

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4370873号
(P4370873)

(45) 発行日 平成21年11月25日(2009.11.25)

(24) 登録日 平成21年9月11日(2009.9.11)

(51) Int.Cl. F I
G06F 17/30 (2006.01) G O 6 F 17/30 2 1 0 D
 G O 6 F 17/30 1 7 0 B

請求項の数 6 (全 12 頁)

<p>(21) 出願番号 特願2003-358081 (P2003-358081) (22) 出願日 平成15年10月17日(2003.10.17) (65) 公開番号 特開2005-122550 (P2005-122550A) (43) 公開日 平成17年5月12日(2005.5.12) 審査請求日 平成18年9月21日(2006.9.21)</p>	<p>(73) 特許権者 000005496 富士ゼロックス株式会社 東京都港区赤坂九丁目7番3号 (74) 代理人 100098084 弁理士 川▲崎▼ 研二 (72) 発明者 加藤 雅弘 神奈川県海老名市本郷2274番地 富士 ゼロックス株式会社内 審査官 鈴木 和樹</p>
---	--

最終頁に続く

(54) 【発明の名称】 文書分類装置、プログラムおよび文書分類方法

(57) 【特許請求の範囲】

【請求項1】

文書の画像を表す画像データを取得する画像取得手段と、
 前記画像データで表される画像を解析することによって、前記文書の各ページを構成する構成要素のレイアウトを表すレイアウト情報を取得し、前記構成要素のうち、各ページ内で文章が空間的に連続している領域である構成要素を、文章領域として抽出するレイアウト解析手段と、
 前記文章領域に含まれる文字列を認識する文字認識手段と、
 前記文字認識手段により認識された文字列から視覚的に強調された文字列を抽出し、抽出された文字列をキーワードとするキーワード抽出手段と、
 前記レイアウト解析手段で抽出された複数の文章領域間の境界で各ページを分割し、各ページを根とし各文章領域を葉とし、前記境界の方向を基準とした木構造を用いて、各ページにおける複数の文章領域を階層的に表す構造データをページ毎に生成し、各文章領域に対応する特徴量をページ間で比較することで文章領域毎の包含関係を求め、前記包含関係に基づいて前記構造データにおける各文章領域の階層をページ間で調整する構造データ生成手段と、
 前記構造データ生成手段で各文章領域の階層が調整された構造データに基づいて前記各キーワードの階層を求め、前記各キーワードの階層及び出現順を前記文書の論理構造として抽出する論理構造抽出手段と、
 前記論理構造抽出手段で抽出された各文書の論理構造に含まれるキーワードの階層及び

10

20

出現順を比較することで前記各文書を分類して記憶する分類手段と
を有する文書分類装置。

【請求項 2】

前記特徴量が、文章領域の位置、文章領域の大きさ、文字の大きさ、段組みの向き、1行あたりの平均的な文字数のうち少なくとも一つを含むことを特徴とする請求項 1 に記載の文書分類装置。

【請求項 3】

前記構造データ生成手段が、垂直方向又は水平方向の少なくとも一方の境界で各ページを分割することを特徴とする請求項 1 に記載の文書分類装置。

【請求項 4】

前記特徴量が、前記レイアウト解析手段により取得されたレイアウト情報を基に各文章領域に対応付けられる情報であることを特徴とする請求項 1 に記載の文書分類装置。

【請求項 5】

コンピュータ装置を、

文書の画像を表す画像データを取得する画像取得手段と、

前記画像データで表される画像を解析することによって、前記文書の各ページを構成する構成要素のレイアウトを表すレイアウト情報を取得し、前記構成要素のうち、各ページ内で文章が空間的に連続している領域である構成要素を、文章領域として抽出するレイアウト解析手段と、

前記文章領域に含まれる文字列を認識する文字認識手段と、

前記文字認識手段により認識された文字列から視覚的に強調された文字列を抽出し、抽出された文字列をキーワードとするキーワード抽出手段と、

前記レイアウト解析手段で抽出された複数の文章領域間の境界で各ページを分割し、各ページを根とし各文章領域を葉とし、前記境界の方向を基準とした木構造を用いて、各ページにおける複数の文章領域を階層的に表す構造データをページ毎に生成し、各文章領域に対応する特徴量をページ間で比較することで文章領域毎の包含関係を求め、前記包含関係に基づいて前記構造データにおける各文章領域の階層をページ間で調整する構造データ生成手段と、

前記構造データ生成手段で各文章領域の階層が調整された構造データに基づいて前記各キーワードの階層を求め、前記各キーワードの階層及び出現順を前記文書の論理構造として抽出する論理構造抽出手段と、

