

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4427342号
(P4427342)

(45) 発行日 平成22年3月3日(2010.3.3)

(24) 登録日 平成21年12月18日(2009.12.18)

(51) Int.Cl. F 1
G 0 6 T 11/60 (2006.01) G 0 6 T 11/60 1 0 0 A

請求項の数 15 (全 33 頁)

(21) 出願番号	特願2004-18221 (P2004-18221)	(73) 特許権者	000006747 株式会社リコー 東京都大田区中馬込1丁目3番6号
(22) 出願日	平成16年1月27日(2004.1.27)	(74) 代理人	100070150 弁理士 伊東 忠彦
(65) 公開番号	特開2004-234656 (P2004-234656A)	(72) 発明者	キャサリン バークナー アメリカ合衆国, カリフォルニア 940 25, メンロ・パーク, サンド・ヒル・ロ ード 2882番, スイート 115 リ コー イノベーション インク内
(43) 公開日	平成16年8月19日(2004.8.19)	(72) 発明者	クリストフ マーレ アメリカ合衆国, カリフォルニア 940 25, メンロ・パーク, サンド・ヒル・ロ ード 2882番, スイート 115 リ コー イノベーション インク内 最終頁に続く
審査請求日	平成19年1月23日(2007.1.23)		
(31) 優先権主張番号	354811		
(32) 優先日	平成15年1月29日(2003.1.29)		
(33) 優先権主張国	米国 (US)		

(54) 【発明の名称】 文書分析情報を使用して文書を再フォーマット化する方法及び製造物

(57) 【特許請求の範囲】

【請求項1】

処理論理部と記憶部とを有するコンピュータにおいて、

前記処理論理部が、前記記憶部に記憶された文書に関するレイアウト分析情報を取得する段階と、

前記処理論理部が、前記レイアウト分析情報を使用して1つ以上のテキストゾーンに文書をセグメント化する段階と、

前記処理論理部が、前記レイアウト分析情報を使用して1つ以上のテキストゾーンの各々についてスケールン及び重要度情報を生成する段階と、

前記処理論理部が、選択されたテキストゾーンの部分がスケールン情報及び異なるシーケンスのラインへのテキストのリフローに基づいてスケールンを受けた後に、重要度及びスケールン情報に基づいて且つ画像表現に適合する選択されたテキストゾーンの一部に基づいて、目標サイズで文書の画像表現に適合する前記文書内の1つ以上のテキストゾーンの一部を選択する段階と、

前記処理論理部が、目標サイズでの画像表現を生成するために、選択されたテキストゾーンに基づいて文書を再フォーマットする段階とを含む方法。

【請求項2】

前記レイアウト分析情報を取得する段階において、前記処理論理部は前記文書を走査することにより前記レイアウト分析情報を取得することを特徴とする、請求項1に記載の方法。

10

20

【請求項 3】

前記レイアウト情報から、キャラクタセットのサイズを決定する段階と、
 キャラクタセットのサイズに基づいて処理されているテキストゾーンのスケールリングファクタを発生する段階とを有する請求項 1 に記載の方法。

【請求項 4】

文書の、再フォーマットされた電子文書及び再フォーマットされていない電子文書を前記記憶部に記憶する段階をさらに有する、請求項 1 に記載の方法。

【請求項 5】

JPMファイル内のコードストリームが、再フォーマットされた電子文書及び再フォーマットされていない電子文書の両方に使用される、請求項 4 に記載の方法。

10

【請求項 6】

オリジナルの電子文書に適用する命令であるレイアウトボックスを有するオリジナルの電子文書を前記記憶部に記憶する段階をさらに有する、請求項 1 に記載の方法。

【請求項 7】

処理論理部と記憶部とを有するコンピュータにおいて、
 前記処理論理部が、前記記憶部に記憶された電子文書についての多重解像度セグメント化画像を生成する段階と、

前記処理論理部が、前記多重解像度セグメント化画像に連結コンポーネント分析を実行して、画像連結コンポーネントと、多重解像度セグメント化画像内のそれらの位置と、多重解像度ビット分布とのリストを生成する段階と、

20

電子文書のレイアウト分析を実行してテキストゾーンの位置を特定する段階と、
 前記電子文書のテキストゾーンへ属性を割当てする段階と、
 前記テキストゾーンに関連するテキストコンポーネントのリストを生成する段階と、
 テキストコンポーネントと前記多重解像度セグメント化画像の画像連結コンポーネントに関連するコンポーネント画像を併合する段階とを含む方法。

【請求項 8】

前記処理論理部が、レイアウト分析の境界ボックスを使用してテキストゾーンのコンポーネント画像を生成する段階をさらに含む、請求項 7 に記載の方法。

【請求項 9】

前記処理論理部が、連結コンポーネントのリスト内の画像コンポーネントの各々について属性を得る段階をさらに含む、請求項 7 に記載の方法。

30

【請求項 10】

前記コンポーネント画像を併合する段階は、テキストコンポーネントを有する各画像コンポーネントの間のオーバーラップを計算する段階を含む、請求項 7 に記載の方法。

【請求項 11】

前記処理論理部が、オーバーラップの量が第 1 のしきい値よりも大きく且つ第 2 のしきい値よりも小さい場合には、画像コンポーネントから少なくとも 1 つのテキストコンポーネントを減ずる段階を含む、請求項 7 に記載の方法。

【請求項 12】

併合は、1 つ以上の属性に基づいている、請求項 7 に記載の方法。

40

【請求項 13】

前記 1 つ以上の属性は、テキストコンポーネントについてのレイアウトの分析からの重要度を含む、請求項 1 2 に記載の方法。

【請求項 14】

文書に関するレイアウト分析情報を取得する手段と、
 前記レイアウト分析情報を使用して 1 つ以上のテキストゾーンに前記文書をセグメント化する手段と、

前記レイアウト分析情報を使用して 1 つ以上のテキストゾーンの各々についてスケールリング及び重要度情報を生成する手段と、

選択されたテキストゾーンの部分がスケールリング情報及び異なるシーケンスのラインへ

50

のテキストのリフローに基づいてスケーリングを受けた後に、重要度及びスケーリング情報に基づいて且つ画像表現に適合する選択されたテキストゾーンの一部に基づいて、目標サイズで文書の画像表現に適合する前記文書内の1つ以上のテキストゾーンの一部を選択する手段と、

目標サイズでの画像表現を発生するために、選択されたテキストゾーンに基づいて文書を再フォーマットする手段とを有する装置。

【請求項15】

電子文書についての多重解像度セグメント化画像を発生する手段と、

多重解像度セグメント化画像内のそれらの位置と多重解像度ビット分布と共に、画像連結コンポーネントのリストを発生するために、前記多重解像度セグメント化画像に、連結コンポーネント分析を実行する手段と、

テキストゾーンの位置を特定するために、電子文書のレイアウト分析を実行する手段と、

前記電子文書のテキストゾーンへ属性を割当てする手段と、

前記テキストゾーンに関連するテキストコンポーネントのリストを生成する手段と、

テキストコンポーネントと前記多重解像度セグメント化画像の画像連結コンポーネントに関連するコンポーネント画像を併合する手段とを有する装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、画像処理の分野に関連し、特に、本発明は、レイアウト分析、文書分析又は、光学式文字認識(OCR)情報を使用する文書の再フォーマット化に関連する。

【背景技術】

【0002】

走査された文書は、しばしば、大きく、典型的には、2百万から2億画素(又は、サンプル)である。ある応用は、制約されたディスプレイとここでは呼ぶ、非常に少ない画素を有するディスプレイ上に文書を表す、更に小型の画像を表示することから利益がある。制約されたディスプレイは、PDA、移動装置、携帯電話、デジタルコピーフロントパネルのような装置等のような、物理的に制限された数の画素を有するディスプレイである。例えば、多くのPDAは、現在100,000より少ない画素を有する。制約されたディスプレイは、より大きな物理的なディスプレイ(例えば、高解像度モニタ、印刷されたページ等)内の領域である。グラフィックユーザインターフェース(GUI)は、文書と関連する(例えば、アイコン、検索結果、等)領域を有する。1つの形式の制約されたディスプレイは、サムネイル画像を表示する領域である。サムネイル画像(又は、サムネイル)は、典型的には、3,000から30,000画素である。制約されたディスプレイは、ディスプレイ内で有効な幅と高さのみが、表示されている文書又は画像と同じ大きさでない。

【0003】

サムネイルは、大きな画像の小さな画像表現であり、通常は、見ることそして大きな画像のグループを管理することを容易に且つ素早くすることを意図されている。多くのサムネイルは、通常は、オリジナルの画像の丁度ダウンサンプル版である。言いかえると、伝統的なサムネイルは、全体の文書を要求される幅と高さにも再スケーリングし、そして、典型的には、アスペクト比を保存する。サムネイル発生処理の高速化に焦点を当てた、ウェブサムネイル生成についてのソフトウェアパッケージが、入手できる。マージンの自動クロッピングを実行するソフトウェア(例えば、UNIX(登録商標)のpnmツール)もある。

【0004】

HTMLフォーマットで利用できる文書のよりよい表現を提供する、“向上されたサムネイル”がある。例えば、非特許文献1を参照する。これらの向上されたサムネイルは、伝統的に生成されたサムネイルのコントラストを低下させることにより、そして、HTML

10

20

30

40

50

しで見つかったキーワードを重ねることにより生成される。

【0005】

他の研究は、非特許文献2に記載のような、更に効果的なサムネールを生成するためになされる。あるサムネール表現は、非特許文献3のような、サムネールの走査又は他の機械入力からのオリジナルの文書の検索を可能とするために、それに符号化される、特別な、機械認識可能な情報を有する。

【0006】

他の研究は、伝統的なサムネールの新たな使用を生成するためになされる。例えば、サマバーは、固定の幅に再フォーマットされた文書であるが、しかし、制限されない幅を有し、そして、HTML文書についてのウェブプログラムで使用される。キーワードが、サマバー内で異なる色のコードで表示される。一般的には、テキストは、判読できない。非特許文献4を参照する。

【0007】

しばしば、アイコンは、内容に関連される代わりに、ファイルの形式(例えば、それを生成したプログラム)を識別する。これらの場合には、サムネール内のオリジナルの文書のテキストの判読性は、目的ではない。サムネール表現は、しばしば、サムネールをみながらもとの文書を検索できる、判読できるテキスト以外の情報を有する。

【0008】

次の10年は、電子文書が好まれるので紙の文書の使用は劇的に減少すると考えるべきである。紙-電子の変化は、企業について、走査される文書ツールの設計を戦略的にしうる。走査された文書の重要な特徴は、オブジェクト特にテキストがファイル内で識別されず且つ認識されないことである。それは、走査された文書のテキスト文字、単語及び線の位置を特定し且つ識別する、光学式文字認識(OCR)(又は、一般的には文書分析)ソフトウェアによりしばしば、後分析を必要とする。OCRの現在の使用は、一般的には、カリフォルニア、マウンテンビューのAdobeのAdobe Acrobat Captureのように、キーワード検索のためのテキストファイル出力として又は、余分な情報として、認識されたテキストを使用し、そして、走査された文書にメタデータとしてテキストとその位置を追加することである。

