

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4014160号  
(P4014160)

(45) 発行日 平成19年11月28日(2007.11.28)

(24) 登録日 平成19年9月21日(2007.9.21)

(51) Int. Cl. F I  
**G06F 17/30 (2006.01)** G O 6 F 17/30 2 1 O D  
**G06F 17/21 (2006.01)** G O 6 F 17/21 5 O 1 T

請求項の数 20 (全 36 頁)

|  |  |
|--|--|
| <p>(21) 出願番号 特願2003-155256 (P2003-155256)<br/>                 (22) 出願日 平成15年5月30日 (2003.5.30)<br/>                 (65) 公開番号 特開2004-355528 (P2004-355528A)<br/>                 (43) 公開日 平成16年12月16日 (2004.12.16)<br/>                 審査請求日 平成16年2月27日 (2004.2.27)</p> | <p>(73) 特許権者 390009531<br/>                 インターナショナル・ビジネス・マシー<br/>                 ズ・コーポレーション<br/>                 INTERNATIONAL BUSIN<br/>                 ESS MASCHINES CORPO<br/>                 RATION<br/>                 アメリカ合衆国10504 ニューヨーク<br/>                 州 アーモンク ニュー オーチャード<br/>                 ロード<br/>                 (74) 代理人 100086243<br/>                 弁理士 坂口 博<br/>                 (74) 代理人 100091568<br/>                 弁理士 市位 嘉宏<br/>                 (74) 代理人 100108501<br/>                 弁理士 上野 剛史</p> |
|--|--|

最終頁に続く

(54) 【発明の名称】 情報処理装置、プログラム、及び記録媒体

(57) 【特許請求の範囲】

【請求項1】

文字列、および、その文字列の表示形式を指定する情報であるタグ情報のそれぞれを文書構成要素として少なくとも含む文書情報について、複数の前記文書構成要素を複数のグループに分類する情報処理装置であって、

前記文書情報における前記複数の文書構成要素のそれぞれを、その文書構成要素がタグ情報ならばその役割を示し、その文書構成要素が文字列ならば文字数およびその文字列に含まれる文字の種類の種類少なくとも一方を示す情報である要素識別情報に変換する構成要素変換部と、

予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組を前記文書情報から選択し、前記文書情報に含まれるこの要素識別情報の組のそれぞれを、この要素識別情報の組の配列パターンを示す他の要素識別情報に変換する文書内配列パターン変換部と、

前記文書内配列パターン変換部により前記文書情報を繰り返し変換した結果得られた文書情報において、前記文書内配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類するグループ分類部と

を更に備える情報処理装置。

【請求項2】

前記構成要素変換部により変換された文書情報における前記要素識別情報を、出現する

10

20

頻度の高い順に順次選択して被選択情報とする構成要素選択部と、

前記構成要素選択部により選択された各被選択情報について、その被選択情報と、前記文書情報においてその被選択情報の次に配列されその被選択情報と同一の種類である前記要素識別情報との間に配列された間隙構成要素を検出する間隙構成要素検出部と

を更に備え、

前記文書内配列パターン変換部は、前記構成要素選択部により順次選択される被選択情報について、その被選択情報及びその間隙構成要素を前記要素識別情報の組として検出する

請求項 1 記載の情報処理装置。

【請求項 3】

前記間隙構成要素検出部は、前記文書情報の最後に配列される被選択情報から、前記文書情報の終端までの間に配列された文書構成要素である終端構成要素を更に検出し、

前記文書内配列パターン変換部は、前記文書情報の最後に配列される被選択情報と、前記終端構成要素とを、前記要素識別情報の組として検出する

請求項 2 記載の情報処理装置。

【請求項 4】

前記文書内配列パターン変換部は、既に検出した前記要素識別情報の組より出現頻度が高いことを更に条件として、繰り返し出現する前記要素識別情報の組を検出する

請求項 2 記載の情報処理装置。

【請求項 5】

前記文書内配列パターン変換部は、前記基準頻度以上で繰り返し出現する前記要素識別情報の組であって、変換前の文書構成要素を当該要素識別情報の組に変換するべく前記文書内配列パターン変換部が繰り返す変換の回数が予め定められた基準回数以下の要素識別情報の組を、前記配列パターンを示す要素識別情報に変換する

請求項 1 記載の情報処理装置。

【請求項 6】

前記文書情報は、前記文書構成要素として、表示画面上に表示する画像を識別する画像識別情報を含み、

前記構成要素変換部は、前記画像識別情報を、前記画像識別情報により識別される画像の形状を示す前記要素識別情報に変換する

請求項 1 記載の情報処理装置。

【請求項 7】

前記文書情報は、利用者に表示する表示情報と、当該表示情報に対する指示に応じて表示すべき他の情報の格納位置を示す格納位置情報とを前記文書構成要素として含み、

前記構成要素変換部は、前記格納位置情報を、前記他の情報が格納される前記格納位置の範囲を示す前記要素識別情報に変換する

請求項 1 記載の情報処理装置。

【請求項 8】

前記文書情報は、前記文書構成要素として、表示画面に表示する表示情報と、当該表示情報の表示形式を指定する情報であるタグ情報とを含むタグ付き文書であり、

前記タグ情報は、当該タグ情報により表示形式を指定する表示情報において更に内側タグ情報を含む、外側タグ情報であり、

前記構成要素変換部による変換後の文書情報において、前記外側タグ情報を根ノードとして生成し、前記内側タグ情報を前記根ノードの葉ノードとして生成した、文書構造情報を生成する文書構造情報生成部と、

複数の文書情報のそれぞれについて、前記文書構造情報生成部により生成された前記文書構造情報を比較することにより、一の文書情報が他の文書情報と同一の構造を有するか否かを出力する文書情報同一性出力部と

を更に備える請求項 1 記載の情報処理装置。

【請求項 9】

10

20

30

40

50

前記文書情報は、表示画面に表示すべき情報を指示する文書情報であり、更に、前記表示画面に表示されないコメント情報を含み、

前記文書内配列パターン変換部は、当該コメント情報を、一の前記要素識別情報の組と、他の前記要素識別情報の組との境界を示す情報として用いる

請求項 1 記載の情報処理装置。

【請求項 10】

前記グループ分類部により分類された各グループについて、当該グループに含まれる文書構成要素に対する目次を示す目次情報を出力する目次情報出力部

を更に備える請求項 1 記載の情報処理装置。

【請求項 11】

前記グループ分類部により分類された各グループについて、そのグループに属する文書構成要素が前記文書情報のどの領域に位置するかを示す配置情報を、前記文書情報とは別体に生成して出力するアノテーション出力部

を更に備える請求項 1 記載の情報処理装置。

【請求項 12】

前記文書情報に応じて表示画面上に情報を出力する表示部と、

前記文書内配列パターン変換部により変換される対象となる要素識別情報の組について、当該要素識別情報の組の変換元である文書構成要素が、前記表示画面に一度に表示可能かどうかを判断し、前記表示画面に表示可能でない場合に、当該要素識別情報の組を、変換を繰り返し行う対象から除外する反復終了判断部と

を更に備える請求項 1 記載の情報処理装置。

【請求項 13】

グループに分類する対象である対象文書情報と予め定められた関係を有する関連文書情報を検出する関連文書検出部と、

前記対象文書情報及び前記関連文書情報から、前記文書内配列パターン変換部により変換された要素識別情報を除外した文書について、当該対象文書情報及び当該関連文書情報の双方において出現する前記要素識別情報の組であって、前記対象文書情報及び前記関連文書情報を併せた文書において予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組を特定し、前記対象文書情報に含まれるこの要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する文書間配列パターン変換部と

を更に備え、

前記グループ分類部は、更に、前記文書間配列パターン変換部により変換された結果得られた前記対象文書情報において、前記文書間配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類する

請求項 1 記載の情報処理装置。

【請求項 14】

前記グループ分類部は、分類したグループのそれぞれについて、当該グループ内の文書構成要素が前記文書情報において果たす役割又は当該グループ内の文書構成要素の内容を示すタイトル情報を、更に生成する

請求項 1 記載の情報処理装置。

【請求項 15】

前記グループ分類部は、当該要素識別情報の組の境界に設けられた、表示画面に表示されないコメント情報に含まれる情報を、前記タイトル情報として生成する

請求項 14 記載の情報処理装置。

【請求項 16】

文字列、および、その文字列の表示形式を指定する情報であるタグ情報のそれぞれを文書構成要素として少なくとも含む文書情報について、複数の前記文書構成要素を、複数のグループに分類する情報処理装置であって、

10

20

30

40

50

グループに分類する対象である対象文書情報と予め定められた関係を有する関連文書情報を検出する関連文書検出部と、

前記対象文書情報及び前記関連文書情報のそれぞれにおいて、前記複数の文書構成要素のそれぞれを、その文書構成要素がタグ情報ならばその役割を示し、その文書構成要素が文字列ならば文字数およびその文字列に含まれる文字の種類の少なくとも一方を示す情報である要素識別情報に変換する構成要素変換部と、

前記対象文書情報及び前記関連文書情報の双方において出現する前記要素識別情報の組であって、前記対象文書情報及び前記関連文書情報を併せた文書において予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組を特定し、前記対象文書情報に含まれるこの要素識別情報の組のそれぞれを、当該文書構成要素の組の配列パターンを示す前記要素識別情報に変換する文書間配列パターン変換部と、

前記文書間配列パターン変換部により変換された結果得られた前記対象文書情報において、前記文書間配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類するグループ分類部と

を備える情報処理装置。

【請求項 17】

前記関連文書検出部は、前記対象文書情報が格納される格納位置から予め定められた範囲内に格納されている文書情報を前記関連文書情報として検出する

請求項 16 記載の情報処理装置。

【請求項 18】

前記対象文書情報は、当該対象文書情報が生成される以前に存在していた既存文書情報を更新することにより生成され

前記関連文書検出部は、前記既存文書情報を前記関連文書情報として検出する

請求項 16 記載の情報処理装置。

【請求項 19】

文字列、および、その文字列の表示形式を指定する情報であるタグ情報のそれぞれを文書構成要素として少なくとも含む文書情報について、複数の前記文書構成要素を、複数のグループに分類する情報処理装置を制御するプログラムであって、

前記情報処理装置を、

前記文書情報における前記複数の文書構成要素のそれぞれを、その文書構成要素がタグ情報ならばその役割を示し、その文書構成要素が文字列ならば文字数およびその文字列に含まれる文字の種類の少なくとも一方を示す情報である要素識別情報に変換する構成要素変換部と、

予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組を前記文書情報から選択し、前記文書情報に含まれるこの要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す他の要素識別情報に変換する文書内配列パターン変換部と

前記文書内配列パターン変換部により前記文書情報を繰り返し変換した結果得られた文書情報において、前記文書内配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類するグループ分類部と

して機能させるプログラム。

【請求項 20】

文字列、および、その文字列の表示形式を指定する情報であるタグ情報のそれぞれを文書構成要素として少なくとも含む文書情報について、複数の前記文書構成要素を、複数のグループに分類する情報処理装置を制御するプログラムであって、

前記情報処理装置を、

グループに分類する対象である対象文書情報と予め定められた関係を有する関連文書情報を検出する関連文書検出部と、

10

20

30

40

50

前記対象文書情報及び前記関連文書情報のそれぞれにおいて、前記複数の文書構成要素のそれぞれを、その文書構成要素がタグ情報ならばその役割を示し、その文書構成要素が文字列ならば文字数およびその文字列に含まれる文字の種類の種類少なくとも一方を示す情報である要素識別情報に変換する構成要素変換部と、

前記対象文書情報及び前記関連文書情報の双方において出現する前記要素識別情報の組であって、前記対象文書情報及び前記関連文書情報を併せた文書において予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組を選択し、前記対象文書情報に含まれるこの要素識別情報の組のそれぞれを、当該文書構成要素の組の配列パターンを示す前記要素識別情報に変換する文書間配列パターン変換部と、

前記文書間配列パターン変換部により変換された結果得られた前記対象文書情報において、前記文書間配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類するグループ分類部として機能させるプログラム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、情報処理装置、プログラム、及び記録媒体に関する。特に本発明は、文書情報の内容を分類する情報処理装置、プログラム、及び記録媒体に関する。

【0002】

【従来の技術】

従来、インターネット等のワールド・ワイド・ウェブ(WWW)システムにおいて、知的障害者・高齢者等に対してウェブページを快適に閲覧させることができる技術として、トランスコーディングが用いられている。トランスコーディングを行う際に、ウェブページに付加されたアノテーション等の情報を用いることにより、ウェブページ中のコンテンツを内容又は種類に応じて並び替えたり、ウェブページ中の目次を作成したりすることができる。これにより、障害者・高齢者等は、ウェブページを利便に閲覧することができる。

【0003】

また、上記トランスコーディングの為のアノテーションを適切に付加するには、ある程度の知識が必要であり、かつアノテーションの付加に要する作業量が大きい。そこで、アノテーションの付加を支援する技術が提案されている(特許文献1、非特許文献2、及び特許文献3参照)。例えば、特許文献3に記載の技術は、ウェブページ内のレイアウトを決定するタグの構造及び特徴に基づいて、ウェブページ内を分類することができる。

【0004】

【特許文献1】

特開2003-85087号公報

【0005】

【非特許文献2】

H. Takagi, C. Asakawa, K. Fukuda, J. Maeda 著、「Site-wide Annotation: Reconstructing Existing Pages to be Accessible」、CSUN2002

【0006】

【特許文献3】

特開2002-245068号公報

【0007】

【発明が解決しようとする課題】

しかしながら、上記技術は、構造が動的に変化するウェブページについて、コンテンツを適切に分類することはできない。例えば、ウェブページが、野球のスコアボード、出場選手、及び出場選手の成績表を含む場合、当該ウェブページは逐次更新され動的に変化する場合がある。このような場合、上記技術は、ウェブページのレイアウトが変化した場合、

10

20

30

40

50

又は、野球の試合が進行してスコアボードや出場選手のレコード等が変化したのかを区別することができない。

【0008】

そこで本発明は、上記の課題を解決することのできる情報処理装置、プログラム、及び記録媒体を提供することを目的とする。この目的は特許請求の範囲における独立項に記載の特徴の組み合わせにより達成される。また従属項は本発明の更なる有利な具体例を規定する。

【0009】

【課題を解決するための手段】

即ち、本発明の第1の形態によると、文字列、および、その文字列の表示形式を指定する情報であるタグ情報のそれぞれを文書構成要素として少なくとも含む文書情報について、複数の前記文書構成要素を複数のグループに分類する情報処理装置であって、文書情報における複数の文書構成要素のそれぞれを、その文書構成要素がタグ情報ならばその役割を示し、その文書構成要素が文字列ならば文字数およびその文字列に含まれる文字の種類の少なくとも一方を示す情報である要素識別情報にそれぞれ変換する構成要素変換部と、予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組を前記文書情報から選択し、前記文書情報に含まれるこの要素識別情報の組のそれぞれを、この要素識別情報の組の配列パターンを示す他の要素識別情報に変換する文書内配列パターン変換部と、文書内配列パターン変換部により文書情報を繰り返し変換した結果得られた文書情報において、文書内配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループとして分類するグループ分類部とを備える情報処理装置、当該情報処理装置を制御するプログラムを提供する。

