を示す符号と開始タグの符号を組み合わせてもよい。

10

【0015】

また、例えば、符号化処理部 40 は、階層ごとに、当該階層に出現する文字列の特性に対応した符号化方式により文字列を符号化する。符号化の際に割り当てる符号は、符号化方式ごとにユニークであればよい。このように、階層ごとの符号化方式により当該階層に出現する文字列を符号化する場合、階層ごとに、符号がユニークであればよい。このため、階層が異なると異なる文字列に同じ符号を割り当てることができる。図 1 の例では、「体温」の階層の文字列「36.0」を符号 B1 に符号化し、「概要」の階層の文字列「XXX・・・」の「XXX」を符号 B1 に符号化している。例えば、符号化処理部 40 は、階層ごとに、出現頻度の高い文字や単語などのパターンを短い符号に変換する符号化方式により文字列を符号化する。これにより、階層ごとに、文字列に含まれる各種のパターンのうち、出現頻度が高いパターンを短い符号に変換できるため、符号化対象ファイル 30 全体を短い符号に変換できる。

20

【0016】

符号化処理部 40 は、符号化方式ごとに、変換した文字列と当該文字列に対応する符号を辞書データ 31 に記憶する。図 1 の例では、文字列「36.0」と符号 B1 とが対応付けて辞書データ 31A に記憶され、文字列「XXX」と符号 B1 とが対応付けて辞書データ 31B に記憶されている。

30

【0017】

符号化処理部 40 は、符号化方式ごとに、符号化した文字列に出現したパターンを示したインデックス 33 を生成する（図 1（3））。インデックスとは、パターンが含まれるファイルを示したデータである。例えば、インデックスには、パターンおよびファイルに 1 つのビットが対応付け、ビットの値により、パターンが出現したか否かを記憶するビットマップ型のインデックスがある。また、インデックスには、パターンおよびファイルに複数のビットが対応付け、複数のビットによりパターンの出現回数の情報を保持するカウントマップ型のインデックスがある。図 1 の例では、符号化処理部 40 は、カウントマップ型のインデックス 33A、33B を生成する。インデックス 33A は、「体温」の階層の文字列に出現したパターンの出現回数の情報を保持する。インデックス 33B は、「概要」の階層の文字列に出現したパターンの出現回数の情報を保持する。図 1 の例では、インデックス 33A、33B には、符号化対象ファイル 30 のファイル番号「1」および符号 B1 に対応付けて、出現回数それぞれ複数ビットで記憶されている。なお、本実施例では、符号化処理部 40 は、符号化の際にインデックス 33A、33B を生成する場合を説明するが、これらに限定されるものではなく、適宜変更可能である。例えば、符号化処理部 40 は、インデックス 33A、33B を生成しなくてもよい。

40

【0018】

[検索処理]

次に、図 2A を用いて、実施例 1 にかかる符号化装置 10 が実施する検索処理の概要に

50

ついて説明する。図 2 A は、検索処理の流れを概略的に示した図である。図 2 A の例では、図 1 により符号化された符号化データ 3 2 と、辞書データ 3 1 A、3 1 B と、インデックス 3 3 A、3 3 B が示されている。なお、図 2 A の例では、圧縮データ 3 2 に符号化された文字列を識別しやすくするため、符号の後に括弧記号「()」で囲んで符号化された文字列を記載している。

【 0 0 1 9 】

符号化装置 1 0 のファイル検索部 5 0 は、検索条件の入力を受け付ける。例えば、図 2 A の例では、ファイル検索部 5 0 は、階層「概要」、文字列「X X X」との検索条件を受け付ける。

【 0 0 2 0 】

ファイル検索部 5 0 は、検索条件を満たすファイルを検索する。例えば、ファイル検索部 5 0 は、階層「概要」の文字列を変換した際の辞書データ 3 1 B を参照して、文字列「X X X」に対応する符号 B 1 を特定する(図 2 A (1))。ファイル検索部 5 0 は、階層「概要」の文字列を変換した際に生成したインデックス 3 3 B を参照して、符号 B 1 が出現したファイルのファイル番号を特定する(図 2 A (2))。図 2 A の例では、インデックス 3 3 B に、符号化対象ファイル 3 0 のファイル番号「1」および符号 B 1 に対応付けて、出現回数が記憶されているため、ファイル番号「1」の符号化対象ファイル 3 0 が検索条件を満たすと検索される。このように、符号化装置 1 0 は、符号化された符号化データ 3 2 に対して文字列の検索を行う場合、符号化データ 3 2 を復号化することなく文字列を検索できるため、活用する際の処理量を減らすことができる。

【 0 0 2 1 】

なお、上述したように、インデックス 3 3 A、3 3 B は、必ずしも生成されなくてもよい。図 2 B は、インデックスが生成されない場合の検索処理の流れを概略的に示した図である。図 2 B の例では、図 1 により符号化された符号化データ 3 2 と、辞書データ 3 1 が示されている。なお、図 2 B の例でも、符号化データ 3 2 に符号化された文字列を識別しやすくするため、符号の後に括弧記号「()」で囲んで符号化された文字列を記載している。

【 0 0 2 2 】

ファイル検索部 5 0 は、検索条件の入力を受け付ける。例えば、図 2 B の例では、ファイル検索部 5 0 は、階層「概要」、文字列「X X X」との検索条件を受け付ける。

【 0 0 2 3 】

ファイル検索部 5 0 は、検索条件を満たすファイルを検索する。例えば、ファイル検索部 5 0 は、共通の符号化方式により符号化されたタグを復号化する。そして、ファイル検索部 5 0 は、ファイル検索部 5 0 は、階層「概要」の文字列を変換した際の辞書データ 3 1 B を参照して、階層「概要」の部分の符号を復号化する(図 2 B (1))。そして、ファイル検索部 5 0 は、復号化された部分から文字列「X X X」を検索する(図 2 B (2))。この場合でも、ファイル検索部 5 0 は、階層「概要」の部分の符号を復号化するのみで検索を行えるため、符号化データ全体を復号化する場合と比較して、活用する際の処理量を減らすことができる。

【 0 0 2 4 】

[装置構成]

次に、符号化装置 1 0 の構成について説明する。図 3 は、符号化装置の構成の一例を示す図である。符号化装置 1 0 は、構造化された文書の圧縮などの符号化を行う装置である。符号化装置 1 0 は、例えば、パーソナルコンピュータ、サーバコンピュータなどのコンピュータや、タブレット端末、スマートフォンなどの情報処理装置である。符号化装置 1 0 は、1 台のコンピュータとして実装してもよく、また、複数台のコンピュータによるクラウドとして実装することもできる。なお、本実施例では、符号化装置 1 0 を 1 台のコンピュータとした場合を例として説明する。図 3 に示すように、符号化装置 1 0 は、記憶部 2 0 と、制御部 2 1 とを有する。なお、符号化装置 1 0 は、コンピュータや情報処理装置が有する上記の機器以外の他の機器を有してもよい。また、本実施例では、符号化装置 1

10

20

30

40

50

0により符号化およびファイル検索を行う場合を例として説明するが、符号化とファイル検索は別な装置で行ってもよい。

【0025】

記憶部20は、ハードディスク、SSD(Solid State Drive)、光ディスクなどの記憶装置である。なお、記憶部20は、RAM(Random Access Memory)、フラッシュメモリ、NVS RAM(Non Volatile Static Random Access Memory)などのデータを書き換え可能な半導体メモリであってもよい。

