

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4046941号  
(P4046941)

(45) 発行日 平成20年2月13日(2008.2.13)

(24) 登録日 平成19年11月30日(2007.11.30)

(51) Int.Cl.

F I

G 0 6 K 9 / 2 0 (2006.01)

G 0 6 K 9 / 2 0 3 4 0 C

請求項の数 11 (全 22 頁)

(21) 出願番号	特願2000-367675 (P2000-367675)	(73) 特許権者	000001007
(22) 出願日	平成12年12月1日(2000.12.1)		キヤノン株式会社
(65) 公開番号	特開2002-170079 (P2002-170079A)		東京都大田区下丸子3丁目30番2号
(43) 公開日	平成14年6月14日(2002.6.14)	(74) 代理人	100076428
審査請求日	平成16年6月11日(2004.6.11)		弁理士 大塚 康德
前置審査		(74) 代理人	100112508
			弁理士 高柳 司郎
		(74) 代理人	100115071
			弁理士 大塚 康弘
		(74) 代理人	100116894
			弁理士 木村 秀二
		(74) 代理人	100130409
			弁理士 下山 治
		(74) 代理人	100134175
			弁理士 永川 行光

最終頁に続く

(54) 【発明の名称】 文書書式識別装置および識別方法

(57) 【特許請求の範囲】

【請求項1】

文書画像の文書書式を識別する文書書式識別装置であって、

文書書式を識別すべき文書画像から抽出された複数のブロックそれぞれの位置座標を含む文書書式データを作成する作成手段と、

前記作成手段により作成された前記識別すべき文書画像の文書書式データと、保存手段に保存されているマスター文書画像の文書書式データとを比較することにより、相似関係があるか否かを判断するとともに、当該判断結果と、相似関係があると判断した前記識別すべき文書画像と前記マスター文書画像との間の変倍率と、を含む相似情報を抽出する相似情報抽出手段と、

前記相似情報抽出手段にて抽出した相似情報及び前記文書書式データに基づいて、前記識別すべき文書画像の前記マスター文書画像に対する類似度を計算することにより、前記識別すべき文書画像の文書書式を識別する識別手段と、を備え、

前記相似情報抽出手段は、

前記識別すべき文書画像から抽出されたブロックの個数と前記マスター文書画像のブロックの個数とが所定数以上でかつ互いに等しい場合に、

前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得たX座標値を所定の順序に基づいて整列させることにより得たX座標値列(この座標値列をY<sub>i</sub>とし、その平均値をY<sub>ave</sub>とする)と、前記マスター文書画像から抽出された複数のブロックそれぞれから得たX座標値を所定の順序に基づいて整列させることにより得たX座標

値列（この座標値列を  $X_i$  とし、その平均値を  $X_{ave}$  とする）との間の相関係数を、  

$$\frac{\{ (X_i - X_{ave}) (Y_i - Y_{ave}) \}}{\{ ( (X_i - X_{ave})^2 ) \times ( (Y_i - Y_{ave})^2 ) \}^{(1/2)}}$$

を用いて求め、

当該求めた相関係数が所定値より大きい場合に、前記識別すべき文書画像に関する X 座標値列と前記マスター文書画像に関する X 座標値列の傾きに基づいて X 座標方向の変倍率を求める一方、当該求めた相関係数が所定値以下の場合には相似関係がないと判断し、

更に、前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得た Y 座標値を所定の順序に基づいて整列させることにより得た Y 座標値列（この座標値列を  $Y_i'$  とし、その平均値を  $Y_{ave}'$  とする）と、前記マスター文書画像から抽出された複数のブロックそれぞれから得た Y 座標値を所定の順序に基づいて整列させることにより得た Y 座標値列（この座標値列を  $X_i'$  とし、その平均値を  $X_{ave}'$  とする）との間の相関係数を、

$$\frac{\{ (X_i' - X_{ave}') (Y_i' - Y_{ave}') \}}{\{ ( (X_i' - X_{ave}')^2 ) \times ( (Y_i' - Y_{ave}')^2 ) \}^{(1/2)}}$$

を用いて求め、

当該求めた相関係数が所定値より大きい場合に、前記識別すべき文書画像に関する Y 座標値列と前記マスター文書画像に関する Y 座標値列の傾きに基づいて Y 座標方向の変倍率を求める一方、当該求めた相関係数が所定値以下の場合には相似関係がないと判断し、

