

【特許請求の範囲】**【請求項 1】**

少なくとも 2 以上のアプリケーションで作成された文書情報を取得する文書情報取得手段と、

取得した前記文書情報から、文書画像を生成する画像生成手段と、

取得した前記文書情報を、アプリケーション別に領域として分割する領域分割手段と、

分割された前記領域毎に、文字コードを抽出可能であるか否か判断する判断手段と、

前記文字コードを抽出可能と判断された前記領域から、第 1 の文字情報を抽出する第 1 の文字情報抽出手段と、

前記文字コードが抽出できないと判断された場合、前記文書画像に対して文字認識処理を施して得られた文字コードを第 2 の文字列情報として抽出する第 2 の文字情報抽出手段と、

前記第 1 の文字情報と、前記第 2 の文字情報と、前記文書情報及び前記文書画像の少なくともいずれか一方と、を対応付けて記憶する記憶手段と、

を備えることを特徴とする文書処理装置。

【請求項 2】

検索条件に基づいて、前記記憶手段の前記第 1 の文字情報及び前記第 2 の文字情報のいずれか一つ以上に対して検索を行う検索手段と、

検索条件を満たした、前記第 1 の文字情報及び前記第 2 の文字情報のいずれか一つ以上と、前記記憶手段で対応付けられている前記文書情報及び前記文書画像の少なくともいずれか一方を出力する出力手段と、

をさらに備えることを特徴とする請求項 1 に記載の文書処理装置。

【請求項 3】

前記第 1 の文字情報抽出手段は、さらに前記第 1 の文字情報の位置を示す第 1 の位置情報を抽出し、

前記記憶手段は、さらに前記第 1 の位置情報を対応付けて記憶すること、

を特徴とする請求項 1 又は 2 に記載の文書処理装置。

【請求項 4】

前記出力処理手段は、前記第 1 の文字情報を、前記記憶手段で対応付けられた前記第 1 の位置情報に配置して出力すること、

を特徴とする請求項 3 に記載の文書処理装置。

【請求項 5】

前記第 2 の文字情報抽出手段は、さらに前記第 2 の文字情報の位置を示す第 2 の位置情報を抽出し、

前記記憶手段は、さらに前記第 2 の位置情報を対応付けて記憶すること、

を特徴とする請求項 1 乃至 4 のいずれか一つに記載の文書処理装置。

【請求項 6】

前記表示処理手段は、前記第 2 の文字情報を、前記記憶手段で対応付けられた前記第 2 の位置情報に配置して出力すること、

を特徴とする請求項 5 に記載の文書処理装置。

【請求項 7】

前記取得した文書情報を読み込むアプリケーションをさらに備え、

前記画像生成手段は、前記文書処理装置内の印刷ドライバ内に含まれ、前記文書情報を読み込んだアプリケーションによる出力命令に従って前記文書画像データを生成すること、

を特徴とする請求項 1 乃至 6 のいずれか一つに記載の文書処理装置。

【請求項 8】

文書情報取得手段が、少なくとも 2 以上のアプリケーションで作成された文書情報を取得する文書情報取得ステップと、

画像生成手段が、取得した前記文書情報から、文書画像を生成する画像生成ステップと

10

20

30

40

50

、
領域分割手段が、取得した前記文書情報を、アプリケーション別に領域として分割する領域分割ステップと、

判断手段が、分割された前記領域毎に、文字コードを抽出可能であるか否か判断する判断ステップと、

第 1 の文字情報抽出手段が、前記文字コードを抽出可能と判断された前記領域から、第 1 の文字情報を抽出する第 1 の文字情報抽出ステップと、

第 2 の文字情報抽出手段が、前記文字コードが抽出できないと判断された場合、前記文書画像に対して文字認識処理を施して得られた文字コードを第 2 の文字列情報として抽出する第 2 の文字情報抽出ステップと、

前記第 1 の文字情報と、前記第 2 の文字情報と、前記文書情報及び前記文書画像の少なくともいずれか一方と、を対応付けて記憶手段に記憶する記憶ステップと、

を有することを特徴とする文書処理方法。

【請求項 9】

検索手段が、検索条件に基づいて、前記記憶手段の前記第 1 の文字情報及び前記第 2 の文字情報のいずれか一つ以上に対して検索を行う検索ステップと、

出力手段が、検索条件を満たした、前記第 1 の文字情報及び前記第 2 の文字情報のいずれか一つ以上と、前記記憶手段で対応付けられている前記文書情報及び前記文書画像の少なくともいずれか一方を出力する出力ステップと、

をさらに有することを特徴とする請求項 8 に記載の文書処理方法。

【請求項 10】

前記第 1 の文字情報抽出ステップは、さらに前記第 1 の文字情報の位置を示す第 1 の位置情報を抽出し、

前記記憶ステップは、さらに前記第 1 の位置情報を対応付けて記憶すること、

を特徴とする請求項 8 又は 9 に記載の文書処理方法。

【請求項 11】

前記出力処理ステップは、前記第 1 の文字情報を、前記記憶手段で対応付けられた前記第 1 の位置情報に配置して出力すること、

を特徴とする請求項 10 に記載の文書処理方法。

【請求項 12】

前記第 2 の文字情報抽出ステップは、さらに前記第 2 の文字情報の位置を示す第 2 の位置情報を抽出し、

前記記憶ステップは、さらに前記第 2 の位置情報を対応付けて記憶すること、

を特徴とする請求項 8 乃至 11 のいずれか一つに記載の文書処理方法。

【請求項 13】

前記表示処理ステップは、前記第 2 の文字情報を、前記記憶手段で対応付けられた前記第 2 の位置情報に配置して出力すること、

を特徴とする請求項 12 に記載の文書処理方法。

【請求項 14】

アプリケーションが、前記取得した文書情報を読み込む読込ステップをさらに有し、

前記画像生成ステップは、前記文書処理装置内の印刷ドライバ内に含まれた前記画像生成手段が、前記文書情報を読み込んだアプリケーションによる出力命令に従って前記文書画像データを生成すること、

を特徴とする請求項 8 乃至 13 のいずれか一つに記載の文書処理方法。

【請求項 15】

コンピュータを、

少なくとも 2 以上のアプリケーションで作成された文書情報を取得する文書情報取得手段と、

取得した前記文書情報から、文書画像を生成する画像生成手段と、

取得した前記文書情報を、アプリケーション別に領域として分割する領域分割手段と、

10

20

30

40

50

分割された前記領域毎に、文字コードを抽出可能であるか否か判断する判断手段と、
前記文字コードを抽出可能と判断された前記領域から、第１の文字情報を抽出する第１
の文字情報抽出手段と、

前記文字コードが抽出できないと判断された場合、前記文書画像に対して文字認識処理
を施して得られた文字コードを第２の文字列情報として抽出する第２の文字情報抽出手段
と、