なお上記の発明の概要は、本発明の必要な特徴の全てを列挙したものではなく、これらの特徴群のサブコンビネーションも又発明となりうる。

【0010】

【発明の実施の形態】

以下、発明の実施の形態を通じて本発明を説明するが、以下の実施形態は特許請求の範囲にかかる発明を限定するものではなく、又実施形態の中で説明されている特徴の組み合わせの全てが発明の解決手段に必須であるとは限らない。

図1は、情報処理装置10の機能ブロック図を示す。情報処理装置10は、HTML等のタグ付き文書である文書情報について、当該文書情報に含まれる複数の文書構成要素をグループに分類する装置である。具体的には、本実施形態において、情報処理装置10は、WWWサーバ等から取得したHTML文書について、当該文書情報に含まれるタグ情報等をグループに分類して表示する。これに代えて、情報処理装置10は、WWWサーバ等のサーバ装置であって、ウェブブラウザ等からのリクエストに応じてHTML文書内の情報を分類し、分類した当該情報をウェブブラウザに返送してもよい。

【0011】

情報処理装置10は、関連文書検出部100と、構成要素変換部105と、構成要素選択部110と、間隙構成要素検出部120と、文書内配列パターン変換部130と、変換指示入力部135と、反復終了判断部140と、文書間配列パターン変換部160と、グループ分類部170と、並替出力部180と、目次情報出力部190と、グループ識別情報生成部200と、アノテーション出力部202と、表示部205と、文書構造情報生成部210と、文書情報同一性出力部220とを備える。

【0012】

関連文書検出部100は、グループ分けを行う対象の対象文書情報を取得すると、当該対象文書情報と予め定められた関係を有する関連文書情報を検出する。例えば、関連文書検出部100は、対象文書情報が格納される格納位置から予め定められた範囲内に格納されている文書情報を、関連文書情報として検出する。対象文書情報が格納される格納位置から予め定められた範囲内とは、例えば、対象文書情報と同一のディレクトリ内であってもよいし、同一のサイト内であってもよい。また、関連文書検出部100は、複数の関連文

10

20

30

40

50

書情報を検出することが好ましい。そして、関連文書検出部 100 は、対象文書情報及び関連文書情報を順次構成要素変換部 105 に送る。

【0013】

構成要素変換部 105 は、関連文書検出部 100 から文書情報を受け取ると、当該文書情報における複数の文書構成要素のそれぞれを当該文書構成要素の種類又は役割を示す要素識別情報にそれぞれ変換し、変換結果の文書情報を構成要素選択部 110 に送る。なお、要素識別情報とは、文書構成要素そのものであってもよい。即ち、構成要素変換部 105 は、箇条書きの始点を示す文書構成要素である < l i > を、同様に箇条書きの始点を示す要素識別情報である < l i > に変換してもよい。また、好ましくは、構成要素変換部 105 は、変換前の文書構成要素及び変換後の要素識別情報を対応付けた情報をメモリ等に格納しておく。

10

【0014】

また、構成要素変換部 105 は、文書内配列パターン変換部 130 又は文書間配列パターン 160 からの通知を受けて、当該構成要素変換部 105 により新たに変換すべき文書構成要素の組合せを示す新規登録要素集合を、当該新規登録要素集合の変換先とするべき要素識別情報に対応付けて、メモリ等から取得する。そして、構成要素変換部 105 は、文書情報における新規登録要素集合のそれぞれを、対応する要素識別情報に変換する。

【0015】

構成要素選択部 110 は、構成要素変換部 105 により変換された文書情報又は反復終了判断部 140 から受け取った文書情報について、当該文書情報における要素識別情報の何れかを、例えば、出現する頻度の高い順に被選択情報として選択し、選択結果を間隙構成要素検出部 120 に送る。ここで、出現する頻度とは、例えば、一の文書情報における要素識別情報の個数である。これに代えて、出現する頻度とは、文書情報のサイズ当たりの、要素識別情報の個数であってもよい。

20

【0016】

間隙構成要素検出部 120 は、構成要素選択部 110 により選択された各被選択情報について、当該被選択情報と、文書情報において当該被選択情報の次に配列され当該被選択情報と同一の種類である情報との間に配列された間隙構成要素を検出する。そして、間隙構成要素検出部 120 は、間隙構成要素を被選択情報毎に文書内配列パターン変換部 130 に送る。また、好ましくは、間隙要素検出部 120 は、文書情報において複数の被選択情報の後に配列される終端構成要素を更に検出する。この場合、間隙要素検出部 120 は、更に、終端構成要素を文書内配列パターン変換部 130 に送る。

30

【0017】

文書内配列パターン変換部 130 は、関連文書検出部 100 により変換された文書情報において、予め定められた基準頻度以上で繰り返し出現する要素識別情報の組のそれぞれを変換元の候補として選択する。そして、文書内配列パターン変換部 130 は、変換元の候補である当該要素識別情報の組の配列パターンを示す要素識別情報を変換先の候補として生成する。そして、文書内配列パターン変換部 130 は、変換元の候補及び変換先の候補を変換指示入力部 135 に送る。

【0018】

文書内配列パターン変換部 130 は、変換指示入力部 135 から変換指示を受けると、変換元の候補である要素識別情報の組のそれぞれを、変換先の候補である要素識別情報に変換し、変換結果を反復終了判断部 140 に送る。要素識別情報の組を検出する具体例として、文書内配列パターン変換部 130 は、以下の処理を行う。まず、文書内配列パターン変換部 130 は、間隙構成要素を被選択情報毎に、間隙要素検出部 120 から受け取る。また、文書内配列パターン変換部 130 は、終端構成要素を、複数の被選択情報のうち文書情報の最後に配列される被選択情報に対応付けて受け取る。

40

【0019】

続いて、文書内配列パターン変換部 130 は、被選択情報及び当該被選択情報に対応した間隙構成要素を、要素識別情報の組として検出し、当該要素識別情報の組が基準頻度以上

50

で繰り返し出現する場合に、各要素識別情報の組を変換する。また、文書内配列パターン検出部 130 は、複数の被選択情報のうち文書情報の最後に配列される被選択情報と、終端構成要素とを、要素識別情報の組として検出する。これにより、文書情報の終端に配列された終端構成要素を、文書内配列パターン変換部 130 による変換の対象に含めることができる。

なお、好ましくは、文書内配列パターン変換部 130 は、変換前の要素識別情報の組及び変換後の要素識別情報を対応付けてメモリ等に記録しておく。

#### 【0020】

更に、文書内配列パターン変換部 130 は、配列パターンを示す要素識別情報に変換した変換元である複数の文書構成要素を、変換先の要素識別情報に対応付けて、構成要素変換部 105 による新たな変換対象である新規登録要素集合としてメモリ等に登録し、構成要素変換部 105 に通知する。また、文書内配列パターン変換部 130 は、当該新規登録要素集合を登録すべきか否かを判断するべく、新規登録要素集合の候補となる情報を利用者等に対して出力してもよい。そして、文書内配列パターン変換部 130 は、利用者からの入力に基づいて、当該新規登録要素集合を登録してもよい。

10

#### 【0021】

変換指示入力部 135 は、文書内配列パターン変換部 130 における変換元の候補である要素識別情報の組及び変換先の候補である要素識別情報を利用者に対して出力することにより、文書内配列パターン変換部 130 による要素識別情報への変換を行うか否かを入力させ、入力結果に応じて変換指示を文書内配列パターン変換部 130 に送る。また、変換指示入力部 135 は、文書内配列パターン変換部 130 による要素識別情報への変換を行うか否かを利用者に入力させなくともよい。この場合、文書内配列パターン変換部 130 は、変換処理の毎に利用者にお問い合わせることなく、予め定めた規則に従い、文書構成要素の組を変換する。

20

#### 【0022】

反復終了判断部 140 は、文書情報を文書内配列パターン変換部 130 から受け取ると、文書内配列パターン変換部 130 によって変換された文書情報に対して、更に、文書内配列パターン変換部 130 による変換を繰り返し行わせるべく、受け取った文書情報を構成要素選択部 110 に送る。また、反復終了判断部 140 は、予め定められた終了条件を満たした場合、例えば、文書内配列パターン変換部 130 によって変換対象の要素識別情報の組が検出されなかった場合に、繰り返し処理を終了するべく、文書情報を文書間配列パターン変換部 160 に送る。

30

#### 【0023】

文書間配列パターン変換部 160 は、文書情報、例えば、対象文書情報及び関連文書情報を反復終了判断部 140 から順次受け取る。そして、文書間配列パターン変換部 160 は、対象文書情報及び関連文書情報から、文書内配列パターン変換部 130 により変換された要素識別情報を除外した文書について、当該対象文書情報及び当該関連文書情報の双方において出現する要素識別情報の組を特定する。続いて、文書間配列パターン変換部 160 は、特定した当該要素識別情報の組のうち、当該対象文書情報及び当該関連文書情報を併せた文書において予め定められた基準頻度以上で繰り返し出現する要素識別情報の組を検出する。そして、文書間配列パターン変換部 160 は、基準頻度以上で繰り返し出現する要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す要素識別情報に変換し、変換結果の対象文書情報及び関連文書情報をグループ分類部 170 に送る。

40

#### 【0024】

なお、文書間配列パターン変換部 160 は、変換前の要素識別情報の組及び変換後の要素識別情報に対応付けてメモリ等に記録しておく。更に、文書間配列パターン変換部 160 は、配列パターンを示す要素識別情報に変換した変換元である複数の文書構成要素を、変換先の要素識別情報に対応付けて、構成要素変換部 105 による新たな変換対象である新規登録要素集合としてメモリ等に登録し、構成要素変換部 105 に通知する。また、文書

50



間配列パターン変換部 160 は、当該新規登録要素集合を登録するべきか否かを判断するべく、新規登録要素集合の候補となる情報を利用者等に対して出力してもよい。そして、文書間配列パターン変換部 160 は、利用者からの入力に基づいて、当該新規登録要素集合を登録してもよい。

【0025】

グループ分類部 170 は、文書間配列パターン変換部 160 から受け取った対象文書情報及び関連文書情報のそれぞれにおいて、文書内配列パターン変換部 130 により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループとして分類する。同様に、グループ分類部 170 は、文書間配列パターン変換部 160 により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループとして分類する。そして、グループ分類部 170 は、分類結果と共に、対象文書情報及び関連文書情報を並替出力部 180、目次情報出力部 190、グループ識別情報生成部 200、アノテーション出力部 202、及び文書構造情報生成部 210 に送る。

10

【0026】

更に、グループ分類部 170 は、分類したグループのそれぞれについて、当該グループ内の文書構成要素が文書情報において果たす役割又は当該グループ内の文書構成要素の内容を示すタイトル情報を生成し、目次情報出力部 190、及びアノテーション出力部 202 に出力してもよい。より詳しくは、当該グループに分類された要素識別情報の組が、コメント情報を境界として検出されたものである場合には、グループ分類部 170 は、当該コメント情報に含まれる情報をタイトル情報として生成してもよい。

20

【0027】

また、グループ分類部 170 は、分類したグループのそれぞれについて、文書情報における当該グループの重要度を示す重要度情報を更に生成し、並替出力部 180、目次情報出力部 190、及びアノテーション出力部 202 に出力してもよい。具体的には、グループ分類部 170 は、当該要素識別情報の組が文書情報のどこに位置するかを示す配置情報、当該要素識別情報の組に変換された文書構成要素の色の情報、当該文書構成要素の大きさの情報、及び当該文書構成要素が文字列である場合の文字列の内容に基づいて、重要度情報を生成してもよい。

【0028】

更に、グループ分類部 170 は、グループに属する文書構成要素を利用者に対して出力し、当該グループに分類するべきか否かを示す指示を受け取ってもよい。この場合、グループ分類部 170 は、当該指示に応じてグループに分類するか否かを判断する。また、グループ分類部 170 は、グループに属する文書構成要素を利用者に対して出力し、当該グループに対応付けて生成するべきタイトル情報を指示する役割指定指示を受け取ってもよい。この場合、グループ分類部 170 は、役割指定指示に応じてグループ情報を生成してもよい。

30

【0029】

並替出力部 180 は、対象文書情報において、複数の文書構成要素を、グループ分類部 170 により分類されたグループ毎に並び替える。更に、並替出力部 180 は、グループの内容に応じて、複数のグループを並び替えてもよい。例えば、並替出力部 180 は、文字や画像を含むグループを重要度の高いグループとしてより先に並び替え、リンクリスト、ヘッダ、フッタ、及び広告を含むグループを重要度の低いグループとしてより後に並び替えてもよい。即ち、並替出力部 180 は、複数のグループ中の文書構成情報を、当該グループの重要度情報が高い順に並び替えてもよい。そして、並替出力部 180 は、並び替えた結果の対象文書情報を表示部 205 に送る。

40

【0030】

処理の具体例としては、並替出力部 180 は、構成要素変換部 105、文書内配列パターン変換部 130、及び文書間配列パターン変換部 160 によりメモリ等に格納されてきた情報に基づいて、各グループに分類された文書構成要素を特定し、当該文書構成要素をグ

50

グループ毎に並び替える。なお、並替出力部 180 は、更に、グループ分類部 170 から受け取った関連文書情報の分類結果に基づいて、対象文書情報及び関連文書情報のそれぞれから、タイトル情報が同一であるグループを選択してもよい。この場合、並替出力部 180 は、選択したこれらのグループに属する文書構成要素を表示部 205 に送ってもよい。

#### 【0031】

目次情報出力部 190 は、グループ分類部 170 により分類された各グループについて、当該グループに含まれる文書構成要素に対する目次を示す目次情報を、受け取ったタイトル情報に基づいて生成し表示部 205 に出力する。また、目次情報出力部 190 は、グループ分類部 170 により分類された各グループについて、当該グループに含まれる文書構成要素が、文書情報のどの部分に位置するかを示す情報を更に出力してもよい。より具体的には、この情報とは、HTML 文書におけるアンカーを用いた情報であってもよい。更に、目次情報出力部 190 は、各グループの位置を示すこの情報に、グループのタイトル又はグループの重要度情報を対応付けて出力してもよい。

#### 【0032】

グループ識別情報生成部 200 は、グループ分類部 170 により分類された各グループに対応付けて、当該グループを識別するグループ識別情報を文書情報内に生成し、表示部 205 に出力する。例えば、グループ識別情報生成部 200 は、文書情報内のグループの境界に、グループの境界であることを明示するための画像を生成してもよい。これらを受けて、表示部 205 は、対象文書情報及び関連文書情報に応じて、表示画面上に情報を出力する。