前記 X 座標方向の変倍率と前記 Y 座標方向の変倍率の両方が求められた場合に相似関係があると判断することを特徴とする文書書式識別装置。

【請求項 2】

前記識別手段は、

前記相似情報に基づいて、前記類似度の計算に用いる文書書式データを補正し、類似度の計算を行うことを特徴とする請求項 1 記載の文書書式識別装置。

【請求項 3】

前記識別手段は、

前記変倍率を前記類似度の計算に用いる文書書式データに乗算することを特徴とする請求項 2 記載の文書書式識別装置。

【請求項 4】

前記相似情報抽出手段は、

前記識別すべき文書画像から抽出されたブロックの個数と前記マスター文書画像のブロックの個数とが所定数以上でかつ互いに等しい場合に、

前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得た X 座標値を所定の順序に基づいて整列させることにより得た X 座標値列（この座標値列を  $Y_i$  とし、その平均値を  $Y_{ave}$  とする）と、前記マスター文書画像から抽出された複数のブロックそれぞれから得た X 座標値を所定の順序に基づいて整列させることにより得た X 座標値列（この座標値列を  $X_i$  とし、その平均値を  $X_{ave}$  とする）との間の相関係数を、

$$\frac{\{ (X_i - X_{ave}) (Y_i - Y_{ave}) \}}{\{ ( (X_i - X_{ave})^2 ) \times ( (Y_i - Y_{ave})^2 ) \}^{(1/2)}}$$

を用いて求め、

当該求めた相関係数が所定値より大きい場合に、前記識別すべき文書画像に関する X 座標値列と前記マスター文書画像に関する X 座標値列の傾き及びその切片に基づいて、X 座標方向の変倍率と X 座標方向の原点ずれ量とを求める一方、当該求めた相関係数が所定値以下の場合には相似関係がないと判断し、

更に、前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得た Y 座標値を所定の順序に基づいて整列させることにより得た Y 座標値列（この座標値列を  $Y_i'$  とし、その平均値を  $Y_{ave}'$  とする）と、前記マスター文書画像から抽出された複数のブロックそれぞれから得た Y 座標値を所定の順序に基づいて整列させることにより得た Y 座標値列（この座標値列を  $X_i'$  とし、その平均値を  $X_{ave}'$  とする）との間の

相関係数を、

$$\frac{\{ (X_i' - X_{ave}') (Y_i' - Y_{ave}') \}}{\{ (X_i' - X_{ave}')^2 \} \times \{ (Y_i' - Y_{ave}')^2 \}^{(1/2)}}$$

を用いて求め、

当該求めた相関係数が所定値より大きい場合に、前記識別すべき文書画像に関する Y 座標値列と前記マスター文書画像に関する Y 座標値列の傾き及びその切片に基づいて、Y 座標方向の変倍率と Y 座標方向の原点ずれ量を求める一方、当該求めた相関係数が所定値以下の場合には相似関係がないと判断し、

前記 X 座標方向の変倍率と前記 Y 座標方向の変倍率の両方が求められた場合に相似関係があると判断し、

前記識別手段は、

前記原点ずれ量を前記類似度の計算に用いる文書書式データに加算することを特徴とする請求項 3 記載の文書書式識別装置。

【請求項 5】

前記相似情報は、前記変倍率に基づいて算出したペナルティを含み、前記識別手段は、該ペナルティを前記類似度の計算に用いることを特徴とする請求項 2 記載の文書書式識別装置。

【請求項 6】

前記相似情報抽出手段は、

前記変倍率の適正を判定する判定手段を更に備え、

前記識別手段は、

前記判定手段にて変倍率が不適正であると判定した場合、前記相似情報を用いずに前記類似度計算を行うことを特徴とする請求項 1 記載の文書書式識別装置。

【請求項 7】

前記相似情報抽出手段は、

前記識別すべき文書画像から抽出されたブロックの個数と前記マスター文書画像から抽出されたブロックの個数とが等しくない場合、相似関係がないと判断することを特徴とする請求項 1 に記載の文書書式識別装置。