前記第１の文字情報と、前記第２の文字情報と、前記文書情報及び前記文書画像の少な
くともいずれか一方と、を対応付けて記憶手段に記憶させる登録手段と、

として機能させることを特徴とする文書処理プログラム。

【請求項１６】

10

検索条件に基づいて、前記記憶手段の前記第１の文字情報及び前記第２の文字情報のい
ずれか一つ以上に対して検索を行う検索手段と、

検索条件を満たした、前記第１の文字情報及び前記第２の文字情報のいずれか一つ以上
と、前記記憶手段で対応付けられている前記文書情報及び前記文書画像の少なくともい
ずれか一方を出力する出力手段と、

をさらに機能させることを特徴とする請求項１５に記載の文書処理プログラム。

【発明の詳細な説明】

【技術分野】

【０００１】

本発明は、文書処理装置、文書処理方法及び文書処理プログラムに関するものである。

20

【背景技術】

【０００２】

近年、コンピュータ関連技術の向上、ネットワーク環境が整備によって文書の電子化が
進んでいる。これによりオフィスのペーパーレス化が促進されている。

【０００３】

そして、文書の電子化が進むにつれ、電子化された文書データを一括管理し、検索を行
えるようにしたいという要求がある。

【０００４】

そこで、特許文献１は、描画コードから、テキスト情報、すなわち文字コードを抽出し
、抽出したテキスト情報と、描画コードから生成された文書画像データとを対応付けてい
る。この描画コードから文書画像データを生成している以上、文字コード等を含む中間デ
ータであると考えられる。このため、描画コードから容易に文字コードを抽出している。

30

【０００５】

【特許文献１】特開平８－２１２３３１号公報

【発明の開示】

【発明が解決しようとする課題】

【０００６】

ところで、文書データに画像データとして図や表を入れたりする場合も多い。このよう
な図や表は画像データ等で貼り付けられている。また、図や表に限らず、ＨＴＭＬで表現
されるＷｅｂのページの場合には見栄えを重視して文字も画像として入れる場合が頻繁に
見受けられる。

40

【０００７】

しかしながら、引用文献１に記載された技術では、描画コードから文字コードを抽出す
るのでは、図や表を示す画像データが埋め込まれている場合に抽出できないという問題が
ある。

【０００８】

一方で、文書データから生成した文書画像データに対して、文字認識を掛けて、図や表
から画像データを抽出するのでは文字の抽出精度が落ちるという問題がある。

【０００９】

本発明は、上記に鑑みてなされたものであって、文字の抽出精度を向上させた文書処理

50

装置、文書処理方法及び文書処理プログラム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】

【0010】

上述した課題を解決し、目的を達成するために、請求項1にかかる発明は、少なくとも2以上のアプリケーションで作成された文書情報を取得する文書情報取得手段と、取得した前記文書情報から、文書画像を生成する画像生成手段と、取得した前記文書情報を、アプリケーション別に領域として分割する領域分割手段と、分割された前記領域毎に、文字コードを抽出可能であるか否か判断する判断手段と、前記文字コードを抽出可能と判断された前記領域から、第1の文字情報を抽出する第1の文字情報抽出手段と、前記文字コードが抽出できないと判断された領域に含まれる第2の文字情報を、前記文書画像に対して文字認識処理を施して得られた文字コードから抽出する第2の文字情報抽出手段と、前記第1の文字情報と、前記第2の文字情報と、前記文書情報及び前記文書画像の少なくともいずれか一方と、を対応付けて記憶する記憶手段と、を備えることを特徴とする。

10

【0011】

また、請求項2にかかる発明は、請求項1にかかる発明において、検索条件に基づいて、前記記憶手段の前記第1の文字情報及び前記第2の文字情報のいずれか一つ以上に対して検索を行う検索手段と、検索条件を満たした、前記第1の文字情報及び前記第2の文字情報のいずれか一つ以上と、前記記憶手段で対応付けられている前記文書情報及び前記文書画像の少なくともいずれか一方を出力する出力手段と、をさらに備えることを特徴とする。

20

【0012】

また、請求項3にかかる発明は、請求項1又は2にかかる発明において、前記第1の文字情報抽出手段は、さらに前記第1の文字情報の位置を示す第1の位置情報を抽出し、前記記憶手段は、さらに前記第1の位置情報に対応付けて記憶すること、を特徴とする。

【0013】

また、請求項4にかかる発明は、請求項3にかかる発明において、前記出力処理手段は、前記第1の文字情報を、前記記憶手段で対応付けられた前記第1の位置情報に配置して出力すること、を特徴とする。

【0014】

また、請求項5にかかる発明は、請求項1乃至4のいずれか一つにかかる発明において、前記第2の文字情報抽出手段は、さらに前記第2の文字情報の位置を示す第2の位置情報を抽出し、前記記憶手段は、さらに前記第2の位置情報に対応付けて記憶すること、を特徴とする。

30

【0015】

また、請求項6にかかる発明は、請求項5にかかる発明において、前記表示処理手段は、前記第2の文字情報を、前記記憶手段で対応付けられた前記第2の位置情報に配置して出力すること、を特徴とする。

【0016】

また、請求項7にかかる発明は、請求項1乃至6のいずれか一つにかかる発明において、前記取得した文書情報を読み込むアプリケーションをさらに備え、前記画像生成手段は、前記文書処理装置内の印刷ドライバ内に含まれ、前記文書情報を読み込んだアプリケーションによる出力命令に従って前記文書画像データを生成すること、を特徴とする。

40

【0017】

また、請求項8にかかる発明は、文書情報取得手段が、少なくとも2以上のアプリケーションで作成された文書情報を取得する文書情報取得ステップと、画像生成手段が、取得した前記文書情報から、文書画像を生成する画像生成ステップと、領域分割手段が、取得した前記文書情報を、アプリケーション別に領域として分割する領域分割ステップと、判断手段が、分割された前記領域毎に、文字コードを抽出可能であるか否か判断する判断ステップと、第1の文字情報抽出手段が、前記文字コードを抽出可能と判断された前記領域から、第1の文字情報を抽出する第1の文字情報抽出ステップと、第2の文字情報抽出手

50

段が、前記文字コードが抽出できないと判断された領域に含まれる第２の文字情報を、前記文書画像に対して文字認識処理を施して得られた文字コードから抽出する第２の文字情報抽出ステップと、前記第１の文字情報と、前記第２の文字情報と、前記文書情報及び前記文書画像の少なくともいずれか一方と、を対応付けて記憶手段に記憶する記憶ステップと、を有することを特徴とする。

【００１８】

また、請求項９にかかる発明は、請求項８にかかる発明において、検索手段が、検索条件に基づいて、前記記憶手段の前記第１の文字情報及び前記第２の文字情報のいずれか一つ以上に対して検索を行う検索ステップと、出力手段が、検索条件を満たした、前記第１の文字情報及び前記第２の文字情報のいずれか一つ以上と、前記記憶手段で対応付けられている前記文書情報及び前記文書画像の少なくともいずれか一方を出力する出力ステップと、をさらに有することを特徴とする。

10

【００１９】

また、請求項１０にかかる発明は、請求項８又は９にかかる発明において、前記第１の文字情報抽出ステップは、さらに前記第１の文字情報の位置を示す第１の位置情報を抽出し、前記記憶ステップは、さらに前記第１の位置情報を対応付けて記憶すること、を特徴とする。