【0009】

文書分析システムは、2つの部分:レイアウト分析と文字認識(光学式文字認識又はOCRとも呼ばれる)を有する。文字認識部は、ASCIIのような記号的な形式で出力を発生するために、文字及び文字のグループを解釈するために言語に特定の情報を使用する。レイアウト分析部は、文字認識を実行する前に必要なステップより構成され、即ち、個々の前景画素をストローク(結合されたインクのしみ)のような文字又は文字要素にグループ化し、テキストを含む画像領域を見つけ、そして、パラグラフ、線、単語及びキャラクタのようなテキスト情報ユニットをグループ化する。これらのユニットは、矩形の境界ボックスにより特徴化される。文字認識は、難しいタスクであり、そして、OCRソフトウェアは、文書上に幾つかの間違いをしうる。あるタイトル、見だし、等の、大きなフォントの少量のテキストは、特に認識するのが困難である。これは、ユーザを困らせそして、アプリケーションに誤りを導く。

【0010】

レイアウト情報は、白スペースを拡張するのに(Chilton, J. K., Cullen, J. F. の "デジタル走査装置のための文書画像内の白スペースの拡張(Expansion of White Space in Document Image for Digital Scanning Devices)"を参照する)、白スペースを減少させるのに(名称"ポータブル電子文書に記載されている単語を識別する方法及び装置(Method and Apparatus for Identifying Words Described in a Portable Electronic Document)"の特許文献1参照)又は、制約されたディスプレイに適用する(非特許文献5参照)ことに、既に使用されている。

10

20

30

40

50

【特許文献1】米国特許番号5,832,530

【非特許文献1】Woodruff, A., Faulring, A., Rosenholtz, R., Morrison, J., Pirolli, P.による、Proc. SIGCHI 2001、pp. 198-205、2001の、"ウェブを検索するためのサムネールの使用(Using thumbnail to search the Web)"

【非特許文献2】Ogden, W., Davis, M., Rice, S.によるTRC、1988、pp. 528-534の、"高速に関連性を判断するための文書サムネールの視覚化(Document thumbnail visualizations for rapid relevance judgments: When do they pay off?)"

10

【非特許文献3】Peairs, M.による第3回ICDAR95のプロシーディングのvol. 2、pp. 1174-1179、1995年の、"アイコンペーパー(Icon Paper)"

【非特許文献4】Graham, J.による、Proc. SIGCHI '99、pp. 481-488、1999年の、"リーダのヘルパ:個人化された文書読出し環境(The Reader's Helper: a personalized document reading environment)"

【非特許文献5】Breuel, T.M., Janssen, W.C., Popat, K., Baird, H.S.による、IEEE 2002、pp. 476-479の、"PDAへのペーパー(Paper to PDA)"

20

【発明の開示】

【発明が解決しようとする課題】

【0011】

Adobeは、テキストの検索性を可能とするために、OCR情報を走査された文書の画像に添付する。OCR情報は、しかしながら、サムネールを生成するのに使用できない。OCRがあるテキストに失敗する場合には、そのテストは検索可能でない。

【0012】

しかしながら、走査されていない文書を再フォーマットする方法は、レイアウト分析に基づいて、2次元の制約されたディスプレイを目標とするためになされた。

【課題を解決するための手段】

30

【0013】

電子文書を再フォーマットする方法及び装置が開示される。一実施例では、テキストゾーンの位置を特定するために文書の電子版のレイアウト分析を実行し、文書の電子版のテキストゾーンヘスケールと重要度についての属性を割当て、画像を生成するために属性に基づいて文書の電子版内のテキストを再フォーマットする。

【0014】

本発明は、本発明の種々の実施例の以下の詳細な説明と添付の図面により更に完全に理解されようが、しかしながら、本発明は特定の実施例に限定されると考えるべきではなく、例示と理解のためのみである。

【発明を実施するための最良の形態】

40

【0015】

走査された文書を再フォーマットする方法と装置が開示される。ここの教示は、制約されたディスプレイ上のより良い文書表現を達成するために、走査された文書を再フォーマットする問題と取り組むのに使用される。走査された文書は、画像又は画像を表す文書であってもよい。それは、スキャナ、カメラ又は、他の装置により捕捉され、又は、レンダリングによりデジタル形式で生成されうる。上述のように、制約されたディスプレイの一例は、サムネールである。一実施例では、結果の再フォーマットされた画像は、できるだけオリジナルの文書に含まれる多くの関連するテキストを、読取可能な方法で、表示する。

【0016】

50

特に、ここに開示された技術は、テキストの語義に関する意味を使用せずに、文書内のテキストの再配置を可能とする要素を提供する。これらの要素は、テキストの境界ボックス、テキスト読取順序、テキスト領域の相対的な重要度の評価、スケーリングの可能性及びテキストのリフローの使用を含む。再フォーマットは、(例えば、以下に詳細に説明するように重要度値のような重要な情報を使用して) ブランク空間の除去、スケーリング、ラインとパラグラフの再整形、及び情報を知的に捨てることを使用して実行される。

【 0 0 1 7 】

以下の説明では、多くの詳細が本発明の徹底的な理解を提供するために述べられる。しかしながら、当業者には、本発明のこれらの特定の詳細なしに実行されうことは、理解されよう。他の例では、良く知られた構造と装置は、本発明を曖昧にすることを避けるために、詳細よりも、ブロック図の形式で示される。

10

【 0 0 1 8 】

以下の詳細な説明のある部分は、アルゴリズム及びコンピュータメモリ内のデータビットに関する動作の記号的な表現により示される。これらのアルゴリズム記載と表現は、他の当業者へ研究の実体を最も効果的に伝えるデータ処理技術の当業者により使用される手段である。アルゴリズムはここでは、そして、一般的には、望ましい結果を導くステップの自己一貫性のあるシーケンスであると考えられる。ステップは物理的な量の物理的な操作を必要とする。通常は、必要ではないが、これらの量は、記憶され、伝送され、結合され、比較されそして操作される、電気又は、磁気信号の形式をとる。これらの信号をビット、値、要素、シンボル、キャラクタ、項、数等と呼ぶことは、共通使用の理由により、原理的にしばしば便利である。

20

【 0 0 1 9 】

これらの全ての又は同様な用語は、適切な物理的な量と関連されそして、これらの量に与えられた単に便利なラベルであることは憶えておくべきである。特に述べない限り以下の説明から明らかなように、この記載を通して、" 処理 " 又は、" 計算 " 又は、" 決定 " 又は、" 表示 " 等のような用語を使用する説明は、コンピュータシステムのレジスタ及びメモリ内の物理的(電子的)量として表現されたデータを、コンピュータシステムのメモリ又はレジスタ又は他のそのような情報記憶装置、伝送又は、表示装置内の物理的な量として同様に表現される他のデータへ、操作又は変換する、コンピュータシステム又は同様な電子計算装置の動作又は処理を指すことは理解されよう。

30

【 0 0 2 0 】

本発明は、ここの動作を実行する装置にも関連する。この装置は、要求された目的のために特に構成され、又は、それは、コンピュータ内に格納されたコンピュータプログラムにより選択的に活性化され又は再構成される汎用コンピュータを含みうる。そのようなコンピュータプログラムは、限定はされないが、フレキシブルディスク、光ディスク、CD-ROM、及び、光磁気ディスクのような任意の形式のディスク、読み出し専用メモリ(ROM)、ランダムアクセスメモリ(RAM)、EPROM、EEPROM、磁気又は光カード、又は、電子的命令を格納するのに適する他の形式の媒体のようなそして、各々はコンピュータシステムバスに接続された、コンピュータ読み出し可能な蓄積媒体に格納されう。

40

【 0 0 2 1 】

ここで示されたアルゴリズムと表示は、特定のコンピュータ又は他の装置に固有に関連はしない。種々の汎用システムは、ここの技術に従ってプログラムと共に使用されえ、又は、要求された方法ステップを実行するために更に特化された装置を構成することが便利であるとわかる。種々のこれらのシステムについての要求された構造は、以下の説明から明らかとなる。更に、本発明は、特定のプログラミング言語を参照して記述されていない。種々のプログラミング言語は、ここに記載の本発明の教示を実行するために使用されうことは、理解されよう。

【 0 0 2 2 】

50

機械読み出し可能な媒体は、機械（例えば、コンピュータ）により読み出し可能な形式で情報を格納し又は伝送する機構を含む。例えば、機械読み出し可能な媒体は、読み出し専用メモリ（"ROM"）、ランダムアクセスメモリ（"RAM"）、磁気ディスク記憶媒体、光記憶媒体、フラッシュメモリ装置、電氣的、光学的、音響的又は、他の形式の伝搬信号（例えば、搬送波、赤外信号デジタル信号とう）、等を含む。

【0023】

概要

ここに記載の技術は、走査された文書画像を制約されたディスプレイ文書表現へ、再フォーマットするために文書分析を実行することにより与えられるレイアウト分析情報を使用する。図1は、制約されたディスプレイ文書表現発生器の一実施例のデータフロー図を示す。この発生器は、ハードウェア（例えば、回路、専用論理）、（汎用コンピュータシステム又は専用機で実行されるような）ソフトウェア又は、その組み合わせを含み得る。

10

【0024】

図1を参照すると、分析段階101は、走査された入力画像100を受信し、そして、走査された画像内のテキストゾーンの組を、テキストゾーンの各々についての属性の組と共に発生する。分析段階101は、レイアウト分析情報を発生する、レイアウト分析器110を有する。一実施例では、レイアウト分析器110は文書分析ソフトウェア110Aと光学的にフィックスアップ機構110Bを使用する。一実施例では、レイアウト分析情報は、読取順序の文書内の見つかるテキストゾーンのリスト、読取順序の文書内の見つかるテキストラインのリスト、各テキストラインの単語の境界ボックスのリスト、各テキストゾーンについてのキャラクタサイズを記述する統計値を含む。例えば、あるゾーンで使用される各キャラクタセットについて、このキャラクタセットの平均の寸法（幅と高さ）は、使用される統計値である。一実施例では、レイアウト分析器110は、テキストラインのアラインメント（例えば、左、右、中央、調整）、フォント情報、通常/太字/斜体、等及び単語の信頼性も提供する。

20

【0025】

フィックスアップ機構110Bは、情報を階層で構造に組織化するパーサを含む。これは、以下に更に詳細に記載される。フィックスアップ機構110Bは、対応するもとの走査された画像にへ対応する傾斜除去された（デスクュー、*deskewing*）画像から出力されたレイアウト分析内の座標情報を調整する機能も有する。

30

【0026】

分析段階101は、以下の詳細に記載するように、属性を割当て、属性発生器111も有する。レイアウト分析情報と属性は、合成段階102へ送られる。一実施例では分析結果130も、出力される。

【0027】

レイアウト分析情報が、読み出し順序情報を有しない場合には、上から下及び/又は、左から右又は、右から左へのような位置順序が使用され得る。

【0028】

合成段階102は、記号フォーマッタ112と画像形成器113を有する。一実施例では、フォーマッタ112は、各テキストゾーンのスケールを選択するスケールセクタ112A、テキストゾーン上のリフローを実行するリフロー計算ユニット112B及び制約されたディスプレイ文書表現又は、制約されたディスプレイ出力画像のレイアウトを発生するレイアウトユニット112Cを有する。テキストのリフローが良く知られている。例えば、Gormish他への、米国特許番号6,043,802の、名称「モニタに文書を表示するための解像度減少技術（Resolution Reduction Technique For Displaying Documents on a Monitor）」を参照し、これは、モニタ上に表示するために走査された文書のテキストをのリフローを開示する。