【請求項 8】

前記相似情報抽出手段は、

前記識別すべき文書画像から抽出されたブロックの個数と前記マスター文書画像から抽出されたブロックの個数とが前記所定数より小さくかつ互いに等しい場合、

前記識別すべき文書画像の幅と高さ、ならびに前記マスター文書画像の幅と高さに基づいて、X 座標方向の変倍率と Y 座標方向の変倍率とを求め、

前記求めた X 座標方向の変倍率と Y 座標方向の変倍率、前記識別すべき文書画像から抽出されたブロックの幅と高さ、ならびに前記マスター文書画像から抽出されたブロックの幅と高さが、所定の関係式をみたすか否かを判断し、

当該所定の関係式を満たさない場合に相似関係がないと判断し、当該所定の関係式を満たす場合には相似関係があると判断することを特徴とする請求項 1 記載の文書書式識別装置。

【請求項 9】

前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得た X 座標値は、各ブロックの左上角の X 座標値であり、

前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得た Y 座標値は、各ブロックの左上角の Y 座標値であることを特徴とする請求項 1 記載の文書書式識別装置。

【請求項 10】

文書画像の文書書式を識別する文書書式識別方法であって、

作成手段が、文書書式を識別すべき文書画像から抽出された複数のブロックそれぞれの位置座標を含む文書書式データを作成する作成工程と、

10

20

30

40

50

相似情報抽出手段が、前記作成工程で作成された前記識別すべき文書画像の文書書式データと、保存手段に保存されているマスター文書画像の文書書式データとを比較することにより、相似関係があるか否かを判断するとともに、当該判断結果と、相似関係があると判断した前記識別すべき文書画像と前記マスター文書画像との間の変倍率と、を含む相似情報を抽出する相似情報抽出工程と、

識別手段が、前記相似情報抽出工程にて抽出した相似情報及び前記文書書式データに基づいて、前記識別すべき文書画像の前記マスター文書画像に対する類似度を計算することにより、前記識別すべき文書画像の文書書式を識別する識別工程と、を備え、

前記相似情報抽出工程では、

前記識別すべき文書画像から抽出されたブロックの個数と前記マスター文書画像のブロックの個数とが所定数以上でかつ互いに等しい場合に、

前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得たX座標値を所定の順序に基づいて整列させることにより得たX座標値列（この座標値列をY<sub>i</sub>とし、その平均値をY<sub>ave</sub>とする）と、前記マスター文書画像から抽出された複数のブロックそれぞれから得たX座標値を所定の順序に基づいて整列させることにより得たX座標値列（この座標値列をX<sub>i</sub>とし、その平均値をX<sub>ave</sub>とする）との間の相関係数を、  

$$\frac{\{ (X_i - X_{ave}) (Y_i - Y_{ave}) \}}{\{ ( (X_i - X_{ave})^2 ) \times ( (Y_i - Y_{ave})^2 ) \}^{(1/2)}}$$

を用いて求め、

当該求めた相関係数が所定値より大きい場合に、前記識別すべき文書画像に関するX座標値列と前記マスター文書画像に関するX座標値列の傾きに基づいてX座標方向の変倍率を求める一方、当該求めた相関係数が所定値以下の場合には相似関係がないと判断し、

更に、前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得たY座標値を所定の順序に基づいて整列させることにより得たY座標値列（この座標値列をY<sub>i</sub>'とし、その平均値をY<sub>ave</sub>'とする）と、前記マスター文書画像から抽出された複数のブロックそれぞれから得たY座標値を所定の順序に基づいて整列させることにより得たY座標値列（この座標値列をX<sub>i</sub>'とし、その平均値をX<sub>ave</sub>'とする）との間の相関係数を、

$$\frac{\{ (X_i' - X_{ave}') (Y_i' - Y_{ave}') \}}{\{ ( (X_i' - X_{ave}')^2 ) \times ( (Y_i' - Y_{ave}')^2 ) \}^{(1/2)}}$$

を用いて求め、

当該求めた相関係数が所定値より大きい場合に、前記識別すべき文書画像に関するY座標値列と前記マスター文書画像に関するY座標値列の傾きに基づいてY座標方向の変倍率を求める一方、当該求めた相関係数が所定値以下の場合には相似関係がないと判断し、