【００２０】

また、請求項１１にかかる発明は、請求項１０にかかる発明において、前記出力処理ステップは、前記第１の文字情報を、前記記憶手段で対応付けられた前記第１の位置情報に配置して出力すること、を特徴とする。

20

【００２１】

また、請求項１２にかかる発明は、請求項８乃至１１のいずれか一つにかかる発明において、前記第２の文字情報抽出ステップは、さらに前記第２の文字情報の位置を示す第２の位置情報を抽出し、前記記憶ステップは、さらに前記第２の位置情報を対応付けて記憶すること、を特徴とする。

【００２２】

また、請求項１３にかかる発明は、請求項１２にかかる発明において、前記表示処理ステップは、前記第２の文字情報を、前記記憶手段で対応付けられた前記第２の位置情報に配置して出力すること、を特徴とする。

30

【００２３】

また、請求項１４にかかる発明は、請求項８乃至１３のいずれか一つにかかる発明において、アプリケーションが、前記取得した文書情報を読み込む読込ステップをさらに有し、前記画像生成ステップは、前記文書処理装置内の印刷ドライバ内に含まれた前記画像生成手段が、前記文書情報を読み込んだアプリケーションによる出力命令に従って前記文書画像データを生成すること、を特徴とする。

【００２４】

また、請求項１５にかかる発明は、コンピュータを、少なくとも２以上のアプリケーションで作成された文書情報を取得する文書情報取得手段と、取得した前記文書情報から、文書画像を生成する画像生成手段と、取得した前記文書情報を、アプリケーション別に領域として分割する領域分割手段と、分割された前記領域毎に、文字コードを抽出可能であるか否か判断する判断手段と、前記文字コードを抽出可能と判断された前記領域から、第１の文字情報を抽出する第１の文字情報抽出手段と、前記文字コードが抽出できないと判断された領域に含まれる第２の文字情報を、前記文書画像に対して文字認識処理を施して得られた文字コードから抽出する第２の文字情報抽出手段と、前記第１の文字情報と、前記第２の文字情報と、前記文書情報及び前記文書画像の少なくともいずれか一方と、を対応付けて記憶手段に記憶させる登録手段と、として機能させることを特徴とする。

40

【００２５】

また、請求項１６にかかる発明は、請求項１５にかかる発明において、検索条件に基づいて、前記記憶手段の前記第１の文字情報及び前記第２の文字情報のいずれか一つ以上に

50

対して検索を行う検索手段と、検索条件を満たした、前記第 1 の文字情報及び前記第 2 の文字情報のいずれか一つ以上と、前記記憶手段で対応付けられている前記文書情報及び前記文書画像の少なくともいずれか一方を出力する出力手段と、をさらに機能させることを特徴とする。

【発明の効果】

【0026】

本願発明によれば、文書情報を作成したアプリケーションに拘わらず、文字情報を抽出可能にすると共に、抽出される文字情報の精度を向上させるという効果を奏する。

【発明を実施するための最良の形態】

【0027】

以下に添付図面を参照して、この発明にかかる文書処理装置、文書処理方法及び文書処理プログラムの最良な実施の形態を詳細に説明する。

【0028】

(第 1 の実施の形態)

図 1 は、第 1 の実施の形態にかかる文書処理装置 100 の構成を示すブロック図である。本図に示すように文書処理装置 100 の内部は、文書メタ DB 101 と、データ格納部 102 と、文書取得部 103 と、出力部 104 と、表示処理部 105 と、検索部 106 と、文書アプリケーション 107 と、データ参照アプリケーション 108 と、判断部 109 と、テキスト情報抽出部 110 と、プリンタドライバ 111 と、文字認識テキスト情報抽出部 112 と、登録部 113 と、入力受付部 115 と、種別判別部 116 と、領域分割部 117 と、を備え、取得した文書データを解析した後、当該文書データの管理を行う。また、当該文書データの検索を可能とする。

【0029】

文書データとは、任意のアプリケーションで作成された、文字コード等を含む電子文書と、ページ毎に画像として表されている文書画像データとを含むものとする。文書画像データとしては、例えば紙文書をスキャナで読み込んだデータ等とする。本実施の形態にかかる文書データは、特に、画像データやオブジェクトを含む文書データとし、換言すれば文書作成用アプリケーション（例えば、文書アプリケーション 107）、並びに画像作成アプリケーション及び表作成アプリケーションの少なくとも一方という、複数のアプリケーションで作成されたデータを意味する。

【0030】

データ格納部 102 は、取得した文書データや、当該文書データから生成されたページ画像データを格納する。当該文書データに関する詳細な情報は、後述する文書メタデータベース 101 で管理されている。

【0031】

文書メタデータベース 101 は、文書管理テーブルと、ページ管理テーブルと、領域管理テーブルと、テキスト管理テーブルとを有している。

【0032】

図 2 は、文書管理テーブルのテーブル構造を示した図である。本図に示すように、文書管理テーブルは、文書 ID と、タイトルと、作成更新日と、ページ数と、ファイルフォーマットと、ファイルパスと、ファイル名とを対応付けて保持する。

【0033】

文書 ID は、文書データ毎に付与されたユニークな ID であり、これにより文書データを特定できる。タイトルは文書データのタイトルである。作成更新日は、文書データの作成日又は最終更新日を保持する。ページ数は文書データのページ数を保持している。ファイルフォーマットは、文書データ毎のフォーマットを保持している。これにより、管理している文書データが、アプリケーションで作成された電子文書、Web ページ及び文書画像データ等のうちいずれかのフォーマットであるか特定することができる。

【0034】

ファイルパスは、文書データが格納された場所を示している。そして、ファイル名は、

10

20

30

40

50

文書データのファイル名を示している。

【 0 0 3 5 】

図 3 は、ページ管理テーブルのテーブル構造を示した図である。本図に示すように、ページ管理テーブルは、ページ ID と、文書 ID と、ページ番号と、ページ画像データパスとを対応付けて保持している。次に各フィールドについて説明する。

【 0 0 3 6 】

ページ ID は、文書データを構成するページ毎に付与されたユニークな ID であり、この ID により、文書データのページを一意に特定できる。文書 ID は、当該ページを含んでいる文書データを特定する ID とする。ページ番号は、当該ページを含んでいる文書データ中における、当該ページのページ番号とする。ページ画像データパスは、当該ページを表したページ画像データが格納されている場所を保持する。

10

【 0 0 3 7 】

図 4 は、領域管理テーブルのテーブル構造を示した図である。本図に示すように、領域管理テーブルは、領域 ID と、文書 ID と、ページ ID と、領域座標と、種別と、タイトルと、テキストと、を対応付けて保持している。次に各フィールドについて説明する。

【 0 0 3 8 】

領域管理テーブルの領域 ID は、文書データに含まれるページを分割した領域毎に付与されたユニークな ID であり、この ID により、文書データの各ページに含まれている領域を特定できる。文書 ID とページ ID は、当該領域を含んでいる文書データ及びページを特定する ID とする。領域座標は、当該領域を特定する座標を保持し、本実施の形態では左上の頂点座標と右下の頂点座標を保持することで当該領域を特定する。