#### 【0033】

また、文書情報が HTML 文書である場合には、グループ識別情報生成部 200 は、以下の処理を行ってもよい。例えば、グループ識別情報生成部 200 は、グループ識別情報を音声で出力させるための出力指示情報を、文書情報内の所定のタグにおける alt オプションのパラメータとして生成してもよい。出力指示情報とは、例えば、透明かつ非常に小さな画像など、表示画面上では視覚的に識別することが困難な画像などである。この結果、通常のブラウザ画面での表示に与える影響を非常に小さなものに抑える一方、音声ブラウザ等は、当該画像ファイルの内容等に基づいて、グループ識別情報を音声として出力する。これにより、通常のブラウザを使用する健常者の利便性を保ちつつ、音声ブラウザを使用する障害者等の利便性を高めることができる。

#### 【0034】

アノテーション出力部 202 は、グループ分類部 170 により分類された各グループを識別する情報を、文書情報とは別体に生成して出力する。例えば、アノテーション出力部 202 は、複数のグループのそれぞれについて、当該グループに属する文書構成要素が対象文書情報のどの領域に位置するかを示す配置情報であるアノテーション情報を、各グループを識別する情報として出力する。一例としては、アノテーション出力部 202 は、XPath 又は XPointer 等の技術により配置情報を生成し出力してもよい。更に、アノテーション出力部 202 は、受け取ったタイトル情報及び重要度情報を配置情報に更に対応付けた情報を、アノテーション情報として出力してもよい。

#### 【0035】

また、情報処理装置 10 は、アノテーション情報の作成者による作業を支援することができる。例えば、アノテーション情報の作成者は、情報処理装置 10 により出力されたアノテーション情報に基づき、当該アノテーション情報を修正又は変更することにより、所望のアノテーション情報を作成することができる。この結果、アノテーション情報の作成者は、初めからアノテーション情報を作成する場合と比較して、所望のアノテーション情報を効率的に作成することができる。

#### 【0036】

文書構造情報生成部 210 は、グループ分類部 170 から受け取った対象文書情報及び関連文書情報のそれぞれについて、当該文書情報の構造を示す文書構造情報を生成し、当該文書構造情報を、受け取った分類結果に対応付けて文書情報同一性出力部 220 に送る。

10

20

30

40

50

具体的には、文書情報は、文書構成要素として、表示画面に表示する表示情報と、当該表示情報の表示形式を指定するタグ情報とを含むタグ付き文書である。そして、タグ情報は、当該タグ情報により表示形式を指定する表示情報において更に内側タグ情報を含む、外側タグ情報である。この場合、文書構造情報生成部 210 は、構成要素変換部 105 による変換後の文書情報において、外側タグ情報を根ノードとして生成し、内側タグ情報を当該根ノードの葉ノードとして生成した情報を、文書構造情報として生成する。

#### 【0037】

文書情報同一性出力部 220 は、文書構造情報生成部 210 から受け取った文書情報のそれぞれについて、文書構造情報生成部 210 により生成された文書構造情報を比較することにより、一の文書情報が他の文書情報と同一の構造を有するか否かを判断し、判断結果を出力する。例えば、文書情報同一性出力部 220 は、文書構造情報の葉ノードのそれぞれに対して予めハッシュ値を定めておき、当該ハッシュ値が異なるか否かを判断する DOMHASH 技術を用いてもよい。更に、文書情報同一性出力部 220 は、分類情報を用いて同一性を判断してもよい。

10

#### 【0038】

このように、情報処理装置 10 は、文書情報に含まれる複数の文書構成要素をグループに分類する。更に、情報処理装置 10 は、分類したグループ毎に文書構成要素を並び替え、目次を作成する等の処理を行うことができる。

#### 【0039】

図 2 は、情報処理装置 10 のフローチャートを示す。関連文書検出部 100 は、グループ分けを行う対象の対象文書情報を取得すると、当該対象文書情報と予め定められた関係を有する関連文書情報を検出する (S200)。例えば、関連文書検出部 100 は、対象文書情報が格納される格納位置から予め定められた範囲内に格納されている文書情報を、関連文書情報として検出する。

20

#### 【0040】

これに代えて、対象文書情報が、当該対象文書情報が生成される以前に存在していた既存文書情報を更新することにより生成された場合には、関連文書検出部 100 は、当該既存文書情報を関連文書情報として検出してもよい。より詳細には、関連文書検出部 100 は、対象文書情報と同一のファイル名、パス名、又は URL を有する文書を、関連文書情報として検出してもよい。

30

また、関連文書検出部 100 は、対象文書情報が格納される格納位置から予め定められた範囲内に格納されている文書情報及び対象文書情報が生成される以前に存在していた既存文書情報の双方を、関連文書情報として検出してもよい。

#### 【0041】

構成要素変換部 105 は、文書情報 (例えば、対象文書情報又は関連文書情報) における複数の文書構成要素のそれぞれを当該文書構成要素の種類又は役割を示す要素識別情報にそれぞれ変換する (S210)。そして、文書内配列パターン変換部 130 は、変換後の文書情報において、予め定められた基準頻度以上で繰り返し出現する要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する (S220)。

40

#### 【0042】

文書間配列パターン変換部 160 は、対象文書情報及び関連文書情報から、文書内配列パターン変換部 130 により変換された要素識別情報を除外した文書について、当該対象文書情報及び当該関連文書情報の双方において出現する要素識別情報の組を特定する (S230)。続いて、文書間配列パターン変換部 160 は、特定した当該要素識別情報の組のうち、当該対象文書情報及び当該関連文書情報を併せた文書において予め定められた基準頻度以上で繰り返し出現する要素識別情報の組を検出する。そして、文書間配列パターン変換部 160 は、基準頻度以上で繰り返し出現する要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する。

#### 【0043】

50

グループ分類部 170 は、文書間配列パターン変換部 160 から受け取った対象文書情報及び関連文書情報のそれぞれにおいて、文書内配列パターン変換部 130 により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループとして分類する (S240)。同様に、グループ分類部 170 は、文書間配列パターン変換部 160 により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループとして分類する。

【0044】

文書構造情報生成部 210 は、文書間配列パターン変換部 160 から受け取った対象文書情報及び関連文書情報のそれぞれについて、当該文書情報の構造を示す文書構造情報を生成する (S250)。そして、文書情報同一性出力部 220 は、対象文書情報及び関連文書情報のそれぞれが同一の構造を有するか否かを、文書構造情報生成部 210 により生成された文書構造情報を比較することにより判断し、判断結果を出力する。

10

【0045】

なお、文書構造情報生成部 210 及び文書構造情報生成部 210 による上記処理は、本実施形態における必須の構成とはならない。即ち、文書構造情報生成部 210 は、文書構造情報を生成しなくともよい。また、文書情報同一性出力部 220 は、同一性の判断結果を出力しなくともよい。

【0046】

また、情報処理装置 10 は、S220 及び S230 のうち一方の処理を行わなくともよい。即ち、文書内配列パターン変換部 130 は、基準頻度以上で繰り返し出現する要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す要素識別情報に変換しなくともよい。また、文書間配列パターン変換部 160 は、対象文書情報及び関連文書情報から、文書内配列パターン変換部 130 により変換された要素識別情報を除外した文書について、当該対象文書情報及び当該関連文書情報の双方において出現する要素識別情報の組を特定しなくともよい。

20

【0047】

図 3 は、図 2 の S210 におけるフローチャートを示す。構成要素変換部 105 は、文書情報における各文書構成要素について以下の処理を繰り返し実行する。文書情報は、文書構成要素として、表示画面上に表示する画像を識別する画像識別情報を有している。構成要素変換部 105 は、文書構成要素がこの画像識別情報である場合に (S300: YES)、当該画像識別情報を、当該画像識別情報により識別される画像の形状、データサイズ、又は格納位置を示す要素種別情報に変換する (S310)。

30

【0048】

例えば、構成要素変換部 105 は、当該画像識別情報により識別される画像における、縦の長さ及び横の長さの比率が所定値より大きい場合に、当該画像識別情報を、文章等の区切り画像を示す要素識別情報である <divider /> に変換してもよい。また、構成要素変換部 105 は、透明又は単一色、かつ所定サイズより小さな画像であって、縦の長さ及び横の長さの比率が所定値より大きい領域に配置された画像を識別する画像識別情報を、<divider /> に変換してもよい。

【0049】

一方、文書構成要素が画像識別情報でない場合に (S300: NO)、構成要素変換部 105 は、文書構成要素がテキストデータであるか否かを判断する (S320)。文書構成要素がテキストデータである場合に (S320: YES)、構成要素変換部 105 は、当該テキストデータを、当該テキストデータの内容又はデータサイズを示す要素種別情報に変換する (S330)。

40

【0050】

文書構成要素がテキストデータでない場合に (S320: NO)、構成要素変換部 105 は、文書構成要素が、利用者からの指示を受けて他の情報を表示させるリンク情報であるか否かを判断する (S340)。文書構成要素がリンク情報である場合に (S340: YES)、構成要素変換部 105 は、当該リンク情報のリンク先に基づく要素識別情報に変

50

換する（S350）。

【0051】

構成要素変換部105は、以上の処理を各文書構成要素に適用した後、次の処理を行う。構成要素変換部105は、文書構成要素が所定の規則に適合するか否かを判断する（S360）。例えば、構成要素変換部105は、文書内の索引を示すリンクリスト情報を検出するべく、当該リンクリスト情報を形成する文書構成要素の配列パターンを予め格納しておき、文書構成要素の配列が当該配列パターンに適合するか否かを判断する。

【0052】

変換後の文書情報において、要素識別情報が所定の規則に適合する場合に（S360：YES）、構成要素変換部105は、当該文書構成要素を、当該規則の内容を示す要素識別情報に変換する（S370）。例えば、構成要素変換部105は、リンクリストを形成する文書構成要素の組を、リンクリストである旨を示す要素識別情報に変換する。他の例としては、構成要素変換部105は、野球等のスポーツのスコアボードを形成する文書構成要素の組を、スコアボードを示す要素識別情報に変換する。また、構成要素変換部105は、所定の広告用サイトへのリンク情報であって、リンク先の情報を表示するべく指示を受ける表示情報が画像である場合に、当該リンク情報及び画像を、広告を示す要素識別情報である<a d />に変換する。

【0053】

また、構成要素変換部105は、互いに異なる複数の文書構成要素を同一の要素識別情報に変換してもよい。例えば、構成要素変換部105は、所定サイズ以下の画像を示す文書構成要素であると、箇条書きの開始点等を示す点である「・」とを、同一の要素識別情報である<bullet />に変換してもよい。このように、構成要素変換部105は、互いに異なる複数の文書構成要素であっても、同一の種類であれば、当該文書構成要素を、当該種類を識別する要素識別情報に変換する。

【0054】

また、構成要素変換部105は、互いに同一な複数の文書構成要素を、互いに異なる複数の要素識別情報に変換してもよい。例えば、構成要素変換部105は、箇条書きを示す文書構成要素である「・今日のニュース」を、それぞれが要素識別情報である<bullet />及び<shorttext />に変換する。一方、構成要素変換部105は、文書構成要素である「今日のテストは算数・理科です」を、要素識別情報である<shorttext />に変換する。即ち、構成要素変換部105は、互いに同一な文書構成要素である点「・」であっても、当該文書構成要素の役割に応じて互いに異なる要素識別情報に変換する。

【0055】

また、構成要素変換部105は、文書内配列パターン変換部130又は文書間配列パターン160からの通知を受けて、新規登録要素集合を、当該新規登録要素集合の変換先とするべき要素識別情報に対応付けて、メモリ等から取得する。そして、構成要素変換部105は、文書情報における新規登録要素集合のそれぞれを、対応する要素識別情報に変換する。即ち、情報処理装置10は、以前に変換されたことのある文書構成要素の組については、文書内配列パターン変換部130又は文書間配列パターン160に代えて、構成要素変換部105により変換処理を行う。これにより、情報処理装置10は、変換処理を効率化できる。

【0056】

このように、構成要素変換部105は、文書情報における複数の文書構成要素のそれぞれを、文書構成要素の種類又は役割を示す要素識別情報にそれぞれ変換する。また、構成要素変換部105は、所定の条件を満たす要素識別情報の組を、当該条件の内容を示す要素識別情報に更に変換する。この結果、情報処理装置10は、例えば、ある文書構成要素の組が、利用者に対して表形式で表示するべく配置されたものであるか、或いは、文書情報のレイアウトを整えるべく表形式を利用したものであるかを区別することができる。更

10

20

30

40

50

に、情報処理装置 10 は、一部の文書構成要素の組については、当該文書構成要素の組の果たす役割、例えば、スコアボードである旨又は広告である旨を特定することができる。

【0057】

また、構成要素変換部 105 による処理の手順は、本図の例に限定されない。例えば、構成要素変換部 105 は、S360 における所定の規則が、要素識別情報でなく文書構成要素を対象とする場合には、画像、テキスト、及びリンクの変換 (S310、S330、及び S350) に先だって、所定の規則に適合するか否かを判断してもよい (S360)。

【0058】

図 4 は、図 2 の S220 におけるフローチャートを示す。構成要素選択部 110 は、構成要素変換部 105 により変換された各文書情報における要素識別情報を、出現する頻度の高い順に、順次被選択情報として選択する (S400)。そして、間隙構成要素検出部 120 は、構成要素選択部 110 により選択された当該被選択情報と、文書情報において当該被選択情報の次に配列され当該被選択情報と同一の種類である情報との間に配列された間隙構成要素を検出する (S410)。

【0059】

ここで、間隙構成要素は、被選択情報と同一の種類の要素識別情報を含まない。即ち、間隙構成要素検出部 120 は、被選択情報の次に配列される要素識別情報を順次検出し、被選択情報と同一の種類の要素識別情報を検出するまでに検出された要素識別情報を、間隙構成要素とする。

【0060】

このように、文書内配列パターン変換部 130 は、被選択情報及び間隙構成要素を、要素識別情報の組として検出する。これに代えて、対象文書情報が、対象文書中に配列され表示画面に表示されないコメント情報を含む場合、文書内配列パターン変換部 130 は、当該コメント情報を、一の要素識別情報の組と、他の要素識別情報の組との境界を示す情報として用いてもよい。例えば、対象文書情報が、コメント情報として、< ! - S T A R T A A > 及び < ! - E N D A A > を含む場合、文書内配列パターン変換部 130 は、これらのコメント情報の間に配列される情報を、要素識別情報の組として検出してもよい。更に、グループ分類部 170 は、当該コメント情報に含まれる情報、例えば、「A A」の部分の文字列に基づいて、グループの果たす役割を特定してもよい。更に、文書内配列パターン変換部 130 は、S T A R T 等の所定の文字列がコメント情報内に含まれていない場合には、当該コメント情報を要素識別情報の組の境界として用いなくともよい。

【0061】

また、更に他の例として、文書内配列パターン変換部 130 は、所定以上のサイズの空白又は改行を示す文書構成要素を、一の要素識別情報の組と、他の要素識別情報の組との境界を示す情報として用いてもよい。