前記X座標方向の変倍率と前記Y座標方向の変倍率の両方が求められた場合に相似関係があると判断することを特徴とする文書書式識別方法。

【請求項11】

請求項10に記載の文書書式識別方法をコンピュータによって実行させるための制御プログラムを格納する記憶媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、大量の帳票を処理する分野で、帳票の書式ごとに分類する装置を構築する際に、自動分類を可能にする帳票識別に関するものである。

【0002】

【従来の技術】

帳票内に記載された情報をOCR等の光学式文字認識装置で読みとるためには、帳票の書式を識別し、帳票内の情報記載位置を正確に把握する必要がある。帳票の書式を識別する方法として、あらかじめ登録されたマスタ帳票と、識別したい帳票の帳票内のテーブルや文字を比較し、識別したい帳票と一致するマスタ帳票を抽出する方法がある。

## 【 0 0 0 3 】

帳票内のテーブルや文字を比較するには、識別したい帳票のテーブルブロックおよびテキストブロックに、最も近い座標位置にあるマスタ帳票のブロックを検出し、テーブルブロックおよびテキストブロックごとの詳細情報のマッチングを取る手法が一般的である。そのブロック位置の検出には、帳票ページの左上角を原点として、各ブロックの左上角位置の座標値を使用していた。

## 【 0 0 0 4 】

## 【発明が解決しようとする課題】

しかしながら、識別したい帳票が F A X 等により送信された場合、F A X 等の給紙制約から、識別したい帳票が拡大または縮小されることがある。そして、図 2 の ( A )、( B ) に示すようにマスタ帳票 B に対して、拡大または縮小した識別したい帳票 A を比較すると、識別したい帳票の各ブロックの左上角位置は変倍されているため、識別したい帳票のブロックに対応するマスタ帳票のブロックを正確に検出することができない。また、ブロックごとの詳細構造のマッチングでも、テーブルブロックのサイズ、テーブル内の罫線的位置情報も、前記と同様に変倍されているので、相似形のテーブルブロックでもテーブル構造のマッチング計算では、異なるテーブルだと識別されることになる。その結果、拡大または縮小した帳票は、類似度が非常に低くなり、異なる帳票として判断されることになる。

## 【 0 0 0 5 】

本発明は、上記課題を鑑みてなされたものであり、異なる変倍率で拡大または縮小された複数の文書が混在する環境でも、文書書式を正しく識別することを目的とする。

## 【 0 0 0 6 】

## 【課題を解決するための手段】

かかる課題を解決するため、例えば本発明の文書書式識別装置は以下の構成を備える。すなわち、

文書画像の文書書式を識別する文書書式識別装置であって、

文書書式を識別すべき文書画像から抽出された複数のブロックそれぞれの位置座標を含む文書書式データを作成する作成手段と、

前記作成手段により作成された前記識別すべき文書画像の文書書式データと、保存手段に保存されているマスター文書画像の文書書式データとを比較することにより、相似関係があるか否かを判断するとともに、当該判断結果と、相似関係があると判断した前記識別すべき文書画像と前記マスター文書画像との間の変倍率と、を含む相似情報を抽出する相似情報抽出手段と、

前記相似情報抽出手段にて抽出した相似情報及び前記文書書式データに基づいて、前記識別すべき文書画像の前記マスター文書画像に対する類似度を計算することにより、前記識別すべき文書画像の文書書式を識別する識別手段と、を備え、

前記相似情報抽出手段は、

前記識別すべき文書画像から抽出されたブロックの個数と前記マスター文書画像のブロックの個数とが所定数以上でかつ互いに等しい場合に、

前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得た X 座標値を所定の順序に基づいて整列させることにより得た X 座標値列 (この座標値列を  $Y_i$  とし、その平均値を  $Y_{ave}$  とする) と、前記マスター文書画像から抽出された複数のブロックそれぞれから得た X 座標値を所定の順序に基づいて整列させることにより得た X 座標値列 (この座標値列を  $X_i$  とし、その平均値を  $X_{ave}$  とする) との間の相関係数を、  

$$\frac{\{ (X_i - X_{ave}) (Y_i - Y_{ave}) \}}{\{ ( (X_i - X_{ave})^2 ) \times ( (Y_i - Y_{ave})^2 ) \}^{(1/2)}}$$
を用いて求め、

