

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4785655号
(P4785655)

(45) 発行日 平成23年10月5日(2011.10.5)

(24) 登録日 平成23年7月22日(2011.7.22)

(51) Int. Cl.		F I	
HO4N	1/387	(2006.01)	HO4N 1/387
GO6T	11/60	(2006.01)	GO6T 11/60 100A
GO6T	3/00	(2006.01)	GO6T 3/00 400J

請求項の数 14 (全 26 頁)

(21) 出願番号	特願2006-190826 (P2006-190826)	(73) 特許権者	000001007 キヤノン株式会社 東京都大田区下丸子3丁目30番2号
(22) 出願日	平成18年7月11日(2006.7.11)	(74) 代理人	100076428 弁理士 大塚 康德
(65) 公開番号	特開2008-22159 (P2008-22159A)	(74) 代理人	100112508 弁理士 高柳 司郎
(43) 公開日	平成20年1月31日(2008.1.31)	(74) 代理人	100115071 弁理士 大塚 康弘
審査請求日	平成21年7月8日(2009.7.8)	(74) 代理人	100116894 弁理士 木村 秀二
		(72) 発明者	高田 智美 東京都大田区下丸子3丁目30番2号 キヤノン株式会社内

最終頁に続く

(54) 【発明の名称】 文書処理装置及び文書処理方法

(57) 【特許請求の範囲】

【請求項1】

文書処理装置であって、
文書画像から複数の物理ページを抽出する第1抽出手段と、
前記第1抽出手段によって抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第2抽出手段と、
 前記第2抽出手段によって抽出された夫々のオブジェクトの有するテキストの特徴を解析し、当該特徴に基づいて少なくとも1つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定手段と、
前記判定手段によって前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合手段と、
を有することを特徴とする文書処理装置。

10

【請求項2】

文書処理装置であって、
文書画像から複数の物理ページを抽出する第1抽出手段と、
前記第1抽出手段によって抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第2抽出手段と、
 前記第2抽出手段によって抽出された夫々のオブジェクトの有する表の特徴を解析し、当該特徴に基づいて少なくとも1つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定手段と、

20

前記判定手段によって前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合手段と、
を有することを特徴とする文書処理装置。

【請求項 3】

文書処理装置であって、
文書画像から複数の物理ページを抽出する第 1 抽出手段と、
前記第 1 抽出手段によって抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第 2 抽出手段と、

前記第 2 抽出手段によって抽出された夫々のオブジェクトの色や形状の特徴を解析し、当該特徴に基づいて少なくとも 1 つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定手段と、

10

前記判定手段によって前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合手段と、
を有することを特徴とする文書処理装置。

【請求項 4】

文書処理装置であって、
文書画像から複数の物理ページを抽出する第 1 抽出手段と、
前記第 1 抽出手段によって抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第 2 抽出手段と、

前記第 2 抽出手段によって抽出された夫々のオブジェクトの位置関係を解析し、当該位置関係に基づいて少なくとも 1 つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定手段と、

20

前記判定手段によって前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合手段と、
を有することを特徴とする文書処理装置。

【請求項 5】

前記結合手段は、前記複数の物理ページの位置又は倍率に基づいて、前記判定手段によって前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合することを特徴とする請求項 1 乃至 4 の何れか 1 項に記載の文書処理装置。

【請求項 6】

前記結合手段は、前記オブジェクトの有するテキストのサイズと位置座標とに応じて、前記判定手段によって前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合することを特徴とする請求項 1 に記載の文書処理装置。

30

【請求項 7】

前記第 2 抽出手段は夫々のオブジェクトに関するメタデータを抽出し、前記オブジェクトと前記抽出されたメタデータとを関連付けて格納する格納手段を更に有することを特徴とする請求項 1 乃至 6 の何れか一項に記載の文書処理装置。

【請求項 8】

前記第 1 抽出手段によって抽出された前記複数の物理ページの夫々のレイアウトを解析するレイアウト解析手段と、

40

前記レイアウト解析手段によって解析されたレイアウトに基づいて前記文書画像の論理構造を解析する論理構造解析手段とを更に有し、

前記第 2 抽出手段は、前記論理構造解析手段によって解析された論理構造とページ構成に基づいてメタデータを抽出することを特徴とする請求項 7 に記載の文書処理装置。

【請求項 9】

オブジェクトを検索するための検索条件を入力するための検索条件入力手段と、前記検索条件入力手段によって入力された検索条件に基づいてオブジェクトに関連付けられたメタデータを検索する検索手段とを更に有することを特徴とする請求項 8 に記載の文書処理装置。

【請求項 10】

50

文書処理装置の文書処理方法であって、
 文書画像から複数の物理ページを抽出する第1抽出工程と、
 前記第1抽出工程において抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第2抽出工程と、
 前記第2抽出工程において抽出された夫々のオブジェクトの有するテキストの特徴を解析し、当該特徴に基づいて少なくとも1つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定工程と、
 前記判定工程において前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合工程と、
 を有することを特徴とする文書処理方法。

10

【請求項11】

文書処理装置の文書処理方法であって、
 文書画像から複数の物理ページを抽出する第1抽出工程と、
 前記第1抽出工程において抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第2抽出工程と、
 前記第2抽出工程において抽出された夫々のオブジェクトの有する表の特徴を解析し、当該特徴に基づいて少なくとも1つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定工程と、
 前記判定工程において前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合工程と、
 を有することを特徴とする文書処理方法。

20

【請求項12】

文書処理装置の文書処理方法であって、
 文書画像から複数の物理ページを抽出する第1抽出工程と、
 前記第1抽出工程において抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第2抽出工程と、
 前記第2抽出工程において抽出された夫々のオブジェクトの色や形状の特徴を解析し、当該特徴に基づいて少なくとも1つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定工程と、
 前記判定工程において前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合工程と、
 を有することを特徴とする文書処理方法。

30

【請求項13】

文書処理装置の文書処理方法であって、
 文書画像から複数の物理ページを抽出する第1抽出工程と、
 前記第1抽出工程において抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第2抽出工程と、
 前記第2抽出工程において抽出された夫々のオブジェクトの位置関係を解析し、当該位置関係に基づいて少なくとも1つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定工程と、
 前記判定工程において前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合工程と、
 を有することを特徴とする文書処理方法。

40

【請求項14】

コンピュータを、請求項1乃至9の何れか1項に記載の文書処理装置の各手段として機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、文書処理装置及び文書処理方法に関する。

50

【背景技術】

【0002】

近年、電子文書の普及に伴い、それらを有効活用したいという需要が高まっている。

【0003】

図1は、電子文書に対する処理の流れの一例を示す図である。図1に示すように、電子文書に対する操作は、一度作成・利用した後に蓄積・保存し、更にこれを編集・加工することによって新しい文書を作成するなど、文書作成のコスト削減のために再利用するのが一般的である。一方、印刷文書もコンピュータに取り込み、その内容を再利用したい、という要求がある。

【0004】

印刷文書や電子文書を効率的に再利用するためには、大量の文書の中から必要な情報を探し出すための検索技術が重要となる。文書の中には、オブジェクトデータとして、文字情報だけでなく、図、表、写真等の画像情報も含まれており、特に利用頻度が高いと考えられる。文書に含まれる文字情報の場合は、指定された検索語と文字情報のマッチングを行うことで容易に検索することができる。しかし、画像情報等の場合は、それ自体は文字情報をもたないため、画像情報等に検索のためのメタデータを付加する技術が提案されている。

【0005】

文書画像を複数の領域に分割し、各領域の特徴量によりテキストや画像等の種類を識別する技術（例えば、特許文献1参照）が提案されている。

【特許文献1】特開2000-293671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、文書では、情報量の多い図表や画像を1ページに記述すると小さくて見難いため、複数のページにまたがって記述することがある。また、雑誌等の書籍では、向かい合った左右の2ページに1つの内容を記載する見開きを使用することがよくあり、これらは物理的には2ページ、論理的には1ページとみなすことができる。

【0007】

このような物理的に複数のページが論理的な1ページである文書画像から、画像やテキスト等のオブジェクトを抽出する場合、次のような問題があった。

【0008】

文書を物理ページ毎に読み込むと、複数の物理ページにまたがって記述されている1つのオブジェクトが分割されて抽出される。

【0009】

また、ページ画像を読み込む際にはページ画像や領域毎に色や濃度を最適化し、ページ画像毎に倍率を調整するため、分割された画像毎に異なった画像処理が施される。また、ページ画像に歪みが発生することもあり、分割された画像を単に結合するだけでは、必ずしも元の画像が得られない。

【0010】

また、論理ページを構成する各物理ページを別々に読み込み、ページ画像として合成した場合、左右の物理ページの間空白が入ったり、位置が上下にずれたりすることがある。このような文書画像から抽出される画像やテキスト等のオブジェクトは、複数の物理ページにまたがって記述されている1つのオブジェクトが分割されて抽出されたものである。

【0011】

また、ページ画像から分割して抽出された全てのオブジェクトについて、分割されたオブジェクトの各領域の違いを解析し、色情報や位置・倍率等の全ての要素を正確に補正して結合することは文書処理装置のCPUに大変負荷のかかる処理である。また、補正して結合した各オブジェクトを文書処理装置の内部に保持しておく、ディスクの負荷が増大

10

20

30

40

50

する。

【0012】

本発明は、文書中の見開きのようなページに含まれる分割されたオブジェクトデータを1つのオブジェクトデータとして有効に利用することを目的とする。

【課題を解決するための手段】

【0013】

本発明は、文書処理装置であって、文書画像から複数の物理ページを抽出する第1抽出手段と、前記第1抽出手段によって抽出された前記複数の物理ページの夫々からオブジェクトを抽出する第2抽出手段と、前記第2抽出手段によって抽出された夫々のオブジェクトの有するテキストの特徴を解析し、当該特徴に基づいて少なくとも1つのオブジェクトが前記複数の物理ページにまたがっているか否かを判定する判定手段と、前記判定手段によって前記複数の物理ページにまたがっていると判定されたオブジェクト同士を結合する結合手段と、を有することを特徴とする。

