

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号
特許第6082306号
(P6082306)

(45) 発行日 平成29年2月15日(2017.2.15)

(24) 登録日 平成29年1月27日(2017.1.27)

(51) Int.Cl.

G06K 9/20 (2006.01)

F I

G06K 9/20 340K

G06K 9/20 340L

請求項の数 25 (全 20 頁)

(21) 出願番号	特願2013-84694 (P2013-84694)	(73) 特許権者	511072895
(22) 出願日	平成25年4月15日 (2013.4.15)		キング・アブドゥルアジズ・シティ・フォー・サイエンス・アンド・テクノロジー (ケイ・エイ・シー・エス・ティ)
(65) 公開番号	特開2013-235574 (P2013-235574A)		KING ABDULAZIZ CITY FOR SCIENCE AND TECHNOLOGY (KACST)
(43) 公開日	平成25年11月21日 (2013.11.21)		サウジアラビア、11442 リヤド、ビィ・オウ・ボックス・6086、ザ・ナショナル・センター・フォー・テクノロジー・ディベロップメント
審査請求日	平成28年2月16日 (2016.2.16)		
(31) 優先権主張番号	13/467, 873	(74) 代理人	110001195
(32) 優先日	平成24年5月9日 (2012.5.9)		特許業務法人深見特許事務所
(33) 優先権主張国	米国 (US)		
早期審査対象出願			
			最終頁に続く

(54) 【発明の名称】 光学式文字認識用に画像を前処理するための方法およびシステム

(57) 【特許請求の範囲】

【請求項 1】

光学式文字認識 (OCR) 用に画像を前処理する方法であって、画像は複数の欄を含み、複数の欄のうちの各欄は文字を含み、前記方法は、

文字に関連付けられる複数の構成要素を定めることを備え、構成要素は一連の接続されたピクセルを含み、前記方法は、さらに、

前記複数の構成要素に関連付けられる行高さおよび欄間を計算することと、
行高さおよび欄間の少なくとも一方に基づき、前記複数の構成要素のうちの少なくとも1つの構成要素を、前記複数の欄のうちの ある欄に関連付けることと、

前記複数の欄のうちの各欄について第1の組の特性パラメータを計算することと、
前記第1の組の特性パラメータに基づき、前記複数の欄のうちの各欄の中の複数の構成要素を合成して、少なくとも1つの副単語および少なくとも1つの単語の少なくとも一方を形成することとを備え、前記複数の欄は、第1の領域と第2の領域とを含み、

前記方法は、
最も近い文字構成要素が第1の領域に入る構成要素を当該第1の領域に合成することと

、
最も近い文字構成要素が第2の領域に入る構成要素を当該第2の領域に合成することとを備え、前記第2の領域の少なくとも部分は前記第1の領域の少なくとも部分と縦方向に重なり、

各欄に関連付けられる語間を計算することをさらに備え、前記語間を計算することは、

各欄に関連付けられた複数の構成要素の連続する構成要素間の間隔のヒストグラムを作成することと、ヒストグラムから頻出間隔を特定することとを含み、前記頻出間隔は行高さによって定められるしきい値範囲内にあり、さらに、前記頻出間隔に基づき語間を計算することを含み、

各欄に関連付けられる行間を計算することをさらに備え、前記行間を計算することは、各欄に関連付けられる複数の構成要素の複数の水平射影のヒストグラムを作成することを含み、前記複数の水平射影のうちのある水平射影は、ラスタスキャンの各掃引に対応して複数の構成要素に関連付けられるピクセルの数を示し、さらに、2つの連続する最大水平射影間の平均距離を計算することと、前記平均距離に基づき行間を計算することを含む、方法。

10

【請求項2】

前記画像は、濃淡画像およびカラー画像の少なくとも一方を2進画像に変換することによって得られる、請求項1に記載の方法。

【請求項3】

前記画像は、ごま塩雑音をふるい落とすことによって得られる、請求項1 または請求項2に記載の方法。

【請求項4】

前記画像は、変形ハフ変換を用いて歪みを修正することによって得られ、前記変形ハフ変換はアラビア文字に適合される、請求項1 ～請求項3のいずれか1項に記載の方法。

【請求項5】

前記複数の構成要素を定めることは、
前記画像に対してラスタスキャンを行なうこと、
前記ラスタスキャンの少なくとも1回の掃引に対応する複数の構成要素のうちの少なくとも1つに関連付けられる複数のピクセルを特定すること、および
複数のピクセル間の相互接続に基づき前記複数のピクセルを統合して、少なくとも1組の接続されたピクセルを形成することを含む、請求項1 ～請求項4のいずれか1項に記載の方法。

20

【請求項6】

ピクセルは当該ピクセルの8個の隣接するピクセルの少なくとも1つと相互接続される、請求項5に記載の方法。

30

【請求項7】

前記行高さを計算することは、
前記複数の構成要素の各々の高さに対応する高さのヒストグラムを作成すること、
高さのヒストグラムから頻出高さを特定すること、および
頻出高さに基づき行高さを計算することを含む、請求項1 ～請求項6のいずれか1項に記載の方法。

【請求項8】

前記欄間は、行高さに基づき計算される、請求項7に記載の方法。

【請求項9】

各欄に関連付けられる前記語間を計算することは、
各欄に関連付けられた前記複数の構成要素の連続する構成要素間の間隔のヒストグラムを作成することと、
前記ヒストグラムから頻出間隔を特定することとを含み、前記頻出間隔は前記行高さによって定められるしきい値範囲内にあり、さらに
前記頻出間隔に基づき前記語間を計算することを含む、請求項1 ～請求項8のいずれか1項に記載の方法。

40

【請求項10】

前記連続する構成要素は、縦方向に重なる構成要素および所定の距離離れている構成要素の少なくとも一方を含み、前記縦方向に重なる構成要素は、縦軸に沿って少なくとも1つの座標を共有する、請求項9に記載の方法。

50

【請求項 1 1】

前記複数の構成要素を合成することは、

前記少なくとも 1 つの副単語および少なくとも 1 つの単語の少なくとも一方を形成するために、語間に基づき、各欄に関連付けられる連続する構成要素を結合することと、

前記第 1 の組の特性パラメータに基づき、アラビア文字に関連付けられる複数の構成要素から、非文字項目に関連付けられる複数の構成要素のうちの少なくとも 1 つの構成要素をふり落とすこととを含む、請求項 9 に記載の方法。

【請求項 1 2】

ある欄に関連付けられる少なくとも 1 つの座標に基づき、複数の欄を分類することをさらに備え、前記少なくとも 1 つの座標は、画像における欄の位置に関連付けられる、請求項 1 1 に記載の方法。

10

【請求項 1 3】

前記方法はさらに、各欄に関連付けられる、各副単語および各単語の少なくとも一方に関連付けられる第 2 の組の特性パラメータを計算することを備え、前記第 2 の組の特性パラメータは、各副単語および各単語の少なくとも一方に関連付けられる行高さ、各副単語および各単語の少なくとも一方に関連付けられる語間、ならびに各副単語および各単語の少なくとも一方に関連付けられる行間の 1 つであり、さらに

第 2 の組の特性パラメータに基づき少なくとも 2 つの副単語をグループ化して、少なくとも 1 つの副単語および少なくとも 1 つの単語の一方を形成することを備える、請求項 1 ~ 請求項 1 2 のいずれか 1 項に記載の方法。

20

【請求項 1 4】

前記方法は、前記少なくとも 1 つの副単語および前記少なくとも 1 つの単語を、各副単語および各単語の少なくとも一方に関連付けられる行高さ、ならびに各副単語および各単語の少なくとも一方に関連付けられる行間の少なくとも一方に基づき、少なくとも 1 本の横行に分割することをさらに備える、請求項 1 3 に記載の方法。