40

【0029】

フォーマッタ112により実行されるこれらのリフロー動作は、例えば、高さとの幅のよ

50

うな、制約された出力ディスプレイ表現のサイズに関して、ディスプレイ制約 120 を受信するのに応じて全て実行される。これらの制約はキャンバスサイズ又は、目標画像差入ずとも呼ばれる。フォーマッタ 112 は、境界ボックス座標のようなテキストゾーンについての記号データに動作する。実際の画像データの処理を必要としない。

【0030】

画像形成器 113 は、合成段階 102 からの出力に応じて、再フォーマットされた出力画像 114 を発生する。クロッピング (cropping)、スケーリング及び貼り付け (pasting) のような画像データの実際の処理が実行される。

【0031】

図 1 を参照して記載の動作とユニットを、以下に詳細に説明する。

10

【0032】

制約されたディスプレイ文書表現は、単一のページに対して又は、全体の文書 (又は、全体の文書のサブセット) に対して生成されうる。これは、幾つかの文書が同一又はほぼ同一のカバーページを有する時に特に有益である。

【0033】

OCR 結果からの情報 (フォントサイズと位置以外) は、テキストの包含を重み付けするのに使用され得る。例えば、検索で使用されるツールと逆文書頻度 (inverse document frequency) のようなテキストの要約は、テキストを含むことの重要度を増加するのに使用される。

【0034】

20

ここに記載の技術は、OCR でなく、レイアウト分析のみを使用して実行されうる。しかしながら、OCR を使用する実施例の利点がある。OCR が応用で要求されそして、レイアウト分析情報が OCR 及び制約された表現発生と共有される場合には、制約された表現発生を発生するのに要求される追加の計算は少ない。他の利点は、制約された表現発生が OCR 結果を使用せず、それゆえに、OCR エラーに対して免れ、そして、OCR が失敗したときに有益な情報を提供できることである。

【0035】

再フォーマット処理の一実施例

図 2 は、文書を再フォーマットする処理を示す。この処理は、ハードウェア (回路、専用論理等)、ソフトウェア (汎用コンピュータシステム又は専用機で実行されるような) 又は、両方を有する処理論理により実行される。

30

【0036】

この処理では、テキストゾーンを使用して、一実施例では、目標は、可能な限り多くのテキストを表示することであり、各テキストゾーンは、最小の読取可能なサイズにスケーリングされている。一実施例では、最小の読取可能なサイズは、スケーリングファクタで示される。利用できるスペースを効率的に使用するために、ゾーン内のテキストは、出力ディスプレイの幅に適合するようにリフローされる。

【0037】

図 2 を参照すると、処理論理は、最初に、レイアウト分析情報を得る為に、レイアウト分析 (及びオプションで OCR も) 実行する (処理ブロック 201)。これは、図 1 の文書分析ソフトウェア 110A により実行される。OCR は、テキストの位置及び境界ボックスを提供するレイアウト分析情報を提供する。テキストをリフローするために、個々の単語の位置が必要である。大きなテキストグループを選択し及び / 又はクロップすることが使用されうる。分析情報は、例えば、罫線、単語の信頼性、フォント記述、キャラクタ境界ボックス等の、他の情報も提供する。レイアウト分析処理の結果は、画像の境界ボックスも提供する。一実施例では、得られたレイアウト分析情報は、ラインへの単語のグループ化、テキストゾーンへのテキストラインのグループ化、テキストの読み出し順序、及びラインのアラインメント形式 (例えば、中央、左又は、右) を含む。

40

【0038】

レイアウト分析情報を得た後に、処理論理は、レイアウト分析情報に必要な調整を実行

50

する（処理ブロック202）。これは、OCR処理中に行われた傾斜除去（デスクュー、`deskewing`）について補償するために、境界ボックスの座標を調整することを含む。これは、レイアウト分析情報を解析することを含む。

【0039】

レイアウト分析情報を得た後に、処理論理は、選択的にゾーンセグメント化を実行する。

【0040】

一旦テキストゾーンが識別されると、処理論理は各テキストサブゾーンについての幾つかの属性を得る（処理ブロック203）。これは、属性発生器111により実行される。これらの属性は、スケーリング及び/又は、重要度情報を含む。一実施例では、スケーリング情報は、スケーリングファクタであり、重要度情報は重要度値又は、等級である。スケーリングファクタと重要度値は各テキストゾーンについて発生される。

10

【0041】

スケーリングファクタ属性は、テキストゾーンがスケーリングされる量を示す変数である。一実施例では、スケーリングファクタは、読取できなくなる前にテキストがスケーリングされうる下限である。一実施例では、（特定の形式のディスプレイについての）テキスト内のキャラクタの平均サイズと低い方のスケーリング限度の間の経験的な関係が、スケーリングファクタを決定するのに使用される。これは次の様である：

$$\text{scaling_limit} = \text{minimal_readable_char_size} / \text{char_size}$$
（スケーリング__限度 = 最小__読み出し可能__キャラクタ__サイズ / キャラクタ__サイズ）。

20

特定のビューアー及びディスプレイに依存して、例えば、画素の最小の読み出し可能なキャラクタサイズは72dpi CRTモニタについては6に等しいがしかし、例えば、LCD、高コントラストディスプレイのような、他の装置又は、他のフォント又は、意図された観測距離等で異なる。この関係を解釈する方法は、`minimal_readable_char_size`画素のキャラクタサイズが、典型的なユーザが快適に特定のディスプレイ上で読むことのできる、最小であると考えられることである。このスケーリング限度ファクタによりスケールリングすることにより、テキストはこの`minimal_readable_char_size`画素寸法に縮小される。

【0042】

30

代わりに、最小の読取可能なキャラクタサイズは、目標ディスプレイ解像度及び読取距離についての調整により決定されうる。合成に使用されるスケーリングファクタは、目標ディスプレイについて調整されるべきである。GUIは、ユーザが、最小のテキストについての望ましいサイズを選択することを可能とする（図12）。各選択は、スケーリングファクタへの異なる調整に対応する。`minimal_readable_char_size`は、ディスプレイ特性、観測条件及び/又は、観測者の嗜好に基づいて選択された望ましいサイズでありそして、読取可能基準へのみに限定されない。

【0043】

テキストサブゾーンについての重要度値は、文書内のテキストサブゾーンの重要度を視覚的に評価するのに使用される、属性である。一実施例では、重要度値は、ゾーン内の最大キャラクタセットサイズとページ内のその位置を使用して決定される。一実施例では、重要度値を発生するのに使用される以下の式が与えられる：

40

【0044】

【数 1】

$$\text{重要度値} = \text{char_size} \left[1 - \frac{\left| X - \frac{W}{2} \right|}{W} \right] \cdot \left[1 - \frac{\max\left(Y, \frac{H}{2}\right)}{3H} \right]$$

ここで、XとYは、それぞれ、テキストゾーンの重心の中心の水平及び垂直座標であり、WとHは、文書の幅と高さであり、そしてW/2とH/2は文書の真中の座標である。(X=0は左、X=Wは右、Y=0は上、Y=Hは下である)。上述の式内の他のファクタは、ページ内の水平位置を考える。特に、ページの右又は左の辺のテキストゾーンは、中心のものとは比べて、不利である。さらに使用される他のファクタは、垂直位置を考える。特に、ページの下部のゾーンは、ページの第1の半分内のものと比較して不利にされる。

【0045】

他の実施例では、重要度重みをページの異なる領域と関連付けするのにテンプレートが使用される。更に他の代替の実施例では、重要度値は、テキスト又は、その一部を圧縮するのに、テキスト符号化器(例えば、JBIG)により費やされた幾つかのビットである。

【0046】

属性の決定後に、処理論理は、所定のテキストゾーンに含まれるテキストのリフローを示すためにリフロー計算を実行する(処理ブロック204)。(リフローは、ここで”記号的再フォーマット”及び”画像生成”の2段階で実行されることに注意する)。リフローの記号的再フォーマットは、画像データを評価することなしに実行される。処理論理は、ユニットの物理的なリフローの後に、テキスト要素の境界ボックスがどのように整列されるかを記述することにより、テキストの再マッピングのためのパラメータを計算する。この点では、実際にリフローされた画像出力データは生成されないことに注意する。実際のリフローの生成は、画像形成段階でのみなされる。再マッピング計算と実際の実行の間の分離は、再マッピング情報を見た後に、処理が、リフローが使用されないことを決定する場合には、計算的な効率を可能とする。

【0047】

リフローの実行後に、処理論理は、キャンバス内に合うゾーンを選択する(処理ブロック205)。これは、図1の表示制約120を使用して実行される。キャンバスは、画素ユニットの形状(例えば、矩形)を含みうる。

【0048】

処理論理は、クロッピングも実行しうる。一実施例では、キャンバスに合うゾーンを選択することは、重要度値を減少させる順序で、テキストゾーン上をループする処理論理を含む。これは、必要なリフローを計算することにより達成され、それにより、スケールングファクタ属性によるスケールング後に、リフローテキストゾーンはキャンバス(目標サイズ)に適合し、スケールングされそしてリフローされたテキストを示す高さを計算する。そして、処理論理は、現在のサブゾーンと前のものを表示するのに十分なスペースがあるかを試験し、ない場合には、ループが抜けられる。ゾーンが合わないのでループが抜けられた場合にはそして、最後のリフローされたゾーンがラインのしきい値数(例えば、10)よりも長い場合には、処理論理はこのゾーンの第1の半分のみを保持し、そして、ループを再開する。しきい値は、ユーザ又はアプリケーションにより設定される。この最後のゾーンが、しきい値数ライン(例えば、10)よりも小さい場合には、処理論理はループを再開することなしに、できる限り多くのラインを保持する。そして、その合計に全てのゾーンを表示する十分なスペースがある場合には、そして、利用できるスペースの

10

20

30

40

50

設定された量（例えば、60%）より小さい量を使用される場合には、テキストのスケールが増加される。この点で、ループは増加されたスケールファクタを使用して再び実行される（処理ブロック206）。一実施例では、スケーリングファクタは、25%、50%、100%又は、任意の割合だけ増加される。

【0049】

処理論理は、表示命令のリストを発生する（例えば、クロップ、スケール、及び/又は、貼り付け命令）（処理ブロック207）。一実施例では、これらの命令は、読取順序である。他の実施例では、命令の出力リストが、クロップ位置（例えば、座標、高さ及び幅）、寸法、スケーリング（例えば、浮動少数点、合理的な数）及び貼り付け位置（例えば、xとy座標）とともに、発生される。

10

【0050】

一旦、スケーリングとリフロー命令を有する選択されたゾーンの組が選択されると、画像形成段階中に、処理論理は、リフローされ且つスケーリングされたテキストゾーンを有する小さな画像オブジェクトを生成する。処理論理は、そして、このオブジェクトをより大きなキャンバスに張りつける。一実施例では、処理論理は最初に、リフローされテキストのためにブランクキャンバスを生成する。処理論理は、そして、一連のクロップと貼り付け動作を実行する。即ち、処理論理は、実際の画像を生成するために、テキストゾーンへの必要なクロッピング、スケーリング及び貼り付け動作を、オリジナルの走査された文書から実行する（処理ブロック209）。即ち、実際の画像は、オリジナルの走査された文書からテキストゾーンをクロッピングし、それらをスケーリングし、そしてそれらをキャンバスの等しい又は、等しくない空間に貼り付けることにより生成される。クロッピングは、全体のパラグラフのようなテキストゾーン内に含まれるもののある部分を取り除く又は、表示されるためにテキストゾーン内に残されるものを識別する動作を含む。