前記論理構造抽出手段で抽出された各文書の論理構造に含まれるキーワードの階層及び出現順を比較することで前記各文書を分類して記憶する分類手段

として機能させるためのプログラム。

【請求項 6】

CPUが文書の画像を表す画像データを取得する画像取得ステップと、

CPUが前記画像データで表される画像を解析することによって、前記文書の各ページを構成する構成要素のレイアウトを表すレイアウト情報を取得し、前記構成要素のうち、各ページ内で文章が空間的に連続している領域である構成要素を、文章領域として抽出するレイアウト解析ステップと、

CPUが前記文章領域に含まれる文字列を認識する文字認識ステップと、

CPUが前記文字認識ステップにより認識された文字列から視覚的に強調された文字列を抽出し、抽出された文字列をキーワードとするキーワード抽出ステップと、

CPUが前記レイアウト解析ステップで抽出された複数の文章領域間の境界で各ページを分割し、各ページを根とし各文章領域を葉とし、前記境界の方向を基準とした木構造を用いて、各ページにおける複数の文章領域を階層的に表す構造データをページ毎に生成し、各文章領域に対応する特徴量をページ間で比較することで文章領域毎の包含関係を求め、前記包含関係に基づいて前記構造データにおける各文章領域の階層をページ間で調整する構造データ生成ステップと、

CPUが前記構造データ生成ステップで各文章領域の階層が調整された構造データに基

10

20

30

40

50

づいて前記各キーワードの階層を求め、前記各キーワードの階層及び出現順を前記文書の論理構造として抽出する論理構造抽出ステップと、

C P Uが前記論理構造抽出ステップで抽出された各文書の論理構造に含まれるキーワードの階層及び出現順を比較することで前記各文書を分類して画像蓄積部に記憶する分類ステップと

を有する文書分類方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、文書の画像を用いて文書を分類し、記憶する技術に関する。

10

【背景技術】

【0002】

文書の種類を識別し、文書の種類を表す情報と文書画像とを対応付けてファイリングする文書管理システムが提案されている（例えば、特許文献1および2）。

特許文献1に記載の技術では、予め文書フォームとそれに対応する文書の種類が登録されており、入力された文書のフォームを識別し、そのフォームに対応する文書の種類を表す情報と当該文書画像とを対応付けて格納する。入力された文書のフォームに該当するフォームが登録されていない場合には、新たにそのフォームを登録する。これによって、予め登録されていないフォームを有する文書が入力された場合でも、その文書を分類してファイリングすることが可能になるとしている。

20

特許文献2に記載の技術では、入力された文書の画像から文章、図、表などの領域を抽出し、各領域のレイアウトを表す情報と当該文書画像とを対応付けて格納する。これによって、非定型、すなわちフォームによって識別できない文書であっても、そのレイアウトを表す情報によって分類してファイリングすることが可能になるとしている。

【特許文献1】特開2002-269126号公報

【特許文献2】特開2002-342343号公報

【発明の開示】

【発明が解決しようとする課題】

【0003】

しかしながら、特許文献1および2の技術では、厳密な定型フォームを持たない文書の場合に登録されるフォームが際限なく増大してしまうおそれがある。例えば論文誌に掲載された論文のように書式が規定されてはいるものの、1件毎にページ数やレイアウトが異なる文書の場合がこれにあたる。

30

本発明は、上述した背景のもとになされたものであり、厳密な定型フォームではない文書を適切に分類することのできる技術の提供を目的とする。

【課題を解決するための手段】

【0004】

上述の課題を解決するために、本発明は、文書の画像を表す画像データを取得する画像取得手段と、前記画像データで表される画像を解析することによって、前記文書の各ページを構成する構成要素のレイアウトを表すレイアウト情報を取得し、前記構成要素のうち、各ページ内で文章が空間的に連続している領域である構成要素を、文章領域として抽出するレイアウト解析手段と、前記文章領域に含まれる文字列を認識する文字認識手段と、前記文字認識手段により認識された文字列から視覚的に強調された文字列を抽出し、抽出された文字列をキーワードとするキーワード抽出手段と、前記レイアウト解析手段で抽出された複数の文章領域間の境界で各ページを分割し、各ページを根とし各文章領域を葉とし、前記境界の方向を基準とした木構造を用いて、各ページにおける複数の文章領域を階層的に表す構造データをページ毎に生成し、各文章領域に対応する特徴量をページ間で比較することで文章領域毎の包含関係を求め、前記包含関係に基づいて前記構造データにおける各文章領域の階層をページ間で調整する構造データ生成手段と、前記構造データ生成手段で各文章領域の階層が調整された構造データに基づいて前記各キーワードの階層を求

40

50

め、前記各キーワードの階層及び出現順を前記文書の論理構造として抽出する論理構造抽出手段と、前記論理構造抽出手段で抽出された各文書の論理構造に含まれるキーワードの階層及び出現順を比較することで前記各文書を分類して記憶する分類手段とを有する文書分類装置を提供する。