【請求項 1 5】

光学式文字認識 (OCR) 用に画像を前処理するためのシステムであって、画像は複数の欄を含み、複数の欄の各欄は、アラビア文字および非文字項目の少なくとも一方を含み、前記システムは、

メモリと、

30

前記メモリに結合されるプロセッサとを備え、前記プロセッサは、

複数の欄の中のアラビア文字および非文字項目の少なくとも一方に関連付けられる複数の構成要素を定めるように構成され、構成要素は一連の接続されたピクセルを含み、前記プロセッサは、さらに、

前記複数の構成要素に関連付けられる行高さおよび欄間を計算することと、

行高さおよび欄間に基づき、複数の構成要素のうちの少なくとも 1 つの構成要素を、複数の欄のうちの欄に関連付けることと、

前記複数の欄のうちの各欄について第 1 の組の特性パラメータを計算することと、

第 1 の組の特性パラメータに基づき、複数の欄のうちの各欄の中の複数の構成要素を合成して、少なくとも 1 つの副単語および少なくとも 1 つの単語の少なくとも一方を形成することを行なうように構成され、前記複数の欄は、第 1 の領域と第 2 の領域とを含み、

40

前記プロセッサは、

最も近い構成要素が第 1 の領域に入る構成要素を当該第 1 の領域に合成することと、

最も近い構成要素が第 2 の領域に入る構成要素を当該第 2 の領域に合成することを行なうように構成され、前記第 2 の領域の少なくとも部分は前記第 1 の領域の少なくとも部分と縦方向に重なり、

前記第 1 の組の特性パラメータは、各欄に関連付けられる行高さ、各欄に関連付けられる語間、各欄に関連付けられる行間、各構成要素に対応するピクセルの数、各構成要素の幅、各構成要素の高さ、各構成要素の座標、各構成要素の密度、または各構成要素のアスペクト比の 1 つを含み、

50

各欄に関連付けられる行間を計算するために、前記プロセッサは、

各欄に関連付けられる前記複数の構成要素の中の複数の水平射影のヒストグラムを作成するように構成され、前記複数の水平射影のうちのある水平射影は、ラスタスキャンの各掃引に対応して、前記複数の構成要素に関連付けられるピクセルの数を示し、さらに前記プロセッサは、

2つの連続する最大水平射影間の平均距離を計算し、

前記平均距離に基づき行間を計算するよう構成されている、システム。

【請求項 16】

前記プロセッサは、

濃淡画像およびカラー画像の少なくとも一方を2進画像に変換すること、

ごま塩雑音をふるい落とすこと、および

変形ハフ変換を用いて歪みを修正すること、のうちの少なくとも一つを行なうようさらに構成されている、請求項15に記載のシステム。

【請求項 17】

複数の構成要素を定めるために、前記プロセッサは、

画像に対してラスタスキャンを行ない、

ラスタスキャンの少なくとも1回の掃引に対応して前記複数の構成要素の少なくとも1つの構成要素に関連付けられる複数のピクセルを特定し、

複数のピクセル間の相互接続に基づき、前記複数のピクセルを統合して少なくとも1組の接続されたピクセルを形成するようさらに構成されている、請求項15または請求項16に記載のシステム。

【請求項 18】

前記行高さを計算するために、前記プロセッサは、

前記複数の構成要素の各々の高さに対応する高さのヒストグラムを作成し、

前記高さのヒストグラムから頻出高さを特定し、

前記頻出高さに基づき行高さを計算するよう構成されている、請求項15～請求項17のいずれか1項に記載のシステム。

【請求項 19】

前記プロセッサは、行高さに基づき欄間を計算するようさらに構成されている、請求項18に記載のシステム。

【請求項 20】

各欄に関連付けられる語間を計算するために、前記プロセッサは、

各欄に関連付けられる前記複数の構成要素のうちの連続する構成要素間の間隔のヒストグラムを作成し、

前記ヒストグラムから頻出間隔を特定するように構成され、前記頻出間隔は行高さによって定められるしきい値範囲内にあり、前記プロセッサは、

前記頻出間隔に基づき語間を計算するようさらに構成されている、請求項15～請求項19のいずれか1項に記載のシステム。

【請求項 21】

前記プロセッサは、

語間に基づき各欄に関連付けられる連続する構成要素を結合して、少なくとも副単語および少なくとも1つの単語の少なくとも一方を形成し、

前記第1の組の特性パラメータに基づき、アラビア文字に関連付けられる複数の構成要素から非文字項目に関連付けられる前記複数の構成要素のうちの少なくとも1つの構成要素をふるい落とすようさらに構成されている、請求項20に記載のシステム。

【請求項 22】

前記プロセッサは、ある欄に関連付けられる少なくとも1つの座標に基づき、複数の欄を分類するようさらに構成されており、前記少なくとも一つの座標は画像における欄の位置に関連付けられる、請求項21に記載のシステム。

【請求項 23】

前記プロセッサは、

各欄に関連付けられる、各副単語および各単語の少なくとも一方に関連付けられる第2の組の特性パラメータを計算するようにさらに構成され、第2の組の特性パラメータは、各副単語および各単語の少なくとも一方に関連付けられる行高さ、各副単語および各単語の少なくとも一方に関連付けられる語間、ならびに各副単語および各単語の少なくとも一方に関連付けられる行間のうちの1つであり、さらに前記プロセッサは、

第2の組の特性パラメータに基づき少なくとも2つの副単語をグループ化して、少なくとも1つの副単語および少なくとも1つの単語の一方を形成するようにさらに構成されている、請求項15～請求項22のいずれか1項に記載のシステム。

【請求項24】

10

前記プロセッサは、少なくとも1つの副単語および少なくとも1つの単語を、各副単語および各単語の少なくとも一方に関連付けられる行高さ、ならびに各副単語および各単語の少なくとも一方に関連付けられる行間の少なくとも一方に基づき、少なくとも1つの横行に分割するようにさらに構成されている、請求項23に記載のシステム。

【請求項25】

コンピュータによって実行されるプログラムであって、請求項1～請求項14のいずれか1項に記載の方法を前記コンピュータに実行させるための、プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

20

関連出願への相互参照

これは、米国特許出願連続番号第12/814,448号の継続出願であり、その全体を引用により援用する。

【0002】

発明の分野

本発明は一般に光学式文字認識(OCR:Optical Character Recognition)用に画像を前処理することに関し、画像はアラビア文字および/または非文字項目を含む。より具体的には、本発明は複数の欄を含む画像の前処理に関し、各欄はアラビア文字および/または非文字項目を含む。

【背景技術】

30

【0003】

背景

スキャンされた文書を編集可能および検索可能な文章に自動的に変換するには、正確かつ安定した光学式文字認識(OCR)システムを使用する必要がある。英文に対するOCRシステムは、さまざまな理由により、高いレベルの精度に達している。その主な理由の1つは、OCRシステムへの入力として、英文を分離された文字に前処理する機能にある。英文の各文字は、印刷された英字が繋がっていない性質により、分離することができる。しかし、スキャンされた繋がっている文字はOCRシステムへの課題であり、ピッチが変化している場合にその精度を落とす。

【発明の概要】

40

【発明が解決しようとする課題】

【0004】

アラビア語のスキャンされた文章は一連の繋がっている文字を含むので、文字に分割するのが難しい。アラビア文字での別の問題は、文字や後に続く母音の発音を示すために、多くの点やアクセント記号が文字の上下に入っていることである。これにより、英語向けに設計された前処理技術が正確にアラビア文字を処理することに用いられるのを妨げることとなる。

【0005】

アラビア文字の別の特徴は、アラビア語の文章は母音を示すアクセント記号を付けてもまたは付けなくても書くことができる点にある。さらに、英文は大文字または小文字の表

50

示を有するが、多くのアラビア語の文字は、その位置が単語の初め、単語の間、もしくは単語の終わりにあるのか、または単独の単語としてあるのかに応じて、3つまたは4つの形を含む。したがって、アクセント記号および単語内での文字の位置により、アラビア文字にはさまざまな組合せが可能であるので、現在のOCR前処理システムでアラビア文字を前処理することは不的確となる。