【0026】

記憶部20は、制御部21で実行されるOS(Operating System)や各種プログラムを記憶する。例えば、記憶部20は、後述する符号化処理や検索処理を行うプログラムを記憶する。さらに、記憶部20は、制御部21で実行されるプログラムで用いられる各種データを記憶する。例えば、記憶部20は、符号化対象ファイル30と、辞書データ31と、符号化データ32と、インデックス33とを記憶する。

【0027】

符号化対象ファイル30は、符号化対象のテキストデータが記憶されたデータである。例えば、符号化対象ファイル30には、XMLにより、構造化された文書が記憶されている。

【0028】

辞書データ31は、データの符号化および復号化に用いる辞書のデータである。

【0029】

ここで、本実施例では、構造化された文書を符号化する際に、構造や属性に応じて符号化方式を切り替える。辞書データ31は、辞書を用いて符号化する符号化方式で用いる辞書のデータである。辞書データ31は、辞書を用いて符号化する符号化方式ごとに設けられる。例えば、辞書データ31は、階層化された文書の階層のうち、辞書を用いて符号化を行う階層ごと、または、辞書を用いて符号化を行う階層でデータ属性が類似する階層ごとに、設けられている。辞書データ31は、静的辞書34と、動的辞書35とを有する。

【0030】

静的辞書34は、文書の構造や属性に応じて出現頻度の高いパターンに対応する符号を保持したデータである。動的辞書35は、文書の構造や属性に応じて出現頻度の低いパターンに対応する符号を保持したデータである。静的辞書34は、予め設けられる。動的辞書35は、必要に応じて動的に生成される。

【0031】

静的辞書34は、対応する階層に出現する文字列の特性に対応して、文字列に対応する符号が記憶されている。例えば、静的辞書34は、対応する階層に標準的に出現する文字列や数字などのパターンに対応する符号が記憶されている。また、静的辞書34は、対応する階層で出現頻度の高いパターンに短い符号が対応付けて記憶されている。例えば、人間の体温は、通常、35.0 ~ 42.0 の範囲に収まり、36.0 前後の頻度が高い。そこで、例えば、体温の階層に対応する静的辞書34には、35.0 ~ 42.0 の数値に対して符号が対応付けて記憶されており、36.0 前後に対して短い符号が割り当てられて記憶されている。また、本実施例では、概要に出現する文字列を単語の単位で符号化する。例えば、本実施例では、単語を、一般的な文書を解析して、出現頻度が相対的に高い高頻度単語と、出現頻度が相対的に低い低頻度単語とに分けている。例えば、出現頻度の高い順に所定の順位までの基礎単語を高頻度単語とし、所定の順位以降の基礎単語を低頻度単語とする。高頻度単語については、短い符号を予め割り当てて、割り当てた符号と高頻度単語を対応付けて静的辞書34に記憶させる。例えば、高頻度単語については、予め2バイト(16ビット)の符号を割り当て、割り当てた符号を静的辞書34に予め記憶させる。低頻度単語については、出現した際に符号を動的に割り当てて、割り当てた符号を動的辞書35に記憶させる。すなわち、符号は、高頻度単語については予め登録され、低頻度単語については動的に割り当てられて動的辞書35に記憶される。なお、概要に出現する文字列や数字などのパターンが特定のパターンに定まる場合は、概要の階層に対

10

20

30

40

50

応する静的辞書 3 4 には、特定のパターンと符号を対応付けて予め記憶させてもよい。

【 0 0 3 2 】

動的辞書 3 5 は、対応する階層に出現する文字列の特性に対応して、動的に割り当てられた符号に関する各種の情報を保持したデータである。例えば、概要の階層に対応した動的辞書 3 5 には、低頻度単語など出現頻度の低いパターンに動的に割り当てられた符号が記憶される。

【 0 0 3 3 】

図 4 は、符号の割当ての一例を示す図である。図 4 には、2 バイト (1 6 ビット) の符号に対する割当ての一例が示されている。上部の横方向の項目は、最初の 1 バイト目を 0 ~ F の 1 6 進表記で示しており、「 * 」は、2 バイト目を示している。例えば、「 1 * h 10
」は、1 バイト目が 2 進数表記で「 0 0 0 0 0 0 0 1 」であることを示す。左側の縦方向の項目は、2 バイト目を 0 ~ F の 1 6 進表記で示しており、「 * 」は、1 バイト目を示している。例えば、「 * 2 h 」は、2 バイト目が 2 進数表記で「 0 0 0 0 0 0 1 0 」であることを示す。

【 0 0 3 4 】

図 4 では、縦方向の項目と横方向の項目に対応する領域に、符号に対応させるパターンを示す。例えば、「 0 * h 」、 「 1 * h 」の符号については、各階層とも、同じ制御コードに同じ符号を対応付けている。また、「 2 * h 」 ~ 「 5 * h 」の符号については、各階層とも、同じタグに同じ符号を対応付けている。また、「 6 * h 」 ~ 「 F * h 」の符号については、各階層でそれぞれ個別にパターンに符号を割当て可能としている。例えば、文字列を単語の単位で符号化する場合、「 6 * h 」 ~ 「 9 * h 」の符号については、予め定めた高頻度単語に対して割当てている。「 A * h 」 ~ 「 F * h 」の符号については、低頻度単語が出現した際に符号を動的に割り当てる。「 E * h 」、 「 F * h 」は、符号の不足に対応するため、3 バイトの符号としている。 20

【 0 0 3 5 】

辞書データ 3 1 は、辞書を用いて符号化を行う階層ごと、または、辞書を用いて符号化を行う階層でデータ属性が類似する階層ごとに設けられ、「 6 * h 」 ~ 「 F * h 」の符号について、階層に出現する文字列の特性に対応して文字列と符号を対応付けて記憶する。

【 0 0 3 6 】

なお、辞書データ 3 1 は、タグに対して動的に符号を割当て可能としてもよい。図 5 は、符号の割当ての一例を示す図である。図 5 の例では、1 バイト目が「 5 * h 」の符号について、特定階層のタグとして動的に符号を割当て可能としている。 30

【 0 0 3 7 】

図 3 に戻り、符号化データ 3 2 は、符号化対象ファイル 3 0 をそれぞれ符号化したデータである。インデックス 3 3 は、符号化した文字列に出現したパターンの出現回数を記憶したデータである。例えば、インデックス 3 3 は、符号化方式ごとに設けられ、符号化した文字列に出現したパターンの出現回数と出現したファイルのファイル番号を対応付けて記憶される。

【 0 0 3 8 】

制御部 2 1 は、符号化装置 1 0 を制御するデバイスである。制御部 2 1 としては、C P U (Central Processing Unit)、M P U (Micro Processing Unit) 等の電子回路や、A S I C (Application Specific Integrated Circuit)、F P G A (Field Programmable Gate Array) 等の集積回路を採用できる。制御部 2 1 は、各種の処理手順を規定したプログラムや制御データを格納するための内部メモリを有し、これらによって種々の処理を実行する。制御部 2 1 は、各種のプログラムが動作することにより各種の処理部として機能する。例えば、制御部 2 1 は、符号化処理部 4 0 と、ファイル検索部 5 0 と、復号化処理部 6 0 とを有する。 40

【 0 0 3 9 】

符号化処理部 4 0 は、符号化対象ファイル 3 0 に記憶された構造化された文書を読み出し、読み出した文書を、文書構造に対応した階層構造に応じた符号化方式により符号化し 50