10

【発明の効果】

【0014】

本発明によれば、文書中の見開きのようなページに含まれる分割されたオブジェクトデータを1つのオブジェクトデータとして有効に利用することができる。

【発明を実施するための最良の形態】

【0015】

以下、図面を参照しながら発明を実施するための最良の形態について詳細に説明する。

20

【0016】

[第1の実施形態]

第1の実施形態では、

図2は、本発明の一実施形態に係る文書処理システムが構築されるコンピュータ装置の基本構成を示すブロック図である。

【0017】

図2において、201はCPUであり、後述するROMやRAMのプログラムに従って第1の実施形態の文書処理装置における各種制御を実行する。また、CPU201自身の機能や計算機プログラムの機構により、複数の計算機プログラムを並列に動作させることができる。202はROMであり、CPU201の制御手順を記憶する計算機プログラムや制御データが格納されている。203はRAMであり、CPU201が処理するための制御プログラムを格納すると共にCPU201が各種制御を実行する際の作業領域を提供する。

30

【0018】

204はアルファベット、ひらがな、カタカナ、句点等を入力する文字記号入力キーや、カーソル移動を指示するカーソル移動キーのような各種機能キーを備えたキーボードであり、ユーザによる各種入力操作環境を提供する。また、マウスのようなポインティングデバイス、タッチパネル、スタイラスペンを含むこともできる。205はシステムバス(アドレスバス、データバスなど)であり、各構成を接続する。106は様々なデータなどを記憶するための外部記憶装置であり、ハードディスク、光ディスク、磁気ディスク、光磁気ディスク、不揮発性のメモリカード等の記録媒体と、記憶媒体を駆動し、情報を記録するドライブなどで構成される。保管された計算機プログラムやデータはキーボードなどの指示や各種計算機プログラムの指示により、必要な時にRAM上に完全或いは部分的に呼び出される。

40

【0019】

207は表示器であり、ディスプレイなどで構成され、各種入力操作の状態をユーザに対して表示する。208は他の通信装置等と通信を行うためのネットワークコントロールユニット(NCU)である。ネットワーク(LAN)などを介して不図示の遠隔地に存在する装置と通信し、プログラムやデータを共有することが可能になる。209は画像を読み取るためのイメージスキャナであり、セットされた紙原稿を1枚ずつ光学的に読み取り

50

、イメージ信号をデジタル信号列に変換する。読み取られた画像データは、外部記憶装置やRAM等に格納される。

【0020】

尚、通信手段としては、有線通信や無線通信など、何でも良く、またアダプタ装置などと接続され、通信を行っても良い。有線通信としては、RS232CやUSB、IEEE1394、P1284、SCSI、モデム、イーサネット（登録商標）などである。また無線通信としては、Bluetooth（登録商標）、赤外線通信、IEEE802.11xなどである。

【0021】

また、画像データは、イメージスキャナ209だけでなく、NCU208に接続されたネットワークスキャナやコピー装置等の入力機器を介して入力されても良い。読み取られた画像データも、外部記憶装置やRAMなどではなく、ネットワークに接続されたサーバやコピー機等の外部記憶装置等に格納しても良い。

10

【0022】

以上説明した構成は、第1の実施形態における一例であり、特にこれに限定されるものでない。

【0023】

図3は、見開きのページ画像に対して領域抽出処理を行った結果を示す図である。このページ画像は、見開きを構成する各物理ページを別々に読み込んだ後、ページ画像として合成したため、左右の物理ページの間に空白があり、また位置が上下にずれている。そのため、異なる物理ページにまたがって記述されている画像等のオブジェクトが分割されて抽出されている。

20

【0024】

図3において、300は見開きを構成する左右の物理ページを一度にスキャンしたページ画像である。316及び317は各々抽出された物理ページの領域である。301～314は抽出されたオブジェクトデータを示す領域である。301、313及び314は、303～312の本文を構成する領域とは空間的に離れているため、それぞれ独立した文字領域又は画像領域として抽出される。

【0025】

302～307と309～310は文字領域である。本実施形態では、文字列の方向が同じで、文字サイズと文字間値・行間値がほぼ均一であり、更に行方向の配置（字下げ、センタリング、揃えなど）が同じ部分が一つの文字領域として抽出される。308～310は画像領域であり、図として識別されている。315は後述する処理で抽出された物理ページの分割位置である。

30

【0026】

尚、詳細は後述するが、309及び310、311及び312は、それぞれ一つのオブジェクトを構成するが、分割されて抽出されている。また、図3は、第1の実施形態における領域抽出結果の一例を示す図であるが、画像と文字の領域が抽出できれば、他の領域抽出結果でも構わない。

【0027】

図4は、見開きページを物理ページ毎に読み込んだページ画像に対して領域抽出処理を行った結果を示す図である。図4において、400及び401はスキャンしたページ画像である。400は見開きの左側のページ画像であり、401は見開きの右側のページ画像であり、この例では位置が上下にずれている。402～412は抽出されたオブジェクトデータを示す領域である。402、403、412は、404～411の本文を構成する領域とは空間的に離れているため、それぞれ独立した文字又は画像領域として抽出される。

40

【0028】

404、405、407、410、411は文字領域である。これらの文字領域は、文字列の方向が同じで、文字サイズと文字間値・行間値がほぼ均一であり、更に行方向の配置（字下げ、センタリング、揃えなど）が同じ部分が一つの文字領域として抽出される。

50

406、408、409は画像領域であり、図として識別されている。413及び414は、抽出された物理ページの領域である。

【0029】

尚、詳細は後述するが、408及び409、410及び411は、一つのオブジェクトであるのに分割されて抽出されている。また、図4は、第1の実施形態における領域抽出結果の一例を示す図であるが、画像と文字の領域が抽出できれば、他の領域抽出結果でも構わない。

【0030】

また、図3、図4では、画像とテキストが混在した文書画像を例に挙げたが、必ずしも複数の種類のオブジェクトが混在する必要はなく、例えば画像のみで構成された文書画像であっても構わない。

10

【0031】

図5は、第1の実施形態における文書入力時の処理の一例を示すフローチャートである。この処理を示すプログラムは、ROM202に格納されており、CPU201によって実行される。

【0032】

尚、図5に示す処理の説明では、一例として、イメージスキャナ209などの入力機器で読み取られた紙文書を対象として説明を行う。しかし、紙文書だけでなく、ワードプロセッサや編集ソフトで作成した文書、HTMLなどで記述された文書、PDFなどの形式の電子文書でも構わない。

20

【0033】

但し、電子文書の場合、ステップS501の入力処理において、フォーマット変換などの処理が必要となる。また、文章を文字コードで保持している文書の場合は、ステップS503の文字認識処理は不要となる。

【0034】

まず、ステップS501において、CPU201は、イメージスキャナ209やネットワークに接続されたコピー機などの入力機器を用いて文書を読み取り、電子化されたページ単位の文書画像を得る。入力機器によって入力される文書画像には、2値画像、カラー画像などがある。ページ画像を読み込む際に、ページ画像や領域毎に色情報等を最適化したり、またページ毎に位置や向き等が異なったりすることがある。

30

【0035】

尚、電子化された文書画像を得た後、各ページ画像について、ノイズ除去処理や向きと傾きの補正処理を行っても良い。ページ画像の向きと傾きを判定し、修正する方法としては、公知のどのような方法を用いても構わない。

【0036】

次に、ステップS502において、CPU201は、ステップS501で読み取った文書の各ページ画像について領域分割を行う。そして、文字、図、表、写真などの画像を内包する矩形領域をその矩形の種類とサイズ、ページ内での位置座標等の物理的な情報と共に抽出する。

【0037】

尚、文字領域については、CPU201は、縦書き・横書きなどの文字列の読み方向と文字サイズを検出し、検出結果に基づいて文字列行と文字を抽出する。ここでは、文字列の方向が同じで、文字サイズと文字間値と行間値がほぼ均一である領域を一まとまりの文字領域として抽出する。尚、文字領域内の行方向の配置（字下げ、センタリング、揃えなど）を検出し、検出結果に基づいて文字領域を行方向に分割することで、更に、行方向の配置が同じ領域を一まとまりとしても良い。

40

【0038】

また、非文字領域については、写真、表、枠や線などを検出し、領域として抽出する。入力された文書画像がカラー画像などの多値の場合は、2値に変換することで同様に領域分割処理を行うことができる。この領域分割方法としては、公知のどのような方法でも構

50

わない。

【 0 0 3 9 】

次に、ステップ S 5 0 3 において、CPU 2 0 1 は、全ての文字領域に対して文字認識処理を行い、その処理結果を全て RAM 2 0 3 や外部記憶装置 2 0 6 などの記憶媒体に格納する。そして、ステップ S 5 0 4 において、文書の各ページ画像から物理ページを抽出する処理を行う。この処理は、自動又は手動で行う。自動で行う場合は、各ページ画像に対する物理ページの構成を判別する。そして、1枚のページ画像が複数の物理ページで構成されていれば、各ページ画像を物理ページ単位に分割する。物理ページの構成の判別は、ページ画像の縦横比率やステップ S 5 0 2 で抽出された領域を利用する。例えば、横長のページ画像において最上部・最下部にヘッダやページ番号と思われる左右(上下)対象の領域が存在するかによって判別される。この物理ページ構成の判別方法は一例であり、他にもいろいろな方法が考えられる。また、文書入力時にユーザが指定しても良い。

10

【 0 0 4 0 】

次に、ステップ S 5 0 5 において、CPU 2 0 1 は、文書の各物理ページにおけるレイアウトを抽出し、テキストや画像などのコンテンツの種類毎に矩形領域で分割する。そして、得られた矩形領域の物理的な情報に従って各物理ページ画像における各矩形領域の空間的な関係を抽出する。例えば、物理ページ画像内の2つの領域に対する空間的な関係を各矩形領域の位置座標やサイズを用いて解析し、判定する。空間的な関係としては、互いの領域が存在する上下左右の方向や、2つの領域が重なっている、接している、含まれているなどの状態、2つの領域の大小関係などである。また、2つの領域が接していない場合には、隣接する各領域間の物理ページ画像全体における距離の比較から遠い又は近いなどを判定する。また、文字領域については、物理ページ画像内の他の文字領域との位置を比較することにより、行方向の配置を抽出しても良い。