【0006】

さらに、アラビア文字および非文字項目の欄を複数有する画像では、各欄に関連付けられるアラビア文字はフォントのサイズ、スタイル、色などが変わり得る。フォントサイズが変わることにより、隣接する欄は行が揃わず、正確に分割できないかもしれない。

【0007】

したがって、アラビア文字および/または非文字項目を含む複数の欄を有する画像を前処理するための方法およびシステムが必要である。

【0008】

添付の図面であって、同じ参照符号は、それぞれの図面において同じまたは機能的に類似した要素を指し、以下の詳細な説明とともに明細書の中に組込まれてその一部をなす図面は、さまざまな実施例を示し、かつ本発明に従うさまざまな原理および利点を説明する役割を果たす。

【図面の簡単な説明】

【0009】

【図1】本発明のさまざまな実施例に従い、アラビア文字および/または非文字項目を有する複数の欄を含む画像の一例を示す図である。

【図2】本発明の一実施例に従い、画像に関連付けられる歪みを修正する際の画像の一例を示す図である。

【図3】本発明の一実施例に従い、2進画像に変換された画像の一例を示す図である。

【図4】本発明の一実施例に従い、光学式文字認識のために画像を前処理する方法のフロー図である。

【図5】本発明の一実施例に従い、複数の構成要素を定める方法のフロー図である。

【図6】本発明の一実施例に従い、行の高さを計算する方法のフロー図である。

【図7】本発明の一実施例に従い、複数の欄のうちのある欄に関連付けられる1つ以上の構成要素を有する画像を示す図である。

【図8】本発明の一実施例に従い、語間を計算する方法のフロー図である。

【図9】本発明の一実施例に従い、行間を計算する方法のフロー図である。

【図10】本発明の一実施例に従い、1つ以上の副単語および/または1つ以上の単語を形成するために、複数の構成要素を合成する方法を示す図である。

【図11】本発明の一実施例に従い、非文字項目が取除かれた画像の一例を示す図である。

【図12】本発明の一実施例に従い、行高さおよび行間に応じて、1つ以上の単語および1つ以上の副単語を1つ以上の横行に分割する一例を示す図である。

【図13】本発明のさまざまな実施例に従い、アラビア文字および/または非文字項目を含む複数の欄を含む画像を前処理するためのシステムのブロック図である。

【図14】発明の一実施例に従い、1つ以上の別の欄または領域に縦方向および/または横方向に重なる欄または領域と関連付けられる複数の構成要素のうち1つ以上の概略図である。

【図15】発明の一実施例に従い、1つ以上の別の欄または領域に縦方向および/または横方向に重なる欄または領域と関連付けられる複数の構成要素のうち1つ以上の概略図である。

【発明を実施するための形態】

【0010】

図面の要素は簡潔におよび明瞭にするために示されており、必ずしも尺度通りに描かれていないことは、当業者なら理解するであろう。たとえば、図面の一部の要素の寸法は、

10

20

30

40

50

本発明の実施例をわかりやすくするために、他の要素に対して拡大されて示されているかもしれない。

【 0 0 1 1 】

詳細な説明

本発明に従う実施例を詳細に説明する前に、実施例は主に光学式文字認識（OCR）用に画像を前処理するための方法およびシステムに関する方法の工程および装置の構成要素の組合せに基づいていることに注意しなければならない。画像は複数の欄を含み、各欄はアラビア文字および／または非文字項目を含む。したがって、装置の構成要素および方法の工程は、図面において適する場合は従来の記号によって示され、ここでの記載が当業者にとって容易に明らかとなる詳細でもって開示を曖昧にしないよう、本発明の実施例の理解に関連する具体的詳細のみが示されている。

10

【 0 0 1 2 】

本明細書では、第1および第2、上および下などのような相関的用語は、あるエンティティまたは動作を別のエンティティまたは動作と区別するためにのみ用いられており、これらのエンティティまたは動作間において実際にこのような関係または順序を必ずしも必要としないまたは意味しない。「含む」、「有する」またはその他のこのような用語の変形は、限定されない含有を網羅するために意図されており、一連の要素を含むプロセス、方法、物品または装置は、これらの要素のみを含むのではなく、プロセス、方法、物品もしくは装置に明記されていない要素、または固有の他の要素をも含み得る。「...を含む」の用語が付いている要素は、それ以外の制限がなければ、その要素を含むプロセス、方法、物品または装置において付加的同一要素の存在を排除するものではない。

20

【 0 0 1 3 】

ここに記載される発明の実施例は、OCR用に画像を前処理する方法の機能の一部、大部分、またはすべてを実施するために、特定の非トランザクション - クライアント回路と併せて、1つ以上の従来のトランザクション - クライアントと、その1つ以上のトランザクション - クライアントを制御する固有の記憶されているプログラム命令とを含み得ることが理解されるであろう。画像は複数の欄を含み、各欄はアラビア文字および／または非文字項目を含む。非トランザクション - クライアント回路は、無線受信装置、無線送信装置、信号ドライバ、クロック回路、電源回路、およびユーザ入力装置を含むことができるが、これらに限定されない。したがって、これらの機能は、OCR用に、アラビア文字および非文字項目を含む画像を前処理する方法の工程として解釈することができる。代替的に、機能の一部またはすべては、プログラム命令が記憶されていない状態マシンによって、または1つ以上の特定用途向け集積回路（ASIC）において実施することができ、各機能または特定の機能の一部の組合せは、カスタム論理として実施される。当然、これら2つのアプローチを組合せて用いることもできる。これらの機能の方法および手段がここに記載される。さらに、当業者なら、ここに開示されている概念および原理により、たとえば利用可能な時間、現行の技術および経済的な点を考慮して動機付けられる著しい努力および多くの設計的選択事項があったとしても、最小限の実験でもってこのようなソフトウェア命令、プログラムおよびICを容易に生成できると考えられる。

30

【 0 0 1 4 】

一般に、さまざまな実施例に従い、本発明は、OCR用に画像を前処理するための方法およびシステムを提供する。画像は複数の欄を含み、各欄はアラビア文字および／または非文字項目を含む。本方法は、複数の欄の中のアラビア文字および／または非文字項目に関連付けられる複数の構成要素を定めることを含む。ここで、構成要素は1組の繋がっているピクセルを含む。つぎに、複数の構成要素に関連付けられる行高さおよび欄間が計算される。その後、複数の構成要素のうちの1つ以上の構成要素は、行の高さおよび／または欄間に基づき、複数の欄のうちの欄に関連付けられる。さらに、複数の欄のうちの各欄に関連付けられる第1の組の特性パラメータが計算される。各欄に関連付けられる複数の構成要素は、第1の組の特性パラメータに基づいて合成されて、1つ以上の副単語および／または1つ以上の単語を形成する。

40

50

【 0 0 1 5 】

図 1 は、本発明のさまざまな実施例に従い、複数の欄を含む画像の一例を示し、複数の欄の各欄はアラビア文字および／または非文字項目を含む。画像は濃淡画像またはカラー画像のいずれかであり得る。さらに、画像はごま塩雑音を含み、歪んでいるかもしれない。OCR用に画像を前処理する前に、画像に関連付けられるごま塩雑音および歪みは取除かれる。さらに、画像は濃淡画像またはカラー画像から 2 進画像に変換される。

【 0 0 1 6 】

画像に関連付けられる歪みは、画像に関連付けられる基線を定めて、基線の配列に基づき画像を正しい位置に置くことによって修正される。基線の配列は、変形ハフ変換によって定められ、水平射影は複数の方向で定められる。水平射影は、画像の前景に関連付けられるピクセルの数を示す。アラビア語の近似単語長さに対応する妥当なランレングスが考慮されて、最も高いピクセル密度を有する方向が決定される。最も高いピクセル密度の方向が基線の配列と一致すると考えられる。その後、画像は基線の配列に基づき正しい位置に置かれる。図 2 は、画像に関連付けられる歪みを修正する際の画像の一例を示す。

