

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4885112号
(P4885112)

(45) 発行日 平成24年2月29日(2012.2.29)

(24) 登録日 平成23年12月16日(2011.12.16)

(51) Int. Cl.	F 1
G06T 7/00 (2006.01)	G06T 7/00 300F
G06F 17/30 (2006.01)	G06F 17/30 170B
	G06F 17/30 340B

請求項の数 18 (全 24 頁)

(21) 出願番号	特願2007-293392 (P2007-293392)	(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19年11月12日(2007.11.12)	(74) 代理人	100089118 弁理士 酒井 宏明
(65) 公開番号	特開2009-122758 (P2009-122758A)	(72) 発明者	大黒 慶久 東京都大田区中馬込1丁目3番6号 株式会社リコー内
(43) 公開日	平成21年6月4日(2009.6.4)	審査官	鹿野 博嗣
審査請求日	平成22年6月3日(2010.6.3)	(56) 参考文献	特開2005-242579 (JP, A)) 特開2003-208433 (JP, A))

最終頁に続く

(54) 【発明の名称】 文書処理装置、文書処理方法及び文書処理プログラム

(57) 【特許請求の範囲】

【請求項1】

文書画像間の照合を行う文書処理装置において、
前記文書画像に含まれた文字画像毎の外接矩形に基づいて、当該外接矩形を連結した文字行を切り出す文字行切出手段と、

前記文字行内における前記外接矩形の特性を表す配置情報を固定段階に量子化する量子化手段と、

前記量子化された配置情報の各々を固定種類のシンボルにシンボル化するシンボル生成手段と、

所定個の前記シンボルの組合せからなるシンボル系列の出現頻度を算出する出現頻度算出手段と、

照合対象の文書画像と、当該文書画像の被照合対象となる複数の文書画像とについて、前記出現頻度算出手段により算出された出現頻度を照合し、より高い相関を有した被照合対象の文書画像を所定数選定する被照合対象選定手段と、

前記照合対象の文書画像と、前記被照合対象選定手段により選定された被照合対象の文書画像の各々で一致した前記シンボル系列に対応する各配置情報に基づいて、当該各配置情報の何れか又は全てが表す外接矩形の出現位置の分布状態を文書画像毎に導出する分布状態導出手段と、

前記分布状態導出手段により導出された前記照合対象の文書画像についての分布状態と、前記被照合対象の文書画像についての分布状態との類似度を判定し、最も高い類似度を

有した被照合対象の文書画像を照合結果として選定する照合結果選定手段と、
を備えたことを特徴とする文書処理装置。

【請求項 2】

前記分布状態導出手段は、前記文書画像の水平方向及び/又は垂直方向について、前記外接矩形の出現位置の分布状態を導出することを特徴とする請求項 1 に記載の文書処理装置。

【請求項 3】

前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を度数分布ヒストグラムとして導出することを特徴とする請求項 1 又は 2 に記載の文書処理装置。

【請求項 4】

前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を正規分布とみなし、当該正規分布の平均、標準偏差、歪度及び尖度を導出することを特徴とする請求項 1 又は 2 に記載の文書処理装置。

【請求項 5】

前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおける前記文字行に含まれた各文字画像の前記外接矩形のサイズを集計し、当該サイズの平均値又は最頻値により前記正規分布を規定する数値を正規化することを特徴とする請求項 4 に記載の文書処理装置。

【請求項 6】

前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおいて一致した前記シンボル系列に対応する配置情報が表す外接矩形のサイズを集計することを特徴とする請求項 5 に記載の文書処理装置。

【請求項 7】

文書画像間の照合を行う文書処理装置で実行される文書処理方法であって、
文字行切出手段が、前記文書画像に含まれた文字画像毎の外接矩形に基づいて、当該外接矩形を連結した文字行を切り出す文字行切出ステップと、

量子化手段が、前記文字行内における前記外接矩形の特性を表す配置情報を固定段階に量子化する量子化ステップと、

シンボル系列生成手段が、前記量子化された配置情報の各々を固定種類のシンボルにシンボル化するシンボル生成ステップと、

出現頻度算出手段が、所定個の前記シンボルの組合せからなるシンボル系列の出現頻度を算出する出現頻度算出ステップと、

被照合対象選定手段が、照合対象の文書画像と、当該文書画像の被照合対象となる複数文書画像とについて、前記出現頻度算出手段により算出された出現頻度を照合し、より高い相関を有した被照合対象の文書画像を所定数選定する被照合対象選定ステップと、

分布状態導出手段が、前記照合対象の文書画像と、前記被照合対象選定ステップで選定された被照合対象の文書画像の各々とで一致した前記シンボル系列に対応する各配置情報に基づいて、当該各配置情報の何れか又は全てが表す外接矩形の出現位置の分布状態を文書画像毎に導出する分布状態導出ステップと、

照合結果選定手段が、前記分布状態導出ステップで導出された前記照合対象の文書画像についての分布状態と、前記被照合対象の文書画像についての分布状態との類似度を判定し、最も高い類似度を有した被照合対象の文書画像を照合結果として選定する照合結果選定ステップと、

を含むことを特徴とする文書処理方法。

【請求項 8】

前記分布状態導出手段は、前記分布状態導出ステップにおいて、前記文書画像の水平方向及び/又は垂直方向について、前記外接矩形の出現位置の分布状態を導出することを特徴とする請求項 7 に記載の文書処理方法。

【請求項 9】

前記分布状態導出手段は、前記分布状態導出ステップにおいて、前記外接矩形の出現位

10

20

30

40

50

置の分布状態を度数分布ヒストグラムとして導出することを特徴とする請求項7又は8に記載の文書処理方法。

【請求項10】

前記分布状態導出手段は、前記分布状態導出ステップにおいて、前記外接矩形の出現位置の分布状態を正規分布とみなし、当該正規分布の平均、標準偏差、歪度及び尖度を導出することを特徴とする請求項7又は8に記載の文書処理方法。

【請求項11】

前記分布状態導出手段は、前記分布状態導出ステップにおいて、前記照合対象の文書画像と、前記被照合対象の文書画像とにおける前記文字行に含まれた各文字画像の前記外接矩形のサイズを集計し、当該サイズの平均値又は最頻値により前記正規分布を規定する数値を正規化することを特徴とする請求項10に記載の文書処理方法。

10

【請求項12】

前記分布状態導出手段は、前記分布状態導出ステップにおいて、前記照合対象の文書画像と、前記被照合対象の文書画像とにおいて一致した前記シンボル系列に対応する配置情報が表す外接矩形のサイズを集計することを特徴とする請求項11に記載の文書処理方法。

【請求項13】

文書画像間の照合を行うコンピュータを、

前記文書画像に含まれた文字画像毎の外接矩形に基づいて、当該外接矩形を連結した文字行を切り出す文字行切出手段と、

20

前記文字行内における前記外接矩形の特性を表す配置情報を固定段階に量子化する量子化手段と、

前記量子化された配置情報の各々を固定種類のシンボルにシンボル化するシンボル生成手段と、

前記シンボル系列内における、所定個のシンボルの組合せからなるシンボル系列の出現頻度を算出する出現頻度算出手段と、

照合対象の文書画像と、当該文書画像の被照合対象となる複数の文書画像とについて、前記出現頻度算出手段により算出された出現頻度を照合し、より高い相関を有した被照合対象の文書画像を所定数選定する被照合対象選定手段と、

前記照合対象の文書画像と、前記被照合対象選定手段により選定された被照合対象の文書画像の各々とで一致した前記シンボル系列に対応する各配置情報に基づいて、当該各配置情報の何れか又は全てが表す外接矩形の出現位置の分布状態を文書画像毎に導出する分布状態導出手段と、

30

前記分布状態導出手段により導出された前記照合対象の文書画像についての分布状態と、前記被照合対象の文書画像についての分布状態との類似度を判定し、最も高い類似度を有した被照合対象の文書画像を照合結果として選定する照合結果選定手段と、

して機能させることを特徴とする文書処理プログラム。

【請求項14】

前記分布状態導出手段は、前記文書画像の水平方向及び/又は垂直方向について、前記外接矩形の出現位置の分布状態を導出することを特徴とする請求項13に記載の文書処理プログラム。

40

【請求項15】

前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を度数分布ヒストグラムとして導出することを特徴とする請求項13又は14に記載の文書処理プログラム。

【請求項16】

前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を正規分布とみなし、当該正規分布の平均、標準偏差、歪度及び尖度を導出することを特徴とする請求項13又は14に記載の文書処理プログラム。

【請求項17】

前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とに

50

おける前記文字行に含まれた各文字画像の前記外接矩形のサイズを集計し、当該サイズの平均値又は最頻値により前記正規分布を規定する数値を正規化することを特徴とする請求項 16 に記載の文書処理プログラム。

【請求項 18】

前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおいて一致した前記シンボル系列に対応する配置情報が表す外接矩形のサイズを集計することを特徴とする請求項 17 に記載の文書処理プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、文書画像間の照合を行う文書処理装置、文書処理方法および文書処理プログラムに関する。