当該求めた相関係数が所定値より大きい場合に、前記識別すべき文書画像に関する X 座標値列と前記マスター文書画像に関する X 座標値列の傾きに基づいて X 座標方向の変倍率を求める一方、当該求めた相関係数が所定値以下の場合には相似関係がないと判断し、

10

20

30

40

50

更に、前記識別すべき文書画像から抽出された前記複数のブロックそれぞれから得た Y 座標値を所定の順序に基づいて整列させることにより得た Y 座標値列 (この座標値列を  $Y_i'$  とし、その平均値を  $Y_{ave}'$  とする) と、前記マスター文書画像から抽出された複数のブロックそれぞれから得た Y 座標値を所定の順序に基づいて整列させることにより得た Y 座標値列 (この座標値列を  $X_i'$  とし、その平均値を  $X_{ave}'$  とする) との間の相関係数を、

$$\frac{\{ (X_i' - X_{ave}') (Y_i' - Y_{ave}') \}}{\{ (X_i' - X_{ave}')^2 \} \times \{ (Y_i' - Y_{ave}')^2 \}^{(1/2)}}$$

を用いて求め、

当該求めた相関係数が所定値より大きい場合に、前記識別すべき文書画像に関する Y 座標値列と前記マスター文書画像に関する Y 座標値列の傾きに基づいて Y 座標方向の変倍率を求める一方、当該求めた相関係数が所定値以下の場合には相似関係がないと判断し、

前記 X 座標方向の変倍率と前記 Y 座標方向の変倍率の両方が求められた場合に相似関係があると判断することとを特徴とする。

【0007】

【発明実施の形態】

[実施形態1]

以下、図面を参照して本発明の実施の形態を詳細に説明する。

【0008】

図1は、本発明の実施の形態に係る帳票書式自動識別装置の概略構成を示すブロック図である。

【0009】

11はスキャナーであり、帳票イメージを光学的に読み取り、帳票イメージデータを出力する。12はプロセッサでありメモリ15に格納された制御プログラム15dを実行することにより、画像特徴量抽出手段12a、書式データ作成手段12b、類似度計算手段12cとして機能する。スキャナー11で読み取った画像は、帳票イメージ15cとしてメモリ15に格納される。帳票イメージ15cは、2値化処理されて画像特徴量抽出手段12aに送られ、黒ドットのヒストグラム法などの手法により、テーブル、テキスト、ピクチャなどブロックごとに属性分類される。テーブルブロックについては、さらに罫線追跡手法などで、テーブルの詳細構造を求める。また、テキストブロックについては、さらに文字コードに変換するなどの処理を行う。

【0010】

このようにして取得した情報から、書式データ作成手段12bにて、図3に示す帳票のページ書式およびテーブル書式を作成し、メモリ15およびディスク14に保存する。図3は、画像特徴量抽出手段12aで抽出した後の帳票サンプル31を示す。テーブル・ブロック3個(311~313)、ピクチャ・ブロック1個(314)が抽出されている。この帳票の書式データ32は、ページ書式321とテーブル書式322に階層化して保存する。ページ書式321は、ヘッダ部321aに帳票ページ幅、帳票ページ高さを所有する。

【0011】

また、データ部322aにはブロック毎に各種情報が記憶される。例えば、ブロック属性がテーブルの場合には、位置情報としてブロック左端位置、ブロック上端位置、大きさ情報としてブロック幅、ブロック高さの情報を所有する。また、比較帳票をピックアップするのに使用するためのページ原点からの距離および類似度の計算に使用するための当該ブロックの面積を全テーブル・ブロックで割った値も所有する。さらに、テーブル詳細情報とリンクするためにテーブルIDを所有する。このテーブルIDにリンクしたテーブルのセルの詳細構造をテーブル書式322で示す。テーブル内のセル個数、セルの位置、大きさ情報を所有する。