た符号化データ32を生成する。符号化処理部40は、特定部41と、符号化部42と、生成部43とを有する。

【0040】

特定部41は、各種の特定を行う。例えば、特定部41は、符号化対象ファイル30に記憶されたXMLの文書の文書構造を特定する。例えば、特定部41は、符号化対象ファイル30に対応してXMLのスキーマが定義されている場合、符号化対象ファイル30に対応するスキーマから文書構造を特定する。

【0041】

図6は、スキーマの概略的構成を示した図である。XMLの文書では、文書構造を示すXMLスキーマ70が定義される。XMLスキーマ70には、XMLの文書の文書構造の定義や、末端要素の型や制約の定義が、スキーマ言語により記述されている。図6の例では、構造定義として、文書構造を示すタグの入れ子の関係や、タグの制約などが記述される。また、図6の例では、末端要素の型や制約として、格納される文字列のデータ型、数値の最大、最小、string(文字列)の長さ、使える文字、stringを、例えば、Male、Female等の選択型として使用しているかが記述される。符号化対象ファイル30には、XMLスキーマ70の定義に対応して、XMLで文書が記憶される。図6の例では、XMLの文書には、文字コードが1行目に記載され、XMLスキーマ70に対応した文書構造で文書が記憶される。なお、XMLスキーマ70は、文書構造を柔軟に定義でき、符号化対象ファイル30ごとにタグの個数を変更可能な定義もできる。例えば、符号化対象ファイル30Aは、Xタグの配下にYタグが10個存在し、符号化対象ファイル30Bは、Xタグの配下にYタグが20個存在する文書構成とすることもできる。

10

20

【0042】

特定部41は、符号化対象ファイル30に対応してXMLスキーマ70が定義されている場合、XMLスキーマ70から文書構造を特定する。なお、特定部41は、符号化対象ファイル30に記憶された文書を解析して文書構造を特定してもよい。

【0043】

符号化部42は、符号化対象ファイル30に記憶された文書の符号化を行う。例えば、符号化部42は、特定部41により文書構造を特定したXMLの文書を符号化対象ファイル30から読み出す。そして、符号化部42は、読み出した文書を、文書構造に対応した階層構造に応じた符号化方式により符号化する。例えば、符号化部42は、読み出した文書に出現するタグに対して順に符号を割当てて、符号化する。なお、文書構造に出現するタグが定まっている場合、タグと符号を対応付けたタグ用の辞書データを予め記憶させ、符号化部42は、タグ用の辞書データを用いて、読み出した文書に出現するタグを符号化してもよい。また、文書構造に出現する頻度の高い一部のタグについてタグ用の辞書データに記憶させ、符号化部42は、頻度の高い一部のタグをタグ用の辞書データを用いて符号化し、それ以外のタグに順に符号を割当てて符号化してもよい。

30

【0044】

ここで、構造化された文書には、タグにより文書の要素を区分けして文書構造を示した文書と、文書の一部にタグによりメタデータを付与した文書がある。

【0045】

図7Aは、タグにより文書構造を示した文書の一例を示す図である。図7Aの例では、例1の文書は、タグにより「概要」、「本文」が定義されている。また、例1の文書は、タグにより「概要」、「本文」と区分けされた部分にそれぞれ内容に応じた文字列(テキスト)が記憶された文書が示されている。例2の文書は、タグにより「特許」が定義され、「特許」の下層に「名称」、「課題」、「効果」が定義されている。また、例2の文書は、タグにより「名称」、「課題」、「効果」と区分けされた部分にそれぞれ内容に応じた文字列が記憶された文書が示されている。

40

【0046】

符号化部42は、タグを共通の符号化方式により符号化する。例1の文書は、「概要」、「本文」のタグを共通の符号化方式により符号化する。例2の文書は、「特許」、「名

50

称」、「課題」、「効果」のタグを共通の符号化方式により符号化する。

【0047】

また、符号化部42は、タグにより分けられた部分の文字列を、それぞれ階層に応じた符号化方式により符号化する。例えば、符号化部42は、タグにより分けられた部分の文字列を、それぞれ階層に対応した辞書データ31を用いて符号化する。例えば、符号化部42は、文字列に出現した単語が、階層に対応した辞書データ31の静的辞書34または動的辞書35に登録されている場合、出現した単語を静的辞書34または動的辞書35に登録された符号に符号化する。また、符号化部42は、文字列に出現した単語が、階層に対応した辞書データ31の静的辞書34または動的辞書35に登録されていない場合、符号を動的に割り当て、出現した単語を割り当てた符号に符号化する。符号化部42は、出現した単語と割り当てた符号を対応付けて動的辞書35に登録する。これにより、以降、動的辞書35に登録された単語は、出現した際に、動的辞書35を用いて同じ符号に符号化される。なお、符号化部42は、データ属性が類似する階層の文字列を同じ符号化方式により符号化してもよい。これにより、符号化部42は、データ属性が類似する階層の文字列を同じ辞書データ31で符号化できる。

10

【0048】

図7Bは、文書の一部にタグによりメタデータを付与した文書の一例を示す図である。図7Bの例では、例3の文書は、「AAAへのリンクはこちら」の文書の「リンク」部分にタグによりリンク先のURLをメタデータとして付与した場合を示している。例4の文書は、「BBBを訴えたため、CCCを疑い、DDDを投与した」の文書の「BBB」部分が病状を示し、「CCC」部分が病名を示し、「DDD」部分が薬名を示すことをタグによりメタデータとして付与した場合を示している。例5の文書は、「2015/3/6に、大阪で鈴木に会う」の文書の「2015/3/6」部分が日時を示し、「大阪」部分が地名を示し、「鈴木」部分が人名を示すことをタグによりメタデータとして付与した場合を示している。

20

【0049】

符号化部42は、タグを共通の符号化方式により符号化する。例3の文書は、「リンク」のタグを共通の符号化方式により符号化する。例4の文書は、「病状」、「病名」、「薬名」のタグを共通の符号化方式により符号化する。例5の文書は、「日時」、「地名」、「人名」のタグを共通の符号化方式により符号化する。また、符号化部42は、タグにより分けられた部分の文字列を、それぞれ階層に応じた符号化方式により符号化する。例えば、符号化部42は、タグにより分けられた部分の文字列を、それぞれ階層に対応した辞書データ31を用いて符号化する。

30

【0050】

図8Aは、符号化の一例を示す図である。図8Aの例は、タグ「A」の下層にタグ「B」が定義された文字列のデータを符号化した一例を示している。図8Aの例では、タグ「A」、タグ「B」の符号の間に、文字列のデータを符号化したコードが記憶される。なお、図8Aの例では、タグ「A」、タグ「B」の終了タグの符号を、タグの終了を示す符号と開始タグの符号を組み合わせたものとしている。

40

【0051】

図8Bは、符号化の一例を示す図である。図8Bの例は、「大阪で鈴木に会う」の文書の「大阪」部分が地名を示し、「鈴木」部分が人名を示すことをタグによりメタデータとして付与された文字列のデータを符号化した一例を示している。図8Bの例では、地名の開始符号「25h」と終了符号「20h」、「25h」の間で、「大阪」が「B0h」と符号されている。また、人名の開始符号「26h」と終了符号「20h」、「26h」の間で、「鈴木」が「B0h」と符号されている。

【0052】

符号化部42は、階層が異なると異なる文字列に同じ符号を割り当てることのできるため、階層ごとに文字列を短い符号に変換できる。例えば、図8Bの例では、「大阪」と「鈴木」が共に同じ「B0h」に変換されている。このように、符号化部42は、階層ごとに