【0005】

また、本発明は、コンピュータ装置を、文書の画像を表す画像データを取得する画像取得手段と、前記画像データで表される画像を解析することによって、前記文書の各ページを構成する構成要素のレイアウトを表すレイアウト情報を取得し、前記構成要素のうち、各ページ内で文章が空間的に連続している領域である構成要素を、文章領域として抽出するレイアウト解析手段と、前記文章領域に含まれる文字列を認識する文字認識手段と、前記文字認識手段により認識された文字列から視覚的に強調された文字列を抽出し、抽出された文字列をキーワードとするキーワード抽出手段と、前記レイアウト解析手段で抽出された複数の文章領域間の境界で各ページを分割し、各ページを根とし各文章領域を葉とし、前記境界の方向を基準とした木構造を用いて、各ページにおける複数の文章領域を階層的に表す構造データをページ毎に生成し、各文章領域に対応する特徴量をページ間で比較することで文章領域毎の包含関係を求め、前記包含関係に基づいて前記構造データにおける各文章領域の階層をページ間で調整する構造データ生成手段と、前記構造データ生成手段で各文章領域の階層が調整された構造データに基づいて前記各キーワードの階層を求め、前記各キーワードの階層及び出現順を前記文書の論理構造として抽出する論理構造抽出手段と、前記論理構造抽出手段で抽出された各文書の論理構造に含まれるキーワードの階層及び出現順を比較することで前記各文書を分類して記憶する分類手段として機能させるためのプログラムを提供する。

【0006】

また、本発明は、CPUが文書の画像を表す画像データを取得する画像取得ステップと、CPUが前記画像データで表される画像を解析することによって、前記文書の各ページを構成する構成要素のレイアウトを表すレイアウト情報を取得し、前記構成要素のうち、各ページ内で文章が空間的に連続している領域である構成要素を、文章領域として抽出するレイアウト解析ステップと、CPUが前記文章領域に含まれる文字列を認識する文字認識ステップと、CPUが前記文字認識ステップにより認識された文字列から視覚的に強調された文字列を抽出し、抽出された文字列をキーワードとするキーワード抽出ステップと、CPUが前記レイアウト解析ステップで抽出された複数の文章領域間の境界で各ページを分割し、各ページを根とし各文章領域を葉とし、前記境界の方向を基準とした木構造を用いて、各ページにおける複数の文章領域を階層的に表す構造データをページ毎に生成し、各文章領域に対応する特徴量をページ間で比較することで文章領域毎の包含関係を求め、前記包含関係に基づいて前記構造データにおける各文章領域の階層をページ間で調整する構造データ生成ステップと、CPUが前記構造データ生成ステップで各文章領域の階層が調整された構造データに基づいて前記各キーワードの階層を求め、前記各キーワードの階層及び出現順を前記文書の論理構造として抽出する論理構造抽出ステップと、CPUが前記論理構造抽出ステップで抽出された各文書の論理構造に含まれるキーワードの階層及び出現順を比較することで前記各文書を分類して画像蓄積部に記憶する分類ステップとを有する文書分類方法を提供する。

【発明の効果】

【0008】

本発明によれば、厳密な定型フォームではない文書を適切に分類することができる。文章領域のレイアウト上の階層に基づいてキーワードの論理レベルを決定し、文書全体の論理構造を求め、異なる文書間で論理構造を比較するから、文書の論理構造に着目した分類が可能となる。また、文書のカテゴリーを表す分類情報を文書画像と対応付けて記憶するから、カテゴリー毎に文書を検索することが可能となる。

【発明を実施するための最良の形態】

【0009】

10

20

30

40

50

以下、図面を参照して、本発明の実施の形態について説明する。

[構成]

図1は、文書分類装置10のハードウェア構成を示す図である。ROM(Read Only Memory)102には、プログラム10Pが書き込まれている。CPU(Central Processing Unit)101は、文書分類装置10に電源(図示省略)が投入されると、ROM102に書き込まれているプログラム10Pを読み出し、RAM(Random Access Memory)103をワークエリアとしてプログラム10Pを実行する。CPU101がプログラム10Pを実行することによって、文書分類装置10には、図8に示すモジュール群が仮想的に形成される。なお、外部の装置にプログラム10Pを記憶させておき、通信網(図示省略)を介してプログラム10PをダウンロードしてROM102に記憶させることとしてもよい。

10

【0010】

画像蓄積部117は、原稿の画像データを蓄積するハードディスクドライブである。画像処理部118はCPU101による制御の下で、画像蓄積部117に蓄積されている画像データを読み出し、各種の処理を行う。CPU101、ROM102、RAM103、画像蓄積部117および画像処理部118はバス115に接続されている。