20

【 0 0 3 9 】

領域管理テーブルの種別は、当該領域のデータの種別を特定する情報を保持する。データの種別としては、例えばテキスト、画像（画像作成アプリケーションで作成された画像データ）、図（組織図、フローチャート、ガントチャート等の図作成用アプリケーションで作成されたオブジェクト等も含む）、写真、表（例えば表作成アプリケーションで作成されたオブジェクト等も含む）、グラフ（円グラフ、棒グラフなど表計算アプリケーションで作成されたオブジェクト等も含む）等とする。これら種別を、文書データを構成する文書構成要素とする。

【 0 0 4 0 】

タイトルは、当該領域を示すタイトルを保持する。テキストは、当該領域に含まれていたテキスト情報を保持する。このテキスト情報は、当該領域の種別がテキストの場合に限るものではなく、画像等であっても文字認識処理を実行して抽出したテキスト情報を保持する。当該フィールドに対して文字列を検索キーとして検索を行うことで、当該文字列を含む領域を特定できる。

30

【 0 0 4 1 】

図 5 は、テキスト管理テーブルのテーブル構造を示した図である。本図に示すように、テキスト管理テーブルは、文書 ID と、ページ ID と、領域 ID と、文字情報と、開始位置座標と、字間と、行間と、文字揃えと、文字組と、フォントサイズとを、を対応付けて保持している。文書 ID、ページ ID 及び領域 ID により、どの文書のどのページのどの領域のテキストを管理しているか特定できる。文字情報は、当該領域から抽出された一つ又は複数の文字コードを格納している。開始位置座標は、抽出元の領域内において、抽出された文字コードの位置座標（例えば左上の頂点座標）を表している。フォントサイズは抽出された文字コードのフォントサイズを示している。

40

【 0 0 4 2 】

上述したデータ格納部 102 及び文書メタデータベース 101 は、HDD (Hard Disk Drive)、光ディスク、メモリカード、RAM (Random Access Memory) などの一般的に利用されているあらゆる記憶装置により構成することができる。

【 0 0 4 3 】

50

文書取得部 103 は、文書データを取得する。図 6 は、本実施の形態で処理対象となる文書データを示す図である。図 6 に示すように、文書データは、表 601 が画像データとして含むとともに、縦組みの文章 603 も画像として含むものとする。これら表 601、文章 603 は、文書処理装置 100 においてそれぞれ別領域として扱われる。つまり、文書データに含まれる画像データやオブジェクトはこれらを作成したアプリケーションに応じて分割される。

【0044】

また、テキストの領域についても、領域 602 など破線で分けられた領域毎に別領域として扱われる。取得した文書データは、フォーマット形式に応じて適切なアプリケーション、又は当該アプリケーションの API を呼び出して、読み込みが行われるものとする。以下にアプリケーションの例について説明する。

10

【0045】

文書アプリケーション 107 は、第 1 のフォーマット形式の文書データを読み込み可能なアプリケーションとする。当該文書アプリケーション 107 では、第 1 のフォーマット形式の文書データの参照、編集が可能とする。

【0046】

データ参照アプリケーション 108 は、第 2 のフォーマット形式の文書データを読み込み可能なアプリケーションとする。データ参照アプリケーション 108 では、第 2 のフォーマット形式の文書データを参照のみ可能とする。

【0047】

20

なお、文書アプリケーション 107 では、第 2 のフォーマット形式の文書データを認識できず、データ参照アプリケーション 108 では、第 1 のフォーマット形式の文書データを認識できないものとする。

【0048】

プリンタドライバ 111 は、画像生成部 114 を備え、文書データの出力命令を受け付ける。本実施の形態においては、出力命令を受け付けた場合、画像生成部 114 を呼び出して、ページ画像データの生成を行う。

【0049】

画像生成部 114 は、プリンタドライバ 111 が文書データの出力命令を受け付けた場合に、当該文書データの全表示領域を表すページ画像データを生成する。本実施の形態にかかる画像生成部 114 は、文書データのページ毎にページ画像データを生成する。

30

【0050】

領域分割部 117 は、文書データの各ページを、当該ページに含まれるオブジェクト及び画像データ別（つまりデータを作成したアプリケーション別）、並びに文書に含まれる段落若しくはコラム別に、領域として分割する。ページに含まれる領域を分割する手法については、周知の手法を問わず、どのような手法を用いても良い。

【0051】

種別判別部 116 は、文書データを構成する各領域に対して、当該領域を表す種別を判別する。本実施の形態にかかる種別判別部 116 は、種別判別の対象となる領域から特徴量を抽出し、抽出された特徴量に基づいてパターン認識処理を行うことで、領域の種別を判断する。その際に用いるパターン認識手法としては、どのような手法を用いても良いが、例えばニューラルネットやサポートベクターマシン手法を用いてもよい。これらニューラルネットやサポートベクターマシン手法を用いることで、学習用のデータセットを作成し学習させることで、より精度の高い領域の識別の判断を行うことができる。

40

【0052】

判断部 109 は、文書データのページ毎に画像データのみ含むページであるか否か判断する。また、判断部 109 は、当該ページが画像データ以外の文書構成要素を含むページと判断した場合、当該ページを構成する領域毎に、判別された種別が文字コードを抽出可能な種別であるか否か判断する。文字コードが抽出可能な種別としては、テキスト要素の他に、例えば表計算アプリケーションで作成されたオブジェクトを表す文書構成要素など

50

が考えられる。

【 0 0 5 3 】

テキスト情報抽出部 1 1 0 は、種別がテキスト要素等の文字コードが抽出可能と判断された領域から、テキスト情報、すなわち文字コード群を抽出する。なお、本実施の形態は、テキスト情報抽出部 1 1 0 がテキスト情報を抽出する領域として、テキスト要素以外で文字コードが抽出可能な種別としては、文字コードが抽出可能なデータ、例えば任意のアプリケーションで作成されたオブジェクトデータ等が考えられる。

【 0 0 5 4 】

文字認識テキスト情報抽出部 1 1 2 は、画像生成部 1 1 4 により生成されたページ画像データの範囲のうち、判断部 1 0 9 により種別がテキスト要素等ではないと判断された領域に対応する範囲に対して文字認識処理を実行して得られた文字コードを、当該領域のテキスト情報として抽出する。この文字認識処理は、O C R (Optical Character Recognition) に用いられている文字認識処理とする。

【 0 0 5 5 】

このように、文字認識テキスト情報抽出部 1 1 2 は、上述した文字認識処理を行う前に、予め文書データの領域と、生成されたページ画像データの範囲とを対応付けておく。これにより、文字コードを抽出できない領域から、テキスト情報を抽出することが可能となる。なお、範囲と領域とを対応付ける手法については、周知の手法を問わず、どのような手法を用いても良い。