【 0 0 1 7 】

歪みを修正する際、画像に関連付けられるごま塩雑音が取除かれる。ごま塩雑音は、任意に起こる白および黒ピクセルを表わし、暗い背景上の白い点または明るい背景上の黒い点を含み得る。一実施例において、ごま塩雑音はメディアンフィルタおよび／または多数フィルタを用いることによって除去することができる。当業者にとって、ごま塩雑音は当該技術分野におけるノイズ除去技術を用いることによっても除去できることは明らかである。

【 0 0 1 8 】

その後、画像は濃淡画像またはカラー画像から 2 進画像に変換される。たとえば画像が濃淡画像の場合、0 から 255 の各ピクセル値を 0 のピクセル値または 1 のピクセル値に変換することにより、画像は 2 進画像に変換される。ある実施例において、ピクセル値 0 は背景値を表わし、ピクセル値 1 は前景値を表わす。代替的に、ピクセル値 0 は前景値を表わし、ピクセル値 1 は背景値を表わしてもよい。ピクセル値 0 は白ピクセルに関連付けられ、ピクセル値 1 は黒ピクセルに関連付けられる。

【 0 0 1 9 】

あるピクセルのピクセル値を変換する前に、濃淡画像にしきい値が定められ、しきい値より上のピクセル値はピクセル値 1 に変換され、しきい値より下のピクセル値はピクセル値 0 に変換される。一実施例において、しきい値は濃淡画像のピクセル値のヒストグラムを作成することによって計算される。ヒストグラムは、各ピクセル値の頻度を表わす。このヒストグラムを作成する際、連続するピクセル値の頻度を加算して、その連続するピクセル値を、連続するピクセル値の結合された頻度を有する単一のピクセル値に置き換えることにより、平滑化されたヒストグラムを生成することができる。考慮される連続するピクセル値の数は、予め定めることができる。後で、平滑化されたヒストグラムの 2 つの最も顕著なピークが選択され、この 2 つの顕著なピーク間の最小の谷が定められる。最も低い谷の中で最も低い頻度を有するピクセル値がしきい値として選択される。図 3 は例示的に 2 進に変換された画像を示す。

【 0 0 2 0 】

別の例であって、画像がカラー画像の場合、カラー画像はまず濃淡画像に変換され、次に上記のように 2 進画像に変換される。一実施例において、カラー画像を濃淡画像に変換するために、全国テレビジョン方式委員会 (NTSC) のデフォルト値を用いることができる。

【 0 0 2 1 】

画像を 2 進画像に変換する際、ピクセル値 0 およびピクセル値 1 の発生数が数えられる。より低いカウントの 2 進値は前景値であると考えられ、より高いカウントを有する 2 進値は背景値であると考えられる。

【 0 0 2 2 】

図4を参照すると、本発明の一実施例に従い、光学式文字認識用に画像を前処理する方法のフロー図が示される。画像は複数の欄を含み、各欄はアラビア文字および/または非文字項目を含む。前記のように、画像の2進画像への変換、ごま塩雑音の除去、および画像に関連付けられる歪みの修正のいずれか1つ以上を行なうことにより、画像が得られる。画像を前処理するために、画像の中のアラビア文字および/または非文字項目に関連付けられる複数の構成要素がステップ402で決定される。ここで、構成要素は1組の繋がっているピクセルを含む。構成要素は、文字が他の文字に繋がらない場合、アラビア文字の1文字を表わす。したがって、複数の文字が他の文字に繋がる場合、繋がっている文字は1つの構成要素であると考えられる。複数の構成要素を決定する方法は、図5と併せてさらに説明される。

10

【0023】

複数の構成要素を決定する際、複数の構成要素に関連付けられる行高さおよび欄間がステップ404で計算される。複数の構成要素に関連付けられる行の高さは、複数の構成要素の各構成要素の高さに対応する高さのヒストグラムを作成することによって計算される。行高さおよび/または欄間は、画像の複数の構成要素のすべての構成要素の平均値に基づき計算される。たとえば、行高さは、複数の構成要素のすべての構成要素に対して平均化された頻出高さである。行高さを計算する方法は、図6と併せてさらに詳しく説明される。欄間は、行高さの関数として動的に計算される。ステップ406において、複数の構成要素のうちの1つ以上の構成要素は、図7で例示的に示されるように、行高さおよび/または欄間に基づき、複数の欄のうちのある欄に関連付けられる。すなわち、複数の構成要素は、図7の702、704および706で例示的に示されるように、複数の構成要素に関連付けられる行高さおよび欄間に基づき、複数の欄に分離される。たとえば、2つの横方向に連続する構成要素間の間隔が欄間よりも小さければ、その構成要素は同じ欄のものであると考えられ、それに応じて分離される。

20

【0024】

複数の構成要素がある欄に関連付けられると、ステップ408において、第1の組の特性パラメータが各欄について計算される。ある実施例において、第1の組の特性パラメータは、各欄に関連付けられる行高さ、各欄に関連付けられる語間、各欄に関連付けられる行間、各構成要素のピクセルの数、各構成要素の幅、各構成要素の高さ、各構成要素の座標、各構成要素の密度、および各構成要素のアスペクト比を含む。行高さ、語間、および行間を計算する方法は、それぞれ図6、図8、および図9と併せて説明される。

30

【0025】

その後、ステップ410において、各欄に関連付けられた複数の構成要素は、第1の組の特性パラメータに基づき合成されて、1つ以上の副単語および/または1つ以上の単語を形成する。複数の構成要素を合成する方法は、図10と併せてさらに説明される。

【0026】

図5は、本発明の一実施例に従い、複数の構成要素を定める方法のフロー図を示す。ステップ502において、ラスタスキャンが画像に対して行なわれる。ラスタスキャンでは複数の掃引を行ない、複数の構成要素に対応する各ピクセル列に対して1回掃引される。ラスタスキャンの1回以上の掃引の実行により、画像の前景に関連付けられる1つ以上のピクセルがステップ504において特定される。画像の前景は、複数の構成要素に対応する。その後、ステップ506において、複数のピクセル間の相互接続に基づき、その複数のピクセルは統合されて、1つ以上の組の接続ピクセルを形成する。一実施例において、複数のピクセルは、8個の隣接ピクセルと1つ以上繋がっている場合に相互接続していると考えられる。こうして、アラビア文字の連続する文字は、連続する文字に関連付けられる1つ以上のピクセルが互いに相互接続されている場合に、単一の構成要素を形成する。

40

【0027】

たとえば、ラスタスキャンの現行の掃引で特定されたピクセルは、当該ピクセルが前回の掃引で特定されたピクセルと繋がる場合には、そのピクセルと統合される。現行の掃引で特定されたピクセルが、前回の掃引で特定された複数のピクセルと繋がる場合、当該ピ

50

クセルはその複数のピクセルと統合される。別の例では、現行の掃引で特定された複数のピクセルが繋がっている場合、その複数のピクセルは統合される。同様に、ラスタスキャンの後続の掃引で特定される1つ以上のピクセルは、その1つ以上のピクセルが互いに繋がる場合にも統合される。統合されたピクセルは、1組の繋がっているピクセルを形成し、複数の構成要素のうちのある構成要素に対応付けられる。その結果、1つ以上の組の繋がっているピクセルは複数の構成要素に関連付けられる。

【0028】

接続するピクセルの組を決定する際に、各構成要素のピクセルの数、各構成要素の幅、各構成要素の高さ、各構成要素の座標、各構成要素の密度、および各構成要素のアスペクト比のいずれか1つ以上は、各構成要素に関連付けられる接続ピクセルをトラッキングすることにより計算される。