20

【0051】

画像形成動作は、目標画像サイズ（制約）の寸法を使用してブランク画像生成を実行する論理を有し（処理ブロック208）、制約されたディスプレイキャンバスの生成となり、そして、処理の結果をテキストネールキャンバスに貼り付け、それにより、制約された表示画像112を生成する。

【0052】

他の実施例では、画像生成を実行するとき、処理論理は、テキストゾーン内の全ての単語のプールを生成し、そして続いて、それを、ラインが所定の幅を満たしそして、他のテキストラインを開始するまで、テキストラインに加える。

30

【0053】

画像の2部分のスケーリングと、領域クロッピングは、画像がJPEG2000で符号化される場合には、JPEG2000復号器により実行される。画像がJPEG2000圧縮された画像は、低解像度ウェーブレット係数データを単純に復号することによってのみ、圧縮のために使用されたウェーブレット変換の各レベルについて、2のべき乗で縮小されたサイズにスケーリングされる。JPEG2000画像は、あるタイル、プレシント、又は、コードブロックを複合することによってのみクロッピングされる。全ての圧縮されたデータを復号しないことにより、処理時間は減少される。例えば、256×256タイルと5ウェーブレットレベルで圧縮された図17に示された1024×1024画像を考える。400、600と900、800の角を有する矩形をクロッピングしそして、両寸法で1/6にスケーリングすることを考える。256、512及び1024、1024の角を有する矩形（タイル9、10、11、13、14及び15）より構成される6タイルと4の最低解像度レベル（6から、1/4のスケーリングとなる）が復号される。192×128の復号された画像では、クロップされた矩形は、角(400-256)/4=36、(600-512)/4=22及び(900-256)/4=161及び(800-512)/4=72を有しそして、従来技術の方法でクロップされる。(1/6)/(1/4)=2/3のスケーリングは、任意の従来技術の方法で実行されうる。このように、192×128の復号された画像サイズについての処理労力でありそして、大きな1

40

50

0 2 4 x 1 0 2 4 全画像サイズではない。

【 0 0 5 4 】

例示のレイアウト分析とOCRシステム

一実施例では、N.Y.のロチェスタのXerox画像システムからのソラリスOCRソフトウェアのScanWorXバージョン2.2は、レイアウト分析とOCRを実行するのに使用される。この実施例では、ここでは、XDOCと呼ばれるテキストファイルフォーマットで、結果が出力される。出力は、対応するパラメータを有する、一連のインターミックスされたマークアップより構成される。多くの形式のレイアウト分析及び/又は、OCRシステムが、使用されそして従来技術で良く知られていることに注意する。

【 0 0 5 5 】

レイアウト分析により出力される情報から、オブジェクトの境界が走査された画像の画素座標内で識別される。オブジェクトは、例えば、ライン、キャラクタ、単語、テキスト、パラグラフ、罫線（例えば、米国特許の発明社名の上の水平線のような、垂直又は水平ライン）、テキストゾーン、画像ゾーン等である。境界ボックスは、典型的には、矩形領域であるが、しかし、画像領域又は領域の任意の記載でもよい。一実施例では、ソフトウェアはOCRを実行するために実行されそして、それは、自身の座標系（例えば、XDOC座標）で、レイアウト分析情報を表現しそして出力する。この座標系は、走査された画像でない測定の異なる単位を使用し、そして、この画像を正確に記述しないが、代わりに、（Xerox画像システム、ソラリスリリースノートのためのScanWorXMotifバージョン2.2に記載されている）デスクューイング変換後の画像を記述する。画像画素座標系で位置情報を表示するために、2つの動作がXDOC座標に適用される。第1は、逆デスクューイング動作が実行されそして、スケーリング動作が実行される。ページの上方左と下方右の座標が、XDOC座標系で与えられるので、XDOC系の文書の幅と高さは、決定される。

【 0 0 5 6 】

一実施例では、上、左、右及び下のテキストゾーンの境界（以下に詳細に説明する）と、画像ゾーンは、直接OCRソフトウェア出力で（又は、サブゾーンの処理出力内で）表現される。これらは、（上述のように）画像画素座標系に変換される。テキストラインとして、OCRソフトウェアは、基線のY座標と、左及び右X座標のみを提供する。フォント情報（大文字の高さと、ディセンダを有する又はディセンダを有しない小文字の高さ）及び認識されたテキスト（大文字の高さと、ディセンダを有する又はディセンダを有しない小文字の高さ）を使用して、ラインの上方及び下方境界が決定される。キャラクタ認識の失敗への強さのために、変形は認識されたキャラクタを考えずそして、上方の境界については大文字の高さを、そして、下方の境界についてはディセンダを有する小文字の高さのみを使用する。一旦、ライン境界が決定されると、座標は変換されそして、画像とテキストゾーンと同じように、矩形が決定されそして描かれる。

【 0 0 5 7 】

一実施例では、情報をさらに容易に走査するために、情報は、図3に示された例示の構造のような、階層データ構造に再組織化される。図3のボックスは、文書301、フォント記述子302、罫線303、テキストゾーン304、テキストライン305、単語306及び画像ゾーン307を有する。

【 0 0 5 8 】

キャラクタサイズの決定

レイアウト分析情報からのキャラクタサイズ情報は、属性割当て前に、使用するために、統計値へ変換される。アプリケーションに依存して、幾つかの変換が使用される。

【 0 0 5 9 】

一実施例では、テキストゾーンについてのスケーリングファクタが、キャラクタセットのサイズ（又は、フォント）を使用して、決定される。一実施例では、処理されているテキストゾーンについての最大のキャラクタセットサイズを決定するために、処理論理は、キャラクタセット（全ての個々のキャラクタ、フォント）の平均幅と高さの算術平均を決

10

20

30

40

50

定する。これは、以下のように表現される：

```
char__size = max_character_set ( ( < height > + < width > ) / 2 )
```

幾何平均は、同様に使用される。他の実施例では、最小のフォントサイズが使用され又は、平均が使用される。一実施例では、そのゾーン内の最大のキャラクタセットサイズのみが使用される。即ち、スケーリングファクタ属性は、そのゾーン内の最大のフォントを有する各ゾーンについて計算される。代わりに、各キャラクタの境界ボックス又は、各単語又はライン又は点の推定されたフォントサイズの境界ボックスの高さが使用されうる。平均の幅と平均の高さも使用されうる。

【 0 0 6 0 】

再フォーマット化のためのゾーンへのセグメント化

レイアウト分析情報を得た後に、処理論理はゾーンセグメント化も実行し得る。これは、OCR処理により識別されるテキストゾーンは、非常に大きい又は、文書の全てのテキストより構成されうる。一実施例では、処理論理は、共通の特徴（フォント又はアラインメント）と空間の近接により関連される数テキストラインより構成される、オブジェクトを生成することによりゾーンセグメント化を実行する。

【 0 0 6 1 】

テキストラインの適切なグループ化を伴うゾーンを決定する処理の一実施例は、図4に示されており、そして、そのようなオブジェクトは、ここではテキストゾーンと呼ぶ。図4を参照すると、*ispc*はインタースペース（境界ボックスを使用して計算されたライン間の間隔）を示しそして、高さは文書の高さを示す。セグメント化処理は、テキストゾーン400のような、テキストゾーンについて、テキストゾーンを分けるかどうかを決定する処理論理で開始する。一実施例では、処理論理はインタースペース（*ispc*）が5で割られた文書の高さ（*height*）よりも大きい場合には、又は、ゾーンが罫線と識別された場合には、テキストゾーンを分ける（処理論理401）。そして、処理論理は、同じフォント又はアラインメントを有するラインのクラスタを作る（処理論理402）。代替の実施例では、処理論理はフォントのみ又はアラインメントのみに基づいてクラスタを作る。

【 0 0 6 2 】

他の代替のもの、テキストラインについての境界ボックスのみを使用することである。次に、処理論理は、インタースペース（*ispc*）が所定の数（例えば、2）とメディアンインタースペース（*ispc*メディアン）の積よりも大きいときに、テキストゾーンを分ける（処理ステップ403）。

【 0 0 6 3 】

同様に、レイアウト分析ソフトウェアが個々のテキストラインをゾーンにグループ化せずに出力する場合には、同様な特性を有するラインがゾーンにグループ化されうる。

【 0 0 6 4 】

アラインメント形式

一実施例では、アラインメント形式は既にOCR処理により決定されている。しかしながら、それは正確でないか又は、全体のページのアラインメントのみを考え、（ここで記載のリフロー処理について有益な）サブゾーンのテキストラインのアラインメントではないので、処理論理は、各サブゾーンについてアラインメント形式を再評価する。一実施例では、サブゾーンの中央、左エッジ及び右エッジの標準偏差を計算することによりそして、アラインメント形式として最も低い標準偏差を有する軸をとることにより実行される。

【 0 0 6 5 】

リフローの一実施例

一実施例では、リフローされるテキストは、テキストラインと呼ばれるクラスのオブジェクトのリストにより示される。テキストラインオブジェクトのこのリストは、画像データではなく、オリジナルの文書のテキストについての情報（ラインとワードの境界、フォントOCRされたテキスト、属性）を含む。テキストラインオブジェクトのリストは、o

10

20

30

40

50

`old_textline` (古い__テキストライン)と呼ぶ。リフロー計算の出力段階は、ここでは、`reflown_textline` (リフローされた__テキストライン)と呼ばれる新たなテキストラインオブジェクトのリストを出力する。新たなリストは、境界ボックスについての新たな位置としてリフローを記述する。加えて、`reflown_textline`内のテキストオブジェクトは、古いラインとリフローされたラインの間のマッピングも含む。このマッピングは、リフローされたラインの部分(ラインのサブユニット)と古いラインの対応する部分の間の一連の対よりなる。`reflown_textline`内の各テキストオブジェクトについて、リフロー命令は、以下の方法のこれらの対の1つを記述する5アプレットのリストである。

【0066】

リフロー命令/マッピング = (`reflown_start`, `reflown_end`, `old_line`, `old_start`, `old_end`)であり、

1) `reflown_start`, `reflown_end`: ワードのリフローされたテキストラインオブジェクトのリスト内の数値位置により与えられる、部分の第1ワードと最後のワード;

2) `old_line`: 部分がくるテキストラインの(`old_textline`リスト内の)数値位置;

3) `old_start`, `old_end`: ワードのオリジナルのテキストラインオブジェクトのリスト内の数値位置により与えられる、部分の第1ワードと最後のワード。

【0067】

例示の分析データフローが図5に示されている。ここでは`current_textline` (現在の__テキストライン)と呼ぶテキストラインのオブジェクトを扱うループであり、それはリフロー後のテキストラインの記載である。図5を参照すると、処理論理は最初に、`old_textline`の第1ラインをコピーすることにより、`current_textline`を生成し、これは古いテキストラインメモリ500に記憶されそして、その境界ボックスと表示制約を比較する(処理ブロック501)。そして、必要ならば、処理論理は、制約された幅に合うように、`current_textline`を分ける(処理ブロック502)。(分けた後に)最後の部分は新たな`current_textline`505となり、(分ける前の)他の部分は、リフローされたテキストラインメモリ503内の`reflown_textline`リストの先頭として記憶される。分けることが要求されない場合には、`current_textline`リフローあれたテキストラインに記憶される。その後、処理論理は、古いテキストラインメモリ500内に`old_textline`が残っていないかどうか決定される(処理ブロック506)。ない場合には、処理論理は、現在のテキストラインをリフローされたテキストラインメモリ503内へ、前リフローされたテキストラインメモリ503の最後に、に記憶されたリフローされたテキストラインの後に、記憶する。そのようであれば、処理論理は、`current_textline`505と古いテキストラインメモリ500内の次のラインを併合する(処理論理507)。結果のラインは、新たな`current_textline`508であり、これは、ループの先頭に帰還される。古いテキストラインメモリ500内のループに帰還すべきそれ以上のラインがなくなるまで、ループは、継続する。これらの動作中に、リフローされたラインのテキストと古いラインの1つの対応は、以下の記載のように、部分リスト内に記録される。