【背景技術】

【0002】

従来、文字列が画像（文字行）として記録された文書画像中から文字列を抽出する方法として、種々の技術が提案されている。例えば、文書画像に含まれた文字行に外接する矩形の形状及び位置に関する特徴（大きさ、間隔等）について、複数の制約を適用することにより文字行を文字列として認識することが可能な技術が提案されている（例えば、特許文献 1、2 参照）。

【0003】

【特許文献 1】特開平 11 - 219407 号公報

【特許文献 2】国際公開第 00 / 62243 号パンフレット

【発明の開示】

【発明が解決しようとする課題】

【0004】

しかしながら、特許文献 1、2 に記載の技術では、文字行の認識を精度よく行うために外接矩形に關する複数の制約を人手によって最適値に調整する必要がある。また、文字行らしさを判定することはできるものの、文字行の内容に関する特徴を認識することはできないため、文書画像間の照合に用いたとしても十分な精度を得ることができない可能性がある。また、複数行間の相対的な位置関係の利用については何等言及されていないため、文書画像の一部分となる部分画像を照合対象とした場合には、対応することができないという問題がある。

【0005】

本発明は、上記に鑑みてなされたものであって、文書画像間の類似性をより効率的且つ高精度に判定することが可能な文書処理装置、文書処理方法及び文書処理プログラムを提供することを目的とする。

【課題を解決するための手段】

【0006】

上述した課題を解決し、目的を達成するために、請求項 1 に係る発明は、文書画像間の照合を行う文書処理装置において、前記文書画像に含まれた文字画像毎の外接矩形に基づいて、当該外接矩形を連結した文字行を切り出す文字行切出手段と、前記文字行内における前記外接矩形の特性を表す配置情報を固定段階に量子化する量子化手段と、前記量子化された配置情報の各々を固定種類のシンボルにシンボル化するシンボル生成手段と、所定個の前記シンボルの組合せからなるシンボル系列の出現頻度を算出する出現頻度算出手段と、照合対象の文書画像と、当該文書画像の被照合対象となる複数の文書画像とについて、前記出現頻度算出手段により算出された出現頻度を照合し、より高い相関を有した被照合対象の文書画像を所定数選定する被照合対象選定手段と、前記照合対象の文書画像と、前記被照合対象選定手段により選定された被照合対象の文書画像の各々とで一致した前記シンボル系列に対応する各配置情報に基づいて、当該各配置情報の何れか又は全てが表す外接矩形の出現位置の分布状態を文書画像毎に導出する分布状態導出手段と、前記分布状

10

20

30

40

50

態導出手段により導出された前記照合対象の文書画像についての分布状態と、前記被照合対象の文書画像についての分布状態との類似度を判定し、最も高い類似度を有した被照合対象の文書画像を照合結果として選定する照合結果選定手段と、を備えたことを特徴とする。

【0007】

また、請求項2に係る発明は、請求項1に係る発明において、前記分布状態導出手段は、前記文書画像の水平方向及び/又は垂直方向について、前記外接矩形の出現位置の分布状態を導出することを特徴とする。

【0008】

また、請求項3に係る発明は、請求項1又は2に係る発明において、前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を度数分布ヒストグラムとして導出することを特徴とする。

10

【0009】

また、請求項4に係る発明は、請求項1又は2に係る発明において、前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を正規分布とみなし、当該正規分布の平均、標準偏差、歪度及び尖度を導出することを特徴とする。

【0010】

また、請求項5に係る発明は、請求項4に係る発明において、前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおける前記文字行に含まれた各文字画像の前記外接矩形のサイズを集計し、当該サイズの平均値又は最頻値により前記正規分布を規定する数値を正規化することを特徴とする。

20

【0011】

また、請求項6に係る発明は、請求項5に係る発明において、前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおいて一致した前記シンボル系列に対応する配置情報が表す外接矩形のサイズを集計することを特徴とする。

【0012】

また、請求項7に係る発明は、文書画像間の照合を行う文書処理装置で実行される文書処理方法であって、文字行切出手段が、前記文書画像に含まれた文字画像毎の外接矩形に基づいて、当該外接矩形を連結した文字行を切り出す文字行切出ステップと、量子化手段が、前記文字行内における前記外接矩形の特性を表す配置情報を固定段階に量子化する量子化ステップと、シンボル系列生成手段が、前記量子化された配置情報の各々を固定種類のシンボルにシンボル化するシンボル生成ステップと、出現頻度算出手段が、所定個の前記シンボルの組合せからなるシンボル系列の出現頻度を算出する出現頻度算出ステップと、被照合対象選定手段が、照合対象の文書画像と、当該文書画像の被照合対象となる複数文書画像とについて、前記出現頻度算出手段により算出された出現頻度を照合し、より高い相関を有した被照合対象の文書画像を所定数選定する被照合対象選定ステップと、分布状態導出手段が、前記照合対象の文書画像と、前記被照合対象選定ステップで選定された被照合対象の文書画像の各々とで一致した前記シンボル系列に対応する各配置情報に基づいて、当該各配置情報の何れか又は全てが表す外接矩形の出現位置の分布状態を文書画像毎に導出する分布状態導出ステップと、照合結果選定手段が、前記分布状態導出ステップで導出された前記照合対象の文書画像についての分布状態と、前記被照合対象の文書画像についての分布状態との類似度を判定し、最も高い類似度を有した被照合対象の文書画像を照合結果として選定する照合結果選定ステップと、を含むことを特徴とする。

30

40

【0013】

また、請求項8に係る発明は、請求項7に係る発明において、前記分布状態導出手段は、前記文書画像の水平方向及び/又は垂直方向について、前記外接矩形の出現位置の分布状態を導出することを特徴とする。

【0014】

また、請求項9に係る発明は、請求項7又は8に係る発明において、前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を度数分布ヒストグラムとして導出すること

50

を特徴とする。

【 0 0 1 5 】

また、請求項 1 0 に係る発明は、請求項 7 又は 8 に係る発明において、前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を正規分布とみなし、当該正規分布の平均、標準偏差、歪度及び尖度を導出することを特徴とする。

【 0 0 1 6 】

また、請求項 1 1 に係る発明は、請求項 1 0 に係る発明において、前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおける前記文字行に含まれた各文字画像の前記外接矩形のサイズを集計し、当該サイズの平均値又は最頻値により前記正規分布を規定する数値を正規化することを特徴とする。

10

【 0 0 1 7 】

また、請求項 1 2 に係る発明は、請求項 1 1 に係る発明において、前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおいて一致した前記シンボル系列に対応する配置情報が表す外接矩形のサイズを集計することを特徴とする。

【 0 0 1 8 】

また、請求項 1 3 に係る発明は、文書画像間の照合を行うコンピュータを、前記文書画像に含まれた文字画像毎の外接矩形に基づいて、当該外接矩形を連結した文字行を切り出す文字行切出手段と、前記文字行内における前記外接矩形の特性を表す配置情報を固定段階に量子化する量子化手段と、前記量子化された配置情報の各々を固定種類のシンボルにシンボル化するシンボル生成手段と、前記シンボル系列内における、所定個のシンボルの組合せからなるシンボル系列の出現頻度を算出する出現頻度算出手段と、照合対象の文書画像と、当該文書画像の被照合対象となる複数の文書画像とについて、前記出現頻度算出手段により算出された出現頻度を照合し、より高い相関を有した被照合対象の文書画像を所定数選定する被照合対象選定手段と、前記照合対象の文書画像と、前記被照合対象選定手段により選定された被照合対象の文書画像の各々とで一致した前記シンボル系列に対応する各配置情報に基づいて、当該各配置情報の何れか又は全てが表す外接矩形の出現位置の分布状態を文書画像毎に導出する分布状態導出手段と、前記分布状態導出手段により導出された前記照合対象の文書画像についての分布状態と、前記被照合対象の文書画像についての分布状態との類似度を判定し、最も高い類似度を有した被照合対象の文書画像を照合結果として選定する照合結果選定手段と、して機能させることを特徴とする。

20

30

【 0 0 1 9 】

また、請求項 1 4 に係る発明は、請求項 1 3 に係る発明において、前記分布状態導出手段は、前記文書画像の水平方向及び/又は垂直方向について、前記外接矩形の出現位置の分布状態を導出することを特徴とする。

【 0 0 2 0 】

また、請求項 1 5 に係る発明は、請求項 1 3 又は 1 4 に係る発明において、前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を度数分布ヒストグラムとして導出することを特徴とする。

【 0 0 2 1 】

また、請求項 1 6 に係る発明は、請求項 1 3 又は 1 4 に係る発明において、前記分布状態導出手段は、前記外接矩形の出現位置の分布状態を正規分布とみなし、当該正規分布の平均、標準偏差、歪度及び尖度を導出することを特徴とする。

40

【 0 0 2 2 】

また、請求項 1 7 に係る発明は、請求項 1 6 に係る発明において、前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおける前記文字行に含まれた各文字画像の前記外接矩形のサイズを集計し、当該サイズの平均値又は最頻値により前記正規分布を規定する数値を正規化することを特徴とする。