【0062】

文書内配列パターン変換部 130 は、構成要素選択部 110 により選択された被選択情報及び間隙構成要素検出部 120 により検出された間隙構成要素の組が、構成要素選択部 110 及び間隙構成要素検出部 120 により既に検出されていた要素識別情報の組より出現頻度が高い場合に (S420 : Y E S)、構成要素選択部 110 により選択された被選択情報及び間隙構成要素検出部 120 により検出された間隙構成要素の組を変換対象の候補として選択する (S430)。一方、既に検出されていた要素識別情報の組より出現頻度が低い場合に (S420 : N O)、情報処理装置 10 は、S440 に処理を移す。

【0063】

間隙要素検出部 120 が検出する間隙構成要素とは、文書構成要素及び要素識別情報の何れも含まない空集合であってもよい。この場合、文書内配列パターン変換部 130 は、同一の種類の被選択情報が複数連続して配列される場合に、当該複数の被選択情報を、変換対象の候補として選択する。

【0064】

構成要素選択部 110 は、文書情報における何れかの要素識別情報について未だ被選択情

10

20

30

40

50

報として選択していない場合に ( S 4 4 0 : N O )、S 4 0 0 に処理を戻すことにより、次に出現頻度の高い要素識別情報を被選択情報として選択する。一方、構成要素選択部 1 1 0 は、文書情報における全ての要素識別情報を被選択情報として選択した場合に ( S 4 4 0 : Y E S )、S 4 3 0 において変換対象の候補として選択した要素識別情報の組を、所定の判断に基づき、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する処理を行う ( S 4 5 0 )。

【 0 0 6 5 】

反復終了判断部 1 4 0 は、予め定められた終了条件を満たしていない場合に ( S 4 6 0 : N O )、変換後の文書情報について更に変換処理を行うべく、S 4 0 0 に処理を戻す ( S 4 7 0 )。ここで、反復終了判断部 1 4 0 は、文書内配列パターン変換部 1 3 0 により変換される対象となる要素識別情報の組について、当該要素識別情報の変換元である文書構成要素の合計サイズが、予め定められた条件を満たす場合、例えば一つの表示画面に一度に表示可能な情報のサイズに達した場合に、当該要素識別情報の組を、以降の処理において文書内配列パターン変換部 1 3 0 による変換を繰り返し行わせる対象から除外する。例えば、反復終了判断部 1 4 0 は、当該要素識別情報の組を、以降の変換を停止する旨を示す所定の要素識別情報に変換してもよい。これにより、各グループにおける文書構成要素の合計サイズを所定サイズ以下に抑え、各グループを表示画面内に表示させることができる。

【 0 0 6 6 】

これに代えて、反復終了判断部 1 4 0 は、文書構成要素の合計サイズが、2つの表示画面に表示可能な情報のサイズに達した場合に、当該要素識別情報の組を、以降の処理において文書内配列パターン変換部 1 3 0 による変換を繰り返し行わせる対象から除外してもよい。即ち、表示画面に表示可能な情報のサイズとは、利用者が情報を閲覧するべく表示画面をスクロールする量が予め定めた量以下となる情報のサイズであってもよい。

【 0 0 6 7 】

一方、予め定められた終了条件を満たした場合に ( S 4 6 0 : Y E S )、反復終了判断部 1 4 0 は、処理を終了する。例えば、反復終了判断部 1 4 0 は、S 4 5 0 において変換対象の候補の何れも変換されなかった場合に、予め定められた終了条件を満たしたと判断してもよい。

【 0 0 6 8 】

図 5 は、図 4 の S 4 5 0 におけるフローチャートを示す。文書内配列パターン変換部 1 3 0 は、変換対象の候補となる要素識別情報の組のそれぞれについて、以下の処理を行う。文書内配列パターン変換部 1 3 0 は、変換対象の候補となる要素識別情報の組が、文書情報において基準頻度以上出現していない場合に ( S 5 0 0 : N O )、当該要素識別情報の組を変換対象から除外するべく処理を終了する。

【 0 0 6 9 】

一方、変換対象の候補となる要素識別情報の組が、文書情報において基準頻度以上出現した場合に ( S 5 0 0 : Y E S )、変換対象の候補となる要素識別情報の組における変換の階層が、所定の値以下であるか否かを判断する ( S 5 1 0 )。ここで、要素識別情報の組における変換の階層とは、例えば、変換前の文書構成要素を当該要素識別情報の組に変換するべく文書内配列パターン変換部 1 3 0 が文書情報に対して繰り返す変換の回数をいう。

【 0 0 7 0 】

変換対象の候補となる要素識別情報の組における変換の階層が、所定の値を超える場合に ( S 5 1 0 : N O )、文書内配列パターン変換部 1 3 0 は、当該要素識別情報の組を変換対象から除外するべく処理を終了する。一方、変換対象の候補となる要素識別情報の組における変換の階層が、所定の値以下である場合に ( S 5 1 0 : Y E S )、文書内配列パターン変換部 1 3 0 は、変換対象の候補となる要素識別情報の組について、当該要素識別情報の組に変換される前の文書構成要素の合計サイズが、予め定められた基準サイズより小さいか否かを判断する ( S 5 2 0 )。

10

20

30

40

50

## 【 0 0 7 1 】

当該要素識別情報の組に変換される前の文書構成要素の合計サイズが、予め定められた基準サイズ以上の場合に（S 5 2 0 : N O）、文書内配列パターン変換部 1 3 0 は、当該要素識別情報の組を変換対象から除外するべく処理を終了する。一方、当該要素識別情報の組に変換される前の文書構成要素の合計サイズが、予め定められた基準サイズより小さい場合に（S 5 2 0 : Y E S）、文書内配列パターン変換部 1 3 0 は、変換対象の候補となる要素識別情報の組を、変換対象と決定する。即ち、文書内配列パターン変換部 1 3 0 は、文書情報に出現する当該要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する（S 5 3 0）。

## 【 0 0 7 2 】

このように、文書内配列パターン変換部 1 3 0 は、基準頻度以上で繰り返し出現する要素識別情報の組のうち、当該要素識別情報の組に変換される前の文書構成要素の合計サイズが予め定められた基準サイズより小さい要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する。これにより、グループ内の文書構成要素のサイズが大きくなりすぎるのを防ぐことができる。例えば、情報処理装置 1 0 が P D A 等の携帯情報通信装置であり、表示画面の大きさが限定される場合であっても、情報処理装置 1 0 は、当該表示画面に表示可能なグループに分類することができる。

## 【 0 0 7 3 】

また、文書内配列パターン変換部 1 3 0 は、基準頻度以上で繰り返し出現する要素識別情報の組であって、変換前の文書構成要素を当該要素識別情報の組に変換するべく文書内配列パターン変換部 1 3 0 が文書情報に対して繰り返す変換の回数が基準回数以下の要素識別情報の組を、配列パターンを示す要素識別情報に変換する。これにより、グループ内の文書構成要素のサイズが大きくなりすぎるのを防ぐことができる。更に、文書内配列パターン変換部 1 3 0 は、分類すべきグループにおける変換の階層が予め特定できている等の場合には、要素識別情報の組を適切に変換することができる。

## 【 0 0 7 4 】

なお、本図に示した処理の順序は一例であり、情報処理装置 1 0 は、本図の例に代えて、他の順序で判断を行ってもよい。また、情報処理装置 1 0 は、本図に示した判断の少なくとも一部を行わなくともよい。

## 【 0 0 7 5 】

例えば、当該要素識別情報の組に変換される前の文書構成要素の合計サイズが、予め定められた基準サイズ以上の場合であっても（S 5 2 0 : N O）、文書内配列パターン変換部 1 3 0 は、当該要素識別情報の組を変換対象としてもよい。この場合、文書情報全体が一のグループとして分類されることとなる。これにより、より精度が高くかつ柔軟なグループ分けを実現することができる場合がある。以下に例を示す。

## 【 0 0 7 6 】

情報処理装置 1 0 は、文書内配列パターン変換部 1 3 0 による変換の経過を示す情報を利用者に出力する。例えば、情報処理装置 1 0 は、変換処理のそれぞれについて、変換元である要素識別情報の組と、変換先である要素識別情報とを対応付けて、変換の経過を示す情報として階層構造により出力する。利用者は、当該階層構造の中から、実際にグループとして分割することを希望する要素識別情報を選択する。

## 【 0 0 7 7 】

これにより、利用者による入力が必要となるものの、基準サイズによりグループを決定する場合と比較して、より適切かつ柔軟にグループ分けを実現できる。また、基準サイズ以下のグループにおいて当該グループのタイトル情報が特定できなかった場合であっても、基準サイズ以上のグループにおいて当該グループのタイトル情報が特定できた場合には、このタイトル情報を利用して基準サイズ以下のグループのタイトル情報を特定し、グループ分けの精度を高めることができる場合がある。

## 【 0 0 7 8 】

また、情報処理装置 1 0 は、本図に示した判断に対して、更に他の判断を追加して行って

10

20

30

40

50



もよい。例えば、文書内配列パターン変換部 130 は、要素識別情報の組として検出するべきでない配列パターンを、利用者からの入力等に基づいて定めてもよい。即ち、文書内配列パターン変換部 130 は、当該配列パターンと一致する要素識別情報の組を、変換対象から除外する。

【0079】

要素識別情報の組として検出するべきでない配列パターンとは、例えば、表示領域の区切りように用いられる所定の形状の画像を含む配列パターンであってもよいし、所定以上の面積を有する空白を含む配列パターンであってもよい。また、文書内配列パターン変換部 130 は、要素識別情報の組として検出するべきでない配列パターンが否かを、変換処理の毎に利用者に問い合わせてもよいし、予め利用者に入力させておいてもよい。

10

【0080】

また、文書内配列パターン変換部 130 は、図 4 及び図 5 に示した方法に代えて、他の方法を用いて、要素識別情報の組を検出し、当該要素識別情報の組の配列パターンを示す要素識別情報に変換してもよい。即ち、文書内配列パターン変換部 130 が要素識別情報の組を検出する方法は、基準頻度以上出現する要素識別情報の組を検出する方法であればよい。

【0081】

図 6 は、図 2 の S230 におけるフローチャートを示す。文書間配列パターン変換部 160 は、対象文書情報及び関連文書情報のそれぞれから、文書内配列パターン変換部 130 により変換された要素識別情報を除外した文書を検出対象とする (S600)。

20

【0082】

より具体的には、文書間配列パターン変換部 160 は、文書内配列パターン変換部 130 により変換された要素識別情報を、その旨を示す要素識別情報である <group /> に変換してもよい。これにより、文書間配列パターン変換部 160 は、以降の処理で、文書内配列パターン変換部 130 により変換された要素識別情報を適切に検出し、除外することができる。

【0083】

続いて、文書間配列パターン変換部 160 は、特定した当該要素識別情報の組のうち、当該対象文書情報及び当該関連文書情報を併せた文書において、要素識別情報を、出現する頻度の高い順に、順次被選択情報として選択する (S610)。ここで、文書間配列パターン変換部 160 は、当該対象文書情報及び当該関連文書情報を併せた文書に対して、当該対象文書情報及び当該関連文書情報の境界を示す情報を挿入する。そして、文書間配列パターン変換部 160 は、被選択情報を含む要素識別情報の組を、変換対象の候補として選択する (S620)。

30

【0084】

具体的には、文書間配列パターン変換部 160 は、被選択情報及び間隙構成要素を、要素識別情報の組として選択する。ここで、文書間配列パターン変換部 160 は、文書内配列パターン変換部 130 により既に要素識別情報に変換された旨を示す <group /> を、一の文書構成要素の組と、他の文書構成要素の組との境界として用いてもよい。また、文書間配列パターン変換部 160 は、対象文書情報及び関連文書情報の境界を示す情報を、一の文書構成要素の組と、他の文書構成要素の組との境界として用いてもよい。

40

【0085】

変換対象の候補である要素識別情報の組が、予め定められた基準頻度以上出現する場合に (S630: YES)、文書間配列パターン変換部 160 は、当該要素識別情報の組を、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する (S640)。続いて、対象文書情報及び関連文書情報を併せた文書において未だ被選択情報として選択していない要素識別情報がある場合に (S650: NO)、文書間配列パターン変換部 160 は、S610 に処理を戻すことにより、次に出現頻度の高い要素識別情報を被選択情報として選択する。対象文書情報及び関連文書情報を併せた文書における全ての要素識別情報を被選択情報として選択した場合に (S650: YES)、文書間配列パターン変換部 1

50

60は、処理を終了する。

【0086】

なお、文書間配列パターン変換部160が用いる基準頻度は、文書内配列パターン変換部130が用いる基準頻度とは異なってもよい。即ち、文書間配列パターン変換部160は、文書内配列パターン変換部130が判断に用いる基準頻度である第1頻度とは異なる第2頻度以上、変換対象の候補である要素識別情報の組が出現する場合に、要素識別情報の組を、当該要素識別情報の組の配列パターンを示す要素識別情報に変換してもよい。

【0087】

図7は、図2のS240におけるフローチャートを示す。グループ分類部170は、文書内配列パターン変換部130により対象文書情報を繰り返し変換した結果得られた文書情報において、文書内配列パターン変換部130により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループとして分類する(S700)。同様に、グループ分類部170は、文書間配列パターン変換部160により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループとして分類する。

10

【0088】

更に、グループ分類部170は、分類したグループのそれぞれについて、当該グループ内の文書構成要素が文書情報において果たす役割を示すタイトル情報を生成し、目次情報出力部190、グループ識別情報生成部200、及びアノテーション出力部202に出力してもよい。より詳しくは、当該グループに分類された要素識別情報の組が、コメント情報を境界として検出されたものである場合には、グループ分類部170は、当該コメント情報に含まれる情報をタイトル情報として生成してもよい。

20

【0089】

また、グループ分類部170は、分類したグループのそれぞれについて、文書情報における当該グループの重要度を示す重要度情報を更に生成し、並替出力部180、グループ識別情報生成部200、及びアノテーション出力部202に出力してもよい。具体的には、グループ分類部170は、当該要素識別情報の組が文書情報のどこに位置するかを示す配置情報、当該要素識別情報の組に変換された文書構成要素の色の情報、当該文書構成要素の大きさの情報、及び当該文書構成要素が文字列である場合の文字列の内容に基づいて、重要度情報を生成してもよい。

30

【0090】

更に、グループ分類部170は、文書情報のうちグループに分類した文書構成要素を除外した部分、即ち、グループとして分類できなかった部分を、その旨を示すグループに分類する。この場合、グループ分類部170は、グループとして分類できなかった部分を、更に、文書情報内の位置、表示画面上での位置、及び表示画面上での背景色に応じて、更に、複数のグループに分類してもよい。

【0091】

並替出力部180は、対象文書情報及び関連文書情報のそれぞれにおいて、複数の文書構成要素を、グループ分類部170により分類されたグループ毎に並び替える(S710)。そして、目次情報出力部190は、グループ分類部170により分類された各グループについて、当該グループに含まれる文書構成要素に対する目次を示す目次情報を生成する(S720)。更に、グループ識別情報生成部200は、グループ分類部170により分類された各グループに対応付けて、当該グループを音声により識別することを容易にする為の情報を文書情報内に生成する(S730)。