【 0 0 5 6 】

つまり、本実施の形態にかかる文書処理装置 1 0 0 では、テキスト情報抽出部 1 1 0 でテキスト情報を抽出すると共に、当該テキスト情報抽出部 1 1 0 で抽出できない領域については、当該領域に対応するページ画像データの表示領域に対して文字認識処理を実行して、テキスト情報を抽出する。これにより、文書データに含まれるテキスト情報は全て抽出できる。ところで、文字認識処理によるテキスト情報の抽出は、文字認識に誤りが生じる可能性がある。そこで、本実施の形態では、文字コードを抽出可能な領域については、文字認識処理を実行せずに、テキスト情報を抽出することとした。これにより、文書データから生成された画像データの全領域に文字認識処理を実行してテキスト情報を抽出した場合より、テキスト情報の抽出精度を向上させることができる。

【 0 0 5 7 】

登録部 1 1 3 は、文書メタデータベース 1 0 1 に対して、文書データ、当該文書データのページ及び領域に関する情報の登録、並びに文書データ及び画像データをデータ格納部 1 0 2 に格納する処理を行う。なお、詳細な登録手順については後述する。

【 0 0 5 8 】

また、上述した種別判別部 1 1 6、判断部 1 0 9、文字認識テキスト情報抽出部 1 1 2、テキスト情報抽出部 1 1 0、登録部 1 1 3 は、文書アプリケーション 1 0 7 やデータ参照アプリケーション 1 0 8 に組み込まれるプラグイン形式のプログラムでも良いし、別アプリケーションとして作成しても良い。

【 0 0 5 9 】

入力受付部 1 1 5 は、図示しない入力デバイスから情報の入力を受け付ける。例えば、入力受付部 1 1 5 は、文書データの検索要求や、その際の文書データの検索条件等の入力を受け付ける。

【 0 0 6 0 】

検索部 1 0 6 は、入力された検索条件に従って、文書管理テーブル、ページ管理テーブル、領域管理テーブル及びテキスト管理テーブルに対して検索を行う。また、検索部 1 0 6 は、領域管理テーブル及びテキスト管理テーブルに対して検索を行うことで、テキスト情報抽出部 1 1 0 により抽出されたテキスト情報、及び文字認識テキスト情報抽出部 1 1 2 により抽出されたテキスト情報のいずれか一つ以上に検索を行うことができる。

【 0 0 6 1 】

表示処理部 1 0 5 は、図示しない表示装置に対して、任意の情報の表示処理を行う。例

10

20

30

40

50

例えば、例えば、表示処理部 105 は、検索部 106 の検索結果として、検索条件に一致するテキスト情報を含むページ画像データ等の一覧の表示処理を行う。なお、本実施の形態とは異なるが、表示処理部 105 が、検索条件に一致するテキスト情報を含む文書データを表示しても良い。この後、入力受付部 115 が、表示されたページ画像データの一覧から、利用したい文書データのページ画像データの選択を受け付ける。

【0062】

出力部 104 は、利用者により選択されたページ画像データを含む文書データを、当該文書データの読み込み可能なアプリケーション（例えば文書アプリケーション 107）に出力する。例えば、入力受付部 115 が、表示処理部 105 により表示処理されたページの画像データに対する選択を受け付けた場合に、出力部 104 が、選択を受け付けた画像データをページとして含む文書データの出力を行う等が考えられる。これにより、利用者は選択した文書を利用することができる。

10

【0063】

次に、以上のように構成された本実施の形態にかかる文書処理装置 100 における文書を取得してから文書メタデータベース 101 に登録するまでの処理について説明する。図 7 は、本実施の形態にかかる文書処理装置 100 における上述した処理の手順を示すフローチャートである。なお、当該フローチャートでは、一般的となっている文書画像認識技術を最大限利用しているものとする。

【0064】

まず、文書取得部 103 が、スキャナ 151 や外部記憶装置 150 等の外部環境から、文書データを取得する（ステップ S701）。

20

【0065】

そして、文書取得部 103 が、取得した文書データを認識可能なアプリケーション（例えば、文書アプリケーション 107 など）の API を呼び出して、当該文書データの読み込み処理を行う（ステップ S702）。これにより、文書データの内容を把握することができる。

【0066】

次に、登録部 113 が、取得した文書データを、データ格納部 102 に格納すると共に、文書管理テーブルに文書データに関する情報の登録を行う（ステップ S703）。文書データに関する情報としては、文書管理テーブルの各フィールドに示されたタイトル、作成更新日、ページ数、ファイルフォーマット、当該文書データの格納先となるファイルパス、当該文書データのファイル名とする。

30

【0067】

そして、判断部 109 が、当該文書データに含まれるページにおいて、全領域が画像データであるか否か判断する（ステップ S704）。なお、ステップ S703 は、最初は 1 ページ目について判断を行い、当該ステップが繰り返される度に次ページを判断するように処理を行う。

【0068】

次に、当該ページの全領域が画像データと判断された場合（ステップ S704：Yes）、登録部 113 が、当該ページを表す画像データと共に、当該ページの情報を対応付けて、ページ管理テーブルに登録する（ステップ S705）。ページに関する情報としては、当該ページを含む文書の文書 ID と、当該ページのページ番号、当該ページを著す画像の格納先の画像データパスとする。なお、当該ページを表す画像データは、登録部 113 により、データ格納部に格納されるものとする。

40

【0069】

そして、領域分割部 117 が、当該ページ画像を領域毎に分割する（ステップ S706）。なお、ページに含まれる領域の分割は、空白領域の幅に基づいて分割するなどの周知の手法を問わず、あらゆる手法で分割することができる。当該分割により、領域毎の領域座標を特定できる。

【0070】

50

次に、種別判別部 1 1 6 が、当該ページの分割された領域毎に、当該領域から抽出された特徴量に基づいて、当該領域の種別を判別する（ステップ S 7 0 7）。この領域からの特徴量の抽出は、種別判別部 1 1 6 が行うものとする。

【 0 0 7 1 】

そして、種別の判別が行われた後、文字認識テキスト情報抽出部 1 1 2 は、画像データの当該領域に対して文字認識処理を実行して、当該領域に含まれているテキスト情報、すなわち文字コード群を抽出する（ステップ S 7 0 8）。なお、文字認識テキスト情報抽出部 1 1 2 は、文字認識処理により文字コードを抽出する際、各文字コードの位置座標も取得する。なお、位置座標のみならず、適切なフォントサイズも取得する。なお、複数文字を一レコードに登録する場合、字間、行間、文字揃え等も取得する。

10

【 0 0 7 2 】

次に、登録部 1 1 3 が、当該領域に関する情報を領域管理テーブルに登録すると共に、当該領域から抽出されたテキストに関する情報を、テキスト管理テーブルに登録する（ステップ S 7 0 9）。当該領域に関する情報としては、当該領域を含むページのページ ID、当該ページを含む文書 ID、当該領域の領域座標、当該領域の種別、テキストなどが考えられる。