表示部105は、CRT(Cathode Ray Tube)あるいは液晶パネルである。操作部107は、ポインティングデバイス(マウスあるいはデジタイザ)およびキーボードである。スキャナ109は、原稿を光学的に読み取り、画像信号を出力する。プリンタ111は、電子写真方式あるいはインクジェット方式のプリンタである。表示部105、操作部107、スキャナ109およびプリンタ111はそれぞれインターフェイス106、108、110、112を介してバス116に接続されており、バス116とバス115とはバスブリッジ104によって接続されている。バス116は、インターフェイス114を介してネットワーク113に接続されており、これによって文書分類装置10と外部の装置との通信が可能となっている。

20

【0011】

次に、CPU101がプログラム10Pを実行することによって文書分類装置10に仮想的に形成されるモジュール群について、図8を用いて説明する。

画像取得手段21は、文書の画像を表す画像データを取得する手段である。文書の画像を表す画像データとは、スキャナ109等の画像入力装置を用いて文書を走査することによって生成された画像データである。なお、画像データは、予め画像蓄積部117に格納されていてもよいし、外部の装置からネットワーク113を介して文書分類装置10が受信することとしてもよい。

30

【0012】

レイアウト解析手段22は、画像データで表される画像を解析することによって文書の各ページを構成する構成要素のレイアウトを表すレイアウト情報を取得し、ページ内で文章が空間的に連続している文章領域を抽出する手段である。ここで構成要素とは、文章、図、表などである。本実施形態においては、文章が空間的に連続している領域を文章領域と呼ぶ。また、図によって占められる領域を図領域、表によって占められる領域を表領域と呼ぶ。画像データは、文書をスキャナ109で走査して得られた画素値によって構成されており、ページ記述言語などで記述された文書データが内包しているようなレイアウト情報を有していない。そのため、レイアウト解析手段では、既知のレイアウト解析技術を用いて、当該画像で表される文書のレイアウト情報を得る。レイアウト解析は、例えば特開2000-90194号公報に記載されている技術を用いて行う。この技術では、文書画像に含まれる文章が縦書きか横書きかを判定し、その判定結果によって文書を分割する境界を設定する。また、文書画像を構成する画素の投影分布を算出し、所定のしきい値に満たない頻度の区間を用いて分割境界を設定する。

40

レイアウト解析手段22はこのようにして文書のレイアウト情報を取得し、各ページを文章領域、図領域および表領域の集合として認識する。

【0013】

50

文字認識手段 2 3 は、レイアウト解析手段 2 2 により抽出された文章領域に含まれる文字列を認識する手段である。

キーワード抽出手段 2 4 は、文字認識手段により認識された文字列から視覚的に強調された文字列を抽出し、抽出された文字列をキーワードとする手段である。キーワードの抽出は、例えば特開平 9 - 2 9 7 7 6 5 号公報に記載されている方法を用いて行う。ここで、キーワードとは、何らかの方法により視覚的に強調されている文字列である。例えば、予め文字サイズの閾値を定めておき、この閾値を超える大きさの文字列を抽出する。あるいは、太字、斜体など、通常と異なるフォントを用いた文字列、枠で囲まれた文字列、下線を引かれた文字列などを抽出してもよい。

【 0 0 1 4 】

構造データ生成手段 2 5 は、レイアウト解析手段 2 2 で抽出された文章領域のレイアウト上の階層構造を表す構造データをページ毎に生成する手段である。構造データは、図 4 に示すように、各ページを根とする木構造によって形成されており、レイアウト解析により抽出された文章領域の各々は、構造データの葉の各々と対応付けられている。葉の各々には、レイアウト解析の際に得られたレイアウト情報、すなわち、各文章領域の位置、大きさ、行の向きなどの情報が対応付けられている。なお、図、表など文章以外の領域については、構造データに含めない。

【 0 0 1 5 】

論理構造抽出手段 2 6 は、構造データ生成手段 2 5 で生成された構造データとキーワード抽出手段 2 4 で抽出されたキーワードとを用いて、文書の論理構造を抽出する手段である。

分類手段 2 7 は、論理構造抽出手段 2 6 で抽出された論理構造を用いて文書を分類して記憶する手段である。

なお、論理構造抽出手段 2 6 および分類手段 2 7 によって行われる処理については、動作の説明において詳述する。

【 0 0 1 6 】

[動作]

上記の構成からなる文書分類装置 1 0 の動作について説明する。ただし、文書分類装置 1 0 は、ハードウェアがソフトウェアを用いることによって動作する装置であるから、これ以降の説明においては、動作の主体を、仮想的に形成されるモジュールではなく、ハードウェアとする。

ここで、文書の例について説明する。図 3 は、文書分類装置 1 0 による処理の対象とされる文書の例を示す図である。この例は、横書き 2 段組を基本とする書式に従って作成された文書（例えば、論文）の例である。