40

【0092】

アノテーション出力部202は、グループ分類部170により分類された各グループを識別するアノテーション情報を出力する(S740)。例えば、アノテーション出力部202は、複数のグループのそれぞれについて、当該グループに属する文書構造情報が、対象文書情報のどの領域に位置するかを示す配置情報を、アノテーション情報として出力する。一例としては、アノテーション出力部は、XPath又はXPointer等の技術に

50

より、配置情報を生成し出力してもよい。

【0093】

これにより、アノテーション情報を入力とする他の装置等は、対象文書情報中の情報の並び替え、目次作成、又はグループ識別情報の出力のみならず、対象文書情報のダイジェストを作成する等の、他のトランスコーディング技術を用いた処理を行うことができる。

【0094】

また、情報処理装置10は、アノテーション情報の作成者による作業を支援することができる。例えば、アノテーション情報の作成者は、情報処理装置10により出力されたアノテーション情報に基づき、当該アノテーション情報を修正又は変更することにより、所望のアノテーション情報を作成することができる。この結果、アノテーション情報の作成者は、初めからアノテーション情報を作成する場合と比較して、効率的に所望のアノテーション情報を作成することができる。

10

【0095】

図8は、構成要素変換部105が文書構成要素を要素識別情報に変換する一例を示す。本図に示すように、構成要素変換部105は、文書構成要素を、当該文書構成要素の内容、データサイズ、又は格納位置等を示す要素識別情報に変換する。具体的には、構成要素変換部105は、文書構成要素が数値である場合に、当該文書構成要素を、数値を示す要素識別情報である<digit />に変換する。なお、好ましくは、構成要素変換部105は、変換先の要素識別情報として、XML(eXtensible Markup Language)の規格に準拠した形式のタグ情報を生成する。これにより、XML文書を操作する既存のプログラムを利用することができる。

20

【0096】

また、構成要素変換部105は、文書構成要素が100文字以上の文字列である場合に、当該文書構成要素を、長い文字列を示す要素識別情報である<longtext />に変換する。同様に、構成要素変換部105は、30から100文字の文字列、2から30文字の文字列、及び1文字の文字のそれぞれを、要素識別情報である<midtext />、<shorttext />、及び<letter />のそれぞれに変換する。

【0097】

これに代えて、構成要素変換部105は、文字列を表すタグ情報に、当該文字列の長さを示す属性を対応付けた情報を、要素識別情報として生成してもよい。例えば、構成要素変換部105は、<longtext />に代えて、所定長より長くかつ数値でない文字列を示す要素識別情報である<text length = "long" is\_\_digit = "no">を生成してもよい。更に、構成要素変換部105は、文字列の長さを示す情報を属性としたタグ情報を、要素識別情報として生成してもよい。

30

【0098】

このように、構成要素変換部105は、文字列である文書構成要素を、当該文字列が数字であるか数字以外であるかに応じて定まる要素識別情報に変換する。更に、構成要素変換部105は、文字列である文書構成要素を、当該文字列に含まれる語句に応じて定まる要素識別情報に変換してもよい。

【0099】

また、構成要素変換部105は、表示画面に表示する画像を識別する画像識別情報であって、当該画像が対象文書情報と同一のサイト内にあり、かつ当該画像のサイズが300×300ピクセル以上である画像識別情報を、要素識別情報である<in-largeimg />に変換する。一例として、対象文書情報がHTML文書である場合の画像識別情報とは、<img src = "AAA.JPG">等のイメージタグ情報である。なお、構成要素変換部105は、画像の表示サイズを特定するべく、表示するべき画像をレンダリング及び/又は解析してもよい。より具体的には、構成要素変換部105は、画像ファイルのヘッダ部分に記録されたサイズ情報を解析してもよい。

40

【0100】

更に、構成要素変換部105は、レンダリング等により、表示画面において画像が表示さ

50

れる位置又は画像の色彩等の情報を解析してもよい。この場合、構成要素変換部105は、画像識別情報を、当該画像の表示位置又は色彩等を示す要素識別情報に変換してもよい。

#### 【0101】

同様に、構成要素変換部105は、画像識別情報であって、当該画像識別情報における画像が対象文書情報と同一サイト内にあり、かつ当該画像のサイズが100×100から300×300ピクセルである画像識別情報を、要素識別情報である<in-midimg />に変換する。また、構成要素変換部105は、画像識別情報であって、当該画像識別情報における画像が対象文書情報と同一サイト内にあり、かつ当該画像のサイズが100×100ピクセル以下である画像識別情報を、要素識別情報である<in-small 10  
img />に変換する。同様に、構成要素変換部105は、対象文書情報と異なるサイト内にある画像を識別する画像識別情報を、当該画像の大きさに応じて、<out-large 10  
img />、<out-midimg>、又は<out-small 10  
img />に変換する。

#### 【0102】

これに代えて、構成要素変換部105は、画像を表すタグ情報に、当該画像の大きさを示す属性を対応付けた情報を、要素識別情報として生成してもよい。例えば、構成要素変換部105は、<in-large 10  
img />に代えて、所定範囲内のデータサイズかつサイト内に存在する画像を示す要素識別情報である<image size = "large" 20  
location = "in" />を生成してもよい。

#### 【0103】

このように、構成要素変換部105は、文書情報における各文書構成要素を、当該文書構成要素のデータサイズに応じた要素識別情報に変換する。また、好ましくは、構成要素変換部105は、文字列又は画像の種類に応じて異なる要素識別情報に変換する。例えば、構成要素変換部105は、文字列の文字数が同一であっても、文字列において数字、漢字、ひらがな、カタカナ、及びアルファベットのそれぞれが占める割合に応じて、異なる要素識別情報に変換してもよい。また、構成要素変換部105は、画像のデータサイズが同一であっても、画像が縦長又は横長であるかに応じて異なる要素識別情報に変換してもよいし、画像のデータ形式に応じて異なる要素識別情報に変換してもよい。

#### 【0104】

また、文書情報は、利用者に表示する表示情報と、当該表示情報に対する利用者からの指示に応じて表示するべき他の情報の格納位置を示す格納位置情報とを文書構成要素として含んでいる。例えば、文書情報がHTML文書である場合の、表示情報及び格納位置情報とは、利用者からの指示を受けて他の情報を表示させるハイパーリンクである。一例として、格納位置情報とは、<a href = "BBB.HTML" >等のタグ情報及び</a >等のタグ情報であり、表示情報とは、当該<a >タグ及び</a >タグの間に配列された情報である。 30

#### 【0105】

構成要素変換部105は、当該格納位置情報を、当該格納位置情報が示す格納位置の範囲を示す要素識別情報に変換する。例えば、構成要素変換部105は、リンク先の他の情報が、対象文書情報と同一のディレクトリ内に格納されている場合に、当該他の情報の格納位置を示す格納位置情報を、格納位置の範囲を示す要素識別情報である<samedir-link >及び</samedir-link >に変換する。より具体的には、構成要素変換部105は、格納位置情報であると共に表示情報の開始点を示すタグ情報である<a href = "BBB.HTML" >等を<samedir-link >に変換する。また、構成要素変換部105は、表示情報の終了点を示す</a >を、</samedir-link >に変換する。 40

#### 【0106】

これに代えて、構成要素変換部105は、リンクを表すタグ情報に、当該リンクのリンク先を示す属性を対応付けた情報を、要素識別情報として生成してもよい。例えば、構成要 50

素変換部105は、`< s a m e d i r - l i n k >`及び`< / s a m e d i r - l i n k >`に代えて、`< l i n k l o c a t i o n = " s a m e " t a r g e t = " n o n e " / >`を生成してもよい。

【0107】

また、構成要素変換部105は、`< a h r e f = " B B B . H T M L " >`及び`< / a >`間に配列された表示情報を変換しなくともよい。また、構成要素変換部105は、当該表示情報に付いては、別途他の種類を示す要素識別情報に変換してもよい。

【0108】

また、構成要素変換部105は、リンク先の他の情報が、対象文書情報と同一サイト内に格納されている場合に、当該他の情報の格納位置を示す格納位置情報を、格納位置の範囲を示す要素識別情報である`< i n - l i n k >`及び`< / i n - l i n k >`に変換する。

10

【0109】

また、構成要素変換部105は、リンク先の他の情報が、対象文書情報と異なるサイト内に格納されている場合に、当該他の情報の格納位置を示す格納位置情報を、格納位置の範囲を示す要素識別情報である`< o u t - l i n k >`及び`< / o u t - l i n k >`に変換する。

【0110】

このように、構成要素変換部105は、文書構成要素の種類又は役割を、タグ情報における属性のパラメータに応じて検出してもよい。そして、構成要素変換部105は、当該文書構成要素を、当該文書構成要素のパラメータの種類等を示す要素識別情報に変換する。

20

【0111】

このように、構成要素変換部105は、文書情報における複数の文書構成要素のそれぞれを、当該文書構成要素の種類又は役割を示す要素識別情報にそれぞれ変換する。

【0112】

図9は、構成要素変換部105が所定の要素識別情報の組を要素識別情報に変換する例を示す。構成要素変換部105は、図8に示した変換に続いて、予め定めた条件に適合した要素識別情報の組を、要素識別情報に変換する。例えば、本図は、構成要素変換部105が、野球等のスポーツにおけるスコアボードを形成する文書構成要素の組を検出し、スコアボードを示す要素識別情報、例えば、`< b a s e b a l l - s c o r e b o a r d / >`に変換する例を示している。

30

【0113】

構成要素変換部105は、文書情報において、縦3列×横13行以上の表を構成し、かつ以下の各条件を満たす文書構成要素の組を検出する。まず、構成要素変換部105は、1行目の2列目から数値を示す要素識別情報である`< d i g i t / >`が9個以上連続し、かつそれ以降に`< s h o r t t e x t / >`又は`< l e t t e r / >`が3個連続することを条件とする。更に、構成要素変換部105は、2行目及び3行目の1列目は、`< s h o r t t e x t / >`若しくは`< l e t t e r / >`又はリンク情報に対応付けられた`< s h o r t t e x t / >`若しくは`< l e t t e r / >`であることを条件とする。

【0114】

更に、構成要素変換部105は、2行目及び3行目のそれぞれにおいて、`< d i g i t / >`が2列目から12個以上連続することを条件とする。但し、`< d i g i t / >`が連続する個数は、1行目において連続する`< d i g i t / >`の個数に3を加えた数であることを条件とする。なお、構成要素変換部105は、野球のワールドゲーム又は途中経過等を示すスコアボードを検出するべく、他の条件を用いてもよい。

40

【0115】

このように、構成要素変換部105は、データの内容の詳細が日々変化する文書構成要素の組であっても、予め定めたパターンに適合するか否かを判定することにより、当該文書構成要素の組が示すデータの種別を適切に判別し、文書構成要素に変換することができる。

【0116】

50

図10(a)は、文書情報により表示された表示画面の一例を示す。図10(b)は、図10(a)に示す表示画面を表示した文書情報に対して、文書内配列パターン変換部130による変換を行った結果を示す。

【0117】

表示部205は、箇条書きのアイテムの開始点を示す画像と、箇条書きの内容を示す文字列とを対応付けて表示する。例えば、表示部205は、社会、スポーツ、金融・経済、及び政治という文字列のそれぞれを、画像に対応付けて表示している。更に、これらの文字列は、ハイパーリンクとなっており、即ち、これらの文字列に対して利用者から指示を受けた場合には、表示部205は、他の情報を表示する。

【0118】

構成要素変換部105は、箇条書きのアイテムの開始点を示す画像を識別する画像識別情報を、箇条書きのアイテムの開始点を示す要素識別情報である<bullet />に変換する。また、構成要素変換部105は、表示画面に表示すべき文字列を、当該文字列が2文字以上20文字以内であることを示す要素識別情報である<shorttext />に変換する。更に、構成要素変換部105は、当該文字列に対して利用者から指示を受けた場合のリンク先を示す格納位置情報を、当該格納位置情報の示す格納位置が対象文書情報と同一サイト内である旨を示す要素識別情報である<in-link>及び</in-link>に変換する。また、構成要素変換部105は、当該文字列に続く文字列を次の行に表示することを示す文書構成要素、例えば、<br>タグを、その旨を示す要素識別情報である<new-line />に変換する。

【0119】

なお、要素識別情報の一例である<bullet />とは、箇条書きのアイテムの開始点を示す情報であればよく、本例に示した所定サイズ以下の画像の他、特定の条件を満たす<letter />又は<digit />等から構成されていてもよい。

【0120】

構成要素変換部105により変換された対象文書情報について、文書内配列パターン変換部130は、要素識別情報の組である<bullet />、<in-link>、<shorttext />、</in-link>、及び<new-line />を検出する。そして、文書内配列パターン変換部130は、当該要素識別情報の組が基準頻度以上で繰り返し出現する場合に、当該要素識別情報の組の配列パターンを示す要素識別情報、例えば、箇条書きのアイテムを示す要素識別情報である<itemizedlink />等に変換する。

【0121】

グループ分類部170は、文書内配列パターン変換部130により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループとして分類する。例えば、グループ分類部170は、<itemizedlink />に変換された文書構成要素を、箇条書きを示すグループとして分類する。好ましくは、グループ分類部170は、分類したグループのそれぞれについて、当該グループの果たす役割を特定し、当該グループに対応付けて出力する。例えば、グループ分類部170は、2文字以上所定の文字数(例えば、30)以下の文字列を示す要素識別情報、即ち、<shorttext />から構成されるリンク情報のグループについて、当該グループがリンクリストを示すと判断する。

【0122】

このように、文書内配列パターン変換部130は、要素識別情報の出現頻度及び配列順に基づいて、要素識別情報の組を変換することができる。

【0123】

なお、本図の例においては、文書内配列パターン変換部130が、基準頻度より高い頻度で出現する文書構成要素の組を検出したが、これに代えて、構成要素変換部105が、これらの文書構成要素の組を、予め定めた条件に適合したと判断して変換してもよい。例えば、構成要素変換部105は、<bullet />、<in-link>、<short

10

20

30

40

50

t t e x t / >、< / i n - l i n k >、及び< b r >がこの順に配列されており、かつこの配列が連続している場合に、予め定めた条件に適合したと判断してもよい。この場合、構成要素変換部 105 は、これらの要素識別情報の組を、リンクリストを示す< s i t e - i n d e x / >に変換する。本条件を含む条件の一例を正規表現として以下に示す。

【 0 1 2 4 】

(リンクリストを示す正規表現)

[ [ < b u l l e t / > ] ? ( < s a m e d i r - l i n k > < s h o r t t e x t / > < / s a m e d i r - l i n k > ) | ( < i n - l i n k > < s h o r t t e x t / > < / i n - l i n k > ) [ < b r > | < p > ] + ]

10

【 0 1 2 5 】

このように、構成要素変換部 105 は、上記の条件を満たす要素識別情報の組を、リンクリストを示す要素識別情報である< s i t e - i n d e x / >に変換する。