【0068】

例示の結果

図6は、特定の走査された文書を示す。図7は、図6の文書についてのテキストゾーン境界を示す。図7を参照すると、各ゾーンについて計算されたスケールングファクタ(アンダーラインされた数)と重要度値(斜体の数)が示されている。テキストゾーン矩形境界の座標は、画素座標でありそして描かれている。これらの座標をOCR情報を解析し、スケールングと逆デスクューイング変換後に得られたXDOC座標から得る為に、実行される。

10

20

30

40

50

【 0 0 6 9 】

図 8 は、テキストライン境界を示す。一実施例では、テキストライン境界を得る為に、処理論理は、OCR出力情報からXDOC座標を得てそして、同じ動作を画素座標を得る為に適用する。

【 0 0 7 0 】

図 9 と 1 0 は、文書内の例示のゾーンとリフローが適用された後のゾーンをそれぞれ示す。リフロー処理は、テキストをリフローするだけでなく、この場合は白色スペースも減少することに注意する。

【 0 0 7 1 】

図 1 1 は、テキストゾーンの選択と除去、位置決め及びリフローを使用する例示の制約されたディスプレイ文書表現を示す。

10

【 0 0 7 2 】

文書のブラウジング

この技術に従って発生された、制約されたディスプレイ文書表現は、文書の組をブラウズしそして、ユーザが検索したい文書を選択することを可能とするために使用され得る。制約されたディスプレイ文書表現は、ユーザへあるキーテキストを提供するアイコンとして機能できる。一実施例では、これらの制約されたディスプレイ文書表現は、文書（例えば、走査された又はPDF文書）を取り出すボタンとして動作する。ユーザが望む文書を取り出すために多くの制約されたディスプレイ文書表現又はサムネールがウィンドウ内に表示される、ブラウジングのシナリオでは、制約されたディスプレイ文書表現は多くの方法で使用され得る。例えば、ユーザの文書を取り出すために、ユーザが制約されたディスプレイ文書表現のみを見る、独立の制約されたディスプレイ文書表現がある。一実施例では、ユーザは、望むならそして、そのような選択が可能ならば、サムネールへ切り換えることが可能である。

20

【 0 0 7 3 】

他の実施例では、制約されたディスプレイ文書表現とサムネールの組合せが使用される。1つのそのような例では、ユーザは、カーソル制御装置（例えば、マウス）が文書についての領域を入力するときに、サムネールと制約されたディスプレイ文書表現の両方を、そして次に互いに、又は、ポップアップとしてのみ制約されたディスプレイ文書表現を見る。

30

【 0 0 7 4 】

更に使用では、制約されたディスプレイ文書表現又は、正規のサムネールを表示するために、自動化された選択がブラウザにより提供される。一実施例では、制約されたディスプレイ文書表現発生処理は、文書が正規のサムネールのほうが良いそのようなリッチ画像レイアウトを有するかを決定する。

【 0 0 7 5 】

更に他の使用では、サムネールブラウザで、ユーザがカーソル制御装置を使用してテキストゾーンをわたりカーソルを移動するときに、ゾーンのテキストがサイドウィンドウ内に現れる。一実施例では、NJのWestCaldwellのリコーポレーションのeCabinetは、文書を識別するために、キーワード検索を実行するためにOCRを使用し得る。しかしながら、キーワード検索が低信頼性値を有する結果を発生した場合には、制約されたディスプレイ文書表現は、ユーザが検索している文書をユーザが識別することを助けるのに使用されうる。

40

【 0 0 7 6 】

同様に、多機能周辺機器（MFP）又は全体の文書を示すことができない小ディスプレイを有する他の装置については、ここに記載の技術は、装置を通して記憶され及び/又はアクセス可能な文書に視覚的な指示を提供するのに使用されうる。

【 0 0 7 7 】

他の分析方法との組み合わせ

図 1 に記載の分析出力 1 3 0 は、他の分析と組み合わせられうる。図 1 3 は、制約された

50

ディスプレイ文書表現発生を他の形式の画像発生と統合するシステムの一実施例のフロー図でありそして、以下に詳細に記載される。

【0078】

図13を参照すると、走査された文書1700がウェーブレット分析1701に入力され、これは、走査された文書1700にウェーブレット分析を実行する。ウェーブレット分析の結果は、画像1703を生成するために、合成及び画像生成1702により処理される。ウェーブレット分析1701と合成及び画像生成1702に関する更なる情報は、2002年に1月10日に出願された、名称「マルチスケール変換を使用して圧縮された画像のヘッダベースの処理 (Header-Based Processing of Images Compressed Using Multi-Scale transforms)」の米国特許出願番号10/044,420及び、2002年に1月10日に出願された、名称「画像の小さな表現の内容及び表示装置依存生成 (Content and Display Device Dependent Creation of Smaller Representation of Images)」の米国特許出願番号10/044,603を参照し、両者は、本発明の譲り受け人に譲渡されそして、参照によりここに組み込まれる。

10

【0079】

走査された文書1700は、レイアウト分析1705にも入力され、これは、例えば、上述の、キャラクタ、単語、ライン、ゾーンの境界ボックスのような、境界ボックスを識別する。この情報はレイアウト分析1705からOCR1706へ出力され、これは、OCR情報1707を発生するためにOCRを実行する。OCR情報1707は、全テキスト検索、自動キーワード抽出等に使用され得る。

20

【0080】

レイアウト分析1705から出力される情報(画像分析出力)は、制約されたディスプレイ文書表現分析1700(テキスト分析出力)への出力され、これは、制約されたディスプレイ画像1712を発生するために合成及び画像生成1771と共に上述のように動作する。

【0081】

両ウェーブレット分析1702と制約されたディスプレイ文書表現分析1710の出力は、画像1716が発生される、併合、合成及び画像生成ブロック1715に入力される。どのように画像1716が発生されるかの一実施例を以下に示す。

30

【0082】

J2Kベースの出力と制約されたディスプレイ画像表現出力の併合

走査された文書についてのウェーブレット分析出力を併合するために、(例えば、MAPアルゴリズムに計算された)多重解像度セグメント化データ及び多重解像度エントロピー分布が有効でなければならない。更なる情報は、2002年に1月10日に出願された、名称「マルチスケール変換を使用して圧縮された画像のヘッダベースの処理 (Header-Based Processing of Images Compressed Using Multi-Scale transforms)」の米国特許出願番号10/044,420を参照し、これは本発明の譲り受け人に譲渡されそして、参照によりここに組み込まれる。

40

【0083】

次に結合されたコンポーネント分析が、多重解像度セグメント化の出力に実行される。これは、結合された近傍を発生するためにMatlab(Mathworks社)の関数呼出し"bwlabel"を使用して実行される。結合されたコンポーネント分析は技術的に良く知られている。出力は、それらの位置と共に結合されたコンポーネントとのリストである。

【0084】

コンポーネント当りの属性が得られる。一実施例では、これは、セグメント化マップで決定される画像コンポーネントの解像度、コンポーネントを含む最小の矩形のxとy位置

50

と x と y 寸法及びその重要度値即ち、その解像度でコンポーネントを符号化するのに使用されたビット数を含む。

【 0 0 8 5 】

一旦制約されたディスプレイ分析出力が得られると、テキストゾーンのコンポーネント画像が生成される。

【 0 0 8 6 】

一実施例では、第 1 に、テキストゾーンについてのコードブロック解像度でのコンポーネントマップが生成される。このテキストゾーンの寸法に対応するコードブロック解像度での矩形の寸法 (x , y) は、

【 0 0 8 7 】

【 数 2 】

$$x_r = \lceil x / x_{cb} \rceil$$

$$y_r = \lceil y / y_{cb} \rceil$$

で与えられ、 x_{cb} と y_{cb} は、コードブロックの寸法である。

【 0 0 8 8 】

次のステップで、新たなコンポーネントリストが、画像とテキストコンポーネントを併合することにより得られる。各画像コンポーネントについて、任意のテキストコンポーネントとのオーバーラップがあるかどうかに関するチェックがなされる。一実施例では、オーバーラップは、画像コンポーネントとテキストコンポーネント内の画素の最大数により割られたテキストと画像コンポーネントの間のオーバーラップする画素の数として計算される。画像コンポーネントとテキストコンポーネントがオーバーラップする場合には、オーバーラップの更に詳細な分析が実行される。

【 0 0 8 9 】

オーバーラップがない場合には、画像コンポーネントは併合されたコンポーネントリストに加えられる。

【 0 0 9 0 】

オーバーラップがある場合には、すべてのテキストコンポーネントを有する画像コンポーネントについてのオーバーラップの和がしきい値 T_1 (例えば、0.3) より小さいか又は、しきい値 T_2 (例えば、0.7) より大きいに関するチェックがなされる。この場合には、画像コンポーネントは重要であると考えられる。画像コンポーネントとすべてのそのオーバーラップするテキストコンポーネントの間の結合の合計領域は、併合されたコンポーネントリストへコンポーネントとして加えられる。その解像度属性は、オリジナルの画像コンポーネントの属性である。

【 0 0 9 1 】

すべてのテキストコンポーネントを有するオーバーラップの和がしきい値 T_1 より大きいしかし、 T_2 より小さい場合には、画像コンポーネントは、その中にテキストの重要な部分とそこに非テキストの重要な部分を有すると考えられる。この場合には、画像コンポーネントと重要なオーバーラップ (しきい値 T_3 (例えば、0.25) より大きい) を有するテキストコンポーネントは、画像領域から抽出される。結果の、差画像は、ホールとの 1 つの結合されたコンポーネント又は、幾つかの小さな結合されたコンポーネントである。差画像内で結合されたコンポーネントの数を決定するために、結合されたコンポーネント分析がその上に実行される。結果は、もはや、テキストコンポーネントとの任意の重要なオーバーラップを有しない画像コンポーネントの集合である。集合は、併合されたコンポーネントリストに加えられる。その解像度属性は、オリジナルの画像コンポーネントの属性である。

【 0 0 9 2 】

10

20

30

40

50

最後のステップで、全てのテキストコンポーネントは、制約されたディスプレイ文書表現分析からのオリジナルの属性（解像度と重要度）を含む、併合されたコンポーネントリストに加えらる。

【0093】

一旦、併合されたコンポーネントが生成されると、属性が割当てられる必要がある。これらの属性は、上述と同じものである。併合されたコンポーネントリストは、画像とテキストコンポーネントと属性の混合である。解像度属性は既にコンポーネント画像の併合中に割当てられているが、重要度値はさらに、併合されたコンポーネントリストに割当てられる必要がある。