【 0 0 2 3 】

また、請求項 1 8 に係る発明は、請求項 1 7 に係る発明において、前記分布状態導出手段は、前記照合対象の文書画像と、前記被照合対象の文書画像とにおいて一致した前記シ

50

ンボル系列に対応する配置情報が表す外接矩形のサイズを集計することを特徴とする。

【発明の効果】

【0024】

本発明によれば、照合対象の文書画像と被照合対象の文書画像とについて、文字行内における外接矩形の特徴を表した配置情報を抽出し、これらを固定段階に量子化してシンボルを生成することにより、文字認識することなく文字行の特徴の抽出が可能となり、被照合対象の文書画像から、照合対象の文書画像と相関の高い被照合対象の文書画像を所定の数だけ選定することができる。また、照合対象の文書画像と、選定された被照合対象の文書画像とについて、一致するシンボル系列の出現位置の分布状態を照合することで、当該シンボル系列の相対的な位置関係の類似性を判定することができるため、照合対象の文書画像と被照合対象の文書画像との類似性を高精度に判定することができる。これにより、文書画像中の部分画像が照合対象の文書画像とされた場合であっても、この部分画像に含まれた文字画像の外接矩形の位置関係に基づいて、当該部分画像と類似する文書画像を高精度に検索することが可能となる。

10

【発明を実施するための最良の形態】

【0025】

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

【0026】

(文書処理装置のハードウェア構成)

20

図1は、本発明の第1の実施形態にかかる文書処理装置100のハードウェア構成を示したブロック図である。図1に示したように、文書処理装置100は、PC(Personal Computer)などのコンピュータであり、文書処理装置100の各部を制御するCPU(Central Processing Unit)1、CPU1を起動するためのプログラムが記憶されるROM(Read Only Memory)2、後述する画像入力部21により入力された文書画像やオペレーティングシステム、種々のプログラム等を記憶するハードディスク3、CPU1のワークエリアとして機能するRAM(Random Access Memory)4、オペレータからの各種入力を受け付けるキーボード5、入力状況等を表示する表示装置6、CD-ROMなどの各種光情報記録メディア(図示せず)に記憶されたプログラム等を読み取る光ディスクドライブ7、インターネットやLAN(Local Area Network)等の電気通信回線を介して文書画像を送受信する通信装置8、原稿画像の光学的な読み取りを行うスキャナ9等から構成されており、これらの各部間で入出力されるデータをバスコントローラ10が調停して動作する。

30

【0027】

文書処理装置100では、オペレータが電源を投入するとCPU1がROM2内のローダーというプログラムを起動させ、ハードディスク3よりオペレーティングシステムというコンピュータのハードウェアとソフトウェアとを管理するプログラムをRAM4に読み込み、このオペレーティングシステムを起動させる。このようなオペレーティングシステムは、オペレータの操作に応じてプログラムを起動したり、情報を読み込んだり、保存を行ったりする。オペレーティングシステムのうち代表的なものとしては、Windows(登録商標)、UNIX(登録商標)等が知られている。これらのオペレーティングシステム上で走る動作プログラムをアプリケーションプログラムと呼んでいる。

40

【0028】

ここで、文書処理装置100は、CPU1が実行するプログラムとして、後述する文書照合処理)にかかる文書処理プログラムをハードディスク3に記憶している。この意味で、ハードディスク3は、文書処理プログラムを記憶する記憶媒体として機能する。

【0029】

また、一般的には、文書処理装置100のハードディスク3にインストールされるプログラムは、CD-ROMなどの各種光情報記録メディアやFD等の磁気メディア等の記憶媒体に記録され、この記憶媒体に記録されたプログラムがハードディスク3にインストー

50

ルされる。このため、CD-ROMなどの各種光情報記録メディアやFD等の磁気メディア等の可搬性を有する記憶媒体も、文書処理プログラムを記憶する記憶媒体となり得る。さらには、文書処理プログラムは、例えば通信装置8を介して外部から取り込まれ、ハードディスク3にインストールされても良い。

【0030】

CPU1は、オペレーティングシステム上で動作する文書処理プログラムが起動すると、この文書処理プログラムとの協働により後述する各機能部を実現させる。以下、文書処理装置100の機能的構成について説明する。

【0031】

(文書処理装置の機能的構成)

図2は、文書処理装置100の機能的構成を示したブロック図である。図2に示したように、文書処理装置100は機能部として、画像入力部21、照合画像選択部22、矩形抽出部23、行切出部24、量子化部25、シンボル生成部26、出現頻度集計部27、候補画像選定部28、出現位置分布導出部29、照合結果選定部30及び表示部31を含み構成される。

【0032】

画像入力部21は、外部から入力される文書画像を受け付け、ハードディスク3に記憶する。具体的に、画像入力部21の機能は、図1に示した光ディスクドライブ7、通信装置8、スキャナ9により実現することができる。

【0033】

照合画像選択部22は、画像入力部21から入力される文書画像や、キーボード5を介して指定されたハードディスク3に記憶された文書画像を、照合対象の文書画像として選択する。以下、照合対象の文書画像を「照合画像」という。なお、照合画像選択部22は、文書画像中の特定の領域がキーボード5を介して指定された場合には、この領域内に含まれる部分的な文書画像(部分画像)を照合画像として選択するものとする。

【0034】

また、照合画像選択部22は、照合画像の照合先となる被照合対象の文書画像を選択する。ここで、被照合対象の文書画像は、例えば、ハードディスク3に予め記憶された一部又は全ての文書画像としてもよいし、キーボード5を介して指定された文書画像を被照合対象の文書画像としてもよい。以下、被照合対象の文書画像を「被照合画像」という。

【0035】

矩形抽出部23は、文書画像に含まれた各文字画像の外接矩形を抽出する。ここで「文字画像」とは、所定の言語からなる文字が画像として表されたものを意味する。行切出部24は、矩形抽出部23で抽出された外接矩形を連結することで文字行の切り出しを行う。以下、文字行に含まれる外接矩形を「行内矩形」という。

【0036】

量子化部25は、行切出部24で切り出された文字行に含まれる各行内矩形の特性を表す配置情報を固定段階に量子化する。ここで、行内矩形の特性とは、各行内矩形に対応する文字画像の黒画素密度や文字行内における行内矩形の高さ、始点位置等のパラメータ群であって、行内矩形に固有の配置状態を表すものである。なお、配置情報の量子化については後述する。

【0037】

シンボル生成部26は、量子化部25により量子化された配置情報の各々を固定種類のシンボルにシンボル化し、文書画像を構成する各文字行に対応する一連のシンボル系列を生成する。以下、文書画像全体についてのシンボル系列を全体シンボル系列という。

【0038】

出現頻度集計部27は、全体シンボル系列内において、所定個のシンボルの組合せからなるシンボル系列が出現する頻度(出現頻度)を算出する。候補画像選定部28は、照合画像と、当該照合画像の照合先となる被照合画像とについて、出現頻度集計部27により算出された出現頻度を照合し、より高い相関を有した被照合画像を所定個数選定する。以

10

20

30

40

50

下、候補画像選定部 28 により選定された被照合画像を「候補画像」という。

【0039】

出現位置分布導出部 29 は、照合画像と候補画像との各文書画像において、両文書画像で一致した各シンボル系列に対応する配置情報の何れか又は全てが表す行内矩形に基づき、当該行内矩形の出現位置の分布状態を文書画像毎に夫々導出する。また、出現位置分布導出部 29 は、照合画像についての分布状態と、候補画像についての分布状態との類似度を算出し、算出した類似度を対応する候補画像と対応付けて RAM 4 等に保持する。

【0040】

照合結果選定部 30 は、出現位置分布導出部 29 により算出された類似度に基づいて、最も高い類似度を有した候補画像を照合結果として選定する。

10

【0041】

表示部 31 は、画像入力部 21 から入力された文書画像や各処理の経過状況等の表示を行うとともに、照合結果選定部 30 により選定された候補画像の表示を行う。なお、表示部 31 の機能は、図 1 に示した表示装置 6 により実現できる。

【0042】

以下、文書処理装置 100 が実行する各種の処理のうち、本実施の形態に特長的な処理である文書照合処理について以下に説明する。

【0043】

図 3 は、文書照合処理の手順を示したフローチャートである。まず、照合画像選択部 22 は、画像入力部 21 から入力される文書画像や、キーボード 5 を介して指定された文書画像を照合画像として選択する（ステップ S1）。次いで、照合画像選択部 22 は、ステップ S1 で選択した照合画像の照合先となる、被照合画像を選択する（ステップ S2）。

20

【0044】

続いて、矩形抽出部 23、行切出部 24、量子化部 25、シンボル生成部 26 及び出現頻度集計部 27 は、ステップ S1、S2 で選択された各文書画像について、出現頻度集計処理を実行する（ステップ S3）。以下、図 4 を参照して、ステップ S3 の出現頻度集計処理について説明する。なお、出現頻度集計処理は、照合画像及び被照合画像の各々について行われるものとするが、以下の説明では「文書画像」と総称して説明する。