50

、文字列を短い符号に変換できるため、符号化対象ファイル 30 全体を短い符号に変換できる。

【0053】

なお、符号化部 42 は、タグにより分けられた部分の文字列の属性や範囲によっては、辞書データ 31 を用いずに当該文字列を符号化してもよい。例えば、タグにより分けられた部分の文字列が「0」～「255」の範囲の数値を示す文字列である場合、符号化部 42 は、「0」～「255」の範囲の数値を示す文字列を 1 バイトの整数型（例えば、int 型）の符号に符号化してもよい。すなわち、符号化部 42 は、文字列が数値を示す場合、当該数値の範囲に対応したデータ型の符号に符号化してもよい。数値を表す文字列を数値のデータ型の符号に符号化すると、符号化した状態でも数値の比較や集計など各種の演算を行うことができる。

10

【0054】

ここで、XML など構造化された文書では、タグによってコンテキストが規定される。構造化された文書は、タグによってコンテキストが規定され、コンテキストによりデータの処理に関わる要素が定まる。例えば、データの型や値の範囲、文書の構成要素（言語であれば日本語の単語、英語の単語、他言語の単語）など辞書に関わる要素が決まる。また、例えば、テキストであれば検索やマイニング、数値であれば平均値や集計値、頻度分布など、データの内容がどのように活用が可能かの活用分野が定まる。また、図 7A に示すような、タグにより文書構造を示した文書では、単一タグだけでなく、階層構造を上位から順に辿ると判明するコンテキストがある。例えば、<A><合計>T</合計>では、「T」が A に関する合計を示している。また、例えば、図 7B に示すような、文書の一部にタグによりメタデータを付与した文書では、単一タグで囲まれた領域で、階層構造に加え、追加されるコンテキストがある。例えば、<地名>大阪</地名>は、「大阪」が地名であることを示している。よって、符号化部 42 は、タグにより分けられた部分の文字列を、タグによって規定されるコンテキストに適した符号化方式で符号化することにより、活用する際の処理量を減らすことができる。

20

【0055】

符号化部 42 は、符号化対象ファイル 30 に記憶された文書の符号化したデータを符号化データ 32 として格納する。

【0056】

生成部 43 は、符号化方式ごとに、符号化した文字列に出現したパターンを示したインデックス 33 を生成する。例えば、生成部 43 は、符号化した符号化対象ファイル 30 に対して、順にファイル番号を付与する。そして、生成部 43 は、符号化した符号化対象ファイル 30 のファイル番号に対応付けて、符号化対象ファイル 30 に出現した数値や単語などのパターンの出現回数を格納したインデックス 33 を生成する。

30

【0057】

ここで、符号化の流れを説明する。図 9 は、符号化の流れを概略的に示した図である。符号化装置 10 の符号化処理部 40 は、符号化対象ファイル 30 に記憶された文書を読み出し、文書の文書構造を特定する。符号化処理部 40 は、読み出した文書を、文書構造に対応した階層構造に応じた符号化方式により符号化する。例えば、符号化処理部 40 は、タグや文字列に出現した単語が、静的辞書 34 または動的辞書 35 に登録されている場合、出現した単語を静的辞書 34 または動的辞書 35 に登録された符号に符号化する。また、符号化処理部 40 は、タグや文字列に出現した単語が、静的辞書 34 または動的辞書 35 に登録されていない場合、符号を動的に割り当て、タグや出現した単語を割り当てた符号に符号化する。符号化処理部 40 は、タグや出現した単語と割り当てた符号を対応付けて動的辞書 35 に登録する。

40

【0058】

符号化処理部 40 は、符号化対象ファイル 30 に記憶された文書の符号化したデータを符号化データ 32 として格納する。図 9 の例では、「概要」と「本文」のタグの階層の文字列がそれぞれ符号化されている。また、符号化処理部 40 は、符号化対象ファイル 30

50

のファイル番号に対応付けて、符号化対象ファイル30に出現した数値や単語などのパターンの出現回数を格納したインデックス33を生成する。図9の例では、「概要」と「本文」のタグの階層に対応して、出現回数の集計結果としてインデックス33A、33Bが生成されている。符号化装置10は、符号化データ32を他の装置へ移動させる場合、符号化データ32に対応して生成された動的辞書35およびインデックス33A、33Bも移動させる。

【0059】

図3に戻り、ファイル検索部50は、指定された検索条件に従い、ファイルを検索する。ファイル検索部50は、受付部51と、検索部52と、出力部53とを有する。以下、ファイル検索部50の各構成について詳細に説明する。

10

【0060】

受付部51は、検索条件を受け付ける。例えば、受付部51は、検索条件とするキーワードや階層の入力を受け付ける操作画面などの入力インタフェースを提供しており、検索条件とする文字列や階層の入力を受け付ける。

【0061】

検索部52は、検索条件を満たすファイルを検索する。例えば、検索部52は、検索条件の階層に対応する辞書データ31の静的辞書34および動的辞書35を参照して、検索条件のキーワードに対応する符号を特定する。そして、検索部52は、検索条件の階層に対応するインデックス33を参照して、特定された符号が出現したファイルのファイル番号を特定する。なお、検索条件のキーワードが単語や数値を複数含む場合、検索部52は、キーワードを単語や数値に分解して符号化し、それぞれの単語や数値ごとに対応する符号を特定する。検索部52は、検索条件の階層に対応するインデックス33を参照して、それぞれの単語や数値ごとに対応する符号が出現したファイルのファイル番号を特定する。ここで、インデックス33では、検索条件の文字列に含まれる複数の単語や数値の出現順が正しいかを確認できない。そこで、例えば、検索部52は、特定したファイル番号の符号化対象ファイル30に検索条件の文字列が含まれるかを検索する。なお、検索部52は、特定したファイル番号に対応する符号化データ32の検索条件の階層を復号化して、検索条件の文字列が含まれるかを検索してもよい。

20

【0062】

図10Aは、検索の一例を示した図である。図10Aの例は、指定されたファイルが「概要」に「XXX」というキーワードを含み、「本文」に「YYY」というキーワードを含むかを検索する場合を示している。検索部52は、「概要」の階層に対応する辞書データ31の静的辞書34および動的辞書35を参照して、「XXX」に対応する符号を特定する。検索部52は、「概要」の階層に対応するインデックス33を参照して、指定されたファイルのファイル番号に「XXX」に対応する符号が出現したことが記録されているかを特定する。また、検索部52は、「本文」の階層に対応する辞書データ31の静的辞書34および動的辞書35を参照して、「YYY」に対応する符号を特定する。そして、検索部52は、「本文」の階層に対応するインデックス33を参照して、指定されたファイルのファイル番号に「YYY」に対応する符号が出現したことが記録されているかを特定する。そして、検索部52は、指定されたファイルのファイル番号に「XXX」に対応する符号と「YYY」に対応する符号が出現した記録がある場合、「概要」に「XXX」というキーワードを含み、「本文」に「YYY」というキーワードを含むかを検索する。

30

40

【0063】

図10Bは、検索の一例を示した図である。図10Bの例は、「概要」に「ZZZ」というキーワードを含むファイルを検索する場合を示している。検索部52は、「概要」の階層に対応する辞書データ31の静的辞書34および動的辞書35を参照して、「ZZZ」に対応する符号を特定する。そして、検索部52は、「概要」の階層に対応するインデックス33を参照して、「ZZZ」に対応する符号が出現したファイルのファイル番号を特定する。