20

【 0 0 4 1 】

以上の解析結果は、物理ページ毎に木構造やネットワーク構造で表現することができる。ここで挙げた各矩形領域間の関係及びその表現方法は、第1の実施形態における一例であり、他の関係が抽出されても良いし、また解析結果を他の方法で表現しても構わない。例えば、レイアウトとして、各矩形領域の物理ページ全体に対する相対的な位置やサイズなどを抽出しても良い。

【 0 0 4 2 】

図6は、ある物理ページ画像における各領域の空間的な関係を抽出した結果の一例を示す図である。図6では、ページ画像内の2つの領域に対する空間的な関係、更に、2つの領域が接していない場合には、隣接する2つの領域間の相対的な距離をネットワーク構造で表現している。例えば、領域1と領域2の空間的な関係は、領域5が領域4の下にあり、接していないが近い距離にあることを示している。

30

【 0 0 4 3 】

図5に戻り、ステップ S 5 0 6 において、文書の全ての物理ページに対して、連続する複数の物理ページが論理的な1ページを構成しているか、或いは物理的な1ページが論理的な1ページであるかを自動又は手動で判別する。複数の物理ページから成る論理ページの判別を自動で行う場合は、文書の方向やステップ S 5 0 5 で抽出したレイアウト、即ち「左のページ」「右のページ」のような見開きページ内で使用される言語表現などを利用する。

40

【 0 0 4 4 】

尚、ページ番号などを利用して物理ページの連続性を判定し、連続する物理ページについてのみ論理ページの組を判別する。そして、不連続な物理ページについては判別を行わないようにすると効率良く判別できる。更に、文書の種類が折り込みページのない書籍の場合は、向かい合うページと背中合わせのページが必ず交互に並ぶことを考慮して論理ページ構成を判別しても良い。論理ページの判別方法はこれに限るものではなく、他にもいろいろなものが考えられる。

【 0 0 4 5 】

50

次に、ステップS507において、CPU201は、ステップS506の判別結果に基づき論理ページを取得する。そして、ステップS508において、ステップS507で取得した論理ページが見開きのように、複数の物理ページから成る論理ページであるか否かを判定する。複数の物理ページから成ると判定した場合はステップS509へ進み、1物理ページから成ると判定した場合はステップS512へ進む。

【0046】

このステップS509では、CPU201は、論理ページの組になる各物理ページ内の矩形領域についてサイズと位置を合わせる処理を行う。例えば、各ページの背景画像や飾り、抽出された矩形領域のレイアウトの規則性、物理ページの結合位置付近にある矩形領域の位置関係や位置座標・サイズなどを利用し、組になるページサイズの比率と、位置のずれを求めて調整する。矩形の範囲には誤差があるので、矩形の位置やサイズを調整しても、内部の画像やテキストが合致するとは限らないし、矩形のサイズや位置座標等の情報にも誤差があるので、ページ内の全ての矩形領域を完全に合致させるのは難しい。従って、完全に合わせる必要はなく、ある程度の誤差の範囲内で調整できれば良い。ページの倍率と位置を合わせる方法としてはこれに限るものではなく、他にもいろいろな方法が考えられる。

10

【0047】

次に、ステップS510において、CPU201は、論理ページの組になる各物理ページの結合位置付近にある2つの領域が物理ページによって分割された1つのオブジェクトか否かを判別する。この判別処理の詳細については、更に後述する。

20

【0048】

次に、ステップS511において、CPU201は、ステップS509、S510の結果に従って、ステップS505のレイアウト抽出結果に対する補正を行う。即ち、見開きなどの論理ページを対象としたレイアウト抽出処理を行い、ステップS505のレイアウト抽出結果に対して、論理ページに対するレイアウト情報を追加する。論理ページ上でのレイアウト情報として、第1の実施形態では、各領域の位置とサイズから各領域が物理ページのどちら側に属するか、或いは両方に属しているかなどの情報を追加する。補正方法や補正する情報はこれに限るものではなく、他にもいろいろなものが考えられる。

【0049】

次に、ステップS512において、全ての論理ページに対して、ステップS507からステップS511までの処理が終了したか否かを判定する。ここで、未処理の論理ページがある場合は、次の論理ページに対してステップS507からステップS511の処理を行う。

30

【0050】

以上、図5を用いて説明した文書の入力処理は、処理の一例であり、他にもいろいろなものが考えられる。これは、文書入力処理の一例であり、処理の順や処理内容は、これに限定されるものではない。

【0051】

また、第1の実施形態では、文書入力時に、各オブジェクトについての分割判別処理を行っているが、分割判別処理のタイミングとしては文書入力時に限定されるものではなく、他のタイミングで行うようにしても良い。

40

【0052】

図7は、ある文書におけるページ画像や各ページ画像から抽出された領域に関する各種物理的な情報の一例を示す図である。この例では、ページ画像に対して、ページサイズや読み込み時の解像度、電子化されたページ画像データの格納位置などの物理的な情報が付与されている。また、各ページ画像から抽出した物理ページについて、位置やサイズなどの情報と、同じ論理ページを構成している物理ページを示す情報が付与されている。

【0053】

また、抽出された各矩形領域に対して、文字領域、画像領域などの領域種別、矩形領域のサイズ、ページ内での位置座標等の物理的な情報が付与されている。更に、1つのオ

50

プロジェクトが分割されている矩形領域の場合には分割された他のオブジェクトを示す情報が付与されている。更に、文字領域については、文字サイズ、文字認識した結果である文字列が付与され、画像領域については、写真、表などの画像種別が付与されている。

【0054】

例えば、ページ画像1は、幅が290mm、高さが210mmで、処理解像度が300dpiであり、領域1と領域2はページ画像1から抽出された物理ページで見開きページである。また、領域6は、X座標20mm、Y座標50mmの位置にある、幅55mm、高さ50mmの文字領域で、文字サイズ9ポイントで記述されている文字列である。また、領域7及び領域9、領域8及び領域10は、異なる物理ページに分割された1つのオブジェクトである。

10

【0055】

図7は、領域の物理的な情報の一例を示しているが、物理的な情報とはこれに限るものではなく、次のステップにおいて、レイアウト抽出ができれば、他の情報が抽出されても良い。例えば、図7では、矩形領域のサイズと位置座標情報を抽出しているが、矩形領域の左上の位置座標と右下の位置座標を抽出するようにしても良い。

【0056】

図8は、図5に示すステップS510における判別処理の詳細を示すフローチャートである。この処理は、2つの物理ページそれぞれの結合位置付近にある2つの領域が、物理ページによって分割された1つのオブジェクトか否かを判別する処理である。この処理のプログラムは、ROM202に格納されており、CPU201によって実行される。

20

【0057】

まず、ステップS801において、CPU201は、位置情報に基づいて、論理ページの組となる2つの物理ページの結合位置付近にある2つの領域を取得する。そして、ステップS802において、領域に含まれるオブジェクトの種類が同じか否かを判定する。判定の結果、オブジェクトの種類が同じと判定した場合はステップS803へ進み、ステップS509で調整した領域の矩形のサイズと位置、ページ内のレイアウトなどを利用して2領域が1つのオブジェクトである可能性を判定する。

【0058】

例えば、図3では、物理ページを水平方向に結合するので、領域309と領域310の調整後の高さやY座標がほぼ同じであれば、1つのオブジェクトの可能性もある。また、例えば右側の物理ページの主な領域の左上X座標位置よりも、領域310は分割位置315に近く、また左側の物理ページの主な領域の右上X座標位置よりも、領域309は分割位置315に近い。これにより、領域309及び領域310は一つのオブジェクトである可能性が高いと言える。また、領域309及び領域410の距離は左右の物理ページ間の距離とほぼ一致することからも、一つのオブジェクトである可能性が高いと言える。

30

【0059】

尚、矩形の範囲やサイズ・位置座標等には誤差があるので、サイズや位置の比較を行う場合には、誤差とみなせる程度の違いであれば完全に一致していなくても良い。矩形領域のサイズと位置関係を利用して判定する方法としては、これに限るものではなく、他にもいろいろな方法が考えられる。

40

【0060】

次に、ステップS803において、CPU201は、2つの領域が1つのオブジェクトであると判定した場合はステップS804の処理へ進む。そして、領域に含まれるオブジェクトの種類毎にその特徴を利用して2領域が1つのオブジェクトである可能性を判定する。矩形の範囲には誤差があるので、矩形の位置やサイズを調整しても、内部の画像やテキストの位置やサイズが合致するとは限らない。そこで、領域内に記述されている各オブジェクトを解析することによって判定を行う。

【0061】

テキスト領域については、テキスト領域の文字の特徴や文字認識した文字列を利用する。例えば、領域内の文字サイズやスタイル、飾りなどが文書内の標準文字のそれと異なり

50

、かつ一致している場合は、1つのオブジェクトの可能性が高いと言える。また、例えば2つのテキスト領域を分割した場合と結合した場合の領域内の各テキスト文字列について、辞書とのマッチングや形態素解析を行い、解析の結果得られる評価値が大きい方が1つのオブジェクトである可能性が高い。また、見出しやキャプション等と思われる領域については、文字列の特徴を利用することによって判定できる。

【0062】

例えば、図3に示す領域311、領域312、領域307は、画像領域と接しており、それぞれ領域309、領域310、領域408のキャプションと識別できる。また、領域311及び領域307のテキスト文字列は「“図”+英数字+“：”」から始まっているのに対して領域312はこのパターンに当てはまらない。これにより、領域312は領域311と組になると考えられる。文字サイズや文字コードなどは、正確に認識できないこともあるので、完全に一致しなくても誤差とみなせる程度の違いであれば良い。

10

【0063】

表領域については、罫線の位置座標やセルのサイズ、マトリクス構造、セル内のテキストや画像などの領域を利用して判定する。これらの情報は、正確に認識できないこともあるので、完全に一致しなくても誤差とみなせる程度の違いであれば良い。

【0064】

写真等については、結合部分の色や形状等の情報を利用して判定する。画像の色等は、スキャン時にページ画像毎にチューニングされている場合があるので、完全に一致していても、誤差の範囲内で判定すれば良い。