【 0 0 7 3 】

なお、テキスト管理テーブルに文字コードを登録する場合、テキストに関する情報として、抽出先の文書 ID、ページ ID 及び領域 ID と共に、当該領域 ID から抽出された文字コード（文字情報）と、当該文字コードの開始位置座標と、フォントサイズとを対応付けて登録する。開始位置座標とは、当該文字コードを含む矩形の左上座標とする。また、テキスト管理テーブルには、一文字単位でレコードとして登録しても良いし、複数の文字列をレコードに登録しても良い。複数の文字列をレコードに格納する場合、それぞれの文字列の位置座標を特定するための情報、例えば字間、行間、文字揃え、文字組などを対応付けて登録する。これにより、領域に対して文字の配置を再現することが可能となる。

20

【 0 0 7 4 】

次に、判断部 1 0 9 が、当該ページに含まれる全ての領域についてステップ S 7 0 6 ~ S 7 0 8 の処理を行ったか否か判断する（ステップ S 7 1 0）。処理を行っていない領域があると判断した場合（ステップ S 7 1 0 : N o）、処理を行っていない領域に対する種別の判別から開始する（ステップ S 7 0 7）。

30

【 0 0 7 5 】

そして、判断部 1 0 9 が、処理を行ったと判断した場合（ステップ S 7 1 0 : Y e s）、ステップ S 7 2 0 へと処理を進ませる。

【 0 0 7 6 】

一方、判断部 1 0 9 が、当該文書データに含まれるページにおいて、全領域が画像データではないと判断した場合（ステップ S 7 0 4 : N o）、文書データを読み込んだアプリケーションの A P I が、プリンタドライバ 1 1 1 に対して、当該ページの出力命令を行う。

【 0 0 7 7 】

そして、プリンタドライバ 1 1 1 は、当該出力命令に従って、画像生成部 1 1 4 が、当該ページの画像データを生成する（ステップ S 7 1 1）。

40

【 0 0 7 8 】

次に、登録部 1 1 3 が、生成された当該ページを表す画像データと共に、当該ページの情報を対応付けて、ページ管理テーブルに登録する（ステップ S 7 1 2）。

【 0 0 7 9 】

領域分割部 1 1 7 が、当該ページを領域毎に分割する（ステップ S 7 1 3）。なお、ページに含まれる領域の分割は、空白領域の幅に基づいて分割するなどの周知の手法を問わず、あらゆる手法で分割することができる。当該分割により、領域毎の領域座標を特定できる。

【 0 0 8 0 】

50

その後、種別判別部 1 1 6 が、分割された領域毎に、当該領域から抽出された特徴量に基づいて、当該領域の種別を判別する（ステップ S 7 1 4）。

【 0 0 8 1 】

次に、判断部 1 0 9 が、判別された種別に基づいて、当該領域がテキスト要素又は文字コードが抽出可能なオブジェクトであるか否か判断する（ステップ S 7 1 5）。

【 0 0 8 2 】

そして、判断部 1 0 9 が、当該領域がテキスト要素又は文字コードが抽出可能なオブジェクトであると判断した場合（ステップ S 7 1 4：Y e s）、テキスト情報抽出部 1 1 0 が、当該領域に含まれているテキスト情報、すなわち文字コードと文字コードの位置座標を抽出する（ステップ S 7 1 6）。なお、文字コードの位置座標とは、当該領域を含むページが印刷された場合に位置する当該文字コードの位置座標とする。この位置座標は、例えば、生成されたページ画像データから特定することができる。

10

【 0 0 8 3 】

一方、判断部 1 0 9 が、当該領域がテキスト要素又は文字コードが抽出可能なオブジェクトではないと判断した場合（ステップ S 7 1 4：N o）、文字認識テキスト情報抽出部 1 1 2 が、画像データの当該領域に対して文字認識処理を実行して、当該領域に含まれているテキスト情報、すなわち文字コード群を抽出する（ステップ S 7 1 7）。なお、文字認識テキスト情報抽出部 1 1 2 は、文字認識処理を実行して文字コードを抽出する際、各文字コードの位置座標も取得する。

【 0 0 8 4 】

20

その後、登録部 1 1 3 が、当該領域に関する情報を領域管理テーブルに登録すると共に、当該領域から抽出されたテキスト情報を、テキスト管理テーブルに登録する（ステップ S 7 1 8）。

【 0 0 8 5 】

そして、判断部 1 0 9 が、当該ページに含まれる全領域について、ステップ S 7 1 5 ~ S 7 1 8 までの処理を行ったか否か判断する（ステップ S 7 1 9）。処理を行っていない領域があると判断した場合（ステップ S 7 1 9：N o）、処理を行っていない領域に対する種別の判別から開始する（ステップ S 7 1 4）。

【 0 0 8 6 】

一方、判断部 1 0 9 が、ページの全領域に対して処理を行ったと判断した場合（ステップ S 7 1 9：Y e s）、ステップ S 7 2 0 へと処理を進ませる。

30

【 0 0 8 7 】

次に、判断部 1 0 9 は、文書データに含まれる全ページに対して、ステップ S 7 0 4 ~ S 7 1 9 までの処理を行ったか否か判断する（ステップ S 7 2 0）。行っていないと判断した場合（ステップ S 7 2 0：N o）、再びステップ S 7 0 4 から処理を開始する（ステップ S 7 0 4）。

【 0 0 8 8 】

一方、判断部 1 0 9 は、文書データに含まれる全ページに対して、ステップ S 7 0 4 ~ S 7 1 9 までの処理を行ったと判断した場合（ステップ S 7 2 0：Y e s）、処理を終了する。

40

【 0 0 8 9 】

上述した処理手順では、電子文書进行处理する場合について具体的に説明したが、文書画像データの場合も同様に処理することができる。この場合、ステップ S 7 0 4 で常に Y e s と判断される。

【 0 0 9 0 】

上述した処理手順により、文書データに関する情報が文書管理テーブルに登録されると共に、当該文書データに含まれる各ページ、各領域に関する情報が、ページ管理テーブル及び領域管理テーブルに登録される。さらには、各領域に含まれているテキスト情報が、テキスト管理テーブルに、位置座標と対応付けて登録されることになる。このように、当該領域がテキスト要素であるか否かにかかわらず登録されることになる。

50

【 0 0 9 1 】

上述した処理手順では、画像データの解析手法と、文書データの解析手法を組み合わせることで容易に精度よく、当該文書データに含まれているテキスト情報を抽出することができる。

【 0 0 9 2 】

また、現在、文書データのフォーマット形式は、アプリケーション毎に様々な形式のフォーマットがある。そして、これら様々な形式のフォーマットに対応するために、本実施の形態では、上述した文書メタデータベース 1 0 1 に、フォーマット形式によらず全ての文書データの画像データ及びテキスト情報を登録することとした。これにより、文書フォーマットに関係なく、文書データの検索を行うことができる。

10

【 0 0 9 3 】

また、文書データの画像データを生成する際に、プリンタドライバ内の画像生成部 1 1 4 を用いることとした。これにより、文書データのフォーマット形式にかかわらず、文書データの画像データを生成できる。