【 0 1 2 6 】

更に他の例としては、構成要素変換部 105 は、< b u l l e t / >、< i n - l i n k >、< s h o r t t e x t / >、及び< / i n - l i n k >から構成されるテーブル又はリストを、リンクリストを示す< s i t e - i n d e x / >に変換してもよい。

【 0 1 2 7 】

図 11 ( a ) は、文書情報により表示された表示画面の他の例を示す。図 11 ( b ) は、図 11 ( a ) に示す表示画面を表示した文書情報に対して、文書内配列パターン変換部 130 による変換を行った結果を示す。

20

【 0 1 2 8 】

表示部 205 は、箇条書きのアイテムの開始点を示す画像と、箇条書きの内容を示す文字列とを対応付けて表示する。例えば、表示部 205 は、新聞記事の見出しを示す文字列のそれぞれを、箇条書きを示す記号(例えば、「・」)に対応付けて表示している。更に、これらの文字列は、ハイパーリンクとなっており、これらの文字列に対して利用者から指示を受けた場合には、表示部 205 は、他の情報を表示する。

【 0 1 2 9 】

構成要素変換部 105 は、箇条書きのアイテムの開始点を示す黒点を、箇条書きのアイテムの開始点を示す要素識別情報である< b u l l e t / >に変換する。また、構成要素変換部 105 は、表示画面に表示すべき文字列を、当該文字列が 30 文字以上 100 文字以内であることを示す要素識別情報である< m i d t e x t / >に変換する。更に、構成要素変換部 105 は、当該文字列に対して利用者から指示を受けた場合のリンク先を示す格納位置情報を、当該格納位置情報の示す格納位置が対象文書情報と同一サイト内である旨を示す要素識別情報である< i n - l i n k >及び< / i n - l i n k >に変換する。また、構成要素変換部 105 は、当該文字列に続く文字列を次の行に表示することを示す文書構成要素、例えば、< b r >タグを、その旨を示す要素識別情報である< n e w - l i n e / >に変換する。

30

【 0 1 3 0 】

構成要素変換部 105 により変換された対象文書情報について、文書内配列パターン変換部 130 は、要素識別情報の組である< b u l l e t / >、< i n - l i n k >、< m i d t e x t / >、< / i n - l i n k >、及び< n e w - l i n e / >を検出する。そして、文書内配列パターン変換部 130 は、当該要素識別情報の組が基準頻度以上で繰り返し出現する場合に、当該要素識別情報の組の配列パターンを示す要素識別情報、例えば、箇条書きのアイテムを示す要素識別情報である< l o n g i t e m i z e d l i n k / >等に変換する。

40

【 0 1 3 1 】

グループ分類部 170 は、文書内配列パターン変換部 130 により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の文書構成要素をグループと

50

して分類する。例えば、グループ分類部 170 は、`<longitemizedlink />` に変換された文書構成要素を、箇条書きを示すグループとして分類する。好ましくは、グループ分類部 170 は、分類したグループのそれぞれについて、当該グループの果たす役割を特定し、当該グループに対応付けて出力する。例えば、グループ分類部 170 は、文字数が所定範囲内（例えば、30 から 100 文字）の文字列を示す要素識別情報、即ち、`<midtext />` から構成されるリンク情報のグループであって、当該文字列の終端に時刻を示す文字列が配列されているグループについて、当該グループが新聞等のヘッドラインを示すと判断する。

#### 【0132】

図 12 は、変換対象である対象文書情報の表示例を示す。図 13 は、図 12 に示す対象文書情報において、グループ分類部 170 による分類の一例を示す。図 12 における対象文書情報は、新聞記事を示す HTML 文書である。図 13 において実線の矩形のそれぞれは、文書構成要素を示している。また、丸印は、要素識別情報を示している。また、点線で囲まれた領域は、グループ分類部 170 により分類されるグループを示している。また、図 13 においては関連文書情報の一部についても、グループ分類部 170 による分類の一例を図示する。

10

#### 【0133】

情報処理装置 10 は、対象文書情報として、新聞の社会面を示すオブジェクトを表示させる文書構成要素の組 1200 と、それぞれが社会面の記事の詳細に対するリンクでありかつ当該記事のヘッドラインを示す文書構成要素の組 1210、文書構成要素の組 1220、文書構成要素の組 1230、文書構成要素の組 1240、及び文書構成要素の組 1250 とを取得する。

20

#### 【0134】

例えば、文書構成要素の組 1210 は、箇条書きを示す文書構成要素である「・」と、記事のヘッドラインを示す文字列である文書構成要素と、当該文字列から新聞記事へのリンク情報を示す文書構成要素、例えば、`<a>` タグと、文字列の改行を示す文書構成要素、例えば、`<br>` タグとを含む。文書構成要素の組 1220 から文書構成要素の組 1250 のそれぞれは、文書構成要素の組 1210 と略同一であるので説明を省略する。

#### 【0135】

更に、情報処理装置 10 は、対象文書情報として、新聞の政治面を示すオブジェクトを表示させる文書構成要素の組 1270 と、政治面の記事の詳細に対するリンクでありかつ当該記事のヘッドラインを示す文書構成要素の組と、新聞の各面へのリンクの集合を表示させる文書構成要素の組 1290 と、当該新聞のロゴを示す文書構成要素 1292 と、当該新聞記事における著作権の帰属を示す文書構成要素 1295 とを取得する。

30

#### 【0136】

構成要素変換部 105 は、文書構成要素の組 1200 を、新聞記事のジャンルを示す要素識別情報に変換する。そして、構成要素変換部 105 は、文書構成要素の組 1210 を、箇条書きの先頭記号を示す要素識別情報である `<bullet />` と、サイズが所定範囲内の文字列を示す要素識別情報である `<midtext />` と、新聞記事へのリンク、例えば、サイト内のリンクを示す要素識別情報である `<in-link />` と、改行を示す要素識別情報である `<new-line />` とに変換する。

40

#### 【0137】

続いて、文書内配列パターン変換部 130 は、これらの要素識別情報の組を、箇条書きのアイテムを示す要素識別情報である `<itemizedlink />` に変換する。また、文書構成要素の組 1210 から文書構成要素の組 1250 が連続して配置されているので、文書内配列パターン変換部 130 は、`<itemizedlink />` の組を、ヘッドラインの集合を示す要素識別情報 1260 に変換する。更に、文書内配列パターン変換部 130 は、要素識別情報 1260 及び新聞記事のジャンルを示す要素識別情報を、ジャンル毎の新聞記事のヘッドラインを示す要素識別情報 1265 に変換する。

#### 【0138】

50



新聞の政治面の記事についても同様に、構成要素変換部105は、「・」、文字列、リンク、及び改行等の文書構成要素を、<bullet />、<midtext >、<in-link />、及び<new-line>等の要素識別情報に変換する。そして、文書内配列パターン変換部130は、これらの要素識別情報を<itemizedlink />に変換する。

【0139】

続いて、文書内配列パターン変換部130は、箇条書きのアイテムを示す要素識別情報、例えば、<itemizedlink />の組を、ヘッドライン記事の集合を示す要素識別情報1280に変換する。また、文書内配列パターン変換部130は、要素識別情報1280及び新聞記事のジャンルを示す要素識別情報を、ジャンル毎の新聞記事のヘッドラインを示す要素識別情報1285に変換する。

10

【0140】

グループ分類部170は、要素識別情報1265及び要素識別情報1285のそれぞれに変換された複数の文書構成要素を、ジャンル毎の新聞記事のヘッドラインを示すグループ1300及びグループ1310に分類する。

【0141】

また、情報処理装置10は、サイト内リンクを表示させる文書構成要素の組1290を、サイト内リンクを示すグループ1320に分類するべく以下の処理を行う。まず、構成要素変換部105は、それぞれが文書構成要素の組を形成する、「社会」、「政治」、及び途中を一部省略して「ひと」を検出する。文書構成要素の組のそれぞれは、文書構成要素である文字列と、文書構成要素であるリンク情報とを有している。そして、構成要素変換部105は、文字列を、所定の文字数以下の文字列を示す<shorttext />に変換し、リンク情報を、同一サイト内へのリンクを示す<in-link>及び</in-link>に変換し、文字列を区切る記号「|」を、<bullet />に変換する。

20

【0142】

文書内配列パターン変換部130は、<in-link>、<shorttext />、</in-link>、及び<bullet />から構成される要素識別情報の組が、対象文書情報において基準頻度以上で繰り返し出現すると判断する。これにより、グループ分類部170は、リンクリストを示す文字列から構成される文書構成要素の組1295、例えば、「社会|政治|...|ひと|」のように配列されたリンク情報の組を、リンクリストを示すグループとして分類することができる。

30

【0143】

また、情報処理装置10は、フッタを示す文字列及び画像を、フッタ部を示すグループ1330aに分類するべく以下の処理を行う。まず、構成要素変換部105は、文書構成要素1295を、所定のキーワード、例えば、「著作権」を含む文字列を示す要素識別情報である<copyright />に変換する。そして、構成要素変換部105は、文書構成要素1292を、サイズが所定範囲内の画像を示す要素識別情報である<midimage />に変換する。

【0144】

40

文書間配列パターン変換部160は、<copyright />及び<midimage />から構成される要素識別情報の組が、対象文書情報及び関連文書情報を併せた文書において基準頻度以上で繰り返し出現すると判断する。そして、文書間配列パターン変換部160は、この要素識別情報の組を、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する。

【0145】

グループ分類部170は、文書間配列パターン変換部160により変換された要素識別情報をグループとして分類する。例えば、グループ分類部170は、対象文書情報における文書構成要素1292及び文書構成要素1295を、フッタ部を示すグループ1330aに分類する。また、グループ分類部170は、関連文書情報における文書構成要素139

50

2及び文書構成要素1395を、フッタ部を示すグループ1330bに分類してもよい。更に、グループ分類部170は、これらのグループのタイトル情報が、フッタ部である旨を、当該グループに含まれる要素識別情報、例えば、<copyright />に基づいて特定してもよい。

【0146】

図14は、情報処理装置10のハードウェア構成の一例を示す。情報処理装置10は、ホストコントローラ1082により相互に接続されるCPU1000、RAM1020、グラフィックコントローラ1075、及び表示装置1080を有するCPU周辺部と、入出力コントローラ1084によりホストコントローラ1082に接続される通信インターフェイス1030、ハードディスクドライブ1040、及びCD-ROMドライブ1060を有する入出力部と、入出力コントローラ1084に接続されるROM1010、フレキシブルディスクドライブ1050、及び入出力チップ1070を有するレガシー入出力部とを備える。

10

【0147】

ホストコントローラ1082は、RAM1020と、高い転送レートでRAM1020をアクセスするCPU1000及びグラフィックコントローラ1075とを接続する。CPU1000は、ROM1010及びRAM1020に格納されたプログラムに基づいて動作し、各部の制御を行う。グラフィックコントローラ1075は、CPU1000等がRAM1020内に設けたフレームバッファ上に生成する画像データを取得し、表示装置1080上に表示させる。これに代えて、グラフィックコントローラ1075は、CPU1000等が生成する画像データを格納するフレームバッファを、内部に含んでもよい。

20

【0148】

入出力コントローラ1084は、ホストコントローラ1082と、比較的高速な入出力装置である通信インターフェイス1030、ハードディスクドライブ1040、及びCD-ROMドライブ1060を接続する。通信インターフェイス1030は、ネットワークを介して他の装置と通信する。ハードディスクドライブ1040は、情報処理装置10が使用するプログラム及びデータを格納する。CD-ROMドライブ1060は、CD-ROM1095からプログラム又はデータを読み取り、RAM1020を介して入出力チップ1070に提供する。

【0149】

また、入出力コントローラ1084には、ROM1010と、フレキシブルディスクドライブ1050や入出力チップ1070等の比較的低速な入出力装置とが接続される。ROM1010は、情報処理装置10の起動時にCPU1000が実行するブートプログラムや、情報処理装置10のハードウェアに依存するプログラム等を格納する。フレキシブルディスクドライブ1050は、フレキシブルディスク1090からプログラム又はデータを読み取り、RAM1020を介して入出力チップ1070に提供する。入出力チップ1070は、フレキシブルディスク1090や、例えばパラレルポート、シリアルポート、キーボードポート、マウスポート等を介して各種の入出力装置を接続する。

30

【0150】

情報処理装置10に提供されるプログラムは、フレキシブルディスク1090、CD-ROM1095、又はICカード等の記録媒体に格納されて利用者によって提供される。プログラムは、記録媒体から読み出され、入出力チップ1070を介して情報処理装置10にインストールされ、情報処理装置10において実行される。

40

【0151】

情報処理装置10にインストールされて実行されるプログラムは、関連文書検出モジュールと、構成要素変換モジュールと、構成要素選択モジュールと、間隙構成要素検出モジュールと、文書内配列パターン変換モジュールと、変換是非入力モジュールと、反復処理モジュールと、文書間配列パターン変換モジュールと、グループ分類モジュールと、並替出力モジュールと、目次情報出力モジュールと、音声出力命令生成モジュールと、表示モジュールと、文書構造情報生成モジュールと、文書情報同一性出力モジュールとを含む。各

50

モジュールが情報処理装置 10 に働きかけて行わせる動作は、図 1 から図 13 において説明した情報処理装置 10 における、対応する部材の動作と同一であるから、説明を省略する。

#### 【0152】

以上に示したプログラム又はモジュールは、外部の記憶媒体に格納されてもよい。記憶媒体としては、フレキシブルディスク 1090、CD-ROM 1095 の他に、DVD や PD 等の光学記録媒体、MD 等の光磁気記録媒体、テープ媒体、IC カード等の半導体メモリ等を用いることができる。また、専用通信ネットワークやインターネットに接続されたサーバシステムに設けたハードディスク又は RAM 等の記憶装置を記録媒体として使用し、ネットワークを介してプログラムを情報処理装置 10 に提供してもよい。

10

#### 【0153】

以上、本実施形態から明らかなように、情報処理装置 10 は、文書情報に含まれる複数の文書構成要素を、文書構成要素の出現頻度に基づいて、グループに分類する。更に、情報処理装置 10 は、関連する文書情報との類似性に基づいて、文書構成要素を分類する。これにより、情報処理装置 10 は、構造が動的に変化する文書情報、例えば、日々更新されるウェブページについて、文書構成要素を適切に分類することができる。

#### 【0154】

以上に示した実施形態によると、以下の各項目に示す情報処理装置、プログラム、及び記録媒体が実現される。

#### 【0155】

(項目 1) 文書情報に含まれる複数の文書構成要素を、複数のグループに分類する情報処理装置であって、前記文書情報における前記複数の文書構成要素のそれぞれを、当該文書構成要素の種類又は役割を示す要素識別情報にそれぞれ変換する構成要素変換部と、前記構成要素変換部により変換された前記文書情報において、予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す前記要素識別情報に変換する文書内配列パターン変換部と、前記文書内配列パターン変換部により前記文書情報を繰り返し変換した結果得られた文書情報において、前記文書内配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類するグループ分類部とを備える情報処理装置。