20

【0065】

これらは、オブジェクトの特徴を利用して判定する方法の一例であり、他にもいろいろな方法が考えられる。

【0066】

次に、ステップS804において、1つのオブジェクトと判定された場合はステップS805へ処理が進む。そして、CPU201は、テキスト以外の領域についてテキスト領域との関係を利用して1つのオブジェクトである可能性を判定する。例えば、図3に示す領域311と領域312は、それぞれ画像領域309と画像領域310のキャプションであり、1つのオブジェクトと判定する(ステップS804)。そこで、領域309と領域310も、1つのオブジェクトと判定することができる。これは、テキスト領域との関係を利用して判定する方法の一例であり、他にもいろいろな方法が考えられる。

30

【0067】

次に、ステップS805において、2つの領域が1つのオブジェクトと判定された場合はステップS806へ処理が進む。CPU201は、ステップS801で取得した2つの領域の領域抽出結果に対して、一つの領域であることを示す情報を追加する。

【0068】

次に、ステップS807において、CPU201は、2つの物理ページの結合位置付近にある全ての領域に対して、ステップS801からステップS806の処理が終了したか否かを判定する。判定の結果、未処理の領域がある場合は、ステップS801に戻り、全領域について処理を終了するまで、上述の処理を繰り返す。

40

【0069】

以上、図8に示す処理は、ステップS510の処理の一例であり、他にも様々なものが考えられる。例えば、ステップS802～S805の全ての判定処理を行う必要はなく、処理内容や処理順序はこの通りでなくても良い。

【0070】

また、この例では、ステップS802～S805の何れかで可能性がないと判定された場合、それらは別の領域であると判定した。しかし、例えばステップS802～S805の何れかで可能性があるとして判定された場合、その確信度等によるポイントを加算し、全てのステップでの判断による総合ポイントによって判定を行っても良い。

【0071】

50

次に、

図9は、第1の実施形態における見開きページの補正・結合処理を示すフローチャートである。この処理プログラムは、ROM202に格納されており、CPU201によって実行される。

【0072】

まず、ステップS901において、CPU201は、処理対象となるオブジェクトを取得する。利用するオブジェクトは、利用する目的やアプリケーションなどによって異なる。また、オブジェクトではなく論理ページ画像を取得しても良い。

【0073】

次に、ステップS902において、CPU201は、ステップS901で取得した利用対象が分割されているか否かを判定する。即ち、利用対象がオブジェクトの場合は、1つのオブジェクトを含む領域が複数に分割されているか否かを判定する。また、利用対象が論理ページ画像の場合は、その論理ページ内に含まれるオブジェクトを含む領域が物理ページによって分割されているか否かにより判定する。分割されていると判定した場合はステップS903へ処理を進め、分割されていないと判定した場合は、この処理を終了する。

10

【0074】

図5を用いて説明したように、文書入力時に各オブジェクトについての分割判別処理を行っているので、ここではその情報を利用する。しかし、分割判別処理のタイミングとしては文書入力時に限るものではなく、ここで行うようにしても良い。

20

【0075】

次に、ステップS903において、CPU201は、分割されたオブジェクトを含む領域或いはページについて、色や濃度等を補正するか否かを判定する。補正すると判定した場合はステップS904へ進み、補正しないと判定した場合はステップS905へ進む。補正するか否かは、オブジェクトの種類や利用目的によって異なる。例えば、背景やページ飾りなどは、再利用性が低いので補正しなくても良い。

【0076】

ステップS904では、CPU201は、分割された領域又はページを色、濃度、倍率、或いは位置などについて正確に補正する処理を行う。色、濃度、倍率、位置の全てについて補正してもよいし、またこの中の何れか一つについて補正してもよいし、また、この中の組み合わせを補正しても良い。どのように補正するかは、オブジェクトの種類や利用目的によって異なる。

30

【0077】

色の補正は、分割されたオブジェクトを含む画像データについて、例えば各画像領域の彩度、明度、色調の分布を利用することで行える。分割された画像オブジェクトの位置や倍率の補正は、例えば各画像領域の境界部分から複数の対応点を抽出し、対応点のずれを利用して画像領域間の変換式を算出することで行える。テキストオブジェクトの位置や倍率の補正は、各テキスト領域部分の画像データについて、領域中の各テキスト行のサイズと位置を利用して補正することができる。尚、領域中の各テキスト行のサイズは、例えば図4や図5に示すように物理ページを結合するのであれば、高さである。

40

【0078】

表オブジェクトの補正は、オブジェクトを含む領域の画像データを補正してもよいし、表の罫線の位置座標やセルのサイズ、マトリクス構造などの情報を利用して補正してもよい。位置と倍率については正確に補正しない場合でも、図5のステップS509で矩形の位置・サイズを調整した際の情報を利用して、大体の位置と倍率を補正しても良い。また、論理ページ画像の場合は、ページ内に含まれるオブジェクトを含む領域を利用して補正する。

【0079】

上述した補正方法はこれに限るものではなく、他にもいろいろなもの考えられる。

【0080】

50

次に、ステップS905において、CPU201は、分割された領域又はページを結合する処理を行う。この結合処理は、分割された領域又はページについて、結合した画像データを生成するが、表オブジェクトの場合は、画像データを生成するのではなく、結合した表データを抽出しても良い。また、テキストオブジェクトの場合は、各テキスト領域部分を結合した画像データを生成し、再度文字認識処理を行って文字サイズや文字コード等の文字情報を抽出する。結合したデータは、利用後は破棄して構わない。結合方法はこれに限るものではなく、他にもいろいろなものが考えられる。

【0081】

図9に示す処理は、第1の実施例における利用時の補正・結合処理の一例であり、処理の順や処理内容は、この通りでなくても良い。

10

【0082】

第1の実施形態によれば、1つのオブジェクトデータが、複数のページに分割して記述されていることを判定できるようにすることにより、複数のページに含まれている分割されたオブジェクトデータを1つの領域として有効に利用することができる。

【0083】

また、分割されたオブジェクトデータを含む複数のページを1つのページとして有効に利用することができる。

【0084】

[第2の実施形態]

次に、図面を参照しながら本発明に係る第2の実施形態について詳細に説明する。第2の実施形態では、見開きのようなページに含まれている分割されたオブジェクトデータを1つのオブジェクトデータとして、必要に応じた精度で表示し、有効に利用する場合を説明する。

20

【0085】

尚、第2の実施形態における文書処理システムの構成は、第1の実施形態の構成と同様であり、その説明は省略する。

【0086】

図10は、第2の実施形態におけるオブジェクト表示時の処理を示すフローチャートである。この処理のプログラムは、ROM202に格納されており、CPU201によって実行される。

30

【0087】

第2の実施形態では、検索アプリケーションで検索を行った結果の一覧表示、検索結果を確認するために一覧の中から選択して拡大表示、一覧の中から選択したものを編集して再利用するための表示を想定している。しかし、検索結果の一覧表示だけでなく、例えば特定のフォルダやディレクトリ内に格納されているもの、又は何らかの方法でグループ化されたものを表示しても構わない。

【0088】

まず、ステップS1001において、CPU201は、表示対象データを取得する。通常、表示対象は、表示を行うアプリケーションなどによって異なるので、各表示プログラムに応じた適切なものを取得する。例えば、画像検索結果を表示する場合は、検索結果の自然画像や写真等の画像オブジェクトを含む領域のデータを取得し、表検索結果を表示する場合は検索結果として得られた表オブジェクトを含む領域のデータを取得する。尚、画像や表以外のオブジェクトを表示対象としても良いし、論理ページ画像を表示対象としても良い。

40

【0089】

次に、ステップS1002において、CPU201は、ステップS1001で取得した表示対象が分割されているか否かを判定する。即ち、表示対象がオブジェクトの場合は、1つのオブジェクトを含む領域が複数に分割されているか否かを判定する。また、表示対象が論理ページ画像の場合は、その論理ページ内に含まれるオブジェクトを含む領域が物理ページによって分割されているか否かにより判定する。分割されていると判定した場合

50

はステップS1003へ処理を進め、分割されていないと判定した場合はステップS1010へ処理を進める。

【0090】

図5を用いて説明したように、文書入力時に各オブジェクトについての分割判別処理を行っているので、ここではその情報を利用する。しかし、分割判別処理のタイミングとしては文書入力時に限るものではなく、ここで行うようにしても良い。

【0091】

次に、ステップS1003～S1005において、表示目的及び表示方法を判定する。また、ステップS1003～S1005に記述されたもの以外にも様々な表示目的及び表示方法がある。

10

【0092】

このステップS1003では、多くの表示対象を同時に表示する一覧表示か否かを判定する。ここで一覧表示と判定された場合にはステップS1006へ進み、上述のステップS1001で取得した表示対象を結合し、一覧表示のためのサムネイル画像データを生成する。一覧表示の場合、各表示対象は小さい画像であり、大体どのようなものが分かればよいので補正する必要はない。より厳密に処理する場合は、位置や倍率については図5に示すステップS509で矩形の位置・サイズを調整した際の情報を利用して補正しても良い。

【0093】

次に、ステップS1004では、CPU201は、ユーザに選択された特定のオブジェクトや論理ページを確認するための拡大表示か否かを判定する。ここで拡大表示と判定された場合にはステップS1007へ進み、ステップS1001で取得した表示対象を結合し、拡大表示するための画像データを生成する。その際、表示するデータの種類と表示の目的に応じて補正を行う。例えば、検索結果を確認する場合は、データの種類と検索アルゴリズムに応じて確認したい要素を補正して結合する。

20

【0094】

例えば、色特徴量による画像検索結果の場合は、ユーザは画像の色情報を確認したいと想定できる。よって、分割された各画像領域部分の色の違いを正確に補正して結合した画像データを生成する。色の補正は、例えば、分割された各画像領域の彩度、明度、色調の分布を利用することで行うことができる。また、形状特徴量による画像検索結果の場合は、ユーザは画像の形状情報を確認したいと思われるので、分割された各画像領域部分の画像データの倍率と位置を正確に補正して結合した画像データを生成する。位置や倍率の補正は、例えば分割された各画像領域の境界部分から複数の対応点を抽出し、対応点のずれを利用して画像領域間の変換式を算出することで行うことができる。補正には、画像特徴量抽出時の補正情報を利用して良い。