【0045】

図 4 は、出現頻度集計処理の手順を示したフローチャートである。まず、矩形抽出部 23 は、文書画像に含まれた各文字画像の黒画素に外接する外接矩形を抽出する（ステップ S31）。続いて、行切出部 24 は、水平方向に隣接する外接矩形同士を連結して文字行に成長させた後、この文字行を夫々切り出す（ステップ S32）。

30

【0046】

ここで、文書画像の行の切り出しについて、図 5 - 1 ~ 図 5 - 3 を参照して説明する。矩形抽出部 23 は、文書画像（図 5 - 1）について、黒画素の連結成分を求め、それと外接する外接矩形 A, B, C... を求める（図 5 - 2）。そして、行切出部 24 は、矩形抽出部 23 により求められた外接矩形を、水平方向に隣接する外接矩形同士を連結して文字行 Z に成長させる（図 5 - 3）。行内矩形の生成及び文字行の切り出しにかかる処理自体は、公知の手法を用いることができるため詳細な説明は省略する。

40

【0047】

なお、文書画像から一つの文字行として切り出す単位は、行単位や段落単位、章単位等で切り出すことが好ましい。一般的に文書画像に含まれる文字画像のサイズは、行単位や段落単位、章単位で均一となるため、このような纏まりで文字行を切り出すことで、当該文字行内に含まれる文字画像のサイズ（文字サイズ）を揃えることが可能となる。また、本実施形態では、外接矩形の成長を水平方向で実施する態様としたが、これに限らず、文字方向等に応じて垂直方向、或いは、水平方向及び垂直方向の両方で実施する態様としてもよい。

【0048】

図 4 に戻り、量子化部 25 及びシンボル生成部 26 は、ステップ S32 で切り出した各

50

文字行について、シンボル生成処理を実行する。以下、図6を参照してステップS33のシンボル生成処理について説明する。

【0049】

図6は、シンボル生成処理の手順を示したフローチャートである。まず、量子化部25は、ステップS32で切り出された各文字行の高さを計測する(ステップS331)。

【0050】

次いで、量子化部25は、各文字行に含まれる各行内矩形の水平方向の始点(X_s)に基づいて、当該行内矩形を昇順にソートすることで配置順序を整列する(ステップS332)。続いて、量子化部25は、整列した各行内矩形の配置状態を表す配置情報を夫々取得し、この配置情報を固定段階に量子化する(ステップS333)。以下、図7-1、図7-2、図8および図9を参照して、ステップS332、S333の処理を説明する。

【0051】

図7-1および図7-2は、行内矩形の配置例を示す説明図である。欧米系文字行は、図7-1に示すように、大文字と小文字とが混在していることに加え、アポストロフィー、アクサンテギュ、ウムラウトなど、記号類の有無が存在するので、行内矩形の始点の高さは、図7-1のaの位置とbの位置との2カ所に集中することは明らかである。つまり、矩形の配置位置は上下に対称ではない。一方、アジア系文字行は、図7-2に示すように、漢字、ひらがな、カタカナ、ハングルなど、文字の構造が複雑であり、行内矩形の始点の高さは、欧米系文字行で見られるような、2カ所への明確な集中はない。しかし、矩形の配置位置が上下左右、対称ではないことは、欧米系行と同じである。

【0052】

図7-1の欧文文字の行内矩形と、図7-2のアジア系文字の行内矩形とを比較してみると、行内矩形の並び方は言語の種類に関わらず、その文字行の内容に応じて変化していることがわかる。そこで、文字の外接矩形を抽出することで、文字の大まかな特徴を捉えることができる。すなわち、文字そのものを特定しなくても、例えば図8に示すように、矩形座標の始点(X_s, Y_s)と終点(X_e, Y_e)を求め、これを利用した文字画像の外接矩形の配置状態を表す特徴を取得するだけで各文字行の画像特徴を捉えることができる。

【0053】

行内矩形の配置位置が同じであっても、欧米系文字は構造が単純なためアジア系文字と比べて矩形内の黒画素密度は低くなる。なお、アジア系文字においても、構造が簡単なひらがな、カタカナの黒画素密度は低く、構造が複雑な漢字の黒画素密度が高くなることは言うまでもない。

【0054】

このように、文字行内における一つの矩形の配置状態は、行内矩形の始点の高さ、矩形サイズ(幅、高さ)行内矩形中の黒画素密度等を計測することによって唯一に定義することができる。ステップS333の処理では、これら計測結果を配置情報として各文字行の行内矩形毎に取得し、固定段階に量子化する。

【0055】

以下では、行内矩形の始点の高さを基準にして行内矩形の配置状態を定義する一例を示す。図9は、行内矩形の配置状態を示す特徴を量子化する方法を示す説明図である。原稿を特定していない状況下では、行高さは可変であり、処理が行高さの値に依存しないように、行内矩形の高さを次式で正規化する。なお、 y_s は行内矩形始点の高さ、 H はステップS332で取得した行高を意味する。

$$YsRate = y_s / H \quad \dots (1)$$

【0056】

ここで、 $0 < YsRate < 1$ であるから、 $YsRate$ を固定段階に量子化することは容易である。例えば、 N 段階に量子化するなら、

$$YsVal = INT(YsRate * (N - 1)) \quad \dots (2)$$

(ただし、 $INT()$: 小数点以下切捨て)

10

20

30

40

50

とすればよい。各段階は、 $0 \sim (N - 1)$ とラベル付けされる。矩形幅 w および矩形高さ h も同様の手順で量子化される。

【0057】

ところで、記憶容量節約および演算量低減のためなどの理由で、画像処理においては原画像そのものではなく圧縮画像を処理対象にする場合が多い。圧縮画像は、画素数が減るために文字画像の細部に関する情報は失われる。本発明は、図9に示すように、文字画像の外接矩形に注目するものであり、画像そのものの詳細な特徴に基づくものではない。したがって、原画像だけでなく、圧縮画像に対しても有効に機能しうる。

【0058】

なお、上記では文字行画像の特徴として行内矩形の始点の高さを基準としたが、これに限定されない。例えば、文字行画像の特徴として行内矩形の高さを用いる場合は、図9において、次のとおりである。

$$\text{HeightRate} = h / H \quad \dots (3)$$

$$\text{HeightVal}$$

$$= \text{INT}(\text{HeightRate} * (N - 1)) + 0.5 \quad \dots (4)$$

(ただし、 $\text{INT}()$ ：小数点以下切捨て)

各段階は、 $0 \sim (N - 1)$ とラベル付けされる。

【0059】

また、文字行画像の特徴として行内矩形の幅を用いる場合は、次のとおりである。

$$\text{WidthRate} = w / H \quad \dots (5)$$

$$\text{WidthVal}$$

$$= \text{INT}(\text{WidthRate} * (N - 1)) + 0.5 \quad \dots (6)$$

(ただし、 $\text{INT}()$ ：小数点以下切捨て)

各段階は、 $0 \sim (N - 1)$ とラベル付けされる。

【0060】

図5に戻り、続いて、シンボル生成部26は、ステップS333で量子化された配置情報の各々を固定種類のシンボルにシンボル化した後(ステップS334)、図4のステップS34の処理に移行する。

【0061】

以下、図10および図11を参照し、ステップS334の処理について説明する。上述したとおり、ステップS333で取得された配置情報は、対応する行内矩形の配置状態を特徴付けるものとなっている。ステップS334の処理では、量子化された配置情報に含まれる複数種類の測定結果を一つにまとめてシンボル化することで、一つの行内矩形を一つのシンボルに対応させる。

【0062】

例えば、矩形の始点の高さ、矩形高さ、矩形幅の3種の情報をまとめる。仮に、前述の処理で、矩形の始点の高さ(y_s / H)を15段階、矩形高さ(h / H)を8段階、矩形幅(w / H)を2段階に量子化するとする。この結果、図10に示すように、各情報は、矩形の始点の高さ(y_s / H)は15段階であるから4bits、矩形高さ(h / H)は8段階であるから3bits、矩形幅(w / H)は2段階であるから1bitで表現することができる。また、

$$4 \text{ bits} + 3 \text{ bits} + 1 \text{ bit} = 8 \text{ bits}$$

であるから、1byteの各ビットに全情報を格納することができる。そして、これらの3種の情報を一つにまとめたシンボルの種類は、

$$15 \text{ 段階} \times 8 \text{ 段階} \times 2 \text{ 段階} = 240 \text{ 種}$$

となる。

【0063】

ところで、矩形の配置状態を表す複数の特徴を多次元ベクトルの各次元とみなせば、矩形は、その各特徴を用いて一つのベクトルデータに変換(ベクトル量子化)できる。ベクトル量子化とは、周知のように、ベクトルデータの多数のバラエティから、それらを代表