【0094】

テキストと画像コンポーネントについての重要度値を併合する目標を有する併合されたコンポーネントの重要度についてのメトリックの例は次の様であり：

V_1 = 矩形を含む中のラベル付けされたコンポーネント画素の割合、

V_2 = 矩形を含む中のコンポーネント画像解像度での累積的なエントロピー、

V_3 = テキストコンポーネントについてのレイアウトの分析からの重要度：画像コンポーネントについては $V_3 = 0$ 。

マージされたコンポーネントの重要度は、

【0095】

【数3】

$$\text{importance_of_merged_component} = \alpha \cdot \frac{V_1 \cdot V_2}{N_1} + (1 - \alpha) \cdot \frac{V_3}{N_3}$$

画像コンポーネントについては、 $\alpha = 1$ 、

テキストコンポーネントについては、 $\alpha = 0$ 。

マージされたコンポーネント代替りの重要度値は、次のように得られる：

【0096】

【数4】

$$\alpha \cdot \frac{V_1 \cdot V_2}{N_1} + \beta \cdot \frac{V_3}{N_3}$$

ここで、 $\alpha = 0.7$ 及び $\beta = 0.5$ であり、

N_1 についての選択：* 合計の累積エントロピー * (合計累積エントロピー) * (コンポーネントのサイズ)

N_3 ：(画像領域) * (全てのテキストコンポーネントの重要度の和に対するテキストコンポーネントの相対的な重要度) * 。

【0097】

しきい値の例は以下を含む：

threshold1 = 0.4

threshold2 = 0.7

threshold3 = 0.04

= 5000。

【0098】

値は、例えば、

= 定数 * (テキストゾーンを有する文書の範囲の割合)

【0099】

10

20

30

40

50

【数5】

$$\lambda = \text{const.} \cdot (\text{percentage_of_coverage_of_document_with_text_zones})$$

のように適応的に計算されうる。

【0100】

論理記述

コンポーネントの併合の数学的な記述は以下のように記載される。

T_m = テキストコンポーネント, $m = 1, \dots, M$,

I_n = 画像コンポーネント, $n = 1, \dots, N$,

$A(C)$ = コンポーネント C 内のラベル付けされた画素の数、

$R(C)$ = コンポーネント C を有する最小の矩形。

オリジナルのテキストボックスは、矩形形状を有するので、 $A(T_m) = A(R(T_m))$ であるが、一般的には $A(I_m) \neq A(R(I_m))$ である。

テキストと画像コンポーネントのオーバーラップは、オーバーラップ $(I_n, T_m) = m \cdot i_n$ として定義され、

【0101】

【数6】

$$\left(\frac{A(I_n \cap T_m)}{A(T_m)}, \frac{A(I_n \cap T_m)}{A(I_n)} \right)$$

画像とテキストコンポーネントの間の差画像は、

【0102】

【数7】

$$I_n(x, y) - T_m(x, y) = \begin{cases} I_n(x, y) & \text{if } T_m(x, y) = 0 \\ 0 & \text{if } T_m(x, y) > 0 \end{cases}$$

のように定義される。

例示の擬似コードは次のように与えられる：

【0103】

10

20

30

【数 8】

```

for n=1,...,N
  for m = 1,...,M
    compute overlap(In,Tm)
  end

  if(∑moverlap(In,Tm)> 0)
    if(∑moverlap(In,Tm) < thresh1 || ∑moverlap(In,Tm) > thresh2 )
      tmpimage = In
      number_of_separated_text_components = 0
      for m = 1,...,M
        if(overlap(In,Tm)>thresh3)
          tmpImage = tmpImage - (In-Tm)
          number_of_separated_text_components++
        end
      end
      if(number_of_separated_text_components > 0)
        perform connected component analysis on tmpImage
        add components to merged list
      end
    else
      add [In ∪ {Tm, m=1,...,M | overlap(In,Tm)> 0}] to merged
      list
    end
  end
end
else
  add original component In to merged list
end
end
end

for m = 1,...,M
  add original component Tm to merged list
end
end

```

図 1 4 は、多重解像度セグメント化データと制約されたディスプレイテキスト表現分析からのボックスを併合する例を示す。図 1 4 を参照すると、画像 1 4 0 1 は、多重解像度セグメント化画像を示す。画像 1 4 0 1 では、多重解像度セグメント化が、黒 = 高から白 = 低解像度 (黒 = レベル 1、ダークグレー = レベル 2、中間グレー = レベル 3、ライトグレー = レベル 4、白 = レベル 5) を有する MAP 推定として示されている。結合されたコンポーネント分析の実行と最大のビット数を含むもの選択後に、画像 1 4 0 2 が発生され、多重解像度セグメント化画像の結合されたコンポーネントを表す。異なるカラーは異なるコンポーネントを表す。別に、制約されたディスプレイテキスト表現分析が実行

されそして、コードブロック解像度での制約されたディスプレイテキスト表現分析からのゾーンを示す、コンポーネント画像 1 4 0 3 を生成する。コンポーネント画像 1 4 0 2 と 1 4 0 3 は、画像 1 4 0 4 を発生するために、併合される。

【 0 1 0 4 】

図 1 5 は、多重解像度画像セグメント化データと制約されたディスプレイテキスト表現分析からのボックスを併合する処理の一実施例のフロー図である。ボックスの各々は処理論理の場合、ハードウェア（例えば、回路、専用論理等）、（汎用プロセッサ又は専用機で実行される）ソフトウェア又は、両方の結合を有する。

【 0 1 0 5 】

図 1 5 を参照すると、処理論理は最初に有効に $n = 1$ を設定する（処理ブロック 1 5 0 1）。次に、処理論理は画像コンポーネント n と全てのテキストコンポーネントのオーバーラップを計算する（処理ステップ 1 5 0 2）。処理論理はそして、全てのテキストコンポーネントとのオーバーラップの和がゼロより大きいかどうかを試験する（処理ブロック 1 5 0 3）。そのようでない場合には、処理論理は、処理コンポーネント属性を含む最小の矩形を、併合されたコンポーネントリストに加え（処理ブロック 1 5 2 0）そして、処理論理は、 n が N 、画像コンポーネントの合計数、よりも小さいかどうかを試験する（処理ブロック 1 5 2 1）。そのようである場合には、処理論理は、 n を増加し（処理ブロック 1 5 0 4）そして、処理は処理ブロック 1 5 0 2 へ遷移して戻る。そうでない場合異は、処理は終了する。

【 0 1 0 6 】

テキストコンポーネントのオーバーラップの和がしきい値よりも大きい場合には、処理論理は、処理は処理ブロック 1 5 0 6 へ遷移し、ここで、処理論理は、テキストコンポーネントを有するオーバーラップが第 1 のしきい値（ $threshold_1$ ）よりも小さくそして、第 2 のしきい値（ $threshold_2$ ）よりも大きいかどうかを試験する。そのようでない場合には、処理論理はコンポーネントを、画像コンポーネントと n のオーバーラップを有する全てのテキストコンポーネントの結合に等しく設定し（処理ブロック 1 5 1 9）、処理論理は、処理ブロック 1 5 2 0 と 1 5 2 0 へ戻って遷移する。

【 0 1 0 7 】

全てのテキストコンポーネントを有するオーバーラップが $threshold_2$ より大きく且つ $threshold_1$ よりも小さい場合には、処理は処理ブロック 1 5 0 5 へ遷移し、ここで、処理論理は、変数 m を 1 に等しく設定する。その後、処理論理は、テキストコンポーネント m とのオーバーラップが他の $threshold_3$ より大きいかどうかを試験する（処理ブロック 1 5 0 9）。そのようでない場合には、処理論理は、処理ブロック 1 5 0 7 へ遷移し、ここで、変数 m がテキストコンポーネントの合計数 M よりも小さいかどうかを試験する。そのようでない場合には、処理は、処理ブロック 1 5 2 0 へ遷移する。そのようである場合には、処理は、処理ブロック 1 5 0 8 へ遷移し、ここで、変数 m は、1 だけ増加されそして、処理は、処理ブロック 1 5 0 9 へ遷移して戻る。そのようでない場合には、処理は処理ブロック 1 5 1 9 へ遷移する。

【 0 1 0 8 】

テキストコンポーネント M とのオーバーラップが、しきい値 $threshold_3$ よりも大きい場合には、処理論理は、処理ブロック 1 5 1 3 へ遷移し、ここで処理論理は、テキストコンポーネント m と属性を出力リストに記憶する。次に処理ブロック 1 5 1 0 で、処理論理は、画像コンポーネント n からコンポーネント m からのテキストを減じる。

【 0 1 0 9 】

そして、処理論理は、新たな画像セグメントの結合されたコンポーネント分析を実行し（処理ブロック 1 5 1 1）そして、各新たなコンポーネントに対して、処理論理は、併合されたコンポーネントリストを記憶するためにコンポーネント属性を有する最小の矩形を加える（処理ブロック 1 5 1 2）。

【 0 1 1 0 】

制約されたディスプレイ文書表現をファイルに記憶する

10

20

30

40

50

制約されたディスプレイ文書表現をJPEGファイルに記憶する

多くのファイルフォーマットは、文書ページの画像と別のアイコンの両方を記憶する方法を有する。例えば、JPEG圧縮された画像は、典型的には、JFIFファイルフォーマット又はExifファイルフォーマットのいずれかに記憶される。両ファイルフォーマットは、主画像から独立に符号化されたアイコンの記憶を可能とする。典型的には、これらのサムネールは、オリジナルの画像をサブサンプリングすることにより生成されるが、しかし、このようにそれらを得る要求はない。従って、制約されたディスプレイ文書表現発生処理の出力は、符号化されそして、JFIF又はExifファイルで記憶されうる。デジタルカメラ又はPDAのような装置は、しばしば、表示のためにサムネールを復号しそして、自動的に制約されたディスプレイ文書表現を表示する。ユーザがファイルを開く又は、ファイルの一部にズームインすることを求めるときには、装置は、全画像を復号しそして、ディスプレイ上に一部を表示する。これは、正確には望ましい応答である。

10

【0111】

サムネールをJPMファイルに記憶する

文書記憶システムについては、PDF及びJPM(JPEG2000パート6に定義されている)複数のページを記憶するファイルフォーマットは、JFIF又はEXIFよりもさらに有益である。幾つかのフォーマットは、サムネール画像を記憶する複数の方法を提供する。これらの幾つかは、画像データの再利用の能力のために、典型的なサムネールよりも更に効率的である。幾つかの方法は、追加の文書能力を提供する。

【0112】

20

JPMファイルは、規定された形式と長さを有するファイル内のバイトの単純な範囲である、“ボックス”より構成される。各ボックスの内容は通常は(形式と長さ情報を有する)ボックスの組か又は、符号化された画像又はメタデータ又はレイアウト情報である“オブジェクト”の組のいずれかである。しばしば、復号器は、興味のボックスを素早く見つけそして不要な又は復号器により理解できないボックスをスキップするために長さ形式情報を使用できる。JPMファイルは、ファイルとページを組織化するために設計されたいくつかのボックスを含む。

【0113】

JPMファイルは、単一のページについてのレイアウト情報を記憶するために設計された他のボックスを含む。ページは、単一のJPEG又はJPEG2000ファイルによりJPM内で定義されているが、更に一般的には、位置決めされそして合成されねばならない画像とマスクオブジェクトのシーケンスとして定義される。最後に、JPMファイルは、符号化された画像データを記憶するボックスを含む。ページを構成するために合成された画像がある。JPMは多くの異なるオブジェクトのためにこの符号化されたデータを共有するボックスを提供する。