【 0 0 1 7 】

1 ページ目と 2 ページ目とでは文書のレイアウトが異なることがわかる。1 ページ目では、最上部に題目が大きな文字サイズで 1 行記載されており、その下に抄録が 6 行記載されている。さらにその下には、本文が 2 段組で左右の各段に 1 3 行記載されている。2 ページ目では、1 ページ目の本文と同様に 2 段組で記載されており、同様のレイアウトのページが L - 2 ページ目まで続く。L - 1 ページ目では、ページ上方にグラフが挿入されている。L ページ目では、2 ページ目と同様のレイアウトとなり、このレイアウトのページが M - 3 ページ目まで続く。M - 2 ページ目では、右の段の下部にクルマの図が挿入され、M - 1 ページ目では、ページ全体に表が記載されている。M ページ目では 2 ページ目と同様のレイアウトとなり、このレイアウトが N - 1 ページ目まで続く。そして、N ページ目では、右の段に著者の顔写真と紹介記事が記載されている。

【 0 0 1 8 】

図 2 は、CPU 1 0 1 がプログラム 1 0 P を実行することによって行われる処理のフローを示す図である。ここでは、文書分類装置 1 0 には電源が投入されており、CPU 1 0 1 によってプログラム 1 0 P が実行されているものとする。

最初に、文書分類装置 1 0 は画像データを取得する（ステップ S 0 1 ）。ここでは、ま

10

20

30

40

50

ず、スキャナ109によって文書が読み取られ、文書の画像を表す画像信号が文書分類装置10に送信される。文書分類装置10は、スキャナ109から送信された画像信号を受信する。すると、CPU101が、受信された画像信号に基づいて画像データを生成し、画像蓄積部117に格納する。そして、CPU101は、画像蓄積部117に格納された画像データをRAM103上に展開する。

【0019】

次に、CPU101は、レイアウト解析を行って、文書を構成する各ページのレイアウトを表すレイアウト情報を取得する(ステップS02)。図4は、図3に示した文書のレイアウト解析例(上段)およびレイアウト解析結果に基づいて作成された構造データの例(下段)を示す図である。各ページ中の矩形で囲まれた領域が文章領域、図領域、表領域を表しており、これらの中で番号が付されているものが文章領域である。本実施形態においては、文章領域のみを処理の対象とし、図領域および表領域は処理の対象としない。

10

次に、CPU101は、文章領域に含まれる文字を認識するための処理を行い(ステップS03)、認識された文字列からキーワードとなり得る文字列を抽出する(ステップS04)。ここでは、特定の文字サイズ、特定のフォント、枠で囲まれた文字列、下線の引いてある文字列などを抽出する。

【0020】

ステップS03、ステップS04の処理と並行して、CPU101は、構造データの生成(ステップS05)、構造データのレベル調整(ステップS06)を行う。

まず、構造データの生成(ステップS05)について説明する。図4に示すように、構造データは、ページの各々を根とする木構造を有している。木構造の階層は、以下のようにして決定される。1ページ目の例では、まず水平方向に3つの領域に分割された後、最も下の領域が垂直方向に2分割されたとみなされる。これを木構造で表すと、根の1つ下の階層には2つの葉と1つの中間接点が存在し、2つの葉は領域1と領域2に対応付けられる。中間接点の1つ下の階層にはさらに2つの葉が存在し、2つの葉は領域3と領域4に対応付けられる。

20

【0021】

一方、M-2ページ目の例では、まず垂直方向に2つの領域に分割された後、右側の領域が水平方向に2分割されたものとみなされるが、本実施形態においては図領域および表領域を処理の対象としないため、領域2の下方に位置する図領域は無視される。従って、M-2ページ目では、領域1と領域2が等しい階層であるとみなされることとなる。M-1ページ目の表領域、Nページ目の図領域(顔写真)も同様に無視され、図4に示す構造データが得られる。

30

【0022】

次に、レベル調整(ステップS06)について説明する。図7は、構造データのレベル調整のフローを示す図である。まず、ステップS601では、文書の1ページ目の構造データを初期値とする。ステップS602では、2ページ目以降のページについて、直前のページとの間でノード間の対応付けを行う。各ノードには、当該ノードに対応する文章領域がレイアウトされている位置、領域の大きさ、領域内の平均的な文字サイズ、段組みが横方向の分割であるか縦方向の分割であるか、1行あたりの平均的な文字数といった、レイアウト解析の際に得られる種々の特徴量が対応付けられて記憶されている。CPU101は、注目ページとその直前のページとの特徴量を比較し、領域毎の包含関係に基づいて、ノード間の対応付けを行う。

40

【0023】

ここで、ノード間の対応付けについて説明する。図4によると、1ページ目の領域3が2ページ目の領域1に包含されている。同様に、1ページ目の領域4が2ページ目の領域2に包含されている(1対1の包含)。従って、2ページ目の領域1および2は、1ページ目の領域3および4と同等以上の階層に位置付けられることになる。しかし、1ページ目の領域1および2は、それぞれが2ページ目の領域1および2の一部を包含しているため、2ページ目の領域1および2は、1ページ目の領域1および2よりも下の階層に位置