【0064】

50

このように、ファイル検索部 50 は、符号化データ 32 を復号化せずに検索を行えるため、検索の際の処理量が減らすことができ、検索の処理時間を短縮できる。

【0065】

なお、インデックス 33 が生成されない場合、ファイル検索部 50 は、指定された階層のみ復号化して、指定された文字列を検索する。この場合でも、ファイル検索部 50 は、指定された階層を復号化するのみで検索を行えるため、符号化データ全体を復号化する場合と比較して、活用する際の処理量を減らすことができ、検索の処理時間を短縮できる。

【0066】

出力部 53 は、検索結果の出力を行う。例えば、検索部 52 によりファイル番号が特定された場合、出力部 53 は、検索結果として、特定されたファイル番号のファイルのファイル名を出力する。一方、検索部 52 によりファイル番号が特定されない場合、出力部 53 は、検索結果として、該当ファイルなしを出力する。

【0067】

図 3 に戻り、復号化処理部 60 は、符号化データ 32 を復号化する。復号化処理部 60 は、受付部 61 と、復号化部 62 とを有する。以下、復号化処理部 60 の各構成について詳細に説明する。

【0068】

受付部 61 は、復号化の指示を受け付ける。例えば、受付部 61 は、復号化する対象の符号化データ 32 の指定を受け付ける操作画面などの入力インタフェースを提供しており、復号化する対象の符号化データ 32 の指定を受け付ける。なお、受付部 61 は、復号化する対象の符号化データ 32 と共に、符号化する階層の指定を受け付けてもよい。

【0069】

復号化部 62 は、指定された符号化データ 32 を復号化する。例えば、復号化部 62 は、復号化部 62 は、符号化データ 32 のそれぞれの階層の符号データを、当該階層の符号化方式により復号化する。例えば、復号化部 62 は、符号化データ 32 のそれぞれの階層の符号データを、当該階層に対応する辞書データ 31 の静的辞書 34 および動的辞書 35 を用いて、文字列に復号化する。例えば、復号化部 62 は、タグの符号データを、共通の符号化方式により復号化する。そして、復号化部 62 は、タグで区切られた各階層の符号データを、当該階層に対応する辞書データ 31 の静的辞書 34 および動的辞書 35 を参照して、文字列に復号化する。なお、受付部 61 で符号化する階層の指定を受け付けた場合、復号化部 62 は、指定された階層の符号データのみを復号化してもよい。

【0070】

[処理の流れ]

本実施例に係る符号化装置 10 が符号化対象ファイル 30 を符号化する符号化処理の流れについて説明する。図 11 は、符号化処理の手順の一例を示すフローチャートである。この符号化処理は、所定のタイミング、例えば、符号化対象ファイル 30 を指定して符号化開始を指示する所定操作が行われたタイミングで実行される。

【0071】

図 11 に示すように、特定部 41 は、符号化対象ファイル 30 に記憶された構造化された文書の文書構造を特定する (S10)。符号化部 42 は、文書構造を特定した文書の各階層の文字列を、当該文書構造に対応した階層構造に応じた符号化方式により符号化する (S11)。例えば、符号化部 42 は、タグを共通の符号化方式により符号化する。また、符号化部 42 は、タグにより分けられた部分の文字列を、それぞれ階層に応じた符号化方式により符号化する。符号化部 42 は、符号化したデータを符号化データ 32 に格納する (S12)。生成部 43 は、符号化方式ごとに、符号化した文字列に出現したパターンを示したインデックス 33 を生成し (S13)、処理を終了する。

【0072】

次に、本実施例に係る符号化装置 10 が検索条件を満たすファイルを検索する検索処理の流れについて説明する。最初に、検索条件に階層が指定されない場合の検索処理の流れを説明する。図 12 は、検索処理の手順の一例を示すフローチャートである。この検索処

10

20

30

40

50

理は、所定のタイミング、例えば、検索条件を指定して検索開始を指示する所定操作が行われたタイミングで実行される。

【0073】

図12に示すように、検索部52は、辞書データ31の静的辞書34および動的辞書35を参照して、検索条件のキーワードに対応する符号が存在するか判定する(S20)。符号が存在しない場合(S20否定)、検索部52は、キーワードを単語や数値に分解してそれぞれ符号化し、それぞれの単語や数値ごとに対応する符号を特定する(S21)。検索部52は、各インデックス33を参照して、それぞれの単語や数値ごとに対応する符号が出現したファイルのファイル番号を特定する(S22)。検索部52は、特定したファイル番号の符号化対象ファイル30に検索条件の文字列が含まれるかを検索する(S23)。

10

【0074】

一方、符号が存在する場合(S20肯定)、検索部52は、インデックス33を参照して、特定された符号が出現したファイルのファイル番号を特定する(S24)。

【0075】

出力部53は、検索結果を出力し、処理を終了する(S25)。例えば、出力部53は、検索条件の文字列を含む符号化対象ファイル30が検索された場合や、検索部52により符号化対象ファイル30のファイル番号が特定された場合、符号化対象ファイル30のファイル名を出力する。

【0076】

次に、検索条件に階層が指定された場合の検索処理の流れを説明する。図13は、検索処理の手順の一例を示すフローチャートである。この検索処理は、所定のタイミング、例えば、検索条件を指定して検索開始を指示する所定操作が行われたタイミングで実行される。

20

【0077】

図13に示すように、検索部52は、辞書データ31の静的辞書34および動的辞書35を参照して、検索条件のキーワードに対応する符号が存在するか判定する(S30)。符号が存在しない場合(S30否定)、検索部52は、キーワードを単語や数値に分解してそれぞれ符号化し、それぞれの単語や数値ごとに対応する符号を特定する(S31)。検索部52は、指定された階層のインデックス33を参照して、それぞれの単語や数値ごとに対応する符号が出現したファイルのファイル番号を特定する(S32)。検索部52は、特定したファイル番号の符号化対象ファイル30に検索条件の文字列が含まれるかを検索する(S33)。

30

【0078】

一方、符号が存在する場合(S30肯定)、検索部52は、指定された階層のインデックス33を参照して、特定された符号が出現したファイルのファイル番号を特定する(S34)。

【0079】

出力部53は、検索結果を出力し、処理を終了する(S35)。例えば、出力部53は、検索条件の文字列を含む符号化対象ファイル30が検索された場合や、検索部52により符号化対象ファイル30のファイル番号が特定された場合、符号化対象ファイル30のファイル名を出力する。

40

【0080】

次に、本実施例に係る符号化装置10が符号化データ32を復号化する復号化処理の流れについて説明する。図14は、復号化処理の手順の一例を示すフローチャートである。この復号化処理は、所定のタイミング、例えば、復号化する対象の符号化データ32を指定して符号化開始を指示する所定操作が行われたタイミングで実行される。

【0081】

復号化部62は、指定された符号化データ32から符号データを読み出す(S40)。復号化部62は、読み出した符号データを、階層に対応する辞書データ31の静的辞書3

50

4 および動的辞書 3 5 を用いて、文字列に復号化する (S 4 1)。復号化部 6 2 は、符号化データ 3 2 の読み出しが完了したか否かを判定する (S 4 2)。読み出しが完了していない場合は (S 4 2 否定)、S 4 0 へ移行する。一方、読み出しが完了した場合は (S 4 2 肯定)、処理を終了する。

【 0 0 8 2 】

[効果]