10

【0029】

図6を参照すると、本発明の一実施例に従い、行の高さを計算する方法のフロー図が示される。本方法はステップ602において、複数の構成要素の各々の高さに対応する高さのヒストグラムを作成することを含む。ヒストグラムは、複数の構成要素の各々の高さの頻出を表わす。ヒストグラムを作成する際、連続する高さ値の頻度を加算して、その連続する高さ値を、連続する高さ値の結合された頻度を有する単一の高さ値と置き換えることにより、平滑化されたヒストグラムを生成することができる。考慮される連続する高さ値の数は予め定められてもよい。たとえば、連続する高さ値の数が3個であると定められたのなら、高さが20ピクセルの頻度は、高さが19ピクセルの頻度プラス高さが20ピクセルの頻度プラス高さが21ピクセルの頻度となる。

20

【0030】

平滑化されたヒストグラムが得られると、頻出高さがステップ604で特定される。頻出高さを特定するために、アラビア文字に対応するアクセント記号や句読点の小さな構成要素の高さは除外される。これは、しきい値高さを設定し、頻出高さを特定するのに、しきい値高さより大きい高さを有する構成要素のみを考慮することによって行なわれる。頻出高さは、画像が複数の文字サイズを有する場合には、画像の主要文字サイズを表わす。

【0031】

頻出高さを特定する際、行高さはステップ606において頻出高さに基づき計算される。行高さは、頻出高さおよび乗率の積として計算される。乗率は頻出高さに依存する。行高さは、1つ以上の単語および/または1つ以上の副単語を、アラビア文字の1つ以上の横行に分割するのに用いることができる。さらに、行高さは、図8と併せて説明したように、語間を計算するために用いられる。

30

【0032】

図8は、本発明の一実施例に従い、語間を計算する方法のフロー図を示す。本方法は、ステップ802において、複数の構成要素のうちの連続する構成要素間の間隔のヒストグラムを作成することを含む。一実施例において、縦方向に重なり、かつ他の構成要素によって分けられていない2つの構成要素はすべて連続する構成要素であると考えられる。2つの構成要素は、縦軸に沿って1つ以上の共通の座標を共有する場合、縦方向に重なる。すなわち、連続する構成要素はアラビア文字1行に属する。代替的に、2つの構成要素が縦方向に重ならない場合、2つの構成要素は予め定められた距離で分けられている場合に、連続する構成要素であると考えられる。

40

【0033】

連続する構成要素間の間隔のヒストグラムを作成する際、平滑化されたヒストグラムは、連続する間隔値の頻度を加算することにより生成できる。連続する間隔値は、連続する間隔値の結合された頻度を有する単一の間隔値と置き換えられる。たとえば、10ピクセルの間隔値の頻度は、9ピクセルの間隔値の頻度と、10ピクセルの間隔値の頻度と、11ピクセルの間隔値の頻度との合計と置き換えられる。

【0034】

ステップ804において、平滑化されたヒストグラムから頻出間隔が特定される。頻出

50

間隔は、行高さによって定められるしきい値範囲内から特定される。たとえば、5分の1の行高さと半分の行高さとの間にある頻出間隔値を対象とすることができる。ステップ806において、語間は頻出間隔に基づき計算される。語間は、アラビア文字の2つの連続する単語の間の間隔である。

【0035】

図9は、本発明の一実施例に従い、行間を計算する方法のフロー図を示す。ステップ902において、前景に対応する複数の構成要素の複数の水平射影のヒストグラムが作成される。水平射影は、ラスタスキャンの掃引に対応する複数の構成要素に関連付けられるピクセルの数を示す。たとえば、ラスタスキャンの掃引が、複数の構成要素に関連付けられる15個のピクセルを特定すると、その掃引に対するピクセル列の水平射影は15である。

10

【0036】

その後、ステップ904において、2つの連続する最大水平射影間の平均距離が計算される。最大水平射影は、最も高い密度の領域を表わす。その後、ステップ906において、行間は、平均距離に基づき計算される。

【0037】

第1の組の特性パラメータを計算する際、複数の構成要素は合成されて、図10と併せて説明されるように、1つ以上の副単語および/または1つ以上の単語を形成する。

【0038】

図10は、本発明の一実施例に従い、1つ以上の副単語および/または1つ以上の単語を形成するために、複数の構成要素を合成する方法を示す。ステップ1002において、ある欄に関連付けられる連続する構成要素間の間隔が語間の係数未満である場合に、当該連続する構成要素が結合される。語間に加えて、連続する構成要素の座標も、連続する構成要素が結合されるか否かを定めることができる。連続する構成要素の語間および/または座標に基づいてある欄に関連付けられる連続する構成要素を結合することは、アラビア文字のある単語または副単語に対応する異なる構成要素の結合を引起す。

20

【0039】

たとえば、アクセント記号に関連付けられる構成要素は、構成要素の語間および位置に基づき、属する単語と結合される。1つの単語は1つ以上の構成要素を含み得る。構成要素の位置は、構成要素の座標によって定められる。ある構成要素に関連付けられる第1の組の特性パラメータのうち1つ以上が、アラビア文字の句読点またはアクセント記号と類似しており、かつアラビア語の文字に対応する構成要素に対して適切に隣接している場合、その構成要素は文字とともにグループ化されて単語または副単語を形成する。さもなければ、構成要素はノイズであると考えられ、除去される。

30

【0040】

アラビア文字に関連付けられる構成要素を結合することに加えて、非文字項目に関連付けられる構成要素も、語間に基づき結合される。非文字項目に関連付けられる構成要素は結合されて、1つ以上のより大きい構成要素を形成する。

【0041】

ステップ1004において、非文字項目に関連付けられる構成要素は、第1の組の特性パラメータに基づき、アラビア文字に関連付けられる構成要素からふり落とされる。たとえば、大きい高さ、大きい幅、および低い密度を有する構成要素は取除かれる。これらの構成要素は、ある欄の周りまたは他の非文字項目の周りの枠またはボーダーに対応し得る。同様に、大きい高さ、小さい幅、および高い密度を有する構成要素は縦線として認識され、除去される。横線は小さい高さ、大きい幅、および高い密度を有するものとして認識される。

40

【0042】

同様に、他の非文字項目も1つ以上のフィルタに基づき除去される。この1つ以上のフィルタは、画像の共通に起こる構成要素の長さ、構成要素の幅、構成要素のアスペクト比、構成要素の密度、および構成要素の合計数を用いて、非文字項目をアラビア文字からふ

50

るい落とす。2つ以上の欄にわたる非文字項目も、非文字項目に関連付けられる構成要素の寸法を、アラビア文字に関連付けられる構成要素の最もよく起こる寸法と比較することによって除去される。図11は、非文字項目が除去された画像を例示的に示す。

【0043】

非文字項目をアラビア文字からふり落とした後、1つ以上の単語および1つ以上の副単語の第2の組の特性パラメータが計算される。第2の組の特性パラメータは、各副単語および各単語の少なくとも一方に関連付けられる行高さ、各副単語および各単語の少なくとも一方に関連付けられる語間、ならびに各副単語および各単語の少なくとも一方に関連付けられる行間を含む。第2の組の特性パラメータは、1つ以上の副単語および/または1つ以上の単語を形成するために、複数の構成要素を結合するプロセスの精度をさらに上げるために、計算される。第2の組の特性パラメータに基づき、1つ以上の副単語をグループ化して、1つ以上の副単語および/または1つ以上の単語を形成する。

10

【0044】

その後、縦方向に重なり、かつ複数の欄のうちのある欄に関連付けられる1つ以上の副単語および1つ以上の単語は、分割されてアラビア文字の横行を形成する。一実施例において、1つ以上の副単語および1つ以上の単語は、行高さおよび/または行間に基づいても分割されてよい。たとえば、互いに重なる縦方向の構成要素を1つ以上有するので2本の横行が一緒に分割されると、その2本の横行は、行高さおよび/または行間に基づき分けられる。図12は行高さおよび行間に依存して、1つ以上の単語および1つ以上の副単語を1つ以上の横行に分割する例を示す。