30

【0114】

最も単純なJPMファイルは、1ページのみとページを埋めるカラーの一樣な矩形を含む。ページについてのJPEG2000圧縮画像を含む例示のファイルは次のようである：

【0115】

【数9】

JPEG 2000 Signature box	
File Type box	
Compound Image Header box	
Page Collection box	
Page Collection Locator box	
Page Table box	
Page Box	10
Page Header box	
Layout Object box	
Layout Object Header box	
Object box	
Object Header box	
JP2 Header box	
Image Header box	
Contiguous Codestream box	

20

上述の例と、ここで与えられる他の例では、くぼみのレベルは、ファイル内のボックスの入れ子を示す。

【0116】

全てのボックスの完全な説明はJPEG 2000パート6で与えられる。(情報技術 - JPEG 2000 画像符号化規格 - パート6 : Component 画像ファイルフォーマット "ISO/IEC FDIS 15444-6")。簡単には、署名ボックスは、ファイルをファイルフォーマットのJPEG 2000ファミリとして識別する。多くのJPEG 2000ファイルフォーマットが互換性があるので、ファイル形式ボックスは、他のフォーマットリーダーがこのファイルから有益なデータを得られるのは何かを示す。コンパウンド画像ヘッダボックスは、(例えば、ファイル内のページ数、ファイルのプロファイル、ファイル内の幾つかの構造ボックスの位置のような)JPM復号器に有益な幾つかのフィールドを含む。幾つかの多ページ集合ボックスもある。これらのボックスは、文書内の全てのページボックスの位置を捜すために許されるポインタを提供する。それらは、文書内のページ間の順序正しいナビゲーションを可能とするキーである。ページコレクションロケータボックスは、本質的に、これがトップレベルのページコレクションボックスでない場合には現在のページコレクションを含む、ページコレクションへ戻るポインタである。このページテーブルボックスは、ページボックスへのポインタを含む。

30

【0117】

ページボックスは、単一ページの情報を含む。ページヘッダボックスは、ページのサイズと向き、オブジェクトの数及び背景色を規定する。ページで合成されるべき各オブジェクト(マスクと画像)のペアについて1つのレイアウトオブジェクトボックスがある。それは、レイアウトオブジェクトヘッダボックスを含み、これは、レイアウトオブジェクトのサイズとオブジェクトをレイアウトする順序を示す識別子番号を提供する。オブジェクトヘッダボックスは、コンティギアスコードストリームボックス(8バイトのオフセットと4バイトの長さフィールドを有する)へのポインタを含む。ポインタは他のファイル内でコードストリームを示すのに使用され得るが、しかし、このファイル内の追加のデータリファレンスボックスが必要である。

40

【0118】

画像データは、JPEG 2000フォーマットのコンティギアスコードストリームボックスに記憶される。

50

【 0 1 1 9 】

全体のファイルに対する 1 つのサムネール

全体の文書についてのアイコンを記憶するために、JP2 ヘッダボックスは、ファイルレベルで加えられる。JP2 ヘッダボックスが追加される際には、ファイル内の第 1 のコンティギュアスコードストリームボックスはサムネールとして使用される。サムネールは、全体のページを示すのに使用されるコードストリームと等価である。代わりに、第 2 コードストリームがアイコンに加えられる。第 2 コードストリーム画追加される場合には、ファイルは、次のように見える（新たなボックスは下線が付されている）。

【 0 1 2 0 】

【 数 1 0 】

JPEG 2000 Signature box

File Type box

JP2 Header box

Image Header box

Compound Image Header box

Contiguous Codestream box (for thumbnail)

Page Collection box

Page Collection Locator box

Page Table box

Page Box

Page Header box

Layout Object box

Layout Object Header box

ID=1

Object box

Object Header box

JP2 Header box

Image Header box

Contiguous Codestream box for Object

10

20

30

主画像サイズに関連するべきサムネールのサイズに対する要求は、なにもない。複数のページを有する文書については、単一ページよりも大きくそして、1 次以上のからの要素を含みうる。

【 0 1 2 1 】

別のレイアウトオブジェクトとしての各ページについてのサムネール

” 文書 ” サムネールのない、しかし、2 ページの各々についてのサムネールを有するファイルを生成することが可能である。これは、ゼロのレイアウト識別子を有するオブジェクトがサムネールとして使用されそしてページに合成されない、JPM 規格内の規定を利用する。これらの 2 つのサムネールに関連する事項は、下線が付されている。

【 0 1 2 2 】

40

【数 1 1】

JPEG 2000 Signature box
File Type box
Compound Image Header box
Page Collection box
 Page Collection Locator box
 Page Table box
Page Box
 Page Header box
 Layout Object box 10
 Layout Object Header box
 ID=0
 Object box
 Object Header box
 OFFSET points to Contiguous Codestream #1
 JP2 Header box
 Image Header box
 Layout Object box
 Layout Object Header box
 ID=1 20
 Object box
 Object Header box
 OFFSET points to Contiguous Codestream #2
 JP2 Header box
 Image Header box
Contiguous Codestream box #1
Page Box
 Page Header box
 Layout Object box
 Layout Object Header box 30
 ID=0
 Object box
 Object Header box
 OFFSET points to Contiguous Codestream #3
 JP2 Header box
 Image Header box
 Layout Object box
 Layout Object Header box
 ID=1
 Object box 40
 Object Header box
 OFFSET points to Contiguous Codestream #4
 JP2 Header box
 Image Header box
Contiguous Codestream box #3
Contiguous Codestream box #2
Contiguous Codestream box #4

識別子 0 内のレイアウトオブジェクトとそれらに関連するコードストリームは、ページに合成されず、代わりに、それらは、全体のページを復号 / レンダリングすることなしにページについての表現として使用される。ページが最大サイズでレンダリングされる時には、0 以外の識別子を有するレイアウトオブジェクトが使用される。サムネールをファイルの先頭の近くに置くために (コードストリーム 1 及び 3)、この例は全ページコードストリーム (コードストリーム 2 及び 4) を最後に移動した。

【 0 1 2 3 】

もちろん、ボックスの配置と、そして、他の目的のために追加のボックスを含めることについての多くの他の可能性がについて存在する。

【 0 1 2 4 】

レイアウトオブジェクトを再使用する別のページとして記憶されたサムネール

(おそらく各テキスト領域又は、各単語が自身のレイアウトオブジェクトを有する) 幾つかのレイアウトオブジェクトより構成されるページについて、幾つかのレイアウトオブジェクトは、選択されそして、サムネールのためにスケールされる。以下のファイルは、3 オブジェクトを有し、別のコードストリームに記憶された、200 dpi で 8 1 / 2 かける 1 1 インチのページを記述する。"サムネール" は、220 かける 170 サンプルの表示サイズを有する別のページとして記憶されている。主ページからの 2 つのオブジェクトは、サムネールページに含まれるが、しかし、他のオブジェクトはスペースの理由のために消去されている。オブジェクトの 1 つが、10 のファクタでスケールされ、そして、これは、オリジナルのページ上でされたように、同じ関連する量のサムネールを埋める。他のオブジェクトは、5 のファクタで縮小され、そして、このように、主ページでなされるようにサムネール上に比較的長く現れる。これは、10 のファクタで減少されると、テキストは読み取れないと予測されるので、なされる。これは、図 16 に示されているが、しかし、図 16 では、サムネールページ 1660 と最大のレンダリングされたページ 1650 は、同じスケールで描かれない。図 16 を参照すると、コードストリーム 1610 は、2 つの異なるページで使用される。例えば、ページ 1 / レイアウト 1 ボックスス 1601 は、コードストリーム 1610 へのポイントと、レンダリングされたページ 1650 上でスケールしそして配置する命令を含む。ページ 2 / レイアウト 1 ボックスス 1604 は、コードストリーム 1610 へのポイントと、レンダリングされたページ 1660 上でスケールしそして配置する命令を含む。同様にボックス 1603 と 1609 は、2 つの異なるページ上で、コードストリームボックス 1630 を使用する。しかしながら、コードストリーム 1620 は、1 つのページ上でのみ使用される。

【 0 1 2 5 】

以下に記載のファイルは、ページ上でそしてサムネール内でオブジェクトの位置を示すために、幾つかのボックスの幾つかのパラメータをリストする。これらのパラメータのリストの定義は、J P E G 2 0 0 0 規格のパート 6 にある。

【 0 1 2 6 】

10

20

30

【数 1 2】

JPEG 2000 Signature box	
File Type box	
Compound Image Header box	
Page Collection box	
Page Collection Locator box	
Page Table box	
Page Box (Main Page)	
Page Header box	
PHeight=2200, PWidth=1700	
Layout Object box	10
Layout Object Header box	
ID=1, LHeight=200, LWidth=1000, LVoff=200, LHoff=350, Style=2	
Object box	
Object Header box	
OFFSET points to Contiguous Codestream #1	
JP2 Header box	
Image Header box	
Layout Object box	
Layout Object Header box	
ID=2, LHeight=1400, LWidth=1200, LVoff=500, LHoff=250, Style=2	
Object box	
Object Header box	
OFFSET points to Contiguous Codestream #2	20
JP2 Header box	
Image Header box	
Layout Object box	
Layout Object Header box	
ID=3, LHeight=100, LWidth=500, LVoff=2000, LHoff=1000, Style=2	
Object box	
Object Header box	
OFFSET points to Contiguous Codestream #3	
JP2 Header box	
Image Header box	
Page Box (Thumbnail Page)	
Page Header box	
PHeight=220 PWidth=170	
Layout Object box	30
Layout Object Header box	
ID=1, LHeight=20, LWidth=100, LVoff=10, LHoff=35, Style=2	
Object box	
Object Header box	
OFFSET points to Contiguous Codestream #1	
Object Scale Box	
VRN=1,VRD=10,HRN=1,HRD=10	
JP2 Header box	
Image Header box	
Layout Object box	
Layout Object Header box	
ID=2, LHeight=20, LWidth=100, LVoff=180, LHoff=40, Style=2	
Object box	40
Object Header box	
OFFSET points to Contiguous Codestream #3	
Object Scale Box	
VRN=1,VRD=5,HRN=1,HRD=5	
JP2 Header box	
Image Header box	
Contiguous Codestream box #1	
Contiguous Codestream box #2	
Contiguous Codestream box #3	

別のページが、主ページからのレイアウトオブジェクトを利用するサムネールについて加えられるので、別のページは、文書の代わりの視野 (view) のために追加される。新たなページは、全ての同じレイアウトオブジェクトを有するが、しかし異なってスケールリングされそして、異なるサイズ及びページ上に異なって配置される。

【0127】

このように、1つのJPMファイルは、図9と10を記憶できる。2つのページボックスと、ページ上の各アイテムについて各ページボックス内のレイアウトボックスがあるが、しかしデータ自身は複製されない。

【0128】

レイアウトボックスの全てからのオーバーヘッドは重要でありそして、圧縮を減少する。しかしながら、幾つかのシステムは、両レイアウトを有しそして、望ましいレイアウトを提供するためにそれを分析する、サーバ上に1つのファイルを記憶するように選択できる。代わりに、異なるレイアウトは、要求が特定の見る幅を有するページについてなされるときに、発生される。

10

【0129】

代替の実施例

一実施例では、テキストゾーンの重要度ランキングの決定は、個々の提供するキーワードにより増加されうる。即ち、重要度ランキングは、キーワードに基づきうる。他の実施例では、個々は、重要度ランキングを実行する方法として、境界ボックスを選択するのを補助し得る。