上述してきたように、本実施例に係る符号化装置 1 0 は、構造化された文書の文書構造を特定する。符号化装置 1 0 は、文書構造を特定した文書中の特定階層の文字列を、当該文書構造に対応した階層構造に応じた符号化方式により符号化する。これにより、符号化装置 1 0 は、特定階層の部分の符号のみを復号化できるため、活用する際の処理量を減らすことができる。

10

【 0 0 8 3 】

また、本実施例に係る符号化装置 1 0 は、文書中の文書構造を規定する文字列を、共通の符号化方式により符号化する。これにより、符号化装置 1 0 は、共通の符号化方式で復号化することで、文書中の文書構造を規定する文字列を同じ符号化方式で復元できるため、文書構造を速やかに特定でき、特定の階層のデータを速やかに抽出できる。

【 0 0 8 4 】

また、本実施例に係る符号化装置 1 0 は、データ属性が類似する階層の文字列を同じ符号化方式により符号化する。これにより、符号化装置 1 0 は、データ属性が類似する階層の文字列を同じ辞書データ 3 1 で符号化できる。

20

【 0 0 8 5 】

また、本実施例に係る符号化装置 1 0 は、特定階層の文字列を、当該特定階層に出現する文字列の特性に対応した符号化方式により符号化する。これにより、符号化装置 1 0 は、特定階層の文字列を特性に対応した符号化方式で符号化できる。

【 0 0 8 6 】

また、本実施例に係る符号化装置 1 0 は、1 またはデータ属性が類似する複数の階層ごとに、出現頻度の高いパターンを短い符号に変換する符号化方式により符号化する。これにより、符号化装置 1 0 は、符号化対象ファイル 3 0 を高い圧縮率で符号化できる。

【 0 0 8 7 】

また、本実施例に係る符号化装置 1 0 は、符号化した文字列に出現したパターンを示したインデックス 3 3 を生成する。これにより、符号化装置 1 0 は、インデックス 3 3 からパターンが出現した符号化対象ファイル 3 0 を特定できる。

30

【 実施例 2 】

【 0 0 8 8 】

さて、これまで開示の装置に関する実施例について説明したが、開示の技術は上述した実施例以外にも、種々の異なる形態にて実施されてよいものである。そこで、以下では、本発明に含まれる他の実施例を説明する。

【 0 0 8 9 】

例えば、上記の実施例では、出現頻度の高いパターンに対応する符号を辞書データ 3 1 の静的辞書 3 4 に予め記憶させる場合について説明したが、これに限定されない。例えば、文書の階層ごとに、文字列で単語や数字など出現するパターンごとの出現頻度を解析により求めて、出現頻度の高いパターンから短い符号を割り当て符号化してもよい。辞書データ 3 1 は、出現したパターンと割り当てた符号を対応付けて記憶させてもよい。

40

【 0 0 9 0 】

また、上記の実施例では、符号を層構造単位で辞書データ 3 1 に記憶する場合について説明したが、これに限定されない。例えば、共通の辞書データ 3 1 を用いてもよい。また、一部の符号を層構造単位の辞書データ 3 1 で共通に登録して管理してもよい。図 1 5 は、符号の割当ての一例を示す図である。図 1 5 には、一部の符号を層構造単位の辞書データ 3 1 で共通に登録して管理する場合の符号の割当ての一例が示されている。「 8 * h 」 ~ 「 A * h 」の符号については、各階層で符号を共通に登録して管理する。例えば、共通

50

の辞書データ31により、ファイル全体で符号を管理した方が効率がよい符号がある。例えば、数値情報のNA（未入力）やnull値（値なし、文字列や数値共通）が、別の値で表現されている場合がある。この場合、共通の辞書データ31で管理することにより、符号を統一して管理できる。なお、符号を統一して管理する場合でも、ある数値では0.0をNAにし、他の値では-99.9をNAに割り当ててもよい。また、文書内全体で登場するような文字列は、符号を統一して管理することが好ましい。例えば、電子書籍の小説において小説の主人公の名前が概要、本文、講評でも登場する場合、主人公の名前の符号を統一して管理することが好ましい。一方、階層構造単位で符号を管理した方が効率がよい符号がある。例えば、階層構造単位で適切な範囲が定まる場合は、階層構造単位で符号を管理した方がよい。適切な範囲から外れた場合は、NAやNULLとして符号化する。例えば、人の体温の辞書として、35.0～42.0の範囲で辞書データ31を用意する。体温として34.8が出現した場合は、NAやNULL、または、動的に符号を割当てて符号化する。また、人の身長 of 辞書として、120.0～222.3の範囲で辞書データ31を用意する。身長として231.2という値が出現した場合は、NAやNULL、または、動的に符号を割当てて符号化する。

10

20

30

40

50

【0091】

また、図示した各装置の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的状態は図示のものに限られず、その全部または一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的または物理的に分散・統合して構成することができる。例えば、符号化装置10の特定部41、符号化部42、生成部43、受付部51、検索部52、出力部53、受付部61および復号化部62の各処理部が適宜統合されてもよい。符号化装置10の上記各処理部の処理が適宜複数の処理部の処理に分離されてもよい。さらに、各処理部にて行なわれる各処理機能は、その全部または任意の一部が、CPUおよび当該CPUにて解析実行されるプログラムにて実現され、あるいは、ワイヤードロジックによるハードウェアとして実現され得る。

【0092】

[符号化プログラム]

また、上記の実施例で説明した各種の処理は、あらかじめ用意されたプログラムをパーソナルコンピュータやワークステーションなどのコンピュータシステムで実行することによって実現することもできる。そこで、以下では、上記の実施例と同様の機能を有するプログラムを実行するコンピュータシステムの一例を説明する。最初に、符号化処理を行う符号化プログラムについて説明する。図16は、符号化プログラムを実行するコンピュータの一例を示す図である。

【0093】

図16に示すように、コンピュータ400は、CPU(Central Processing Unit)410、HDD(Hard Disk Drive)420、RAM(Random Access Memory)440を有する。これら400～440の各部は、バス500を介して接続される。

【0094】

HDD420には上記の符号化装置10の特定部41、符号化部42および生成部43と同様の機能を発揮する符号化プログラム420aが予め記憶される。なお、符号化プログラム420aについては、適宜分離してもよい。

【0095】

また、HDD420は、各種情報を記憶する。例えば、HDD420は、OSや符号化に用いる各種データを記憶する。

【0096】

そして、CPU410が、符号化プログラム420aをHDD420から読み出して実行することで、実施例の各処理部と同様の動作を実行する。すなわち、符号化プログラム420aは、特定部41、符号化部42および生成部43と同様の動作を実行する。

【0097】

なお、上記した符号化プログラム 420 a については、必ずしも最初から HDD 420 に記憶させることを要しない。

【0098】

[検索プログラム]

次に、符号化データ 32 を検索する検索プログラムについて説明する。図 17 は、復号化プログラムを実行するコンピュータの一例を示す図である。なお、図 16 と同一の部分については同一の符号を付して、説明を省略する。

【0099】

図 17 に示すように、HDD 420 には上記の符号化装置 10 の受付部 51、検索部 52 および出力部 53 と同様の機能を発揮する検索プログラム 420 b が予め記憶される。なお、検索プログラム 420 b については、適宜分離してもよい。