10

20

30

40

50

する少数のベクトルデータを求めることである。求められた代表ベクトルに順にラベル付けすれば、ベクトルデータの系列を単なる一次元のシンボルデータの系列に変換することができる。ベクトル量子化に関しては、「ベクトル量子化と情報圧縮」(コロナ社) Allen Gersho, Robert M. Gray 著、田崎三郎ほか訳、に詳しい。

【0064】

なお、まとめる情報の種類及びその格納のための記憶エリアは、記憶サイズは固定ではなく、識別対象である文字行を特定するのに好適な情報を適宜選択し、決定することが可能であることは言うまでもない。また、図10では、矩形の始点の高さ、矩形高さ、矩形幅についてシンボル化する例を示したが、これに限らず、上述した黒画素密度などの配置情報を含めてシンボル化する態様としてもよい。

10

【0065】

以上の作業を経ることによって、シンボル生成部26は、各文字行に含まれる行内矩形を、固定個のシンボル(ラベル)に変換することができる。したがって、実際の行内矩形の配置は、図11に示すような単なるシンボルの並びとみなすことができる。これで、シンボル系列の並び傾向を記録することができ、行内矩形の並び傾向を記録することと等価となる。

【0066】

図4に戻り、出現頻度集計部27は、ステップS34でシンボル化した各配置情報に対して、所定個のシンボルの組合せからなるシンボル系列の出現頻度を照合画像及び被照合画像の各々について夫々算出、集計し(ステップS34)、図3のステップS4の処理に移行する。

20

【0067】

以下、ステップS34の処理について説明する。配置情報がシンボル化された後には、テキスト検索と同様に、一般的な検索手法によって検索することが可能になる。つまり、照合画像と被照合画像についてシンボル系列間の完全一致を求めればよい。ただし、文字行画像の読み取り誤差によって、文字矩形の特徴の計測結果は異なるので、文字行が同一であっても、そのシンボル変換結果が同一にならない場合もある。よって、シンボル系列の完全一致を求めるのみでは、同一文字行画像を検索できない虞がある。

【0068】

そこで、ステップS34の処理では、シンボル系列の完全一致ではなく、シンボル系列の並び傾向の相関を求める。具体的には、照合画像及び被照合画像について生成された全シンボル系列の各々における、所定個のシンボルの組みからなるシンボル系列の出現頻度を算出し集計する。

30

以下、詳述する。

【0069】

並びの傾向を記録する手段としては、 n -gramモデルがある。 n -gramモデルは、クロード・エルウッドシャノンによって提案された言語モデルである。このモデルでは、系列中のシンボルの出現が、直前の n 個(n は自然数)のシンボルに影響されるとしている。現在の状態が n 個前の入力に依存して決まる確率プロセスを n 重マルコフ過程と呼び、 n -gramモデルは($n-1$)重マルコフモデルとも呼ばれる。特に、 $n=3$ の場合をtrigramと呼び、広く使用されている。

40

【0070】

具体的には、下記式(7)で示されるモデルである。さらに、式(8)にしたがって、照合画像及び被照合画像の各全シンボル系列から3つのシンボルの組みからなるシンボル系列(trigram)の出現頻度を夫々算出する。

【数1】

$$P(W) = \prod_{i=1}^n P(w_i | w_{i-2}, w_{i-1}) \quad \dots(7)$$

$P(w_i | w_{i-2}, w_{i-1})$: w_{i-2}, w_{i-1} の後に w_i が出現する
条件付き確率

$$P(w_i | w_{i-2}, w_{i-1}) = C(w_{i-2}, w_{i-1}, w_i) / C(w_{i-2}, w_{i-1}) \quad \dots(8)$$

$C(w_{i-2}, w_{i-1})$: 系列 w_{i-2}, w_{i-1} の出現頻度
 $C(w_{i-2}, w_{i-1}, w_i)$: 系列 w_{i-2}, w_{i-1}, w_i の出現頻度

10

【0071】

一方で、trigramの出現頻度順位を求め、出現頻度の高い順にtrigramを集計する。表1に、trigram集計結果の一例を示す。

【表1】

表1

順位	trigram	出現頻度
1	[s013,s045,s032]	324
2	[s013,s064,s033]	312
3	[s015,s005,s221]	302
4	[s016,s145,s203]	287
5	[s134,s002,s102]	279
6	[s134,s095,s244]	269
7	[s137,s187,s105]	256
8	[s138,s076,s228]	244
9	[s140,s097,s003]	235
10	[s141,s045,s013]	209
...		

20

30

【0072】

表1において、出現頻度はtrigramに示した3つ組みのシンボル系列、即ち三つの行内矩形を表す配置情報が、全体シンボル系列中にこの順序で出現する頻度を表している。例えば、trigram[s013, s045, s032]では、s013, s045の後ろにs032が出現する頻度が324であり、trigram[s013, s064, s033]では、s013, s064の後ろにs033が出現する頻度が312であることを示している。このように、文書画像の全シンボル系列に関して表1に示したようなtrigram集計結果を求めることが、各文書画像の特徴を求めること(学習)に相当する。

40

【0073】

以上の動作を経ることによって、出現頻度集計部27は、照合画像および被照合画像の各文書画像について、表1に示したようなtrigramの出現確率の集計結果を導出する。

【0074】

50

続いて、候補画像選定部 28 は、ステップ S 3 の処理で導出された照合画像に対応する集計結果（照合画像集計結果）と、被照合画像に対応する集計結果（被照合画像集計結果）とを照合し、より高い相関を有した上位 n 個の被照合画像を候補画像として選定する（ステップ S 4）。ここで、「n」は 1 以上の整数であって、任意の値を設定することが可能であるものとする。

【0075】

照合画像集計結果と、被照合画像集計結果とを照合する場合、一つの文字行に含まれる行内矩形の個数は同値とならないことが多いため、出現頻度そのものを比較することは有意ではない。そのため、ステップ S 4 では、下記式（9）に示した順位相関係数を用いることで、照合画像集計結果と、被照合画像集計結果との相関を判定する。

$$R_{xy} = 1 - (6 * (R_{xi} - R_{yi})^2) / (n * (n^2 - 1)) \dots (9)$$

【0076】

ここで、n はデータ数、 R_{xi} は照合画像集計結果の順位毎の出現頻度、 R_{yi} は被照合画像集計結果の順位毎の出現頻度を意味しており、各順位について R_{xi} と R_{yi} との差を二乗した値の総和が n により演算されるようになっている。なお、順位相関係数に関しては、「ノンパラメトリック法」（培風館）柳川堯著に詳しい。

【0077】

候補画像選定部 28 は、照合画像集計結果と、被照合画像集計結果とに含まれる各出現頻度について、順位相関係数 R_{xy} を算出し、被照合画像のうち、 R_{xy} の値が“1”に近いものから n 個分の被照合画像を候補画像として選定する。なお、順位相関係数を統計的に検定し、最大の順位相関係数が有意な値を示さない場合には、照合画像に類似する被照合画像はない、と判断することとしてもよい。

【0078】

ここで、図 12 を参照して、上述したステップ S 1 ~ S 4 迄の処理の概要を説明する。ステップ S 1、S 2 の処理において、照合画像 X と、複数の被照合画像 Y とが選択されると、ステップ S 3 の処理では、これら文書画像を構成する各文字行に含まれた行内矩形の各々が、配置情報に基づいてシンボル化され、照合画像 X についての全シンボル系列 X 1 と、各被照合画像 Y についての全シンボル系列 Y 1 とが夫々生成される。そして、全シンボル系列中における、*trigram* の出現頻度が集計されることで照合画像 X に対応する照合画像集計結果 X 2 と、被照合画像 Y の夫々に対応する被照合画像集計結果 Y 2 とが導出される。続いて、照合画像集計結果 A 2 に含まれた順位毎の出現頻度と、被照合画像集計結果 B 2 の夫々に含まれた順位毎の出現頻度と、に基づいて順位相関係数 R_{xy} が算出される。

【0079】

続くステップ S 4 において、ステップ S 3 で算出された被照合画像 Y 毎の順位相関係数 R_{xy} の値に基づいて、この値が“1”に近いものから n 個分の被照合画像 Y が候補画像として選定されることになる。

【0080】

図 3 に戻り、出現位置分布導出部 29 は、出現位置分布照合処理を実行する（ステップ S 5）。以下、ステップ S 5 の出現位置分布照合処理について説明する。

【0081】

図 13 は、出現位置分布照合処理の手順を示したフローチャートである。まず、出現位置分布導出部 29 は、ステップ S 4 の処理で選定された n 個の候補画像から、本処理の対象とする候補画像を一つ選択する（ステップ S 51）。

【0082】

続いて、出現位置分布導出部 29 は、ステップ S 51 で処理対象とした候補画像の被照合画像集計結果と、照合画像の照合画像集計結果とに基づいて、両文書画像の間で一致する *trigram*、即ち三つのシンボルの組みからなるシンボル系列を選択する（ステップ S 52）。ここで、選択する *trigram* の個数は特に問わないものとするが、より