【0012】

帳票書式識別装置は、キーボードから帳票の登録、帳票の識別などの命令が入力されると

10

20

30

40

50

、各々の命令に対応する処理をプロセッサ 12 が上記の書式データ 32 を使用して行う。そして、その識別結果をディスプレイ 16 に表示する。

【0013】

図4を参照して、本実施形態の帳票書式識別装置、特に図1のプロセッサ12が実行する各種制御処理の動作を説明する。

【0014】

図4は、帳票書式識別装置による書式識別処理の概略フローチャートを示す。ステップS101にて、識別したい帳票をスキャナーで読み取り、ステップS103で、テーブルブロック、テキストブロックの座標値等の特徴量データを抽出する。ステップS105では、これらの特徴量データを類似度を計算するための書式データに変換する。この書式データを元にして、ステップS107にてマスター帳票の中から、当該識別したい帳票の書式データと同じ帳票である可能性のある帳票を絞り込む。ステップS109で絞り込んだ帳票のすべてについて書式の類似度を計算する(ステップS111)。計算の結果、類似度の高い方の所定数のマスター帳票を類似帳票の候補とし、その識別コードと類似度を出力する(ステップS113)。

10

【0015】

本発明の帳票レイアウトの相似形チェックは、ステップS108にて行う。図5～図8を使用して、詳細にその処理を説明する。

【0016】

ステップS108で帳票レイアウト相似形チェック処理が開始され、ステップS203で識別したい帳票およびマスタ帳票のフレームおよびテーブルブロックを整列する。本実施形態では、ブロック左上角のX座標の昇順にブロック情報を並べる。但し、図3に示すような、ブロック左上角座標のX成分がほとんど同じ位置にある場合には、X座標の誤差により識別したい帳票とマスタ帳票のブロック情報の並びを必ずしも対応付けることができない。識別したい帳票のブロック情報がテーブル1(311) テーブル2(312) ピクチャブロック(314) テーブルブロック3(313)と整列しても、マスタ帳票のブロック情報がテーブル2(312) テーブル1(311) ピクチャブロック(314) テーブルブロック3(313)のように整列する可能性は十分にある。そこで、X座標の位置が5ピクセル以内の差のブロックに関しては、別途Y成分の昇順に整列させる。この結果、識別したい帳票およびマスタ帳票のブロック情報をテーブル1(311) テーブル2(312) ピクチャブロック(314) テーブルブロック3(313)の順に整列することが保証される。

20

30

【0017】

ステップS205では、比較するブロックの個数が同じであるかをチェックしている。ブロックの個数が異なる場合は、相似形でないと判断して、帳票レイアウト相似形チェック処理を終了し、ステップS109へ戻る。

【0018】

ステップS205でブロック個数が同じであった場合には、ステップS207にて、ブロック個数が3個以上かをチェックしている。3個以上の場合と、2個以下の場合では、相似形判定プログラムが異なるからである。

40

【0019】

3個以上の場合、ステップS208\_\_1に進み、ブロック左上角X座標の比較処理を行う。ステップS208\_\_1の詳細を図6のフローチャートを用いて説明する。

【0020】

ステップS209では、ブロック情報のX成分の相似形チェック処理をする。すなわち、識別したい帳票の左上角X座標を縦軸に、マスタ帳票の左上角X座標を横軸にして、下式より相関係数を計算する。

【0021】

【数1】

$$\text{相関係数} = \frac{\sum (X_i - X_{ave})(Y_i - Y_{ave})}{\sqrt{\sum (X_i - X_{ave})^2 \times \sum (Y_i - Y_{ave})^2}}$$

$X_i, Y_i$  : XおよびY座標値

$X_{ave}, Y_{ave}$  : XおよびY座標の平均値

10

#### 【0022】

ここで、相関係数の算出にあたっては、上式の相関係数の分子の値をチェックし（ステップS211）、相関係数の分子が14以下であれば、別プログラムで変倍率  $X$  を求めている。これは、以下に述べるように相関係数の誤差が大きくなるからである。

#### 【0023】

図3の帳票では、テーブルブロック1（311）、テーブルブロック2（312）、ピクチャブロック314の各左上角のX座標はほとんど同じ位置にあるため、 $X_i$ 、 $Y_i$ ともに $X_{ave}$ 、 $Y_{ave}$ に近い値となる。従って、 $X_i$ 、 $Y_i$ が誤差の為に変動した場合、相関係数の変動も大きくなる。このため、 $X_i$ 、 $Y_i$ の誤差が大きいと考えられる環境では、相関係数の信頼度が落ちる。誤差の要因としては、スキャナで読み取るときに発生する誤差、傾斜補正等の画像処理を行ったときに発生する誤差、マッチング（すべての画像を100dpiに正規化してマッチングを行っている）の為に解像度変換を行ったときに発生する誤差などがあげられる。それらの誤差を考慮して、上式の相関係数の分子が14以下の場合には、信頼できないとして（識別したい帳票のページ幅）/（マスタ帳票のページ幅）= 変倍率  $X$  としている（ステップS211\_2）。ただし、帳票をはさみ等で切り取った場合にできるズレによる誤差を考慮して、帳票ページ幅の差分が10ピクセル以内であれば、変倍率  $X$  を1.0とする（ステップS211\_1、S211\_3）。