30

【0095】

また、表の場合は、分割された各表領域部分の画像データの位置や倍率などを補正して結合した画像データを生成しても良い。更に抽出された表の情報、即ち表に関する罫線の位置座標やセルのサイズ、マトリクス構造などを利用して、結合した表示用データを生成しても良い。

40

【0096】

また、論理ページ画像の場合は、ある程度のページ内容が分かればよいと思われるので、論理ページ内に含まれるオブジェクトを含む領域を利用して、物理ページ画像の位置や倍率をある程度補正して結合した画像データを生成する。

【0097】

拡大表示は、検索結果の確認の場合だけとは限らないので、拡大表示する目的に応じて補正する内容は異なる。補正方法と結合方法はこれに限るものではなく、他にもいろいろなものが考えられる。結合した画像データや表示用データは、表示を行った後は破棄して構わない。

【0098】

50

次に、ステップS1005では、CPU201は、ユーザに選択された特定のオブジェクトや論理ページを編集するための表示か否かを判定する。ここで編集のための表示と判定された場合にはステップS1008へ進み、ステップS1001で取得した表示対象を編集のために補正して結合する処理を行う。例えば、画像オブジェクトの場合、分割された各画像領域部分の画像データの色と倍率と位置を正確に補正して結合した画像データを生成する。また、表オブジェクトの場合、分割された各表領域の位置や倍率等を正確に補正して結合し、表の情報、即ち表に関する罫線の位置座標やセルのサイズ、マトリクス構造等を抽出して、表示用データを生成する。

【0099】

また、論理ページ画像の場合、論理ページ内に含まれる各オブジェクトをそれぞれ上述した方法で補正して結合した画像データ・表示用データを合成し、論理ページ画像データを生成する。その際、分割されているテキスト領域は、各テキスト領域部分の画像データについて、領域中の各テキスト行のサイズ（例えば、図3や図4に示すように物理ページを結合するのであれば、高さ）と位置が合うように補正して結合した画像データを生成する。そして、再度文字認識処理を行って文字サイズや文字コード等の文字情報を抽出し、結合した表示用データを生成する。また、背景やページ飾りなどは、再利用性が低いので補正しなくても良い。補正方法と結合方法はこれに限るものではなく、他にもいろいろなものが考えられる。

【0100】

次に、ステップS1003～S1005の何れにも該当しない場合はステップS1009の処理へ進む。そして、CPU201は、分割されている表示対象を表示対象の種類と表示目的及び表示方法に応じて、色、濃度、位置、サイズ等を補正・結合する処理を行う。

【0101】

次に、ステップS1010において、CPU201は、ステップS1006～S1009の何れかで結合された表示対象、又は分割されていない表示対象を各画面に表示する処理を行う。そして、ステップS1011において、全ての表示対象に対して、処理を終了したか否かを判定する。未処理の表示対象がある場合は、ステップS1001に戻り、表示対象がなくなるまで上述の処理を繰り返す。

【0102】

尚、図10に示す処理は、第2の実施形態における表示処理の一例であり、処理の順や処理内容は、この通りでなくても良い。

【0103】

図11は、検索結果や特定のフォルダ内に格納されているオブジェクト及び論理ページを一覧表示した画面例を示す図である。図11は、専用のアプリケーションでウィンドウシステムを利用した場合の画面の例であるが、Webブラウザなどによって同様の機能が提供されるのでも構わない。

【0104】

図11において、1101はタイトルバーと呼ばれるもので、このウィンドウのタイトル表示と、例えば移動や大きさの変更など全体の操作を行う部分である。1102、1103はこのウィンドウに関する機能を提供するボタンで、ヘルプの表示やこのウィンドウを閉じる操作などを指示するためのものである。

【0105】

次に、矩形領域1104、1105は、オブジェクトや論理ページのサムネイル画像を表示する領域である。矩形領域1104に表示されているオブジェクト及び論理ページは分割されており、矩形領域1105は分割されていないことを示している。また矩形領域1104が太枠となっているのは、この領域がユーザによって選択されていることを示しており、1105は選択されていない領域を示している。

【0106】

1106は、この一覧表示画面に表示することができないオブジェクトや論理ページの

10

20

30

40

50

表示を指示する部分である。「前画面」ボタンが選択されたことを検出すると、この画面に表示された一覧の前の一覧を表示し、「次画面」ボタンが押下されると、次の一覧を表示する。

【0107】

1107は選択された領域内に表示されたオブジェクトや論理ページを拡大表示することを指示するためのボタンであり、このボタンが選択されたことを検出すると、拡大表示のための画面へ移行する。

【0108】

1108は選択された領域内に表示されたオブジェクトや論理ページを編集することを指示するためのボタンであり、このボタンが選択されたことを検出すると、編集のための画面へ移行する。そして、1109の「終了」ボタンが選択されたことを検出すると、一覧表示画面を終了する。

【0109】

図12は、第2の実施形態において、あるオブジェクト及び論理ページを拡大表示した画面例を示す図である。図12は、専用のアプリケーションでウィンドウシステムを利用した場合の画面の例であるが、Webブラウザなどによって同様の機能が提供されるので

10

【0110】

図12において、1201はタイトルバーと呼ばれるもので、このウィンドウのタイトル表示と、例えば移動や大きさの変更など全体の操作を行う部分である。1202、1203はこのウィンドウに関する機能を提供するボタンで、ヘルプの表示やこのウィンドウを閉じる操作などを指示するためのものである。

20

【0111】

次に、矩形領域1204は、図11の一覧表示画面等を利用して選択されたオブジェクト及び論理ページを表示する領域であり、ここでは分割されたオブジェクト及び論理ページが表示されている。そして、1205の「終了」ボタンが押下されると、拡大表示画面を終了する。

【0112】

図13は、第2の実施形態において、あるオブジェクト及び論理ページを編集する画面例を示す図である。これは、専用のアプリケーションでウィンドウシステムを利用した場合の画面の例であるが、Webブラウザなどによって同様の機能が提供されるので

30

【0113】

図13において、1301はタイトルバーと呼ばれるもので、このウィンドウのタイトル表示と、例えば移動や大きさの変更など全体の操作を行う部分である。1302、1303はこのウィンドウに関する機能を提供するボタンで、ヘルプの表示やこのウィンドウを閉じる操作などを指示するためのものである。

【0114】

次に、矩形領域1304は、図11の一覧表示画面等を利用して選択されたオブジェクト及び論理ページを表示する領域であり、分割されたオブジェクト及び論理ページができるだけ正確に補正された状態で表示されている。1305は編集を行うためのメニューを表示する部分である。ここでは、例として、「コピー」「切り取り」「貼付」等の項目が表示されているが、編集のメニュー項目としては、これに限るものではなく、他にもいろいろなものと考えられる。

40

【0115】

1306は編集された結果を保存することを指示するためのボタンであり、このボタンが選択されたことを検出すると、編集されたオブジェクト及び論理ページを保存するための画面へ移行する。そして、1307の「終了」ボタンが選択されたことを検出すると、編集画面を終了する。

【0116】

50

第2の実施形態によれば、複数ページに含まれている分割されたオブジェクトデータを、必要に応じて色、濃度、位置座標、倍率の少なくとも何れか一つ又はこれらの組み合わせを補正して結合するか、或いは補正しないで結合する。これにより、文書処理装置のCPUとメモリに負担をかけずに、必要に応じた精度で、1つのオブジェクトデータとして表示することができる。

【0117】

また、分割されたオブジェクトデータを含む領域が存在する複数のページを、必要に応じて色、濃度、位置座標、倍率の少なくとも何れか一つ又はこれらの組み合わせを補正して結合するか、或いは補正しないで結合する。これにより、文書処理装置のCPUとメモリに負担をかけずに必要に応じた精度で、1つのページとして表示することができる。従って、文書中のオブジェクトデータ又はページを有効に再利用することができる。

10

【0118】

[第3の実施形態]

次に、図面を参照しながら本発明に係る第3の実施形態について詳細に説明する。第3の実施形態では、見開きのようなページに含まれている分割されたオブジェクトデータから、文書に関するメタデータを精度良く抽出する場合を説明する。

【0119】

尚、第3の実施形態における文書処理システムの構成は、第1の実施形態の構成と同様であり、その説明は省略する。

【0120】

20

図14は、第3の実施形態における検索用メタデータ抽出時の処理を示すフローチャートである。この処理のプログラムは、ROM202に格納されており、CPU201によって実行される。この処理は、ある一つの検索エンジンのための検索メタデータを抽出する際の処理であり、例えばこのシステムに複数の検索エンジンが実装されている場合には、この処理が複数回実行される。

【0121】

まず、ステップS1401において、CPU201は、検索対象となる写真、図、表などのオブジェクトデータが含まれる領域情報を取得する。検索対象となるオブジェクトデータの種類の種類は、検索の種類によって異なるので、各検索エンジンに応じた適切なオブジェクトデータを取得する。例えば、画像検索の場合は検索対象として自然画像や写真などの画像オブジェクトを取得し、表検索の場合は検索対象として表オブジェクトを取得する。尚、画像や表以外のオブジェクト領域を検索対象として取得しても良い。

30

【0122】

次に、ステップS1402～S1404において、CPU201は、この後抽出されるメタデータを利用する検索エンジンの種類を判定する。尚、ステップS1402～S1404に記述されたもの以外にも様々な検索方法がある。

【0123】

ステップS1402では、CPU201は、言語情報による検索か否かを判定する。ここで、言語情報による検索と判定された場合にはステップS1405へ進む。

【0124】

40

ステップS1403では、CPU201は、色特徴量による画像検索か否かを判定する。ここで、色特徴量による画像検索と判定された場合にはステップS1407へ進む。

【0125】

ステップS1404では、CPU201は、形状特徴量による画像検索か否かを判定する。ここで、形状特徴量による画像検索と判定された場合にはステップS1408へ進む。

【0126】

ステップS1402からステップS1404のいずれにも該当しない場合は、ステップS1409へ進む。

【0127】

50

ステップS1405では、CPU201は、文書中から言語メタデータの抽出対象となる全てのテキスト領域を取得する。ここでは、検索対象となる写真、図、表などの画像オブジェクトに関連付けられているテキスト領域を取得するが、他のテキスト領域をメタデータの抽出対象として取得しても良い。画像オブジェクトとテキスト領域の関連付けは、後述する図15に示すステップS1502で行われる。