10

20

30

40

50

出現頻度の高い `trigram` を選択することが好ましい。また、`trigram` を構成する三つのシンボルのうち、何れかのシンボルを選択する態様としてもよい。

【0083】

次いで、出現位置分布導出部 29 は、照合画像と処理対象の候補画像とについて、文書画像の水平方向および垂直方向における、ステップ S52 で選択したシンボル系列に対応する行内矩形の出現位置の分布状態をヒストグラム（度数分布ヒストグラム）として導出する（ステップ S53）。

【0084】

図 8 に示したように、行内矩形は始点（ X_s 、 Y_s ）と終点（ X_e 、 Y_e ）との 2 点により表現される。そのため、水平方向（ X 軸）に関して分布をとる場合、始点 X_s についてヒストグラムを生成すればよく、垂直方向（ Y 軸）に関しては分布をとる場合、始点 Y_s についてヒストグラムを生成すればよい。

【0085】

図 14 は、行内矩形の存在位置の分布状態をヒストグラムで表現した一例を示した図である。同図に示したように、照合画像と被照合画像との両文書画像の間で一致した行内矩形（図中 K）について、文書画像の水平方向と垂直方向でのヒストグラムを夫々導出する。ヒストグラム集計にあたっての集計幅は、特に問わないものとするが、例えば、ステップ S3 の処理で切り出した各文字行の高さの平均値程度とすることとしてもよい。

【0086】

図 13 に戻り、次に出現位置分布導出部 29 は、ステップ S53 で求めた両ヒストグラムを照合し、その類似度を算出する（ステップ S54）。なお、本実施形態では両ヒストグラムの照合方法として、メジアン（中央値）、モード（最頻値）、平均の各々が属するデータ区間のヒストグラム値を、両ヒストグラムの間で比較するものとする。

【0087】

具体的には、ヒストグラムのデータ区間を座標の小さいものから順次番号付けし、メジアン、モード、平均の所属するデータ区間の番号を求める。ここで、メジアン、モード、平均の所属するデータ区間番号を（`MedianClassNo`、`ModeClassNo`、`AvClassNo`）と表現すれば、以下の 4 種の組が求められる。

（`MedianClassNoXaxQuery`、`ModeClassNoXaxQuery`、`AvClassNoXaxQuery`）・・・（10）

（`MedianClassNoYaxQuery`、`ModeClassNoYaxQuery`、`AvClassNoYaxQuery`）・・・（11）

（`MedianClassNoXaxDB`、`ModeClassNoXaxDB`、`AvClassNoXaxDB`）・・・（12）

（`MedianClassNoYaxDB`、`ModeClassNoYaxDB`、`AvClassNoYaxDB`）・・・（13）

【0088】

なお、「`XaxQuery`」は、照合画像の水平方向のヒストグラムを意味するものであり、上記（10）式は、照合画像の水平方向のヒストグラムにおける、該当するデータ区間番号のヒストグラム値を夫々意味する。また、「`YaxQuery`」は、照合画像の垂直方向のヒストグラムを意味するものであり、上記（11）式は、照合画像の垂直方向のヒストグラムにおける、該当するデータ区間番号のヒストグラム値を夫々意味する。また、「`XaxDB`」は、被照合画像の水平方向のヒストグラムを意味するものであり、上記（12）式は、被照合画像の水平方向のヒストグラムにおける、該当するデータ区間番号のヒストグラム値を夫々意味する。また、「`YaxDB`」は、被照合画像の垂直方向のヒストグラムを意味するものであり、上記（13）式は、被照合画像の垂直方向のヒストグラムにおける、該当するデータ区間番号のヒストグラム値を夫々意味する。

【0089】

出現位置分布導出部 29 は、上記 4 種の組の値を算出した後、下記（14）～（16）式を用いて、垂直方向についての照合画像のヒストグラムと、被照合画像のヒストグラムとの形状の類似度を算出する。

`MedianClassNoXaxDB` + CA = `MedianClassNoXaxQuery`・・・（14）

ModeClassNoXaxDB + CA=ModeClassNoXaxQuery . . . (1 5)

AvClassNoXaxDB + CA=AvClassNoXaxQuery . . . (1 6)

【 0 0 9 0 】

上記 (1 4) ~ (1 6) 式において、「 C A 」は定数であって、最初に処理する 1 式 (例えば (1 4) 式) から求まる値である。出現位置分布導出部 2 9 は、この定数 C A の値が残りの 2 式にて成立するか否か、つまり、残り 2 式での定数 C A からのずれの度合いを、照合画像のヒストグラムと、被照合画像のヒストグラムとの形状の類似度として算出する。なお、定数 C A からのずれの度合いは、例えば、 $C A' / C A$ を算出することで導出できる。ここで、 $C A'$ は、 $C A +$ (は定数 C A からのずれ値) であり、完全一致する際のずれの度合い、即ち類似度は “ 1 ” となる。

10

【 0 0 9 1 】

また、同様に出現位置分布導出部 2 9 は、下記 (1 7) ~ (1 9) 式を用いて、垂直方向についての、照合画像のヒストグラムと、被照合画像のヒストグラムとの形状の類似度を算出する。

MedianClassNoYaxDB + CB=MedianClassNoYaxQuery . . . (1 7)

ModeClassNoYaxDB + CB=ModeClassNoYaxQuery . . . (1 8)

AvClassNoYaxDB + CB=AvClassNoYaxQuery . . . (1 9)

【 0 0 9 2 】

上記 (1 7) ~ (1 9) 式において、「 C B 」は定数であって、上述した C A と同様、最初に処理する 1 式 (例えば (1 7) 式) から求まる値である。出現位置分布導出部 2 9 は、この定数 C B の値が残りの 2 式にて成立するか否か、つまり、残り 2 式での定数 C B からのずれの度合いを、照合画像のヒストグラムと、被照合画像のヒストグラムとの形状の類似度として算出する。なお、定数 C B からのずれの度合いは、上述した定数 C A についてと同様に導出することができる。

20

【 0 0 9 3 】

出現位置分布導出部 2 9 は、上記の手続きにより算出した水平方向および垂直方向での類似度を、処理対象の被照合画像と対応付けて R A M 4 等に保持する。ここで、水平方向 (又は垂直方向) に対して導出されるずれの度合いの個数は、2 式 (或いは 3 式) 分となるが、これらを個別に類似度として保持する態様としてもよいし、これらの平均値を類似度として保持する態様としてもよい。

30

【 0 0 9 4 】

なお、本実施形態では、文書画像の水平方向および垂直方向の両方向について、ヒストグラムの形状の類似度を算出したが、何れか一方向のみについて算出する態様としてもよい。また、本実施形態では、行内矩形の出現位置の分布状態をヒストグラムで表すものとしたが、これに限らず、例えば正規分布を用いて表すものとしてもよい。

【 0 0 9 5 】

図 1 5 は、行内矩形の存在位置の分布状態を正規分布で表現した一例を示した図である。同図に示したように、正規分布を用いて表す場合には、各行内矩形の始点に基づいた集計結果から、水平方向 (X 軸) に関して、平均 μ_x 、標準偏差 σ_x 、歪度、尖度を算出し、また同様に垂直方向 (Y 軸) に関して、平均 μ_y 、標準偏差 σ_y 、歪度、尖度を算出すればよい。

40

【 0 0 9 6 】

この場合、平均値については、照合画像と被照合画像とで画像サイズが異なる可能性があるため、直接比較することは有意ではない。正規分布の形状が一致しているか否かを求めるには、標準偏差、歪度、尖度が類似しているかを判定すればよい。例えば、検索画像の標準偏差、歪度、尖度と、被検索画像の標準偏差、歪度、尖度との各々を比較し、比率が 1 に近いものほど正規分布の形状が類似するものと判断することができる。

【 0 0 9 7 】

なお、照合画像の解像度と、被照合画像の解像度とが一致している場合には、同一文字を構成するドット数は同じになるが、解像度が異なる場合にはドット数は同じにならない

50

。つまり、ヒストグラムや正規分布の形状の一致を評価する場合にも、解像度が同じ場合には両者の数値をそのまま利用しても構わないが、解像度が異なる場合には、ドット数に基づく数値をそのまま利用することができない。

【 0 0 9 8 】

そこで、両文書画像の解像度が異なる場合、或いは解像度自体が未知の場合には、数値の正規化を行う必要がある。一般的な文書画像においては段落単位では文字のサイズは同一であるため、同じ段落に属する文字行は行高さが等しくなる。また、照合画像が被照合画像の一部であれば、同じ行高さになる可能性が高いことは明らかである。よって、被検索画像および検索画像において、各文字行の行高さを集計し、最頻出となる行高さについて、ヒストグラムを規定する数値（平均、モード、メジアン）を除算する。なお、正規分布の場合も同様である。また、最頻出の行高さではなく、各文字行の行高さの平均値で除算してもよい。いずれを選択するかは設計事項であり、使用する環境に応じて決定すればよい。