20

#### 【0024】

ステップS211にて相関係数の分子が14以上であり、かつ相関係数が0.9996以上あれば、X成分に関しては相似形と判断する（ステップS213）。

30

#### 【0025】

一方、ステップS211にて相関係数の分子が14以上であっても、相関係数が0.9996より小さい場合には、相似形でないと判断して相似形チェック処理を終了し、ステップS109へもどる。

#### 【0026】

相関係数が1に近いほど、前記の情報が直線上に並んでいるとみなすことができる。つまり、変倍されている可能性はあるが、X成分の並びは相似形であると考ええる。

#### 【0027】

次に、これらのデータが直線上に並んでいるので、その傾きを計算する。この傾きがX成分の変倍率  $X$  を示す（ステップS215）。傾き=1の時は、識別したい帳票とマスタ帳票のX成分は同じであり、傾きが1より小さいと、マスタ帳票のX成分の増加率が識別したい帳票の比べて大きいことになり、マスタ帳票の幅が識別したい帳票の幅に比べて拡大しているといえる。つまり、傾き=（識別したい帳票のページ幅）/（マスタ帳票のページ幅）の関係が成り立つ。

40

#### 【0028】

ステップS217では、前記の誤差を考慮して、変倍率  $X$  が $1 \pm 0.028$ 以内であれば変倍がないとみなし、変倍率  $X$  を1にリセットする。

#### 【0029】

次に、ステップS219でこの変倍率  $X$  が、テーブルブロックのサイズまで適用できる

50



かをチェックしている。すなわち、ステップS 2 0 3で整列した帳票のブロック情報を順に1個ずつ、「(識別したい帳票のブロック幅)/(マスタ帳票のブロック幅)<変倍率  $X + 0.027 + 1.9 / (\text{マスタ帳票のブロック幅})$ 」を満たすかどうかをチェックし、満たさない場合には、ブロックレイアウトは相似でないとして当該処理を終了し、ステップS 1 0 9にもどる。また、満たす場合には、変倍率  $X$ のテーブルブロックサイズへの適用可と判断し、処理を終了する。

【0030】

ステップS 2 0 8 \_\_ 1と同様に、ステップS 2 0 8 \_\_ 2ではブロック左上角Y座標の比較処理を行う。ステップS 2 0 8 \_\_ 2の詳細フローチャートを図7に示す。

【0031】

ステップS 2 1 0では、ブロック情報のY成分の相似形チェック開始する。すなわち、識別したい帳票の左上角Y座標を縦軸に、マスタ帳票の左上角Y座標を横軸にして、相関係数を計算する。

【0032】

ここで、相関係数の算出にあたっては、相関係数の分子の値をチェックし(ステップS 2 1 2)、相関係数の分子が14以下であれば、別プログラムで変倍率  $Y$ を求めている。相関係数の分子が14以下の場合、信頼できないとして(識別したい帳票のページ高さ)/(マスタ帳票のページ高さ)=変倍率  $Y$ としている(ステップS 2 1 2 \_\_ 2)。ただし、帳票をはさみ等で切り取った場合にできるズレによる誤差を考慮して、帳票ページ幅の差分が10ピクセル以内であれば、変倍率  $Y$ を1.0とする(ステップS 2 1 2 \_\_ 1、S 2 1 2 \_\_ 3)。

【0033】

ステップS 2 1 2にて相関係数の分子が14以上であり、かつ相関係数が0.9996以上あれば、Y成分に関しては相似形と判断する(ステップS 2 1 4)。

【0034】

一方、ステップS 2 1 2にて相関係数の分子が14以上であっても、相関係数が0.9996より小さい場合には、相似形でないとして判断して相似形チェック処理を終了し、ステップS 1 0 9へもどる。