50

付けられることとなる。よって、図5に示すように、2ページ目の領域1および2は、1ページ目の領域3および4と等しい階層（Level_2）に位置付けられる。

一方、Mページ目とNページ目の例では、Mページ目の領域2がNページ目の領域2、3および4を包含している（1対多の包含）。従って、Nページ目の領域2、3および4は、Mページ目の領域2よりも下の階層に位置付けられることとなる。よって、図5に示すように、Nページ目の領域2、3および4は、Mページ目の領域1および2の1つ下の階層（Level_3）に位置付けられる。

【0024】

他方、上述のようなノード間の対応付けの結果、直前のページの構造データにおいてLevel_0を下位の階層へ変更する必要が生じた場合、すなわち、それ以前のページに存在しない上位のノードが注目ページに出現した場合（例えば、図4において1ページ目と2ページ目が入れ替わっている場合）には（ステップS603：YES）、ステップS604で構造データの階層の変更が必要な先頭ページまでさかのぼり、ステップS605で当該先頭ページ以降、未処理ページまでの構造データの階層を変更する。具体的には、仮想ノードを最上位に挿入し、全体の階層を下位にシフトさせる。一方、ステップS603で階層の変更が必要でないと判定された場合には、ステップS606に進み、注目ページの構造データの階層の変更が必要かどうかを判定し、変更が必要であるならばステップS607にて注目ページの構造データの階層を変更する。ステップS608では、全ページについてノード間の対応付けが完了したか否かを判定し、完了していないならばステップS602に戻る。

【0025】

ノード間の対応付けが完了したならば、図2のステップS07に進み、キーワードの階層を求める。このキーワードはステップS04で文章領域から抽出されたキーワードである。このキーワードが属する文章領域の階層は、ステップS06で求められた当該文章領域の階層である。

ステップS08では、特定の階層のキーワードを用いて、文書の論理構造を抽出する。図6は、文書の論理構造を抽出した例を示す図である。この例では、文書Aおよび文書Bはともに公開特許公報である。図6には、「要約」、「特許請求の範囲」、「技術分野」、...と続く特許公報の記載項目の見出しが出現順に示されている。公開特許公報は、帳票のように全ページにおいて文字枠が厳密に既定されているものではなく、ページ数やレイアウトが1件1件異なるものである。また、数式や表が記載された公報と記載されていない公報が存在する。また、各項目毎の文章の分量も1件毎に異なる。このような違いを含んだ文書の場合、従来知られているような厳密なフォーム認識では異なる種類の文書とみなされることになる。しかしながら、文書Aと文書Bの論理構造に着目すれば、両者は明らかに同じ種類の文書である。このように、文書の論理構造を比較することによって、フォームが厳密には一致していない文書や、ページ数が1件毎に異なる文書であっても、同種の文書として分類することができる。ステップS09では、このようにして分類された文書のカテゴリーを表す情報と当該文書の画像データとを対応付けて画像蓄積部117に格納する。

【0026】

以上説明したように、本発明によれば、厳密な定型フォームではない文書を適切に分類することができる。文章領域のレイアウト上の階層に基づいてキーワードの論理レベルを決定し、文書全体の論理構造を求め、異なる文書間で論理構造を比較するから、文書の論理構造に着目した分類が可能となる。また、文書のカテゴリーを表す分類情報を文書画像と対応付けて記憶するから、カテゴリー毎に文書を検索することが可能となる。

【0027】

[変形例]

以上説明した形態に限らず、本発明は種々の形態で実施可能である。例えば、上述の実施形態を以下のように変形した形態でも実施可能である。

スキャナが接続された1または複数のパーソナルコンピュータをネットワークを介して

10

20

30

40

50

文書分類装置 10 に接続し、スキャナで読み込まれた文書の画像データを文書分類装置 10 に送信し、文書分類装置 10 において文書の分類および格納を行うようにしてもよい。このようにすれば、例えば、オフィス内の別々の場所に分散して保管されている文書を文書分類装置 10 で集中管理することが可能となる。