20

【0045】

こうして、OCR用に画像を前処理する方法が開示される。画像は複数の欄を含み、各欄はアラビア文字および/または非文字項目を有する。この方法は、アラビア文字および非文字項目に関連付けられる複数の構成要素を定めることを含む。複数の構成要素のうちのある構成要素は、アラビア語の1つ以上の文字または1つ以上の非文字項目を表わす。複数の文字が相互接続されるのなら、構成要素は2つ以上の文字を表わす。

【0046】

複数の構成要素を決定する際、複数の構成要素に関連付けられる行高さおよび欄間が計算される。行高さおよび欄間は、すべての欄にわたるすべての構成要素の平均値を表わす。複数の構成要素は、平均行高さおよび平均欄間に基づき、1つ以上の欄に分離される。後で、各欄の複数の構成要素に関連付けられる第1の組の特性パラメータが計算される。各欄に関連付けられる複数の構成要素は、後で第1の組の特性パラメータに基づき合成されて、1つ以上の副単語および/または1つ以上の単語を形成する。

30

【0047】

ここに開示されている方法は、繋がっている文字を含み、かつ複数の欄を含むアラビア文字を正確に前処理して分割することを可能にする。本方法は、構成要素がノイズであるのかアラビア文字の一部であるのかを判断する場合に、アラビア文字に関連付けられるアクセント記号および句読点を考慮する。さらに、本方法は画像が複数の欄を含むか否かを特定し、それらを分離する。

【0048】

40

図13は本発明の一実施例に従い、光学式文字認識(OCR)のために、アラビア文字および/または非文字項目を含む画像を前処理するためのシステム1300のブロック図を示す。画像は複数の欄を含み、各欄はアラビア文字および/または非文字項目を含む。画像は濃淡画像およびカラー画像のどちらかであり得る。さらに、画像はごま塩雑音を含み、歪んでいるかもしれない。図13に示されるように、システム1300は、メモリ1302と、メモリ1302に結合されるプロセッサ1304とを含む。OCR用に画像を前処理する前に、プロセッサ1304は、図1と併せて説明したように、変形ハフ変換を用いて画像に関連付けられる歪みを除去する。その後、プロセッサ1304はごま塩雑音を除去し、濃淡画像またはカラー画像を2進画像に変換する。一実施例において、ごま塩雑音は、メディアンフィルタおよび/または多数フィルタを用いて除去され得る。ここで

50

は、画像を前処理するために、プロセッサ 1304 はアラビア文字および / または非文字項目に関連付けられる複数の構成要素を定める。構成要素は接続されたピクセルの組を含む。構成要素は、文字が他の文字と繋がらない場合、アラビア語文字の 1 つの文字を表わす。したがって、複数の文字が他の文字と繋がる場合、繋がっている文字は 1 つの構成要素であると考えられる。

【0049】

一実施例において、複数の構成要素を定めるために、プロセッサ 1304 は画像に対してラスタスキャンを行なう。ラスタスキャンでは複数の掃引を行ない、複数の構成要素に対応する各ピクセル列に対して 1 回掃引される。ラスタスキャンの 1 回以上の掃引の実行により、画像の前景に関連付けられる 1 つ以上のピクセルが特定される。画像の前景は、
10 複数の構成要素に対応する。その後、プロセッサ 1304 は、複数のピクセル間の相互接続に基づき、複数のピクセルを統合して、1 つ以上の組の接続ピクセルを形成する。統合されたピクセルは 1 組の接続ピクセルを形成し、複数の構成要素のうちのある構成要素に関連付けられる。

【0050】

こうしてプロセッサ 1304 によって定められた複数の構成要素は、メモリ 1302 に記憶することができ、プロセッサ 1304 によって用いられて、複数の構成要素に関連付けられる行高さおよび欄間が計算される。行高さおよび欄間を用いて、複数の構成要素のうち
20 の 1 つ以上の構成要素を、複数の欄のある欄に関連付ける。すなわち、複数の構成要素が行高さおよび / または欄間を満たすのなら、複数の構成要素はプロセッサ 1304 によって複数の欄に分離される。たとえば、2 つの縦にまたは横に連続する構成要素間の間隔が欄間よりも小さければ、それらの構成要素は同じ欄のものであると考えられて、分離される。その後、複数の構成要素に関連付けられる第 1 の組の特性パラメータが計算される。ある実施例において、第 1 の組の特性パラメータは、各欄に関連付けられる行高さ、各欄に関連付けられる語間、各欄に関連付けられる行間、各構成要素のピクセルの数、各構成要素の幅、各構成要素の高さ、各構成要素の座標、各構成要素の密度、および各構成要素のアスペクト比を含む。その後、プロセッサ 1304 は第 1 の組の特性パラメータに基づき、複数の構成要素を合成する。合成された構成要素は、1 つ以上の副単語および / または 1 つ以上の単語を形成する。

【0051】

一実施例において、プロセッサ 1304 は、複数の構成要素の各々の高さに対応する高さのヒストグラムを作成することによって行高さを計算する。ヒストグラムから頻出高さがプロセッサ 1304 によって特定される。その後、プロセッサ 1304 は頻出高さおよび乗率の積として、行高さを計算する。乗率は頻出高さに依存する。行高さを用いて、1 つ以上の単語および / または 1 つ以上の副単語を、アラビア文字の 1 つ以上の横行に分割することができる。さらに、プロセッサ 1304 は行高さを用いて語間を計算する。
30

【0052】

次に、プロセッサ 1304 は、複数の構成要素の連続する構成要素間の間隔のヒストグラムを作成することにより、語間を計算する。プロセッサ 1304 は、ヒストグラムから頻出間隔を特定する。頻出間隔はしきい値範囲内から特定され、そのしきい値範囲は行高さに基づいている。その後、語間はプロセッサ 1304 によって頻出間隔に基づき計算される。語間は、アラビア文字の 2 つの連続する単語の間の間隔である。
40

【0053】

プロセッサ 1304 は、複数の構成要素の複数の水平射影のヒストグラムを作成することによって行間を計算するようにも構成されている。水平射影は、ラスタスキャンの各掃引に対応する複数の構成要素に関連付けられるピクセルの数を示す。次に、2 つの連続する水平射影間の平均距離は、プロセッサ 1304 によって計算される。その後、プロセッサ 1304 は平均距離に基づき、行間を計算する。

【0054】

さらに、プロセッサ 1304 は、各構成要素のピクセル数、各構成要素の幅、各構成要
50

素の高さ、各構成要素の座標、各構成要素の密度、および各構成要素のアスペクト比を定める。

【 0 0 5 5 】

前述のように、プロセッサ 1 3 0 4 は、第 1 の組の特性パラメータに基づき、各欄に関連付けられる複数の構成要素を合成する。これを行なうため、プロセッサ 1 3 0 4 は、構成要素間の間隔がその欄に関連付けられる語間の係数よりも小さい場合に、連続する構成要素を結合する。各欄に関連付けられる語間に加えて、連続する構成要素の座標も連続構成要素が結合されるか否かを定めることができる。さらに、プロセッサ 1 3 0 4 は、図 1 0 と併せて説明されたように、第 1 の組の特性パラメータに基づき、アラビア文字に関連付けられる構成要素から、非文字項目に関連付けられる構成要素をふるい落とす。非文字項目をふるい落とすことは、1 つ以上の副単語および / または 1 つ以上の単語をもたらす。