20

(項目 2) 前記構成要素変換部により変換された前記文書情報において、複数の前記要素識別情報の何れかを選択する構成要素選択部と、前記構成要素選択部により選択された被選択情報と、前記文書情報において前記被選択情報の次に配列される前記被選択情報と同一の種類である前記要素識別情報との間に配列された間隙構成要素を検出する間隙構成要素検出部とを更に備え、前記文書内配列パターン検出部は、前記被選択情報及び前記間隙構成要素を前記要素識別情報の組として検出する項目 1 記載の情報処理装置。

30

#### 【0156】

(項目 3) 前記間隙構成要素検出部は、前記文書情報において複数の前記被選択情報の後に配列される終端構成要素を更に検出し、前記文書内配列パターン検出部は、前記複数の被選択情報のうち前記文書情報の最後に配列される被選択情報と、前記終端構成要素とを、前記要素識別情報の組として検出する項目 2 記載の情報処理装置。

40

(項目 4) 前記構成要素選択部は、前記構成要素変換部により変換された文書情報における前記要素識別情報を、出現する頻度の高い順に、順次前記被選択情報として選択し、前記間隙構成要素検出部は、前記構成要素選択部により選択された各被選択情報について、前記間隙構成要素を検出し、前記文書内配列パターン変換部は、前記構成要素選択部により順次選択される被選択情報について、前記被選択情報及び前記間隙構成要素を前記要素識別情報の組として検出する項目 2 記載の情報処理装置。

(項目 5) 前記文書内配列パターン変換部は、既に検出した前記要素識別情報の組より出現頻度が高いことを更に条件として、繰り返し出現する前記要素識別情報の組を検出する項目 4 記載の情報処理装置。

50

## 【 0 1 5 7 】

(項目6) 前記文書内配列パターン変換部は、前記配列パターンを示す前記要素識別情報に変換した変換元である複数の文書構成要素を、変換先の前記要素識別情報に対応付けて、前記構成要素変換部による新たな変換対象である新規登録要素集合として登録し、前記構成要素変換部は、更に、前記文書情報における前記新規登録要素集合のそれぞれを、前記文書内配列パターン変換部による変換先の前記要素識別情報に変換する項目1記載の情報処理装置。

(項目7) 前記文書内配列パターン変換部は、前記基準頻度以上で繰り返し出現する前記要素識別情報の組であって、変換前の文書構成要素を当該要素識別情報の組に変換するべく前記文書内配列パターン変換部が繰り返す変換の回数が予め定められた基準回数以下の要素識別情報の組を、前記配列パターンを示す要素識別情報に変換する項目1記載の情報処理装置。

10

(項目8) 前記文書内配列パターン変換部は、前記基準頻度以上で繰り返し出現する前記要素識別情報の組のうち、当該要素識別情報の組に変換される前の文書構成要素の合計サイズが予め定められた基準サイズより小さい前記要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す前記要素識別情報に変換する項目1記載の情報処理装置。

## 【 0 1 5 8 】

(項目9) 前記文書内配列パターン変換部における変換元の要素識別情報の組及び変換先の要素識別情報を利用者に対して出力することにより、前記文書内配列パターン変換部による要素識別情報への変換を行うか否かを入力させる変換指示入力部を更に備える項目1記載の情報処理装置。

20

(項目10) 前記構成要素変換部は、前記文書構成要素を、当該文書構成要素のデータサイズに応じた前記要素識別情報に変換する項目1記載の情報処理装置。

(項目11) 前記文書情報は、前記文書構成要素として、表示画面上に表示する画像を識別する画像識別情報を含み、前記構成要素変換部は、前記画像識別情報を、前記画像識別情報により識別される画像の形状を示す前記要素識別情報に変換する項目1記載の情報処理装置。

(項目12) 前記文書情報は、利用者に表示する表示情報と、当該表示情報に対する指示に応じて表示するべき他の情報の格納位置を示す格納位置情報とを前記文書構成要素として含み、前記構成要素変換部は、前記格納位置情報を、前記他の情報が格納される前記格納位置の範囲を示す前記要素識別情報に変換する項目1記載の情報処理装置。

30

## 【 0 1 5 9 】

(項目13) 前記文書情報は、前記文書構成要素として、表示画面に表示する表示情報と、当該表示情報の表示形式を指定するタグ情報とを含むタグ付き文書であり、前記タグ情報は、当該タグ情報により表示形式を指定する表示情報において更に内側タグ情報を含む、外側タグ情報であり、前記構成要素変換部による変換後の文書情報において、前記外側タグ情報を根ノードとして生成し、前記内側タグ情報を前記根ノードの葉ノードとして生成した、文書構造情報を生成する文書構造情報生成部と、複数の文書情報のそれぞれについて、前記文書構造情報生成部により生成された前記文書構造情報を比較することにより、一の文書情報が他の文書情報と同一の構造を有するか否かを入力する文書情報同一性出力部とを更に備える項目1記載の情報処理装置。

40

(項目14) 前記文書情報は、表示画面に表示するべき情報を指示する文書情報であり、更に、前記表示画面に表示されないコメント情報を含み、前記文書内配列パターン検出部は、当該コメント情報を、一の前記要素識別情報の組と、他の前記要素識別情報の組との境界を示す情報として用いる項目1記載の情報処理装置。

(項目15) 前記グループ分類部により分類された各グループについて、当該グループに含まれる文書構成要素に対する目次を示す目次情報を出力する目次情報出力部を更に備える項目1記載の情報処理装置。

## 【 0 1 6 0 】

50

(項目16) 前記複数の文書構成要素を、前記グループ分類部により分類されたグループ毎に並び替える並替出力部を更に備える項目1記載の情報処理装置。

(項目17) 前記文書内配列パターン変換部は、前記基準頻度以上で繰り返し出現する要素識別情報を、当該要素識別情報の組の配列パターンを示す要素識別情報として、当該文書情報における当該要素識別情報の組の重要度を示す重要度情報に変換し、前記並替出力部は、更に、複数の前記グループ中の文書構成情報を、当該グループに含まれる重要度情報の高い順に並び替える項目16記載の情報処理装置。

(項目18) 前記グループ分類部により分類された各グループに対応付けて、当該グループを識別する情報を、前記文書情報内に生成するグループ識別情報生成部を更に備える項目1記載の情報処理装置。

10

(項目19) 前記グループ分類部により分類された各グループを識別する情報を、前記文書情報とは別体に生成して出力するアノテーション出力部を更に備える項目1記載の情報処理装置。

(項目20) 前記文書情報に応じて表示画面上に情報を出力する表示部と、前記文書内配列パターン変換部により変換される対象となる要素識別情報の組について、当該要素識別情報の変換元である文書構成要素の合計サイズが、前記表示画面に表示可能な情報のサイズに達した場合に、当該要素識別情報の組を、変換を繰り返し行う対象から除外する反復終了判断部とを更に備える項目1記載の情報処理装置。

#### 【0161】

(項目21) グループに分類する対象である対象文書情報と予め定められた関係を有する関連文書情報を検出する関連文書検出部と、前記対象文書情報及び前記関連文書情報から、前記文書内パターン変換部により変換された要素識別情報を除外した文書について、当該対象文書情報及び当該関連文書情報の双方において出現する前記要素識別情報の組であって、前記対象文書情報及び前記関連文書情報を併せた文書において予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す要素識別情報に変換する文書間配列パターン変換部とを更に備え、前記グループ分類部は、更に、前記文書間配列パターン変換部により変換された結果得られた前記対象文書情報において、前記文書間配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類する項目1記載の情報処理装置。

20

30

(項目22) 前記グループ分類部は、分類したグループのそれぞれについて、当該グループ内の文書構成要素が前記文書情報において果たす役割又は当該グループ内の文書構成要素の内容を示すタイトル情報を、更に生成する項目1記載の情報処理装置。

(項目23) 前記グループ分類部は、当該要素識別情報の組の境界に設けられた前記コメント情報に含まれる情報を、前記タイトル情報として生成する項目22記載の情報処理装置。

#### 【0162】

(項目24) 前記グループ分類部は、分類したグループのそれぞれについて、前記文書情報における当該グループの重要度を示す重要度情報を更に生成する項目1記載の情報処理装置。

40

(項目25) 文書情報に含まれる複数の文書構成要素を、複数のグループに分類する情報処理装置であって、グループに分類する対象である対象文書情報と予め定められた関係を有する関連文書情報を検出する関連文書検出部と、前記対象文書情報及び前記関連文書情報のそれぞれにおいて、前記複数の文書構成要素のそれぞれを、当該文書構成要素の種類又は役割を示す要素識別情報に変換する構成要素変換部と、前記対象文書情報及び前記関連文書情報の双方において出現する前記要素識別情報の組であって、前記対象文書情報及び前記関連文書情報を併せた文書において予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組のそれぞれを、当該文書構成要素の組の配列パターンを示す前記要素識別情報に変換する文書間配列パターン変換部と、前記文書間配列パターン変換部により変換された結果得られた前記対象文書情報において、前記文書間配列パターン変換

50

部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類するグループ分類部とを備える情報処理装置。

【0163】

(項目26) 前記関連文書検出部は、前記対象文書情報が格納される格納位置から予め定められた範囲内に格納されている文書情報を前記関連文書情報として検出する項目25記載の情報処理装置。

(項目27) 前記対象文書情報は、当該対象文書情報が生成される以前に存在していた既存文書情報を更新することにより生成され前記関連文書検出部は、前記既存文書情報を前記関連文書情報として検出する項目25記載の情報処理装置。

10

(項目28) 文書情報に含まれる複数の文書構成要素を、複数のグループに分類する情報処理装置を制御するプログラムであって、前記情報処理装置を、前記文書情報における前記複数の文書構成要素のそれぞれを、当該文書構成要素の種類又は役割を示す要素識別情報にそれぞれ変換する構成要素変換部と、前記構成要素変換部により変換された前記文書情報において、予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組のそれぞれを、当該要素識別情報の組の配列パターンを示す前記要素識別情報に変換する文書内配列パターン変換部と、前記文書内配列パターン変換部により前記文書情報を繰り返し変換した結果得られた文書情報において、前記文書内配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類するグループ分類部として機能させるプログラム。

20

【0164】

(項目29) 文書情報に含まれる複数の文書構成要素を、複数のグループに分類する情報処理装置を制御するプログラムであって、前記情報処理装置を、グループに分類する対象である対象文書情報と予め定められた関係を有する関連文書情報を検出する関連文書検出部と、前記対象文書情報及び前記関連文書情報のそれぞれにおいて、前記複数の文書構成要素のそれぞれを、当該文書構成要素の種類又は役割を示す要素識別情報に変換する構成要素変換部と、前記対象文書情報及び前記関連文書情報の双方において出現する前記要素識別情報の組であって、前記対象文書情報及び前記関連文書情報を併せた文書において予め定められた基準頻度以上で繰り返し出現する前記要素識別情報の組のそれぞれを、当該文書構成要素の組の配列パターンを示す前記要素識別情報に変換する文書間配列パターン変換部と、前記文書間配列パターン変換部により変換された結果得られた前記対象文書情報において、前記文書間配列パターン変換部により変換された要素識別情報のそれぞれについて、当該要素識別情報に変換された複数の前記文書構成要素をグループとして分類するグループ分類部として機能させるプログラム。

30

(項目30) 項目28又は項目29記載のプログラムを記録した記録媒体。

【0165】

以上、本発明を実施形態を用いて説明したが、本発明の技術的範囲は上記実施形態に記載の範囲には限定されない。上記実施形態に、多様な変更または改良を加えることができる。そのような変更または改良を加えた形態も本発明の技術的範囲に含まれ得ることが、特許請求の範囲の記載から明らかである。

40

【0166】

【発明の効果】

上記説明から明らかなように、本発明によれば文書中の情報を適切に分類することができる。

【図面の簡単な説明】

【図1】図1は、情報処理装置10の機能ブロック図を示す。

【図2】図2は、情報処理装置10のフローチャートを示す。

【図3】図3は、図2のS210におけるフローチャートを示す。

【図4】図4は、図2のS220におけるフローチャートを示す。

【図5】図5は、図4のS450におけるフローチャートを示す。

50

【図 6】図 6 は、図 2 の S 2 3 0 におけるフローチャートを示す。

【図 7】図 7 は、図 2 の S 2 4 0 におけるフローチャートを示す。

【図 8】図 8 は、構成要素変換部 1 0 5 が文書構成要素を要素識別情報に変換する一例を示す。

【図 9】図 9 は、構成要素変換部 1 0 5 が所定の要素識別情報の組を要素識別情報に変換する例を示す。

【図 1 0】図 1 0 ( a ) は、文書情報により表示された表示画面の一例を示す。図 1 0 ( b ) は、図 1 0 ( a ) に示す表示画面を表示した文書情報に対して、文書内配列パターン変換部 1 3 0 による変換を行った結果を示す。

【図 1 1】図 1 1 ( a ) は、文書情報により表示された表示画面の他の例を示す。図 1 1 ( b ) は、図 1 1 ( a ) に示す表示画面を表示した文書情報に対して、文書内配列パターン変換部 1 3 0 による変換を行った結果を示す。

【図 1 2】図 1 2 は、変換対象である対象文書情報の表示例を示す。

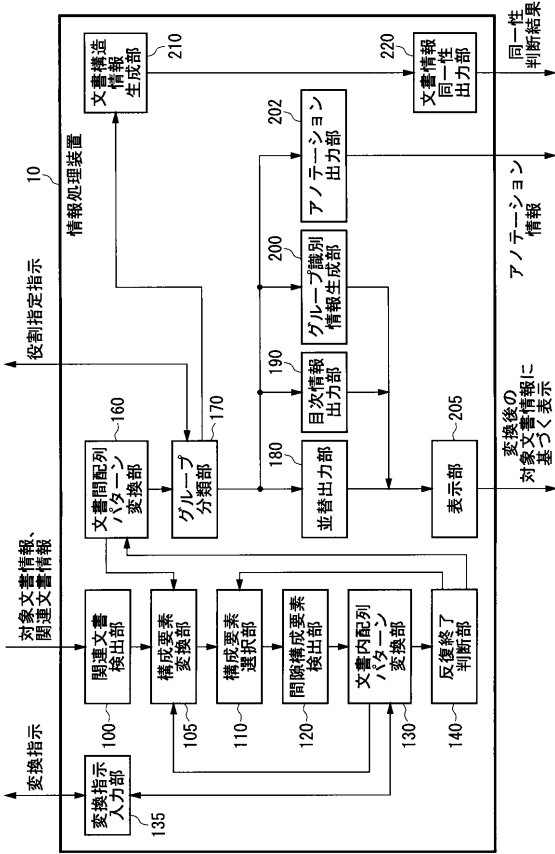
【図 1 3】図 1 3 は、図 1 2 に示す対象文書情報において、グループ分類部 1 7 0 による分類の一例を示す。