10

【0100】

また、HDD 420 は、各種情報を記憶する。例えば、HDD 420 は、OS や検索に用いる各種データを記憶する。

【0101】

そして、CPU 410 が、検索プログラム 420 b を HDD 420 から読み出して実行することで、実施例の各処理部と同様の動作を実行する。すなわち、検索プログラム 420 b は、受付部 51、検索部 52 および出力部 53 と同様の動作を実行する。

【0102】

なお、上記した検索プログラム 420 b についても、必ずしも最初から HDD 420 に記憶させることを要しない。

20

【0103】

[復号化プログラム]

次に、検索条件を満たすファイルを復号化する復号化プログラムについて説明する。図 18 は、復号化プログラムを実行するコンピュータの一例を示す図である。なお、図 16 および図 17 と同一の部分については同一の符号を付して、説明を省略する。

【0104】

図 17 に示すように、HDD 420 には上記の符号化装置 10 の受付部 61 および復号化部 62 と同様の機能を発揮する復号化プログラム 420 c が予め記憶される。なお、復号化プログラム 420 c については、適宜分離してもよい。

30

【0105】

また、HDD 420 は、各種情報を記憶する。例えば、HDD 420 は、OS や検索に用いる各種データを記憶する。

【0106】

そして、CPU 410 が、復号化プログラム 420 c を HDD 420 から読み出して実行することで、実施例の各処理部と同様の動作を実行する。すなわち、復号化プログラム 420 c は、受付部 61 および復号化部 62 と同様の動作を実行する。

【0107】

なお、上記した復号化プログラム 420 c についても、必ずしも最初から HDD 420 に記憶させることを要しない。

40

【0108】

また、例えば、符号化プログラム 420 a、検索プログラム 420 b および復号化プログラム 420 c は、コンピュータ 400 に挿入されるフレキシブルディスク (FD)、CD-ROM、DVD ディスク、光磁気ディスク、IC カードなどの「可搬用の物理媒体」に記憶させてもよい。そして、コンピュータ 400 がこれらからプログラムを読み出して実行するようにしてもよい。

【0109】

さらには、公衆回線、インターネット、LAN、WANなどを介してコンピュータ 400 に接続される「他のコンピュータ (又はサーバ)」などにプログラムを記憶させておく。そして、コンピュータ 400 がこれらからプログラムを読み出して実行するようにして

50

もよい。

【符号の説明】

【0110】

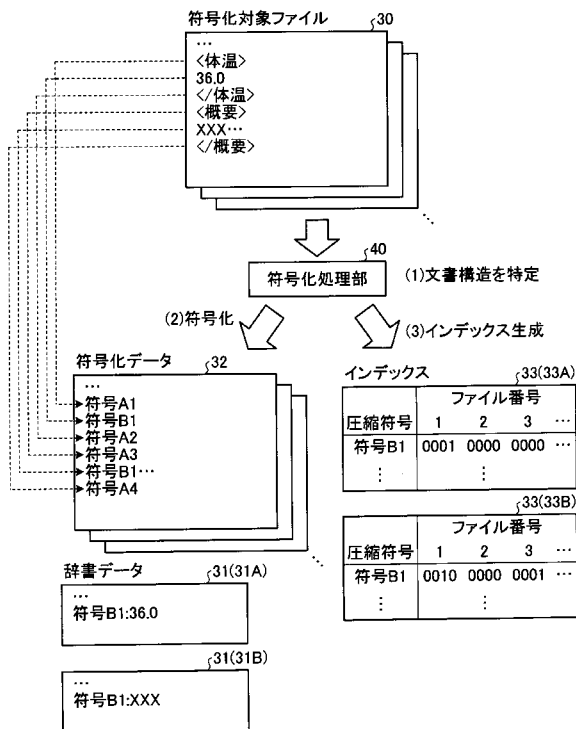
- 10 符号化装置
- 20 記憶部
- 21 制御部
- 30 符号化対象ファイル
- 31 辞書データ
- 32 符号化データ
- 33 インデックス
- 34 静的辞書
- 35 動的辞書
- 40 符号化処理部
- 41 特定部
- 42 符号化部
- 43 生成部
- 50 ファイル検索部
- 51 受付部
- 52 検索部
- 53 出力部
- 60 復号化処理部
- 61 受付部
- 62 復号化部
- 70 XMLスキーマ

10

20

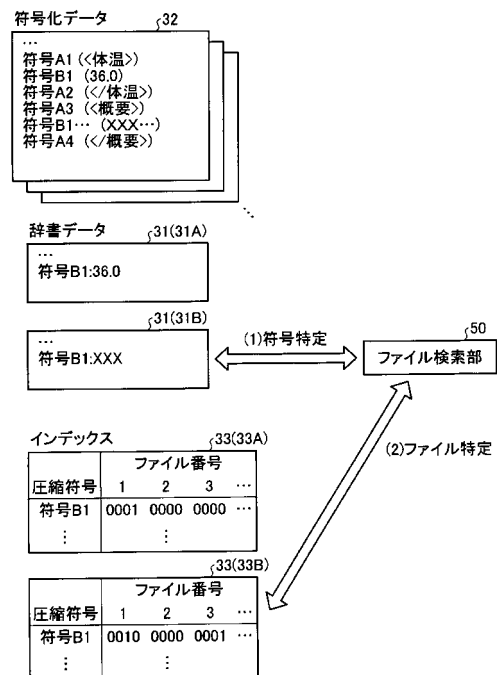
【図1】

符号化処理の流れを概略的に示した図



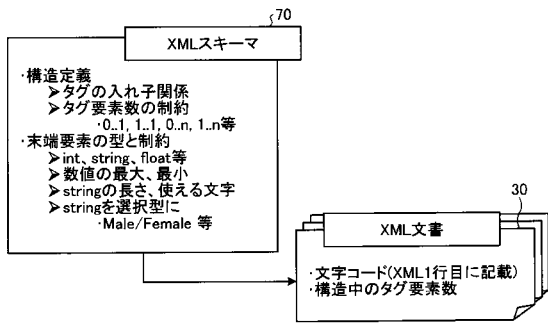
【図2A】

検索処理の流れを概略的に示した図



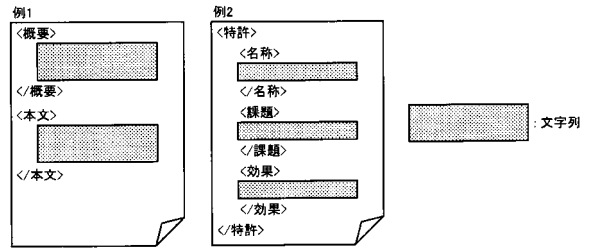
【 図 6 】

スキーマの概略的構成を示した図



【 図 7 A 】

タグにより文書構造を示した文書の一例を示す図



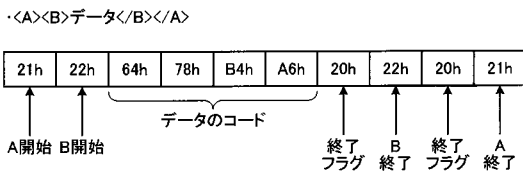
【 図 7 B 】

文書の一部にタグによりメタデータを付与した文書の一例を示す図

例3: AAAへの[リンク](リンク先URL)はこちら
 例4: <症状>BBB</症状>を訴えたため、<病名>CCC</病名>を疑い、<薬名>DDD</薬名>を投与した
 例5: <日時>2015/3/6</日時>に、<地名>大阪</地名>で<人名>鈴木</人名>に会う