【0128】

次に、ステップS1406において、CPU201は、ステップS1405で取得したテキスト領域のうち、分割されているテキスト領域を結合し、結合した領域からテキストを取り出す。その際、各テキスト領域部分の画像データについて、領域中の各テキスト行のサイズと位置が合うように補正して結合した画像データを生成し、再度文字認識処理を行って文字情報を抽出する。領域中の各テキスト行のサイズは、例えば図3や図4に示すように物理ページを結合するのであれば、高さである。テキスト領域の補正方法と結合方法はこれに限るものではなく、他にもいろいろなものが考えられる。また、文字情報を抽出した後は、結合した画像データは破棄して構わない。

10

【0129】

また、ステップS1401で取得した検索対象のオブジェクトについては、分割されていても1つのオブジェクトであることと他の領域との位置関係が分かれば良いので、補正も結合もする必要はない。

【0130】

次に、ステップS1410において、CPU201は、テキスト領域のテキスト情報から検索対象となる写真、図、表などの画像オブジェクトに関連する言語メタデータを抽出する。その際、後述する論理構造解析を利用して画像オブジェクトについて説明している文字列を言語メタデータとして抽出しても良い。例えば、キャプションと思われるテキスト領域の文字情報から画像番号(「図1」と画像名(「システム構成図」)を抽出し、画像名をメタデータとする。また、段落と思われるテキスト領域の文字情報から画像番号を含む文を抽出し、メタデータする。また、例えば「上(の)」のような画像の方向を示す語と画像を示す語を含む文をメタデータとして抽出し、その語が示す画像の方向とステップS511で抽出された論理ページ内での各領域の空間的な関係を照合し、画像と言語メタデータを関連付けても良い。以上は、言語メタデータを抽出する処理方法の一例であり、他にも様々な方法が考えられる。

20

30

【0131】

ステップS1407では、CPU201は、ステップS1401で取得した検索対象オブジェクトが分割されていれば画像領域を結合する。その際、画像オブジェクトの色特徴量が正確に抽出できるように分割された各画像領域部分の色の違いを正確に補正して結合した画像データを生成する。色の補正は、例えば分割された各画像領域の彩度、明度、色調の分布などを利用することで行うことができる。位置や倍率については、検索アルゴリズムが精度をそれほど要求しない場合は、ある程度調整してあれば正確でなくても良いので、図5に示すステップS509で矩形の位置・サイズを調整した際の情報を利用して補正すれば良い。画像領域の補正方法と結合方法はこれに限るものではなく、他にもいろいろなものが考えられる。

40

【0132】

次に、ステップS1411において、CPU201は、検索対象となる画像オブジェクトの画像特徴を解析し、色特徴量を抽出する。色特徴量としては、例えば画像全体や画像を格子状に分割したブロックにおける色分布のヒストグラムや平均色などがある。画像オブジェクトが分割されていた場合に、結合した画像データは、色特徴量を抽出した後は破棄して構わない。また、検索結果表示時に利用できるよう、補正のための情報を保持しておくようにしても良い。

【0133】

ステップS1408では、CPU201は、ステップS1401で取得した検索対象オブジェクトが分割されていれば画像領域を結合する処理を行う。その際、画像オブジェク

50

トの形状特徴量が正確に抽出できるように、分割された各画像領域部分の画像データの倍率と位置を正確に補正して結合した画像データを生成する。位置や倍率の補正は、例えば各画像領域の境界部分から複数の対応点を抽出し、対応点のずれを利用して画像領域間の変換式を算出することで行うことができる。色や濃度については、検索アルゴリズムが精度をそれほど要求しない場合は、補正しなくても良い。画像領域の補正方法と結合方法はこれに限るものではなく、他にもいろいろなものが考えられる。

【0134】

次に、ステップS1412において、CPU201は、検索対象となる画像オブジェクトの画像特徴を解析し、形状特徴量を抽出する。形状特徴量としては、例えば画像全体や画像を格子状に分割したブロックにおける輝度勾配方向の離散化された強度分布などがある。画像オブジェクトが分割されていた場合に、結合した画像データは形状特徴量を抽出した後は破棄して構わない。また、検索結果表示時に利用できるように、補正のための情報を保持しておくようにしても良い。

10

【0135】

ステップS1409では、検索エンジンの種類がステップS1402～S1404の何れにも該当しない場合に、分割されているオブジェクトを検索エンジンに応じて、色、濃度、位置、サイズなどを補正・結合する。例えば、表を検索する検索エンジンで使用するメタデータを抽出する場合は、分割された表を含む領域について、位置座標やサイズなどを補正して結合する。

【0136】

図5を用いて説明したように、文書入力時に各オブジェクトについての分割判別処理を行っており、ステップS1407～S1409では、その結果を利用している。しかし、分割判別処理のタイミングとしては文書入力時に限るものではなく、ステップS1407～S1409の前に行うようにしても良い。

20

【0137】

次に、ステップS1413において、各検索エンジンに応じた方法で検索用メタデータを抽出する。例えば、表検索エンジンの場合は、分割された表に関する罫線の位置座標やセルのサイズ、マトリクス構造をメタデータとして抽出する。オブジェクトが分割されていた場合に結合した画像データは、検索用メタデータを抽出した後は破棄して構わない。また、検索結果表示時に利用できるように、補正のための情報を保持しておくようにしても良い。

30

【0138】

そして、ステップS1414において、各検索エンジンに応じた全ての検索対象オブジェクトに対して、ステップS1401～S1413の処理が終了したか否かを判定する。未処理の検索対象オブジェクトがあると判定した場合はステップS1401に戻り、次の検索対象オブジェクトに対してステップS1401～S1413の処理を行う。

【0139】

図14に示す処理は、第3の実施形態におけるメタデータ抽出処理の一例であり、処理の順や処理内容は、この通りでなくても良い。

【0140】

次に、図15を用いて、第3の実施形態における文書登録時の動作について詳細に説明する。図15は、第3の実施形態における文書登録処理を示すフローチャートである。この処理のプログラムは、ROM202に格納されており、CPU201によって実行される。

40

【0141】

まず、ステップS1501において、CPU201は、画像及び文字情報が混在した1ページ以上で構成される文書画像を入力し、その文書画像を解析し、次の論理構造抽出処理のための前処理を行う。ステップS1501の処理については、図5を用いて説明した通りである。

【0142】

50

次に、ステップS1502において、CPU201は、各領域に関する各種情報、レイアウト抽出結果、及び文字領域に含まれる文字情報の特徴などに基づき、論理構造解析規則に従って解析を行い、文書の論理構造を抽出する。論理構造とは、図7に示すように、ステップS1501で抽出された領域やページに対して論理的な意味属性を抽出して付与したもの、及びそれらの論理的な関係を推定し構造化したものである。論理構造解析規則には、上述の論理ページを処理対象とする規則と物理ページを処理対象とする規則がある。

【0143】

次に、ステップS1503において、検索用メタデータの抽出処理を行う。ステップS1503の処理については、図14を用いて説明した通りである。

10

【0144】

そして、ステップS1504において、ステップS1503で抽出された画像とメタデータを関連付けてDBに格納する。

【0145】

次に、第3の実施形態において、抽出されたメタデータを利用して文書に含まれる写真、図、表などのオブジェクトを検索する時の動作について説明する。

【0146】

第3の実施形態では、写真、図、表などのオブジェクトデータに関連付けられているメタデータを利用して検索を行う。検索は、まずユーザが指定した検索キーワードやキーワードのリストなどの検索条件と各オブジェクトデータに関連付けられたメタデータを対比する。そして、その検索条件と適合するメタデータが付与されているオブジェクトデータをピックアップして検索結果として表示する。

20

【0147】

検索条件と各オブジェクトデータに関連付けられたメタデータを対比する方法は、各検索エンジンによって異なる。また、検索時に、検索条件とピックアップした各オブジェクトデータのメタデータとの類似度を計算して求めても良い。ここで言う類似度とは、ユーザが入力した検索条件が、各オブジェクトデータに付与されたメタデータとの関係を示す表現としてどの程度適切であるかを示すものである。これは、検索方法の例であり、検索方法としてはこれに限るものではなく、どのような方法でも構わない。

【0148】

また、メタデータを利用することにより、文書及び文書中のオブジェクトデータを蓄積する時に、効率的に分類・整理・管理することができるようになる。例えば、メタデータとして付与されている語を分析し、関連するカテゴリでオブジェクトデータを分類することができ、分類するカテゴリはユーザが与えても良いし、クラスタリング等の統計的手法によって自動的に分類するようにしても良い。また、分類時に、カテゴリと各オブジェクトデータのメタデータの類似度を計算して求め、分類に利用しても良い。これは、分類方法、文書管理方法の一例であり、文書管理方法としてはこれに限るものではなく、どのような方法でも構わない。

30

【0149】

第3の実施形態によれば、複数ページに含まれている分割されたオブジェクトデータから、文書に含まれるオブジェクトデータに関するメタデータを抽出する時に、色、濃度、位置座標、倍率の少なくとも何れか一つ又はこれらの組み合わせを補正する。そして、分割されたオブジェクトデータを結合することにより、文書処理装置のCPUとメモリに負担をかけずに、メタデータを精度良く抽出することができる。

40

【0150】

また、複数ページに含まれている分割されたオブジェクトデータから、文書に関するメタデータを抽出する時に、色、濃度、位置座標、倍率の少なくとも何れか一つ又はこれらの組み合わせを補正する。そして、分割されたオブジェクトデータを結合することにより、文書処理装置のCPUとメモリに負担をかけずにメタデータを精度良く抽出することができる。

50

【 0 1 5 1 】

従って、文書中のオブジェクトデータを効率的に再利用できる。また、メタデータを利用することにより、文書及び文書中のオブジェクトデータを蓄積する時に、効率的に分類・整理・管理することができる。

【 0 1 5 2 】

尚、本発明は複数の機器（例えば、ホストコンピュータ、インターフェース機器、リーダー、プリンタなど）から構成されるシステムに適用しても、1つの機器からなる装置（例えば、複写機、ファクシミリ装置など）に適用しても良い。