10

【 0 0 9 9 】

また、照合画像が被照合画像の一部であっても、その一部分の特異な部分だけが照合画像となった場合には、全体画像において最頻出する行高さが、部分画像において最頻出となる行高さとは一致しないことが考えられる。例えば、本文行と見出し行とは行高さが大きく異なる文書画像において、全体画像の行数としては本文行が圧倒的に多いと予想される。その文書の部分画像には見出し行だけしか含まれていない場合には、最頻出行は見出し行となり、全体画像の最頻出行から推定した行高さとは一致しないため、この一致しない結果に基づいて正規化しても正しい比較結果を得ることができないのは明らかである。

20

【 0 1 0 0 】

このような場合、照合画像と被照合画像との両文書画像内において、一致した行内矩形（シンボル系列）だけを対象に矩形サイズの集計を行い、最頻出した矩形サイズのドット数に基づいて、数値（平均、モード、メジアン）を正規化することで対応することができる。

【 0 1 0 1 】

図 1 3 に戻り、出現位置分布導出部 2 9 は、ステップ S 5 4 の処理で求めた類似度を、処理対象の候補画像に対応付けてハードディスク 3 又は R A M 4 に保持する（ステップ S 5 5）。続いて、出現位置分布導出部 2 9 は、ステップ S 4 の処理で選定された n 個の候補画像の全てに対して、本処理の処理対象としたか否かを判定する（ステップ S 5 6）。ここで、本処理の対象としていない未処理の候補画像が存在すると判定した場合には（ステップ S 5 6 ; N o）、ステップ S 5 1 へと再び戻り、未処理の候補画像のうち一つを処理対象として選択する。

30

【 0 1 0 2 】

一方、ステップ S 5 6 において、全ての候補画像を処理対象としたと判定した場合（ステップ S 5 6 ; Y e s）、図 3 のステップ S 6 の処理に移行する。

【 0 1 0 3 】

図 3 に戻り、照合結果選定部 3 0 は、ステップ S 5 の処理により R A M 4 等に保持された n 個の候補画像の類似度に基づいて、最も高い類似度を有した候補画像、即ち類似度の値が“ 1 ”に最も近かった候補画像を照合結果として選定する（ステップ S 6）。

40

【 0 1 0 4 】

続いて、表示部 3 1 は、ステップ S 6 の処理で照合結果に選定された文書画像を、照合画像に対する照合結果として表示装置 6 に表示し（ステップ S 7）、本処理を終了する。

【 0 1 0 5 】

図 1 6 は、上記文書照合処理の動作を説明するための図である。同図において、D 1 1 は照合画像であって、特定の文書画像中の一部分となる部分画像が照合画像に選択された場合を示している。また、D 2 1 ~ D 2 4 は、ステップ S 4 までの処理により選定された 4 つの候補画像を示している。なお、照合画像 D 1 1 は、候補画像 D 2 4 の部分画像となっている。即ち、候補画像 D 2 4 が照合画像 D 1 1 に最も類似する文書画像となっている

50

【0106】

上述したようにステップS4の処理では、行内矩形の配置情報に対応するシンボル系列を照合することで、照合画像D11と相関関係にある文書画像として、候補画像D21～D24までを絞り込むことが可能である。なお、照合画像D11、候補画像D21～D24中矩形Kで表した部分が、各文書画像で一致したシンボル系列（或いはシンボル）の行内矩形を意味している。

【0107】

しかしながら、ステップS4の処理ではシンボル系列の出現頻度に基づいて類似度を判断するのみであるため、候補画像D24が照合画像D11に最も類似する文書画像であること、即ち、照合画像D11が候補画像D24の一部分であることまでを判断することはできない。そのため、ステップS5の処理では、各文書画像で一致したシンボル系列の相対的な位置関係、即ち出現位置の分布状態を照合することで、候補画像D24が照合画像D11に最も類似する文書画像であることを特定することが可能となる。

【0108】

以上のように、本実施形態によれば、照合画像と被照合画像とについて、文字行内における外接矩形の特徴を表した配置情報を抽出し、これらを固定段階に量子化してシンボルを生成することにより、文字認識することなく文字行の特徴の抽出が可能となり、被照合画像から、照合画像と相関の高い被照合画像を所定の数だけ候補画像として選定することができる。また、照合画像と候補画像とについて、一致するシンボル系列の出現位置の分布状態を照合することで、当該シンボル系列の相対的な位置関係の類似性を判定することができるため、照合対象画像と候補画像との類似性を高精度に判定することができる。これにより、文書画像中の部分画像が照合対象の文書画像とされた場合であっても、この部分画像に含まれた文字画像の外接矩形の位置関係に基づいて、当該部分画像と類似する文書画像を高精度に検索することが可能となる。

【0109】

なお、本発明は、上記実施の形態そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化することができる。また、上記実施の形態に開示されている複数の構成要素の適宜な組み合わせにより、種々の発明を形成することができる。例えば、実施の形態に示される全構成要素からいくつかの構成要素を削除してもよい。さらに、異なる実施の形態にわたる構成要素を適宜組み合わせても良い。

【0110】

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

【0111】

また、本実施形態で実行される文書照合処理にかかるプログラムを、ROM等の記憶媒体に予め組み込んで提供するように構成してもよい。

【0112】

また、上記実施形態では、図2に示した各機能部をCPU1とROM2に記憶された所定のプログラムとの協働により実現する態様としたが、これに限らず、ハードウェア構成により実現する態様としてもよい。具体的には、リアルタイム性が重要視される場合には、処理を高速化するため、論理回路（図示せず）を別途設け、論理回路の動作により各種の演算処理を実行するようにすることが好ましい。

【0113】

また、上記実施形態では、文字行よりも小さな単位として行内矩形に着目したが、これに限らず、他の単位でも適用可能である。例えば、文字（文字画像）単位や単語単位の画像特徴でも数値化し量子化することで、上記と同様にシンボル化することが可能であり、

10

20

30

40

50

照合することが可能である。この場合、黒画画素に基づいて文字画像を切り出したのち、当該文字画像の外接矩形を文字単位又は単語単位で用いることで対応することが可能である。なお、文字単位又は単語単位での分割は、OCR (Optical Character Recognition) 等で用いられる公知の文字切り出し手法を用いればよい。

【0114】

代表的な文字切り出し手法として、射影を利用する方法がある。この方法では、水平行について、垂直方向に黒画素数を集計し、その分布を求め、ある黒画素数がしきい値以下の部分を分割位置候補とする。また、分割位置候補に対しては、行高さから推定した文字幅、隣接する分割位置との距離、行全体に亘る分割位置の周期性等の観点から妥協点を評価し、適当な分割位置の選択を行う(垂直行も同様)。

10

【0115】

また、単語単位に分割する他の方法としては、欧文等分かち書きの習慣のある言語については、単語間の空白に基づいて容易に実現することが可能である。このように、文字単位、単語単位等の単位で分割された場合であっても、その範囲の画像に外接する矩形を求めることが可能であり、その外接矩形の開始位置、終点位置を用いることで行内矩形に対する場合と同様な手順で量子化を行うことができる。

【産業上の利用可能性】

【0116】

以上のように、本発明に係る文書処理装置、文書処理方法および文書処理プログラムは、文書画像間を照合する文字処理装置に有用であり、特に、文書画像の一部となる部分画像を照合対象とし、この部分画像に類似する文書画像の検索を行う文書処理装置に適している。

20

【図面の簡単な説明】

【0117】

【図1】文書処理装置のハードウェア構成を示したブロック図である。

【図2】文書処理装置の機能的構成を示したブロック図である。

【図3】文書照合処理の手順を示したフローチャートである。

【図4】出現頻度集計処理の手順を示したフローチャートである。

【図5-1】文字行の切り出しを説明するための図である。

【図5-2】文字行の切り出しを説明するための図である。

30

【図5-3】文字行の切り出しを説明するための図である。

【図6】シンボル生成処理の手順を示したフローチャートである。

【図7-1】行内矩形の配置例を示した図である。

【図7-2】行内矩形の配置例を示した図である。

【図8】行内矩形に対する座標の設定例を説明するための図である。

【図9】行内矩形の配置状態を説明するための図である。

【図10】配置情報の量子化を説明するための図である。

【図11】量子化された配置情報をシンボル化した一例を示した図である。

【図12】文書照合処理の概要を説明するための図である。

【図13】出現位置分布照合処理の手順を示したフローチャートである。

40

【図14】行内矩形の存在位置の分布状態をヒストグラムで表現した一例を示した図である。

【図15】行内矩形の存在位置の分布状態を正規分布で表現した一例を示した図である。

【図16】文書照合処理の概要を説明するための図である。

【符号の説明】

【0118】

100 文書処理装置

1 CPU

2 ROM

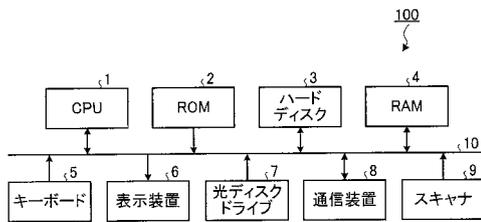
3 ハードディスク

50

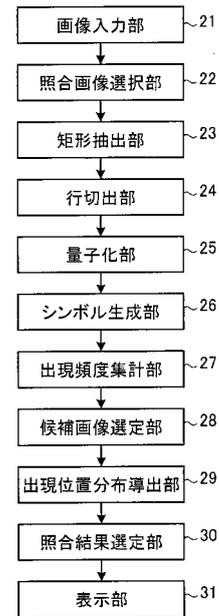
- 4 RAM
- 5 キーボード
- 6 表示装置
- 7 光ディスクドライブ
- 8 通信装置
- 9 スキャナ
- 10 バスコントローラ
- 21 画像入力部
- 22 照合画像選択部
- 23 矩形抽出部
- 24 行切出部
- 25 量子化部
- 26 シンボル生成部
- 27 出現頻度集計部
- 28 候補画像選定部
- 29 出現位置分布導出部
- 30 照合結果選定部
- 31 表示部