文書の論理構造は、特定の階層のキーワードではなく、すべての階層のキーワードをその階層を表す情報とともに表したものであってもよい。

【図面の簡単な説明】

【0028】

【図1】本発明の一実施形態に係る文書分類装置の構成を示す図である。

【図2】CPUがプログラムを実行することによって行われる処理のフローを示す図である。

10

【図3】文書分類装置による処理の対象とする文書の例を示す図である。

【図4】構造データの例を示す図である。

【図5】構造データのレベル調整の例を示す図である。

【図6】文書の論理構造の抽出例を示す図である。

【図7】構造データのレベル調整のフローを示す図である。

【図8】CPUがプログラムを実行することによって形成される仮想的モジュールを示す図である。

【符号の説明】

【0029】

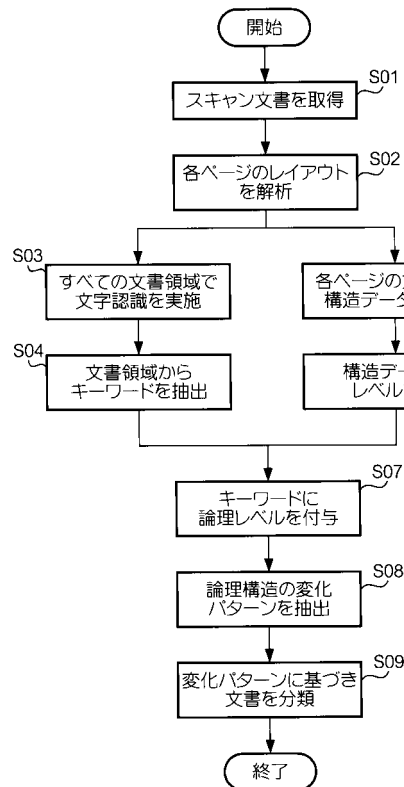
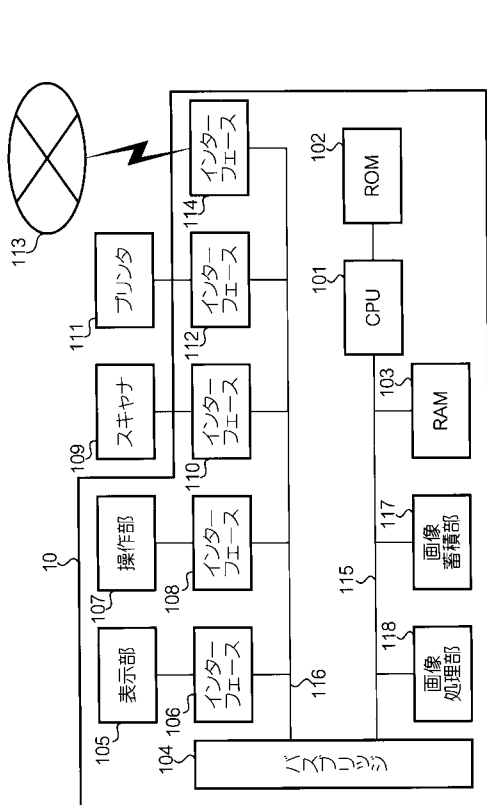
10 ... 文書分類装置、101 ... CPU、102 ... ROM、103 ... RAM、105 ... 表示部、107 ... 操作部、109 ... スキャナ、111 ... プリンタ、113 ... ネットワーク、117 ... 画像蓄積部、118 ... 画像処理部、

20

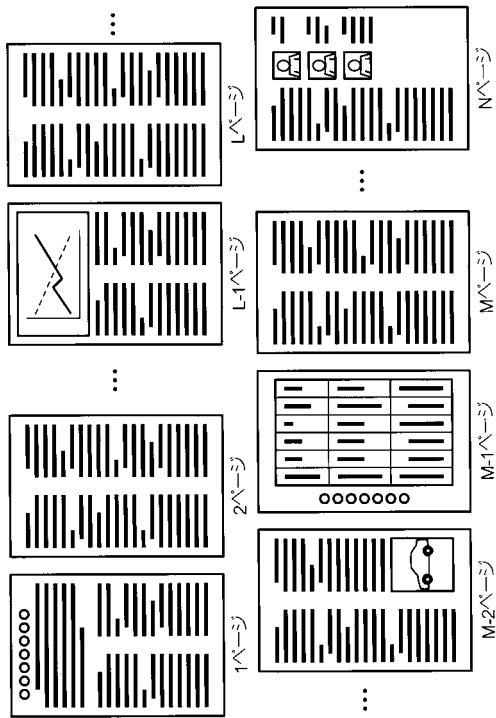
21 ... 画像取得手段、22 ... レイアウト解析手段、23 ... 文字認識手段、24 ... キーワード抽出手段、25 ... 構造データ生成手段、26 ... 論理構造抽出手段、27 ... 分類手段。