【 0 1 5 3 】

また、前述した実施形態の機能を実現するソフトウェアのプログラムコードを記録した記録媒体を、システム或いは装置に供給し、そのシステム或いは装置のコンピュータ（CPU若しくはMPU）が記録媒体に格納されたプログラムコードを読み出し実行する。これによっても、本発明の目的が達成されることは言うまでもない。

10

【 0 1 5 4 】

この場合、記録媒体から読み出されたプログラムコード自体が前述した実施形態の機能を実現することになり、そのプログラムコードを記憶した記録媒体は本発明を構成することになる。

【 0 1 5 5 】

このプログラムコードを供給するための記録媒体として、例えばフレキシブルディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、磁気テープ、不揮発性のメモリカード、ROMなどを用いることができる。

20

【 0 1 5 6 】

また、コンピュータが読み出したプログラムコードを実行することにより、前述した実施形態の機能が実現されるだけでなく、次の場合も含まれることは言うまでもない。即ち、プログラムコードの指示に基づき、コンピュータ上で稼働しているOS（オペレーティングシステム）などが実際の処理の一部又は全部を行い、その処理により前述した実施形態の機能が実現される場合である。

【 0 1 5 7 】

更に、記録媒体から読み出されたプログラムコードがコンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書込む。その後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPUなどが実際の処理の一部又は全部を行い、その処理により前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

30

【 図面の簡単な説明 】

【 0 1 5 8 】

【 図 1 】 電子文書に対する処理の流れの一例を示す図である。

【 図 2 】 本発明の一実施形態に係る文書処理システムが構築されるコンピュータ装置の基本構成を示すブロック図である。

【 図 3 】 見開きのページ画像に対して領域抽出処理を行った結果を示す図である。

【 図 4 】 見開きページを物理ページ毎に読み込んだページ画像に対して領域抽出処理を行った結果を示す図である。

40

【 図 5 】 第 1 の実施形態における文書入力時の処理の一例を示すフローチャートである。

【 図 6 】 ある物理ページ画像における各領域の空間的な関係を抽出した結果の一例を示す図である。

【 図 7 】 ある文書におけるページ画像や各ページ画像から抽出された領域に関する各種物理的な情報の一例を示す図である。

【 図 8 】 図 5 に示すステップ S 5 1 0 における判別処理の詳細を示すフローチャートである。

【 図 9 】 第 1 の実施形態における利用時の補正・結合処理を示すフローチャートである。

【 図 1 0 】 第 2 の実施形態におけるオブジェクト表示時の処理を示すフローチャートであ

50

る。

【図11】検索結果や特定のフォルダ内に格納されているオブジェクト及び論理ページを一覧表示した画面例を示す図である。

【図12】第2の実施形態において、あるオブジェクト及び論理ページを拡大表示した画面例を示す図である。

【図13】第2の実施形態において、あるオブジェクト及び論理ページを編集する画面例を示す図である。

【図14】第3の実施形態における検索用メタデータ抽出時の処理を示すフローチャートである。

【図15】第3の実施形態における文書登録処理を示すフローチャートである。

10

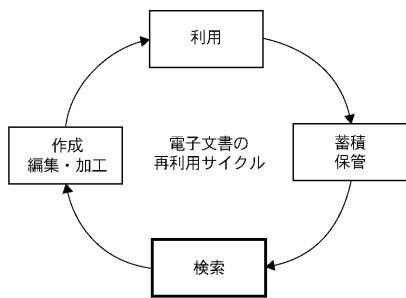
【符号の説明】

【0159】

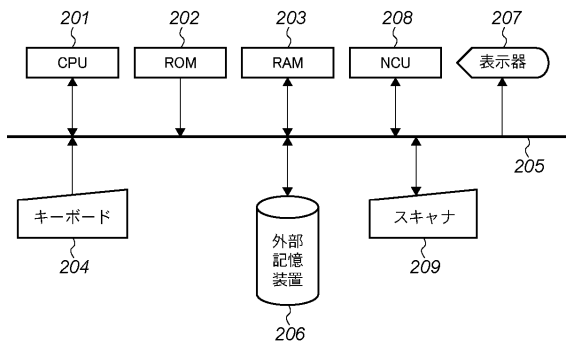
- 201 CPU
- 202 ROM
- 203 RAM
- 204 キーボード
- 205 システムバス
- 206 外部記憶装置
- 207 表示器
- 208 NCU
- 209 スキャナ