10

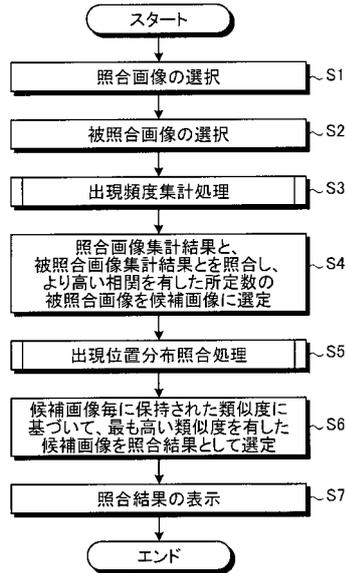
【図1】



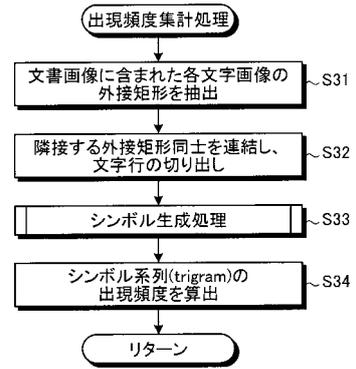
【図2】



【図3】



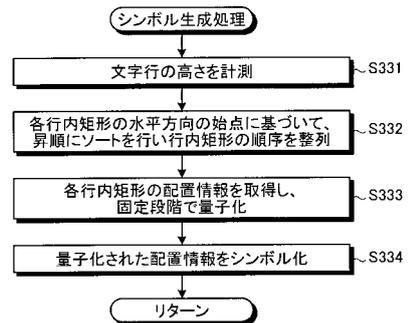
【図4】



【図5-1】

文字認識技術は、紙データを電子化する手段の一つです。オフィス業務のペーパーレス化が進んだ今日でも紙が使用される局面は多く、文字認識技術に対する期待は小さくなることはないようです。

【図6】



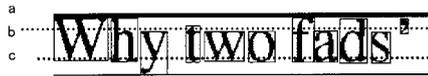
【図5-2】

A B C
文字認識技術は、紙データを電子化する手段の一つです。オフィス業務のペーパーレス化が進んだ今日でも紙が使用される局面は多く、文字認識技術に対する期待は小さくなることはないようです。

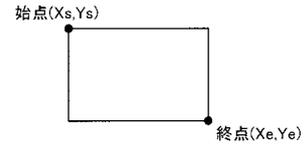
【図5-3】

文字認識技術は、紙データを電子化する手段の一つ
です。オフィス業務のペーパーレス化が進んだ今日
でも紙が使用される局面は多く、文字認識技術に
対する期待は小さくなることはないようです。

【図7-1】



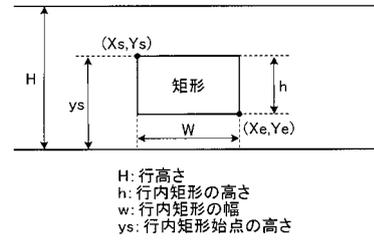
【図8】



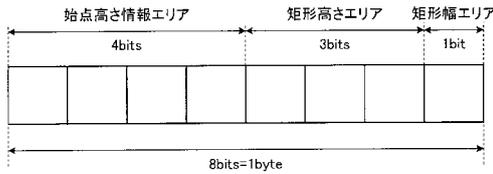
【図7-2】



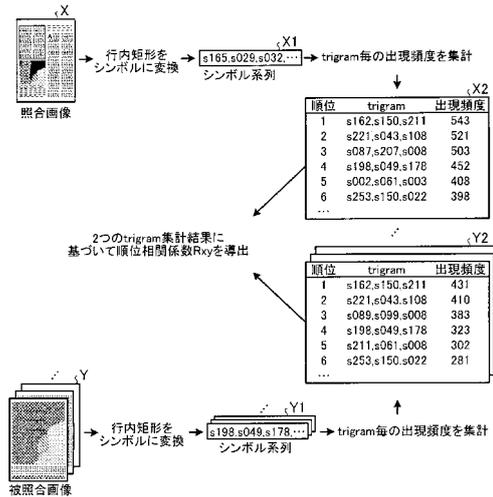
【図9】



【図10】



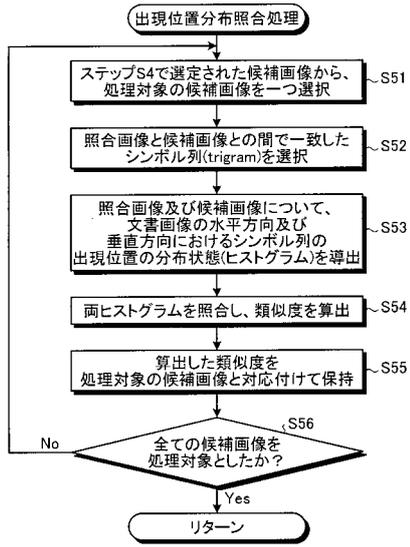
【図12】



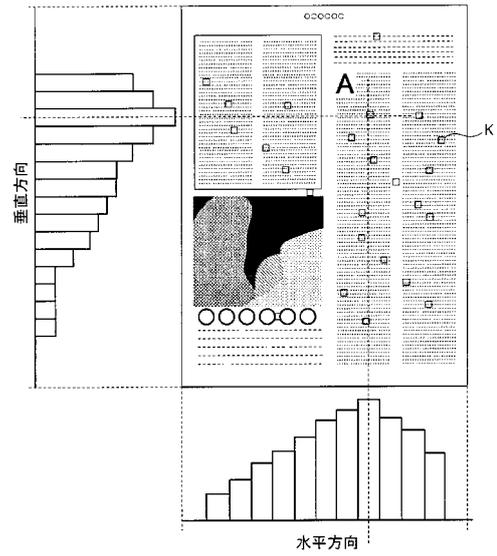
【図11】



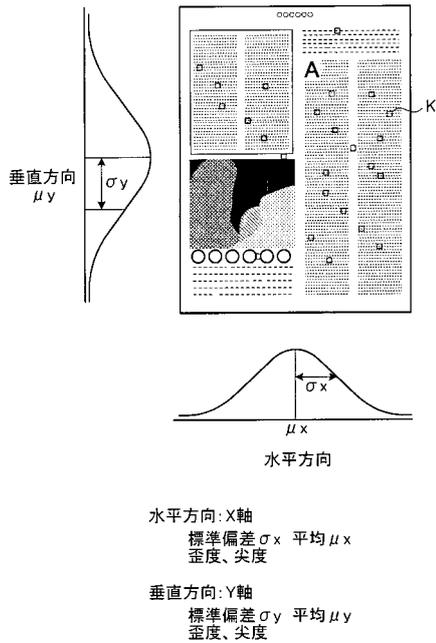
【図13】



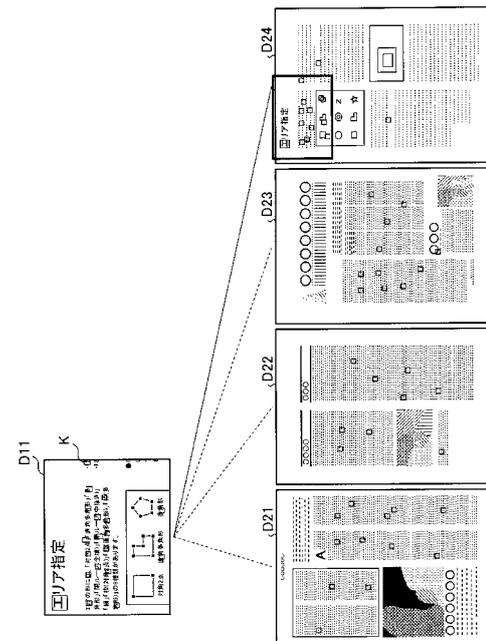
【図14】



【図15】



【図16】



フロントページの続き

(58)調査した分野(Int.Cl. , DB名)

G 0 6 T 7 / 0 0

G 0 6 F 1 7 / 3 0