【図 1 4】図 1 4 は、情報処理装置 1 0 のハードウェア構成の一例を示す。

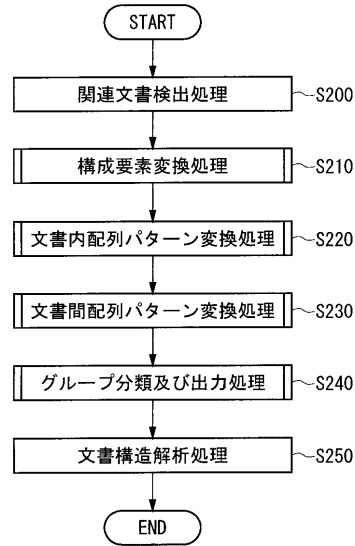
#### 【符号の説明】

|       |              |    |
|-------|--------------|----|
| 1 0   | 情報処理装置       |    |
| 1 0 0 | 関連文書検出部      |    |
| 1 0 5 | 構成要素変換部      | 20 |
| 1 1 0 | 構成要素選択部      |    |
| 1 2 0 | 間隙構成要素検出部    |    |
| 1 3 0 | 文書内配列パターン変換部 |    |
| 1 3 5 | 変換指示入力部      |    |
| 1 4 0 | 反復終了判断部      |    |
| 1 6 0 | 文書間配列パターン変換部 |    |
| 1 7 0 | グループ分類部      |    |
| 1 8 0 | 並替出力部        |    |
| 1 9 0 | 目次情報出力部      |    |
| 2 0 0 | グループ識別情報生成部  | 30 |
| 2 0 2 | アノテーション出力部   |    |
| 2 0 5 | 表示部          |    |
| 2 1 0 | 文書構造情報生成部    |    |
| 2 2 0 | 文書情報同一性出力部   |    |

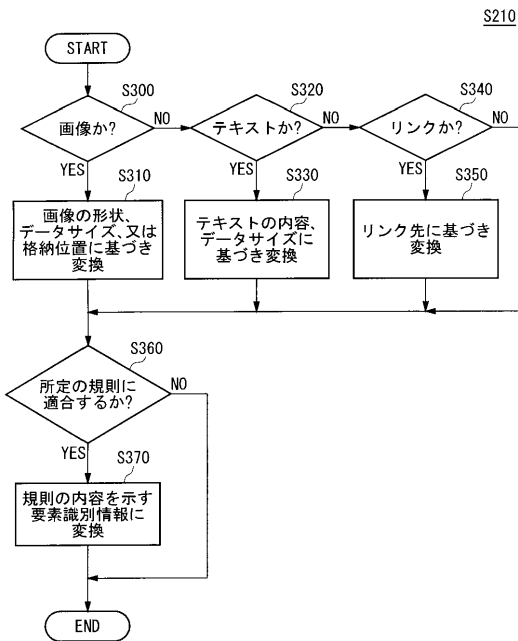
【 図 1 】



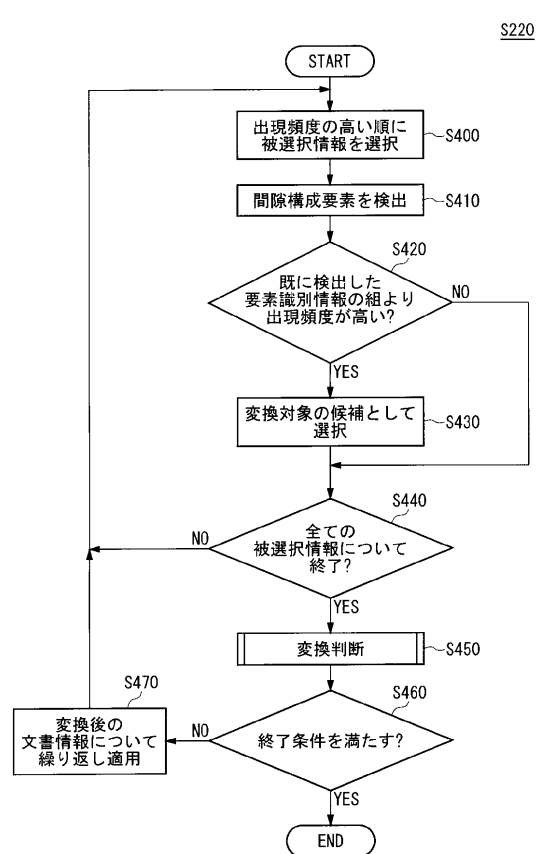
【 図 2 】



【 図 3 】

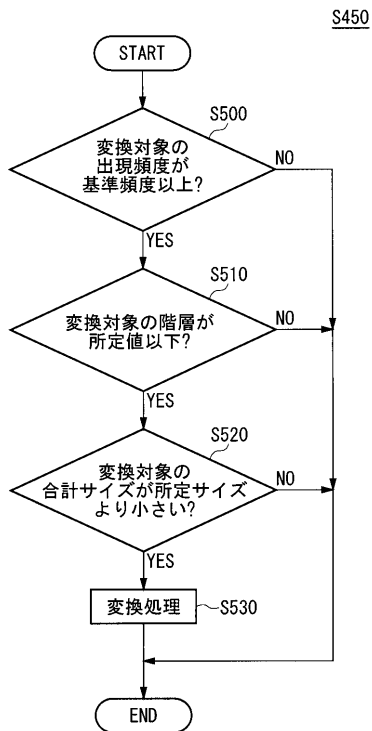


【 図 4 】

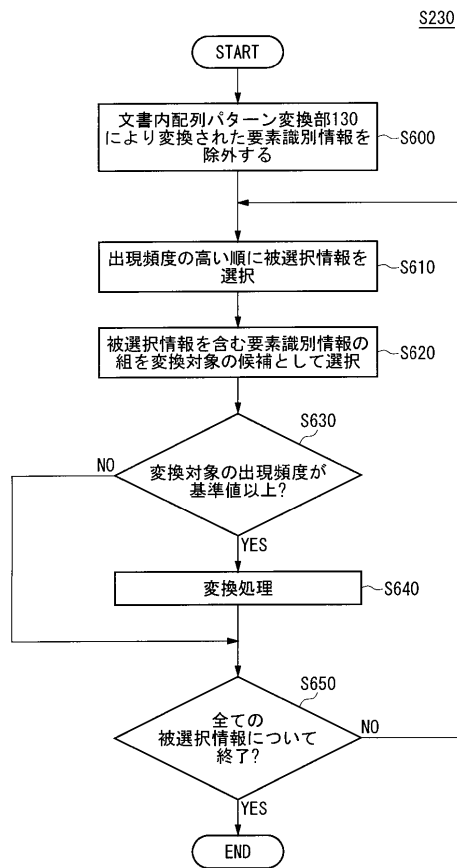




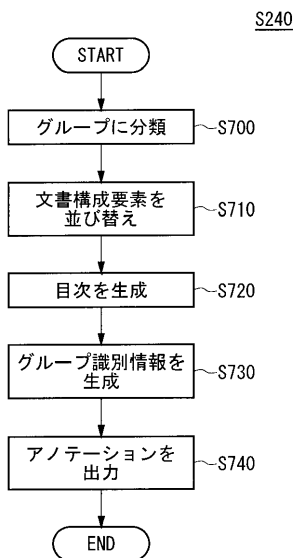
【 図 5 】



【 図 6 】



【 図 7 】



【 図 8 】

| 要素識別情報                           | 文書構成要素                               |
|----------------------------------|--------------------------------------|
| <digit />                        | 数値                                   |
| <longtext />                     | 100文字以上の文字列                          |
| <midtext />                      | 30~100文字の文字列                         |
| <shorttext />                    | 2~30文字の文字列                           |
| <letter />                       | 1文字の文字                               |
| <in-largeimg />                  | 対象文書情報と同一サイト内にあり、300×300ピクセル以上の画像    |
| <in-midimg />                    | 対象文書情報と同一サイト内にあり、100×100~300×300の画像  |
| <in-smallimg />                  | 対象文書情報と同一サイト内にあり、100×100ピクセル以下の画像    |
| <out-largeimg />                 | 対象文書情報と異なるサイト内にあり、300×300ピクセル以上の画像   |
| <out-midimg />                   | 対象文書情報と異なるサイト内にあり、100×100~300×300の画像 |
| <out-smallimg />                 | 対象文書情報と異なるサイト内にあり、100×100ピクセル以下の画像   |
| <samedir-link>...</samedir-link> | 対象文書情報と同一のディレクトリ内へのリンク               |
| <in-link>...</in-link>           | 対象文書情報と同一サイトへのリンク                    |
| <out-link>...</out-link>         | 対象文書情報と異なるサイトへのリンク                   |

【 図 9 】

|      |  |                                       |   |   |   |
|------|--|---------------------------------------|---|---|---|
| 9個以上 | <<letter /><shorttext />又は<br>{<←-link><br><<letter /><shorttext /><br></←-link> | <digit /><br>{<digit /><br></digit /> | <letter /><shorttext /><br>{<letter /><shorttext /><br></letter /><shorttext /> | <letter /><shorttext /><br>{<letter /><shorttext /><br></letter /><shorttext /> | <letter /><shorttext /><br>{<letter /><shorttext /><br></letter /><shorttext /> |
|      | <digit /> ... <digit />  | <digit /> ... <digit />               | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             |
|      | <digit /> ... <digit />  | <digit /> ... <digit />               | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             |
|      | <digit /> ... <digit />  | <digit /> ... <digit />               | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             |
|      | <digit /> ... <digit />  | <digit /> ... <digit />               | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             |
|      | <digit /> ... <digit />  | <digit /> ... <digit />               | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             |
|      | <digit /> ... <digit />  | <digit /> ... <digit />               | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             |
|      | <digit /> ... <digit />  | <digit /> ... <digit />               | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             |
|      | <digit /> ... <digit />  | <digit /> ... <digit />               | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             | <letter /><shorttext /> ... <letter /><shorttext />                             |

【 図 10 】

- (a)
- 社会
  - スポーツ
  - 金融・経済
  - 政治
  - ...
- (b)
- ```

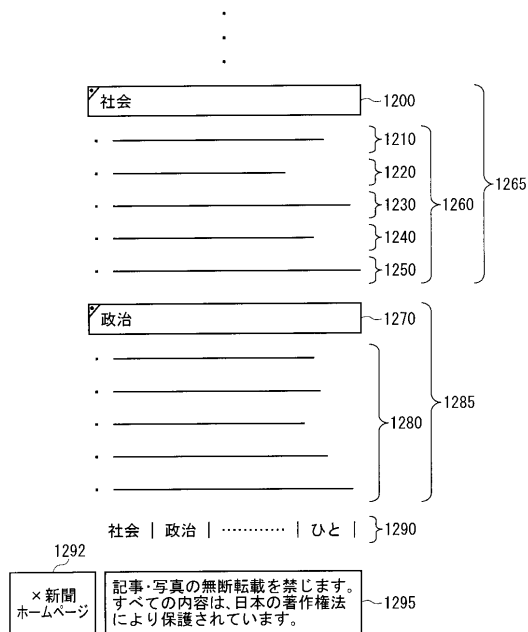
<bullet /><in-link><shorttext /></in-link><new-line />
<bullet /><in-link><shorttext /></in-link><new-line />
<bullet /><in-link><shorttext /></in-link><new-line />
<bullet /><in-link><shorttext /></in-link><new-line />
  
```

【 図 11 】

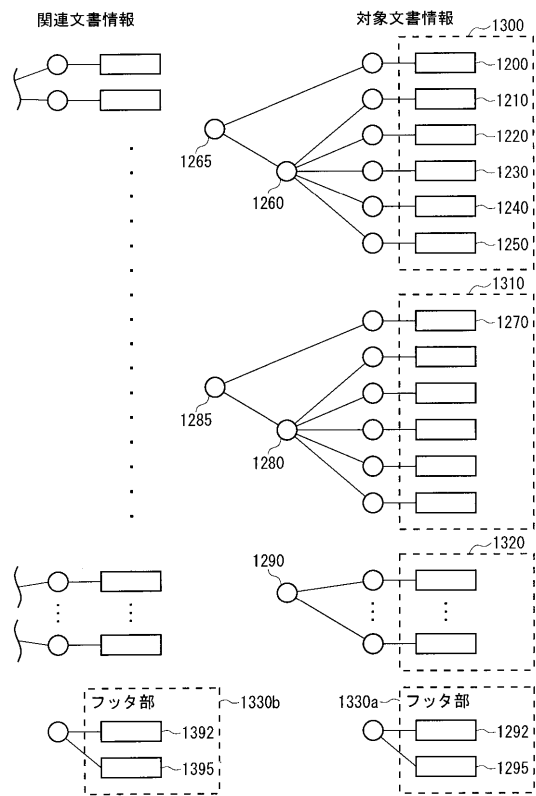
- (a)
- ・ ～道路で死亡事故、××町—〇〇市間が一時通行止めに(21:56)
  - ・ AA市内の川に邦人遺体 事件・事故の両面で捜査(21:18)
  - ・ BB知事、××社への貸付金〇〇〇億円の債権放棄(20:13)
- (b)
- ```

<bullet /><in-link><midtext /></in-link><new-line />
<bullet /><in-link><midtext /></in-link><new-line />
<bullet /><in-link><midtext /></in-link><new-line />
  
```

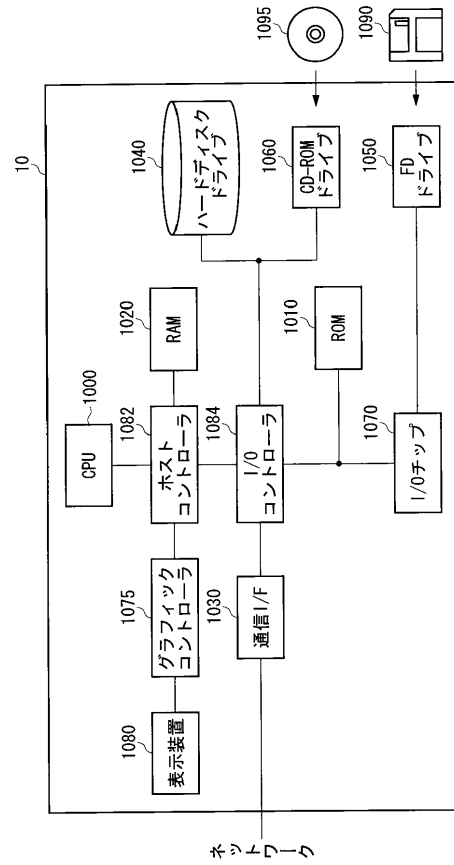
【 図 12 】



【図13】



【図14】



---

フロントページの続き

(72)発明者 福田 健太郎  
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内

審査官 深津 始

(56)参考文献 特開2002-324077(JP,A)  
特開2002-312379(JP,A)  
米国特許出願公開第2002/0022956(US,A1)  
特開2002-334070(JP,A)

(58)調査した分野(Int.Cl., DB名)  
G06F 17/30  
G06F 17/21