10

【 0 0 5 6 】

プロセッサ 1 3 0 4 は、1 つ以上の副単語および / または 1 つ以上の単語に関連付けられる第 2 の組の特性パラメータを計算するようさらに構成されている。第 2 の組の特性パラメータは、各副単語および / または各単語に関連付けられる行高さ、各副単語および / または各単語に関連付けられる語間、ならびに各副単語および / または各単語に関連付けられる行間を含む。次に、2 つ以上の副単語は、第 2 の組の特性パラメータに基づきプロセッサ 1 3 0 4 によってグループ化されて、1 つ以上の副単語および / または 1 つ以上の単語を形成する。すなわち、2 つ以上の副単語は第 2 の組の特性パラメータに基づきグループ化されて、完全な単語またはより大きい副単語を形成する。

20

【 0 0 5 7 】

1 つ以上の副単語および 1 つ以上の単語を形成する際、プロセッサ 1 3 0 4 は、縦方向に重なりかつ複数の欄のうちのある欄に関連付けられる 1 つ以上の副単語および 1 つ以上の単語を分割して、アラビア文字の横行を形成する。一実施例において、1 つ以上の副単語および 1 つ以上の単語は、行高さおよび / または行間に基づき、プロセッサ 1 3 0 4 によって分割されてもよい。

【 0 0 5 8 】

図 1 4 および図 1 5 は、文字欄（または領域）の一部が他の欄内に入るまたは他の欄と重なるある事例を図示する。たとえば、文字の 1 つの矩形のまたは円形の区域が文字のより大きな矩形のまたは円形の区域内に入ってもよい。

30

【 0 0 5 9 】

図 1 4 は、文字領域「テキスト 2」が文字領域「テキスト 1」内にあり、文字領域「テキスト 3」が「テキスト 4」と部分的に縦方向に重なる 1 つの例を示す。これらの文字領域は、文章の 1 つ以上の行または列を有してもよい。以上で述べた技術と同様に、図 1 4 および図 1 5 の例が図示する技術は分割および前処理エンジンの使用に係る。複数の欄、文字および非文字、起こり得る歪み、起こり得るノイズを有する、濃淡またはカラーの完全未加工スキャン文書が分割および前処理エンジンに与えられる。この例では、このシステムの出力は分類された欄のシーケンスであり、かつ 2 進または濃淡形式の、最適な行認識エンジンに渡すべき、正しい側が上になった、歪んでいない、ノイズのない各欄内の行画像の一覧である。

40

【 0 0 6 0 】

文書は、分割および前処理エンジンがページ構造もしくは内容またはフォントまたはフォントのサイズの事前の情報を有していなくても、適切に正しい位置に置かれ、ノイズのないきれいな分割された文字行画像に分割される。

【 0 0 6 1 】

上述のように、濃淡画像は、それが元々であってもまたはカラー画像から変換されたものであっても、ページ毎に異なる動的しきい値化を用いて 2 進化画像に変換される。異なる背景および / またはシェーディングを有する領域が特定され、別々に 2 進化される。正確さのために、各々の領域は別々に 2 進化されてもよい。以上のように、文字ページの前

50

景が背景よりも少ないピクセルを有する各領域に対して、前景／背景のチェックが行なわれる。各領域の2進化は濃淡のヒストグラムを用いて行なわれる。

【0062】

繋がっている構成要素は、単一ランレングスのアルゴリズムを用いて2進画像から計算される。実現例に従うと、スキャンは、2進画像に対して左上から右へ文書の下方に移動して行なわれる。各々のランレングス上でセグメントが見出され、長いセグメントは不純物として無視される。このプロセスを次の行以降について繰返して、上述のように、すべての新しい構成要素が形成されるまで行同士の間で接続するセグメントを合成する。古い列に対して新しい列がチェックされる。あるセグメントまたは未完了の構成要素に何も接続していなければ、これは完了済み構成要素に変換される。1つ以上の未完了構成要素に接続している新しいセグメントは、それらに加えられるか、または1つよりも多くの未完了構成要素を1つに合成するように用いられる。

10

【0063】

1つ以上の寸法、アスペクト比、密度、および周などの構成要素パラメータを用いて、繋がっている構成要素に対して統計的分析を行なう。明白な非文字構成要素をふるい落とすように規則が加えられる。

【0064】

次に残余の文字状の構成要素に対して統計的分析を行なって、ページが他のサイズの文字および領域を有するかもしれないとしても、支配的な文字サイズの単語または行の共通の高さを見出す。各々の単一の欄に属する構成要素に対して再び統計的分析を行なって、異なる欄が異なるサイズの文字を有する場合は、新たな欄特有の単語の共通の高さを用いてさらに文字および非文字を分ける。語間および欄間はこの共通高さをを用いて動的に計算される。このように、変更できない値は存在せず、ページ計算は動的であり、ページサイズが拡大または縮小されたとしても変化しない。これは、サイズとともに変化する無次元パラメータを用いて達成される。計算された距離が計算された語間よりも大きくかつ欄間よりも小さければ、単語を合成してもよい。

20

【0065】

構成要素を欄に分離するのに用いられる単語高さおよび欄間に基づいて距離が計算される。欄のうちいくつかは、「テキスト1」領域内に入る「テキスト2」領域などの他の欄内、または「テキスト3」領域の部分の下に入る「テキスト4」領域などの文字欄が横部分の下または上にある任意の方向のL字型欄内に入り得るので、各々の欄を規定する矩形は横および／または縦方向に重なる。区別される重なる欄または領域に構成要素を分けることは、各々の構成要素を分析すること、およびこの構成要素に最も近い文字構成要素がその欄内に入る場合は構成要素を欄と合成することによってなされ、その距離は以上で計算された余白よりも小さい。最も近い欄と構成要素とを合成すれば、1つの大きな欄となるすべての重なり合う欄となるため、合成は、現在の構成要素に最も近い構成要素を包含する欄に対してなされる。

30

【0066】

図15は、図14に示すような、領域「テキスト3」として処理されるように次に合成される領域「テキスト3a」-「テキスト3e」などの1つ以上の構成要素を分析する例示的な技術を示す。同様のフォントサイズおよび近似した寸法を有する互いの上にある層である「テキスト3a」-「テキスト3e」などの欄は、高さまたは幅などの寸法が異なる1つの欄または領域に合成されてもよい。欄または領域は、ページ上のそれらの場所（たとえばそれらの座標）に基づいて分類されて、認識された結果的に得られた文字を正しい順序で流すのを助ける。文書の上から下へかつ左から右へ文字を処理して、ページの幅または高さ全体にわたり得るいくつかの広い欄を利用して、正しい位置に置くことおよび文書の文字をどのように処理するかについての視覚的の手がかりを与えるように、後処理アルゴリズムを構成してもよい。当該アルゴリズムは、処理のために欄を分けるのに非文字行またはボーダーも利用してもよい。平均行高さを求めるのに統計的分析をさらに用いてもよい。行の厚みが平均行高さよりも数倍大きければ、行は連続した行同士の間の最も重

40

50

なっていない横方向の場所でいくつかの行に分けられる。

【 0 0 6 7 】

本発明の多様な実施例は、OCR用に画像を前処理するための方法およびシステムを提供する。画像は複数の欄を含み、各欄はアラビア文字および／または非文字項目を含む。本発明は、アラビア文字を、OCRシステムによって正確に処理することができる分割された行の副単語および単語に分割する。本方法は、構成要素がノイズであるのかアラビア文字の一部であるのかを判断する場合に、アラビア文字に関連付けられるアクセント記号および句読点も考慮に入れる。

【 0 0 6 8 】

当業者は、ここに記載される上記認識される利点および他の利点は一例であって、本発明のさまざまな実施例の利点すべてを含むことは意図されていないと認識するであろう。

10

【 0 0 6 9 】

以上で、本発明の具体的実施例が説明された。しかし、当業者なら、さまざまな変形および変更が、添付の請求項に記載されている本発明の範囲から逸脱することなく行なうことができるとう理解するであろう。したがって、明細書本文および図面は限定するのではなく例示するものであり、変形はすべて本発明の範囲内に含まれることが意図される。利益、利点、問題の解決、および利益、利点または解決を引起すまたは顕著にする要素は、請求項のいずれかまたはすべてにおける重大な、必要な、または必須の特徴もしくは要素であると考えべきではない。本発明は、本願の係属中になされた補正を含む添付の請求項およびこれら請求項の均等物すべてによってのみ規定される。