20

【0130】

他の実施例では、テキストがOCRを受けそして、その結果にOCRエラーを発生する場合には、オリジナルの画像からのビットマップは、OCR結果の代わりに使用される。これは、結合された結果が、エラーを含まないことを保証する。

【0131】

一実施例では、レイアウト分析ソフトウェアが、文書内の単語のサイズを見つけそしてそれらの単語が発生する場所を見つけるために文書を走査するために、JBIG2 (情報技術 - 損失のある / 無損失の2値画像の符号化、ISO / IEC 14492 : 2001、2001年12月15日) のような、辞書ベースの方法で置き換えられる。

【0132】

一実施例では、リフロー処理は、全てのテキストが一様なサイズになる。これは、異なるスケールリングファクタを必要としそして、各テキストゾーンは同じサイズのテキストを含むことに注意する。他の実施例では、リフロー処理の結果は、全てのテキストがテキストの残りに対してその相対的なサイズを維持することである。言いかえると、リフロー前のテキストサイズの比は、リフロー後と同じである。他の単調なマッピングも使用されうる。

30

【0133】

既にOCRを実行するシステムについては、一旦OCRが実行されると、レイアウト情報は単純に捨てられる。OCR情報の発生は、計算時間に関してコストがかかる。しかしながら、ここに記載の技術に従って、分析レイアウト情報を使用することは、文書分析で既に実行されている仕事の量と比較して、少量の余分な仕事のみを必要とする。

40

【0134】

本発明の多くの変更と修正が、前述の記載を読んだ後に当業者に明らかとなるが、説明により示されそして開示された任意の特定の実施例は、限定するものではないことは理解されよう。従って、種々の実施例への参照は、本発明に必須であると考えられる特徴のみを列挙する請求項の範囲を限定するものではない。

【図面の簡単な説明】

【0135】

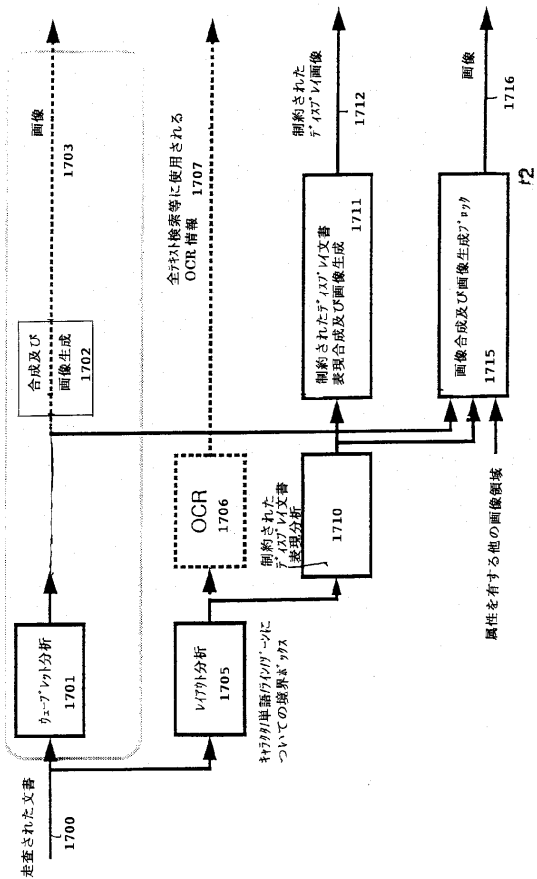
【図1】再フォーマットされた文書発生器の一実施例のデータフロー図である。

【図2】再フォーマットされた文書を発生する処理の一実施例のフロー図である。

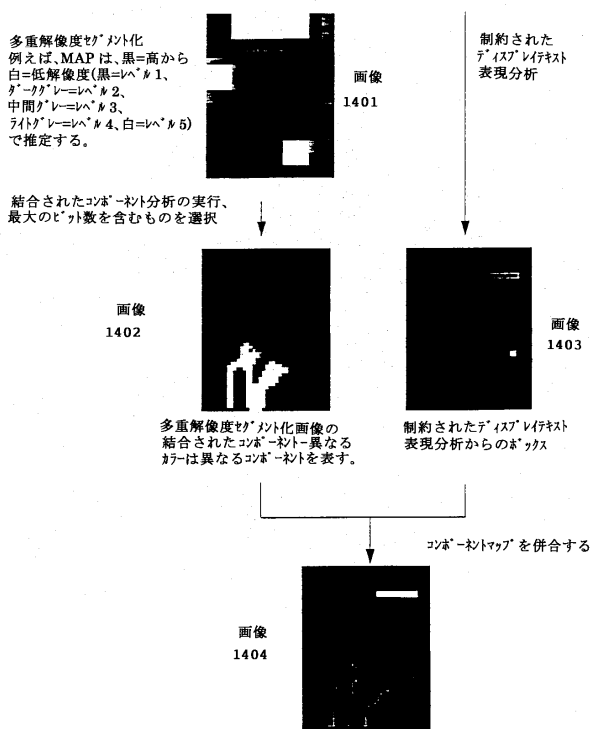
50

- 【図3】例示のデータ構造を示す。
- 【図4】ゾーンセグメント化処理の一実施例のフロー図である。
- 【図5】リフロー処理の一実施例のフロー図である。
- 【図6】オリジナルの走査された文書を示す。
- 【図7】図6の走査された文書のテキストゾーン境界を示す。
- 【図8】図6の走査された文書のテキストライン境界を示す。
- 【図9】オリジナルの文書内のゾーンの例を示す。
- 【図10】リフロー後の図9のゾーンを示す。
- 【図11】制約されたディスプレイ文書表現の例を示す。
- 【図12】最小の読取可能なテキストについてのユーザ選択 " m i n _ s i z e " を有する GUI の例を示す。 10
- 【図13】制約されたディスプレイ文書表現と多重解像度セグメント化画像を発生できるシステムのブロック図である。
- 【図14】多重解像度画像セグメント化データと制約されたディスプレイテキスト表現分析からのボックスを併合する例を示す。
- 【図15】多重解像度画像セグメント化データと制約されたディスプレイテキスト表現分析からのボックスを併合する処理の一実施例のフロー図である。
- 【図16】同じコードストリームからの異なるページ画像の復号を示す。
- 【図17】 J P E G 2 0 0 0 ベースのクロッピングとスケーリングを示す例示の画像である。 20
- 【符号の説明】
- 【 0 1 3 6 】
- 1 0 0 走査された入力画像
- 1 0 1 分析段階
- 1 0 2 合成段階
- 1 1 0 レイアウト分析器
- 1 1 0 A 文書分析ソフトウェア
- 1 1 0 B フィックスアップ機構
- 1 1 1 属性発生器
- 1 1 2 記号フォーマッタ 30
- 1 1 2 A スケールセクタ
- 1 1 2 B リフロー計算ユニット
- 1 1 3 画像形成器
- 1 1 4 再フォーマットされた出力画像
- 5 0 0 古いテキストラインメモリ
- 5 0 3 リフローされたテキストラインメモリ
- 1 7 0 0 走査された文書
- 1 7 0 1 ウェーブレット分析
- 1 7 0 2 合成及び画像生成
- 1 7 0 5 レイアウト分析 40
- 1 7 0 6 O C R
- 1 7 0 7 O C R 情報
- 1 7 1 0 制約されたディスプレイ文書表現分析
- 1 7 1 2 制約されたディスプレイ画像
- 1 7 1 5 併合、合成及び画像生成ブロック
- 1 7 1 6 画像
- 1 7 7 1 合成及び画像生成

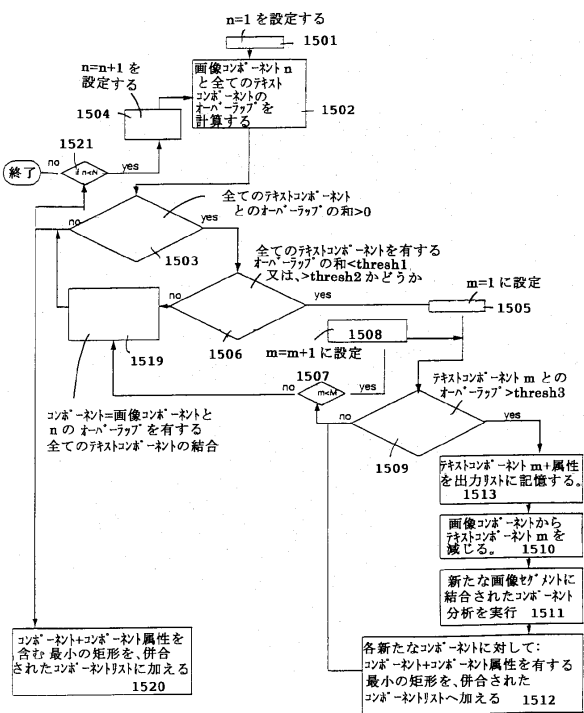
【図13】



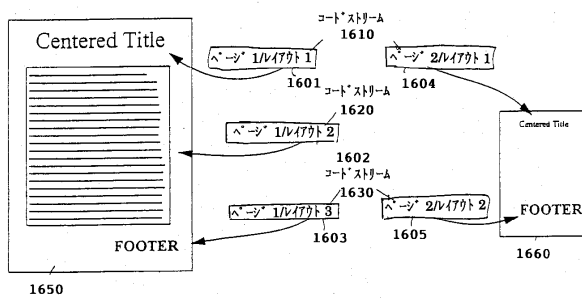
【図14】



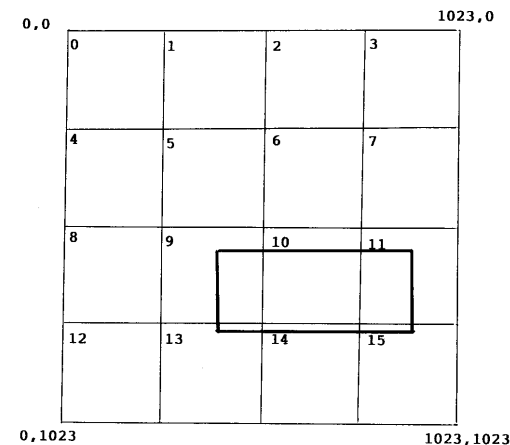
【図15】



【図16】



【図17】



フロントページの続き

- (72)発明者 エドワード エル シュワルツ
アメリカ合衆国, カリフォルニア 94025, メンロ・パーク, サンド・ヒル・ロード 288
2番, スイート 115 リコー イノベーション インク内
- (72)発明者 マイケル ゴーミッシュ
アメリカ合衆国, カリフォルニア 94025, メンロ・パーク, サンド・ヒル・ロード 288
2番, スイート 115 リコー イノベーション インク内

審査官 千葉 久博

- (56)参考文献 特表2002-538638(JP, A)
特開2002-351861(JP, A)
特開2002-185776(JP, A)
特開2001-101164(JP, A)
特開2000-306103(JP, A)
特開2000-224579(JP, A)
特開平11-242654(JP, A)
特開平10-178541(JP, A)
特開平10-162003(JP, A)
特開平10-116065(JP, A)
特開平10-105694(JP, A)
福原隆浩, " 静止画に加えて動画にも対応する国際標準規格 JPEG2000/Motion
- JPEG2000の技術概要と応用 後編", Interface, 日本, CQ出版株式会社
, 2002年12月 1日, 第28巻, 第12号, p.137-147

(58)調査した分野(Int.Cl., DB名)

G06T 11/60
G06F 17/20 - 17/26
G06F 17/30
G09G 5/00 - 5/40