20

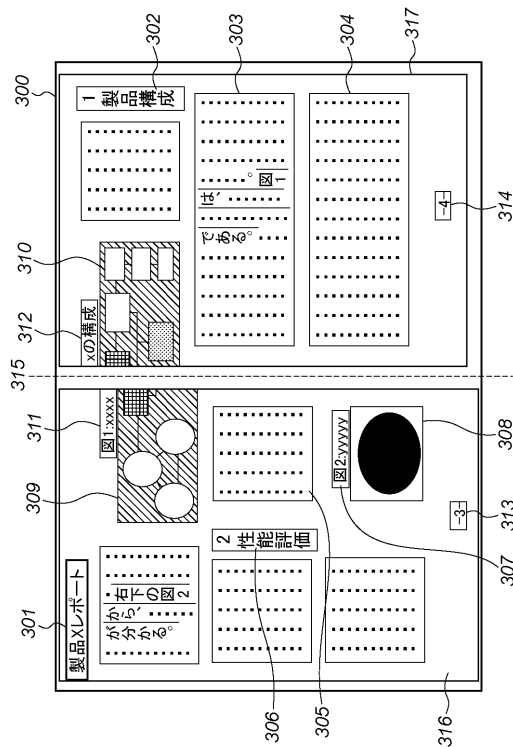
【図1】



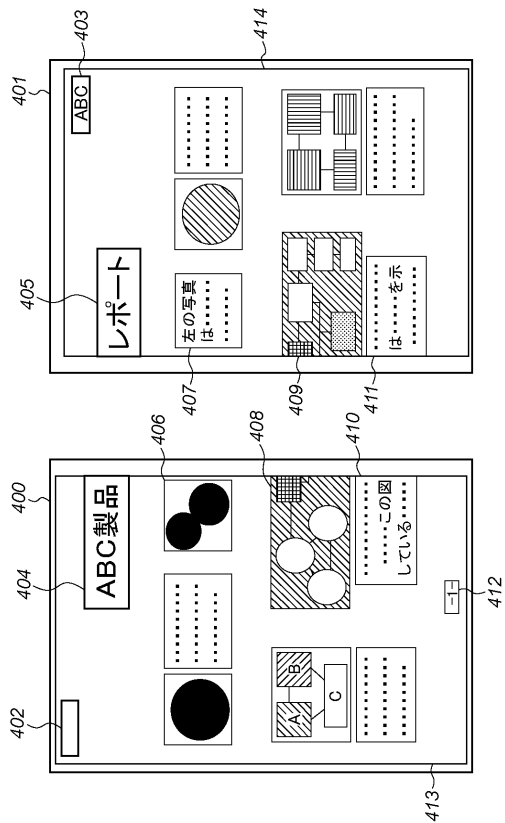
【図2】



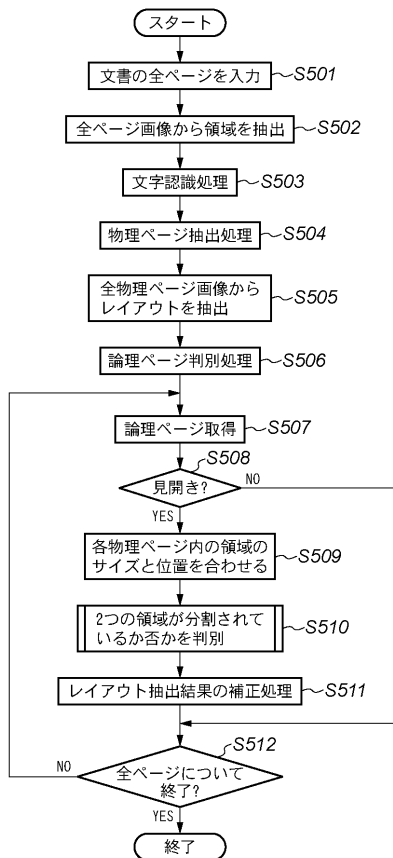
【図3】



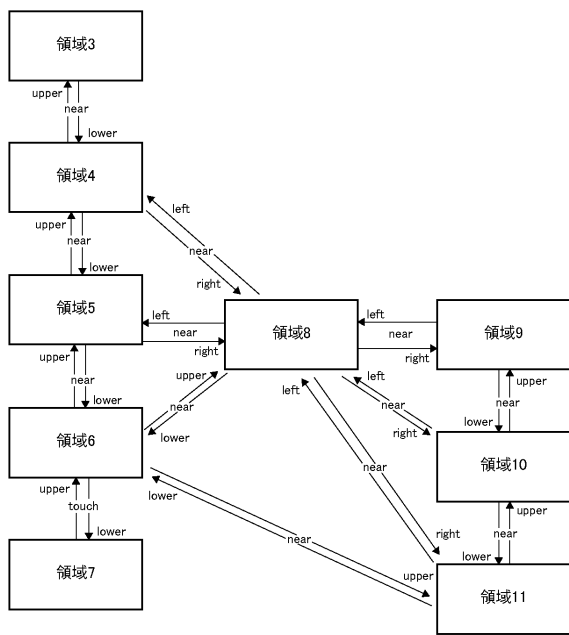
【図4】



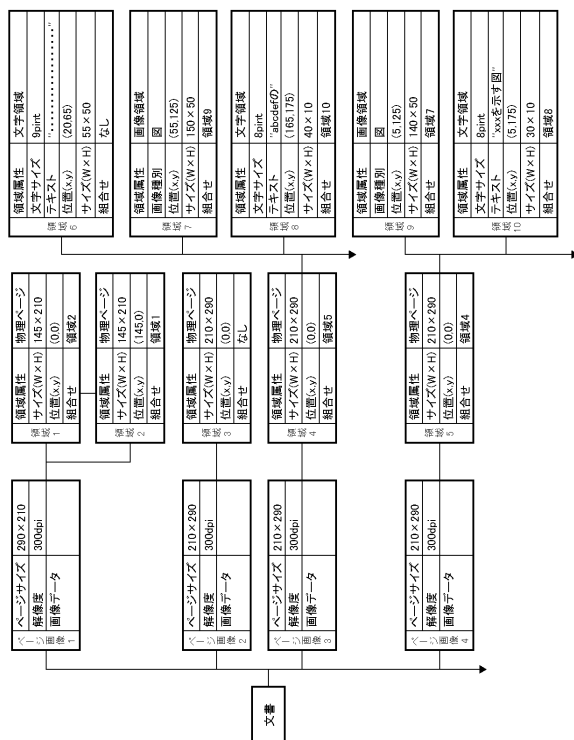
【図5】



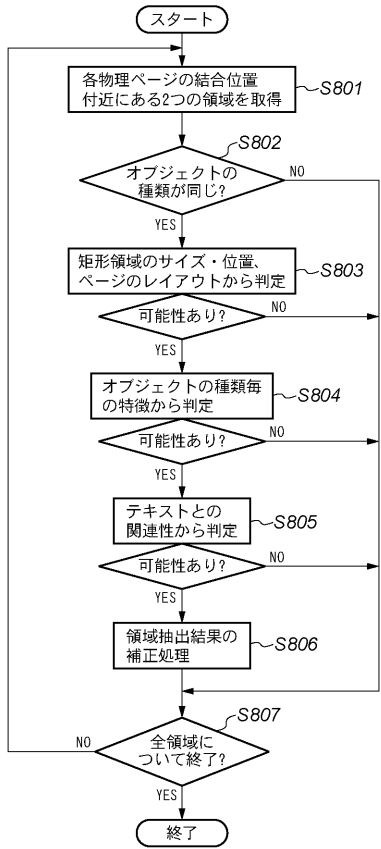
【図6】



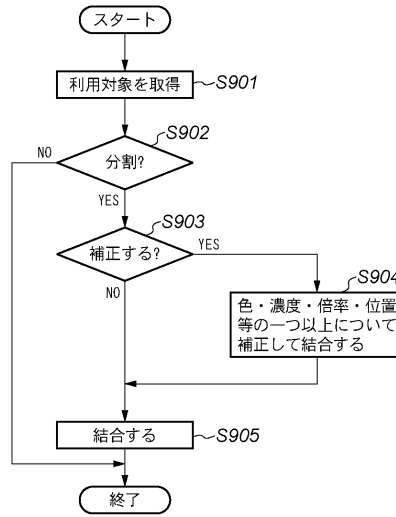
【図7】



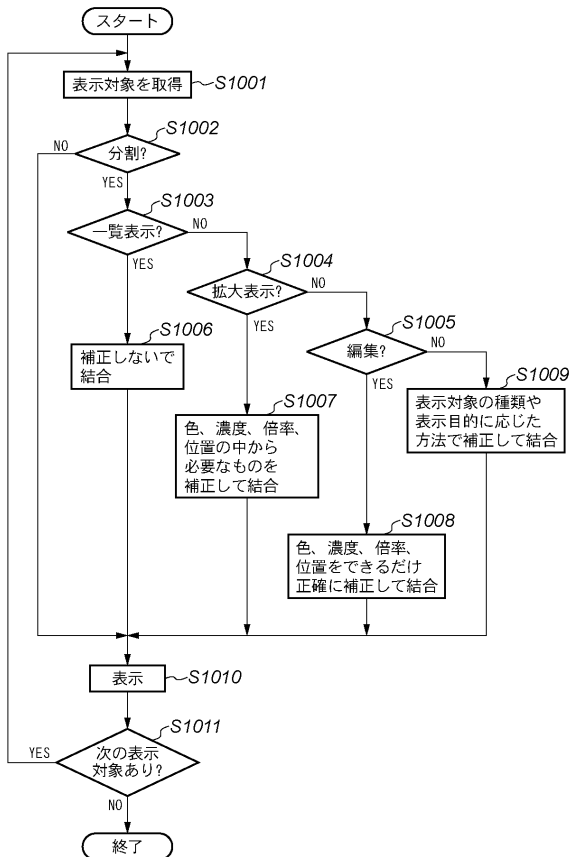
【図 8】



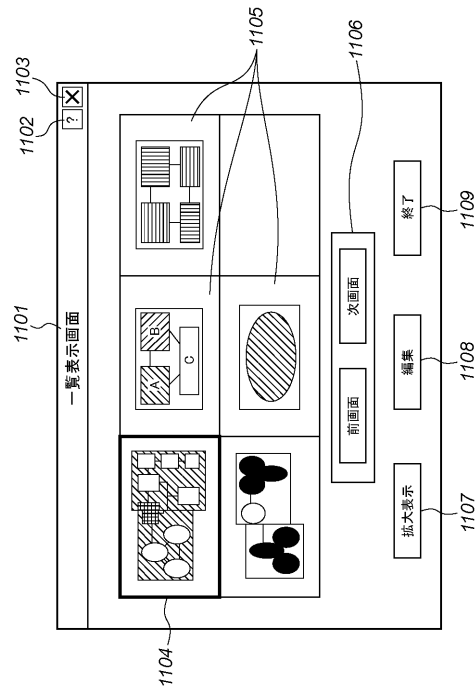
【図 9】



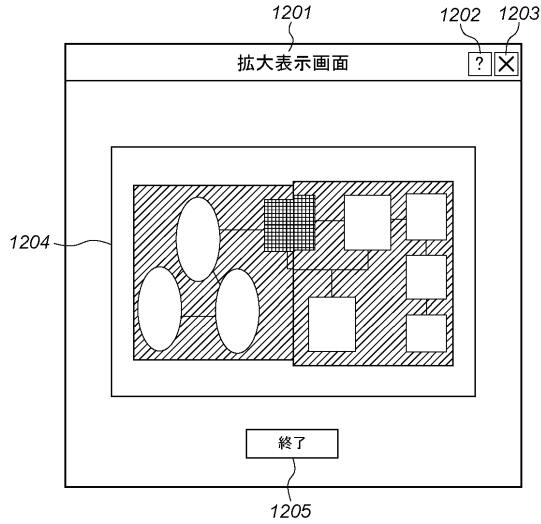
【図 10】



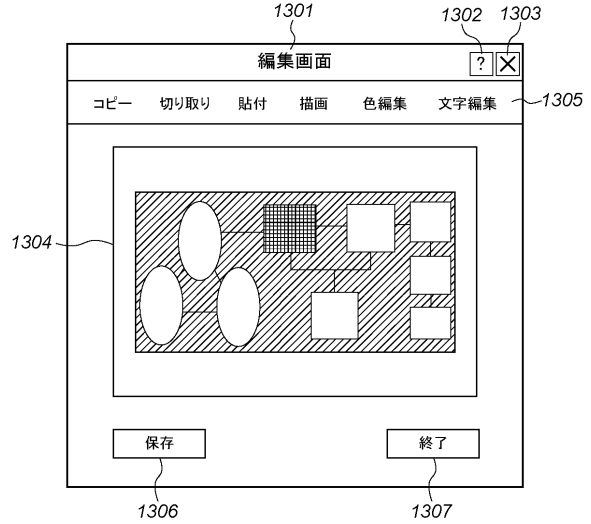
【図 11】



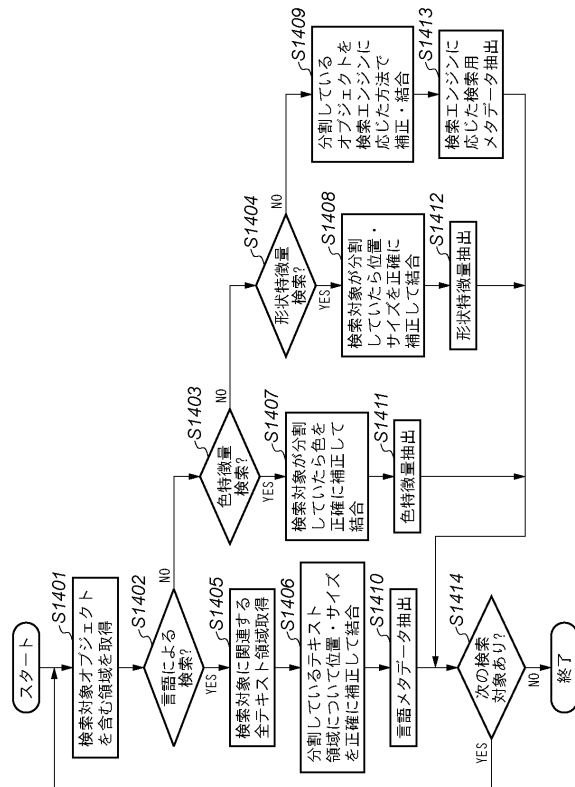
【図12】



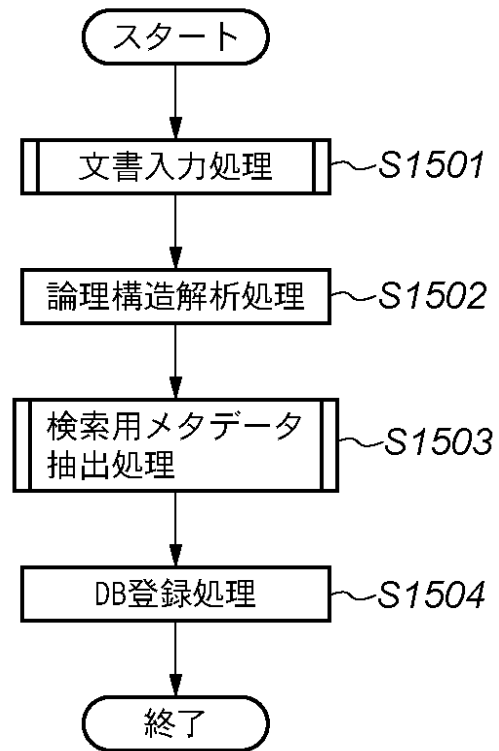
【図13】



【図14】



【図15】



フロントページの続き

審査官 秦野 孝一郎

(56)参考文献 特開2000-293671(JP,A)

(58)調査した分野(Int.Cl., DB名)

H04N	1/387
G06T	3/00
G06T	11/60