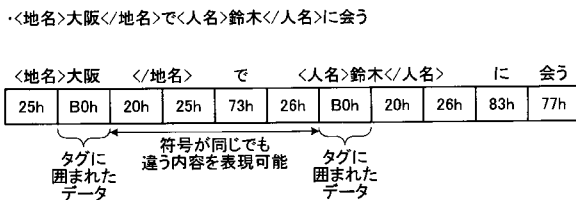
【 図 8 A 】

符号化の一例を示す図



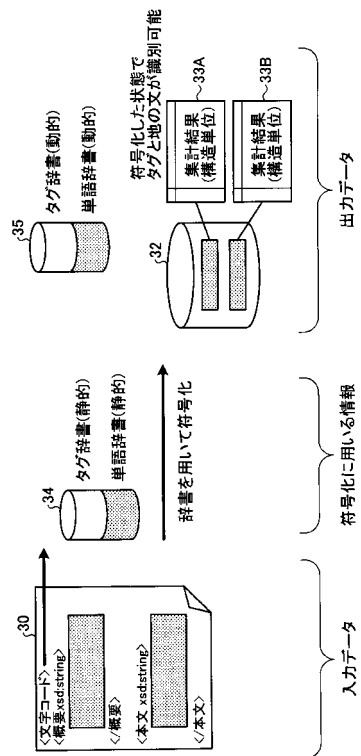
【 図 8 B 】

符号化の一例を示す図



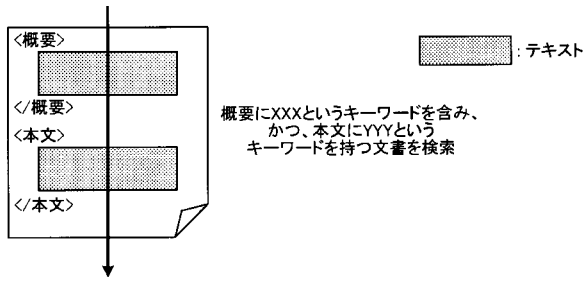
【 図 9 】

符号化の流れを概略的に示した図



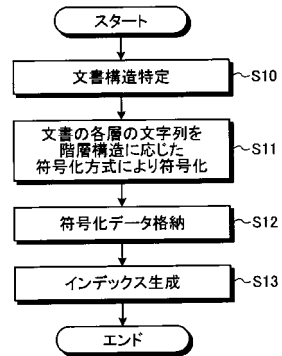
【図10A】

検索の一例を示した図



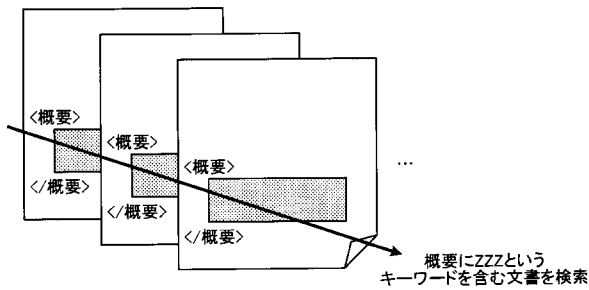
【図11】

符号化処理の手順の一例を示すフローチャート



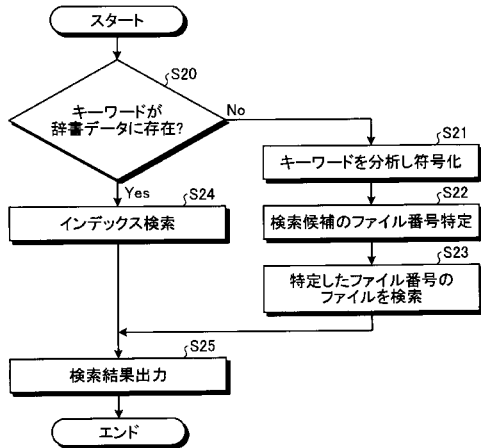
【図10B】

検索の一例を示した図



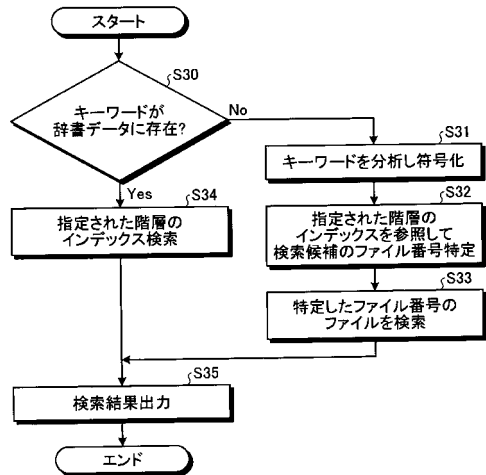
【図12】

検索処理の手順の一例を示すフローチャート



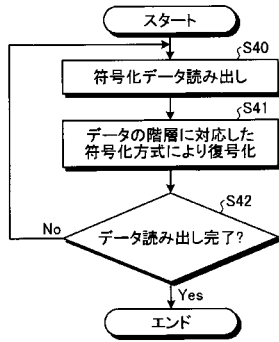
【図13】

検索処理の手順の一例を示すフローチャート



【 図 1 4 】

復号化処理の手順の一例を示すフローチャート



【 図 1 5 】

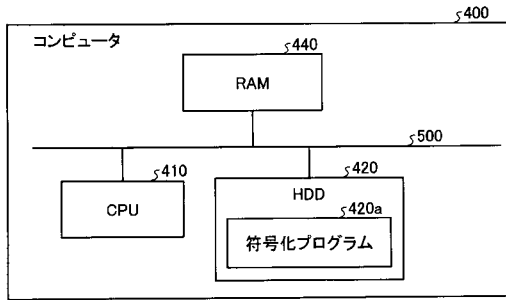
符号の割当ての一例を示す図

	0*h	1*h	2*h	3*h	4*h	5*h	6*h	7*h
*0h	NUL	DLE						
*1h	SOH	DC1						
*2h	STX	DC2						
*3h	ETX	DC3						
*4h	EOT	DC4						
*5h	ENQ	NAK						
*6h	ACK	SYN						
*7h	BEL	ETB						
*8h	BS	CAN						
*9h	HT	EM						
*Ah	LF/NL	SUB/EOF						
*Bh	VT	ESC						
*Ch	FF/NP	FS						
*Dh	CR	GS						
*Eh	SO	RS						
*Fh	SI	US						

	8*h	9*h	A*h	B*h	C*h	D*h	E*h	F*h
*0h				低頻度単語 (動的共通)	低頻度単語 (動的共通)	低頻度単語 (動的共通)	低頻度単語 (動的共通)	低頻度単語 (動的個別) 2バイト
*1h								
*2h								
*3h								
*4h								
*5h								
*6h								
*7h								
*8h								
*9h								
*Ah								
*Bh								
*Ch								
*Dh								
*Eh								
*Fh								

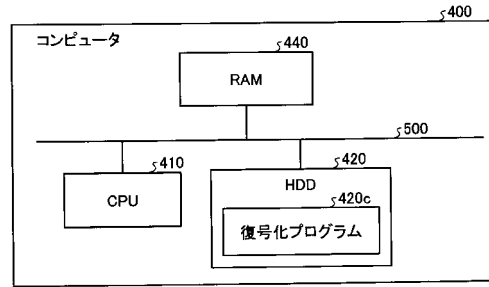
【 図 1 6 】

符号化プログラムを実行するコンピュータの一例を示す図



【 図 1 8 】

復号化プログラムを実行するコンピュータの一例を示す図



【 図 1 7 】

検索プログラムを実行するコンピュータを示す図