20

【 符号の説明 】

【 0 0 7 0 】

1 3 0 0 システム

1 3 0 2 メモリ

1 3 0 4 プロセッサ

4 0 2 アラビア文字および／または非文字項目に関連付けられる複数の構成要素を定める

4 0 4 複数の構成要素に関連付けられる行高さおよび欄間を計算する

4 0 6 複数の構成要素のうちの1つ以上の構成要素をある欄に関連付ける

4 0 8 各欄について、第1の組の特性パラメータを計算する

30

4 1 0 第1の組の特性パラメータに基づいて各欄の複数の構成要素を合成して、1つ以上の副単語および1つ以上の単語のうちの少なくとも1つ以上を形成する

【図 1】



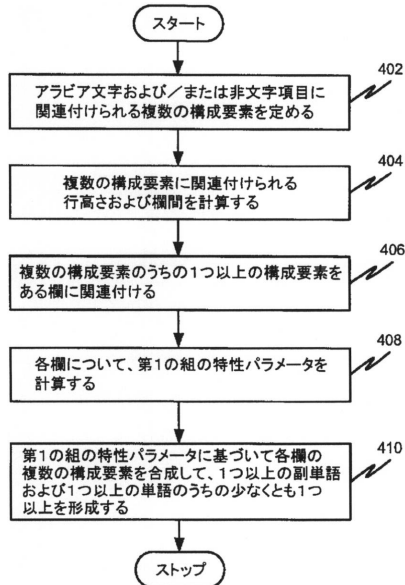
【図 2】



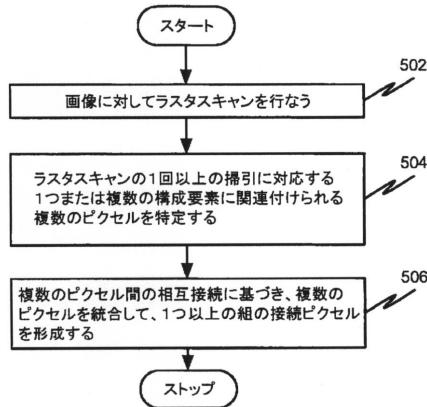
【図 3】



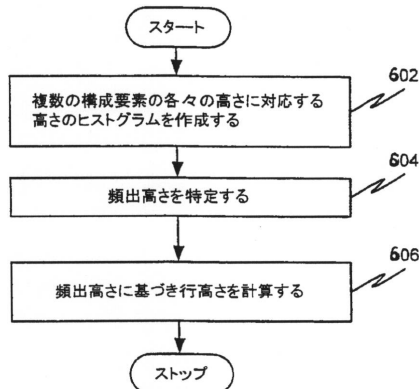
【図 4】



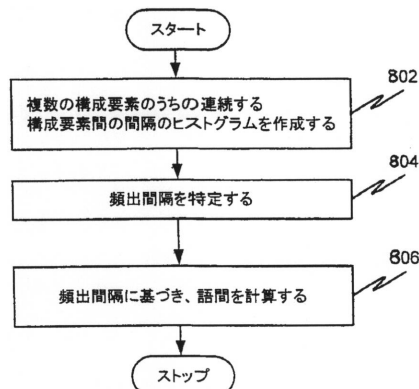
【図 5】



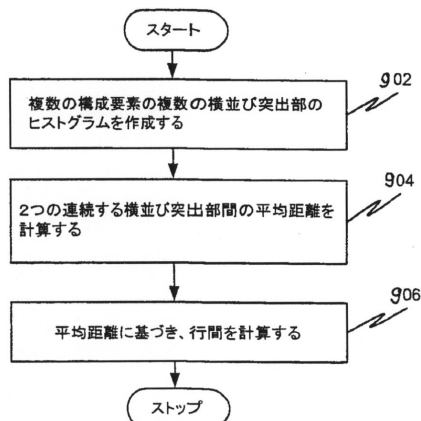
【図 6】



【図 8】



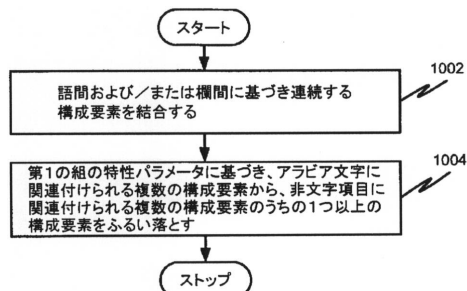
【図 9】



【図 7】



【図 10】



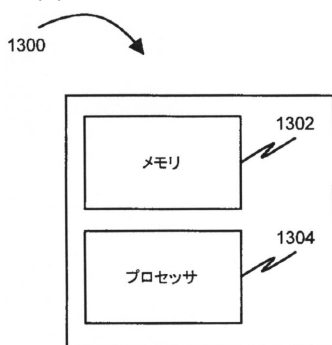
【 ㄨ 1 1 】【

[illegible]

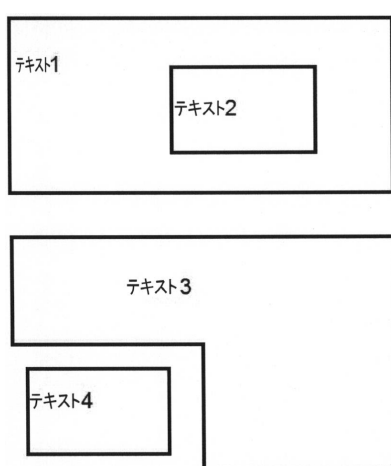
【 ㊦ 1 2 】



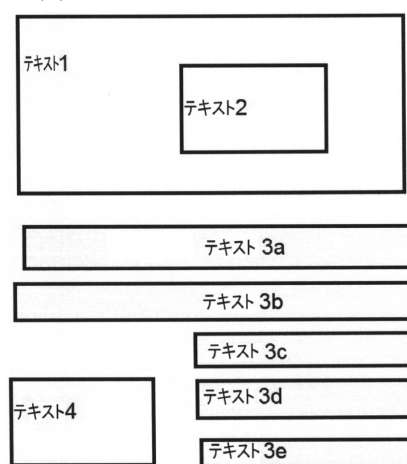
【 図 1 3 】



【 図 1 4 】



【 図 1 5 】



フロントページの続き

- (72)発明者 フセイン・ハリド・アル - オマリ
ヨルダン、 1 1 9 5 3 アンマン、 ピィ・オウ・ボックス・ 2 7 5 0
- (72)発明者 モハメド・スレイマン・ホルシード
サウジアラビア、 1 1 4 4 2 リヤド、 ピィ・オウ・ボックス・ 6 0 8 6、 ビルディング・ 1 7、
ブロック・エフ、 オフィス・ 4 1 5、 キング・アブドゥルアジズ・シティ・フォー・サイエンス・
アンド・テクノロジー

審査官 松浦 功

- (56)参考文献 特開 2 0 1 2 - 0 0 3 7 5 6 (J P , A)
特表 2 0 0 9 - 5 4 5 8 0 7 (J P , A)
国際公開第 2 0 0 8 / 1 3 8 3 5 6 (W O , A 1)
米国特許出願公開第 2 0 1 0 / 0 2 4 6 9 6 3 (U S , A 1)
米国特許出願公開第 2 0 1 0 / 0 2 7 2 3 6 1 (U S , A 1)
米国特許出願公開第 2 0 1 2 / 0 0 8 7 5 8 4 (U S , A 1)

- (58)調査した分野(Int.Cl. , D B 名)
G 0 6 K 9 / 2 0 - 9 / 6 0