【図1】

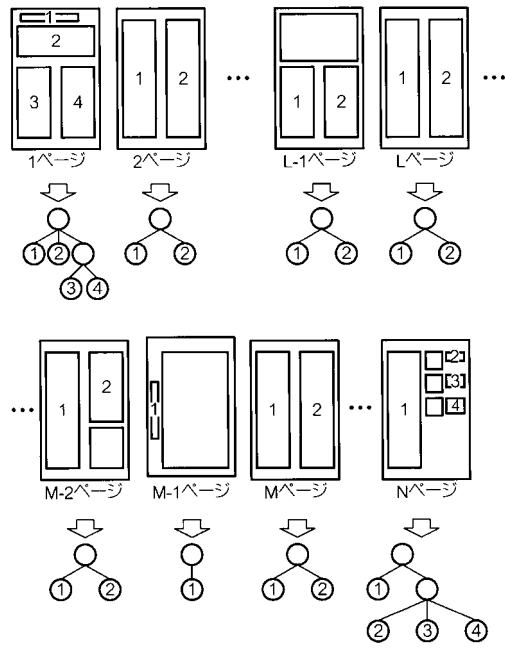
【図2】



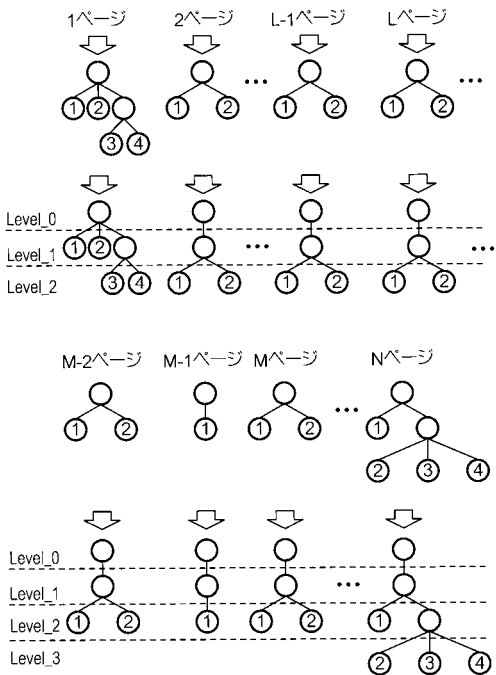
【図3】



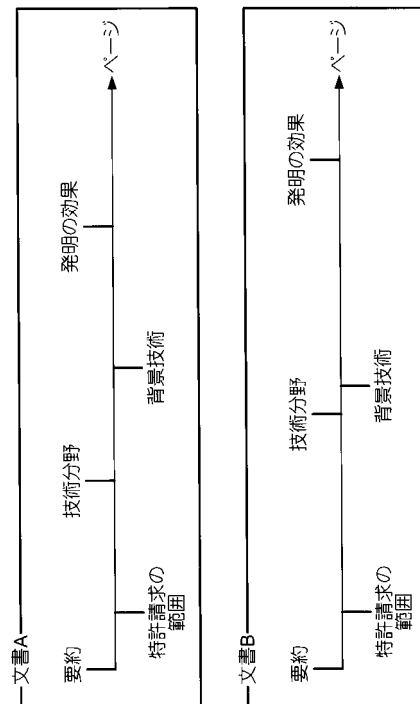
【図4】



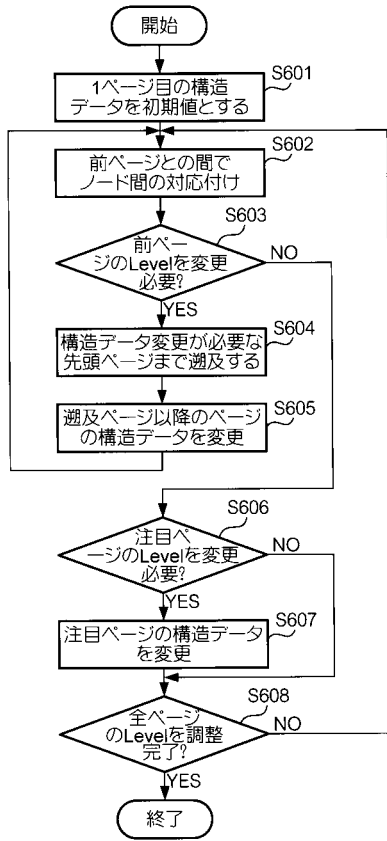
【図5】



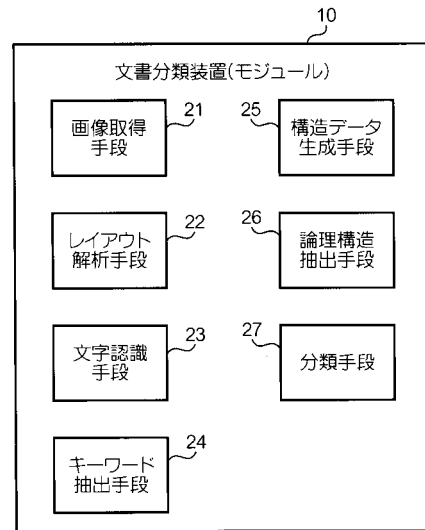
【図6】



【図7】



【図8】



フロントページの続き

- (56)参考文献 特開平06 - 214983 (JP, A)
特開2000 - 067080 (JP, A)
特開2003 - 288334 (JP, A)
特開2002 - 342343 (JP, A)
特開平09 - 297765 (JP, A)
特開平01 - 183784 (JP, A)
黄瀬浩一, 文書画像理解の目指すもの, 電子情報通信学会技術研究報告 (PRMU97 - 240
~ 250), 日本, 社団法人電子情報通信学会, 1998年 2月20日, 第97巻, 第559
号, p. 55 - 62

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