【 0 0 9 4 】

次に、以上のように構成された本実施の形態にかかる文書処理装置 1 0 0 における文書データの検索処理について説明する。図 8 は、本実施の形態にかかる文書処理装置 1 0 0 における上述した処理の手順を示すフローチャートである。

【 0 0 9 5 】

まず、入力受付部 1 1 5 が、利用者が操作する入力デバイスによる、検索条件の入力を受け付ける（ステップ S 8 0 1 ）。

20

【 0 0 9 6 】

次に、検索部 1 0 6 が、入力された検索条件に従って、文書メタデータベース 1 0 1 に対して検索を行う（ステップ S 8 0 2 ）。この検索としては、例えば、利用者により入力された文字列を検索キーとして、領域管理テーブルのテキストフィールドや、テキスト管理テーブルの文字情報フィールドに対して検索を行うことが考えられる。または、ページ管理テーブルや、文書管理テーブルに対して検索条件を設定して検索を行ってもよい。そして、検出された領域、テキスト情報等からページ ID、文書 ID 等を特定できる。

【 0 0 9 7 】

そこで、表示処理部 1 0 5 が、検索条件を満足する領域やテキスト情報を含むページ画像データの一覧を、表示装置に表示する（ステップ S 8 0 3 ）。なお、当該一覧に表示されたページ画像データ毎に、ページ ID が対応付けられているものとする。これにより、利用者が、ページ画像データを選択した場合に、ページ ID が入力されることになる。

30

【 0 0 9 8 】

次に、入力受付部 1 1 5 が、表示されたページ画像データの一覧から、利用者が選択したページ画像データを示すページ ID の入力を受け付ける（ステップ S 8 0 4 ）。

【 0 0 9 9 】

そして、出力部 1 0 4 が、利用者が選択したページ画像データを示すページ ID から、当該ページを含む文書データを特定し、当該文書データを読み込み可能なアプリケーションに出力する（ステップ S 8 0 5 ）。

40

【 0 1 0 0 】

上述した処理手順により、利用者は検索条件に該当する文書データを利用することが可能となる。また、文書メタデータベース 1 0 1 のテーブル群に対して検索を行うので、文書データのフォーマット形式を意識することなく、検索を行うことを可能としている。

【 0 1 0 1 】

また、本実施の形態においては、検索終了後に検索条件に一致するページ画像データを表示したが、ページ画像データに制限するものではなく、文書データを表示しても良い。

【 0 1 0 2 】

ところで、従来、任意のアプリケーションで作成された文書データに画像データが埋め込まれている場合、当該画像データに含まれている文字列を検索できないという問題が生

50

じていた。そこで本実施の形態にかかる文書処理装置 100 では、文書データに画像データが含まれている場合であっても適切に抽出して、文書メタデータベース 101 に登録するので、検索可能となる。

【0103】

上述した実施の形態にかかる文書処理装置 100 によれば、文書データの各ページの領域に応じて、テキストの抽出処理を変更しているので、解析精度を向上させることができる。また、一般的な文字認識処理（OCR）と、文書データのテキスト要素又はオブジェクトからの文字コードの抽出を組み合わせることで、簡易に精度よい解析が可能となる。

【0104】

また、本実施の形態にかかる文書処理装置 100 によれば、文書データの画像データで表現された図や表に対して、文字認識処理を実行することで、テキスト情報を抽出している。これにより、文書データに図や表が画像データとして貼り付けられている場合でも、解析が可能となる。さらに図や表に含まれているテキスト情報で検索することが可能となる。

【0105】

さらに、文書処理装置 100 では、文書画像データであっても、テキスト情報の抽出が可能である。このように文書画像データ、文書データを問わず検索対象とすることが可能となる。

【0106】

上述した文書処理装置 100 では、画像データを生成し、当該画像データからテキスト情報を抽出することで、当該領域に含まれているオブジェクトのフォーマットと依存せず解析が可能である。

【0107】

上述した文書処理装置 100 では、文書管理テーブル、ページ管理テーブル、領域管理テーブル及びテキスト管理テーブルに対して検索を行うことができるので、文書、ページ領域、テキスト情報単位での検索が可能となる。

【0108】

（変形例）

また、上述した各実施の形態に限定されるものではなく、以下に例示するような種々の変形が可能である。

【0109】

上述した第 1 の実施の形態にかかる文書処理装置 100 では、検索結果としてページ画像データの一覧を表示していた。しかしながら、ページ画像データの一覧の表示に制限するものではない。

【0110】

例えば、表示処理部 105 が、ページ画像データの代わりに、領域管理テーブル及びテキスト管理テーブルが保持する情報を利用して各ページを再現して表示しても良い。具体的には、領域管理テーブル及びテキスト管理テーブルが保持する情報を取得して、検索された領域を含むページについて、当該ページを構成する文字コードを、当該文字コードと対応付けられた位置座標に配置することで、当該ページを再現できる。

【0111】

上述した処理手順により、文書データのフォーマット形式にかかわらず、当該文書データのページを再現して表示することができる。さらに、ページ画像データではなく文字コードを用いてページを再現するので、文書処理装置 100 の処理負担を軽減させることができる。

【0112】

図 9 は、文書処理装置 100 の機能を実現するためのプログラムを実行した PC のハードウェア構成を示した図である。本実施の形態の文書処理装置 100 は、CPU（Central Processing Unit）901 と、ROM（Read Only Memory）902 や RAM（Random Access Memory）903 と、HDD（Hard Disk Drive）、CD（Compact Disk）

10

20

30

40

50

ドライブ装置等の外部記憶装置 905 と、ディスプレイ装置等の表示装置 906 と、キーボードやマウス等の入力デバイス 907 と、通信 I/F 904 と、これらを接続するバス 908 を備えており、通常のコンピュータを利用したハードウェア構成となっている。

【0113】

本実施形態の文書処理装置 100 で実行される文書処理プログラムは、インストール可能な形式又は実行可能な形式のファイルで CD-ROM、フレキシブルディスク (FD)、CD-R、DVD (Digital Versatile Disk) 等のコンピュータで読み取り可能な記録媒体に記録されて提供される。

【0114】

また、本実施形態の文書処理装置 100 で実行される文書処理プログラムを、インターネット等のネットワークに接続されたコンピュータ上に格納し、ネットワーク経由でダウンロードさせることにより提供するように構成しても良い。また、本実施形態の文書処理装置 100 で実行される文書処理プログラムをインターネット等のネットワーク経由で提供または配布するように構成しても良い。

【0115】

また、本実施形態の文書処理装置 100 で実行される文書処理プログラムを、ROM等に予め組み込んで提供するように構成してもよい。

【0116】

上述した実施の形態の文書処理装置 100 で実行される文書処理プログラムは、文書処理装置 100 において上記記録媒体から読み出して実行することにより RAM 902 上にロードされ、上記ソフトウェア構成で説明した各部が RAM 903 上に生成されるようになっている。

【産業上の利用可能性】

【0117】

以上のように、本発明にかかる文書処理装置、文書処理方法及び文書処理プログラムは、文書データの処理に有用であり、特に、画像を含む文書データを検索可能に格納する技術に適している。

【図面の簡単な説明】

【0118】

【図 1】実施の形態にかかる文書処理装置の構成を示すブロック図である。

【図 2】文書管理テーブルのテーブル構造を示した図である。

【図 3】ページ管理テーブルのテーブル構造を示した図である。

【図 4】領域管理テーブルのテーブル構造を示した図である。

【図 5】テキスト管理テーブルのテーブル構造を示した図である。

【図 6】実施の形態で処理対象となる文書データを示す図である。

【図 7】実施の形態にかかる文書処理装置における、文書を取得してから文書メタデータベースに登録するまでの処理の手順を示すフローチャートである。

【図 8】本実施の形態にかかる文書処理装置における、文書データの検索処理の手順を示すフローチャートである。

【図 9】本実施の形態にかかる文書処理装置の機能を実現するためのプログラムを実行した PC のハードウェア構成を示した図である。

【符号の説明】

【0119】

- 100 文書処理装置
- 101 文書メタデータベース
- 102 データ格納部
- 103 文書取得部
- 104 出力部
- 105 表示処理部
- 106 検索部

10

20

30

40

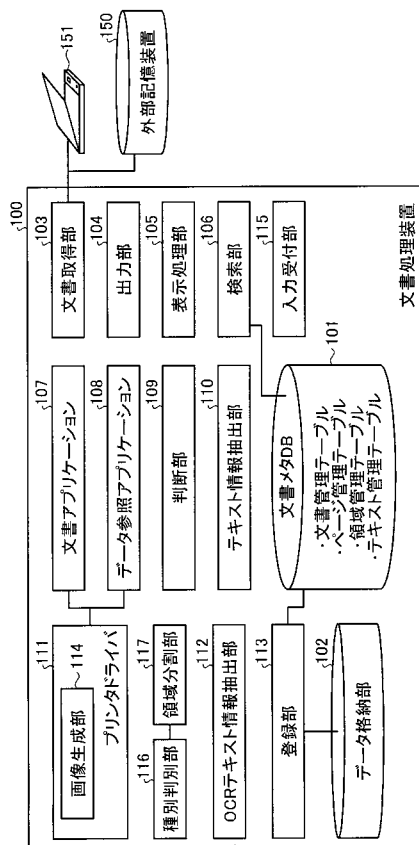
50

- 1 0 7 文書アプリケーション
- 1 0 8 データ参照アプリケーション
- 1 0 9 判断部
- 1 1 0 テキスト情報抽出部
- 1 1 1 プリントドライバ
- 1 1 2 文字認識テキスト情報抽出部
- 1 1 3 登録部
- 1 1 4 画像生成部
- 1 1 5 入力受付部
- 1 1 6 種別判別部
- 1 5 0 外部記憶装置
- 1 5 1 スキャナ
- 9 0 1 C P U
- 9 0 2 R O M
- 9 0 3 R A M
- 9 0 4 通信 I / F
- 9 0 5 外部記憶装置
- 9 0 6 表示装置
- 9 0 7 入力デバイス
- 9 0 8 バス

10

20

【 図 1 】



【 図 2 】

文書ID	タイトル	作成更新日	ページ数	ファイルフォーマット	ファイルパス	ファイル名
DOC0001	イメージについて	2005/11/19	22	doc	/doc/image.doc	image.doc
DOC0002	文書の書き方	2007/1/14	4	tiff	/doc/doc.tiff	doc.tiff
⋮	⋮	⋮	⋮	⋮	⋮	⋮

【 図 3 】

ページID	文書ID	ページ番号	画像データベース
P000001	DOC0001	1
⋮	⋮	⋮	⋮

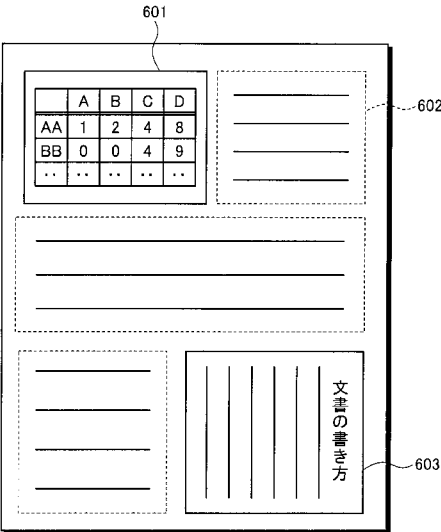
【 図 4 】

領域ID	文書ID	ページID	領域座標	種別	タイトル	テキスト
R000001	DOC0001	P000001	(0,0)-(500,200)	テキスト	イメージ	今...
⋮	⋮	⋮	⋮	⋮	⋮	⋮

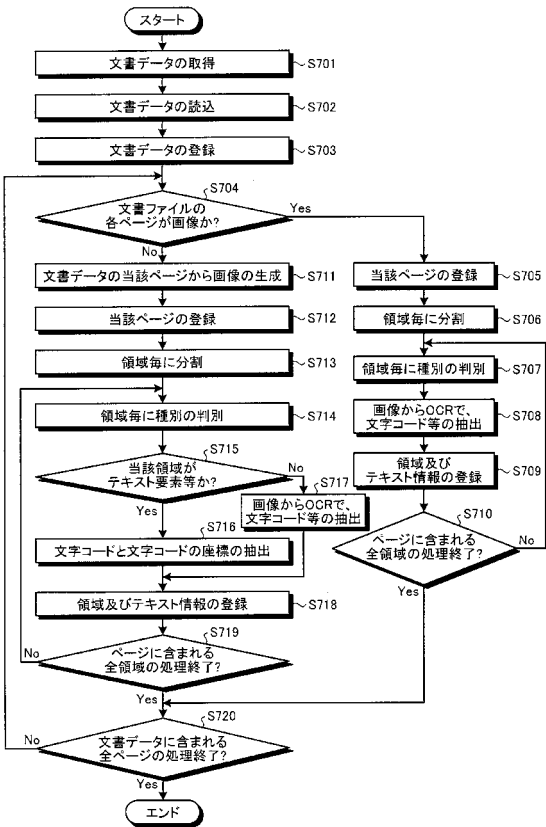
【 図 5 】

	文書ID	ページID	領域ID	文字情報	開始位置座標	字間	行間	文字揃え	文字組	フォントサイズ
	DOC0001	P000001	R000001	今.....	10,15	0	4	左寄せ	縦組み	21p
	DOC0002	P002001	R002001	秋	30,20	-	-	-	-	21p
										..

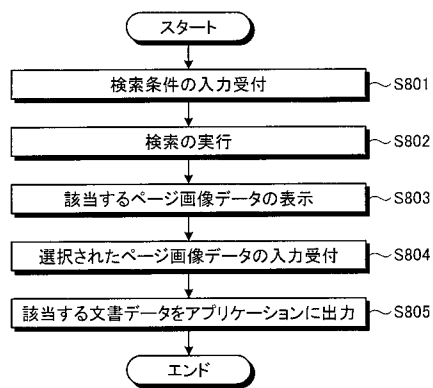
【 図 6 】



【 図 7 】



【 図 8 】



【 図 9 】