【0035】

相関係数が1に近いほど、前記の情報が直線上に並んでいるとみなすことができる。つまり、変倍されている可能性はあるが、Y成分の並びは相似形であると考える。

【0036】

次に、これらのデータが直線上に並んでいるので、その傾きを計算する。この傾きがY成分の変倍率  $Y$ を示す(ステップS 2 1 6)。傾き=1の時は、識別したい帳票とマスタ帳票のY成分は同じであり、傾きが1より小さいと、マスタ帳票のY成分の増加率が識別したい帳票の比べて大きいことになり、マスタ帳票の幅が識別したい帳票の幅に比べて拡大しているといえる。つまり、傾き=(識別したい帳票のページ高さ)/(マスタ帳票のページ高さ)の関係が成り立つ。

【0037】

ステップS 2 1 8では、誤差を考慮して、変倍率  $Y$ が $1 \pm 0.028$ 以内であれば変倍がないとみなし、変倍率  $Y$ を1にリセットする。

【0038】

次に、ステップS 2 2 0でこの変倍率  $Y$ が、テーブルブロックのサイズまで適用できるかをチェックしている。すなわち、ステップS 2 0 3で整列した帳票のブロック情報を順に1個ずつ、「(識別したい帳票のブロック高さ)/(マスタ帳票のブロック高さ)<変倍率  $Y + 0.027 + 1.9 / (\text{マスタ帳票のブロック高さ})$ 」を満たすかどうかをチェックし、満たさない場合には、ブロックレイアウトは相似でないとして当該処理を終了し、ステップS 1 0 9にもどる。また、満たす場合には、変倍率  $Y$ のテーブルブロックサイズへの適用可と判断し、処理を終了する。

【0039】

10

20

30

40

50

ステップS 2 2 1では、X成分、Y成分両方ともに相似形であり、変倍率( X、 Y )を取得できた帳票のみ、レイアウトが相似形であると判断し、それ以外の場合には、相似でないとして当該処理を終了し、ステップS 1 0 9にもどる。

【 0 0 4 0 】

さて、ステップS 2 0 7でブロックの個数が2個以下の場合は、図8のブロック個数が2個以下の場合の処理を行う(ステップS 3 0 1)。

【 0 0 4 1 】

すなわち、ステップS 3 0 1でブロック個数が2個以下の場合の処理が開始され、ステップS 3 0 2にて識別したい帳票のページ幅とマスター帳票のページ幅を比較し、差分が10ピクセルより大きい場合には、変倍率  $X = (\text{識別したい帳票のページ幅}) / (\text{マスター帳票のページ幅})$  とし、差分が10ピクセル以内の場合には変倍率  $X = 1.0$  とする(ステップS 3 0 2、S 3 0 3、S 3 0 4)。

10

【 0 0 4 2 】

さらにステップS 3 0 5にてこの変倍率 Xが、テーブルブロックのサイズまで適用できるかをチェックしている。すなわち、ステップS 2 0 3で整列した帳票のブロック情報を順に1個ずつ、「 $(\text{識別したい帳票のブロック幅}) / (\text{マスター帳票のブロック幅}) < \text{変倍率 } X + 0.027 + 1.9 / (\text{マスター帳票のブロック幅})$ 」を満たすかどうかをチェックし、満たさない場合には、ブロックレイアウトは相似でないとして当該処理を終了し、ステップS 1 0 9にもどる。また、満たす場合には、変倍率 Xのテーブルブロックサイズへの適用可と判断する。

20

【 0 0 4 3 】

同様に、ステップS 3 0 6にて識別したい帳票のページ高さとのマスター帳票のページ高さを比較し、差分が10ピクセルより大きい場合には、変倍率  $Y = (\text{識別したい帳票のページ高さ}) / (\text{マスター帳票のページ高さ})$  とし、差分が10ピクセル以内の場合には変倍率  $Y = 1.0$  とする(ステップS 3 0 6、S 3 0 7、S 3 0 8)。

【 0 0 4 4 】

さらにステップS 3 0 9にてこの変倍率 Yが、テーブルブロックのサイズまで適用できるかをチェックしている。すなわち、ステップS 2 0 3で整列した帳票のブロック情報を順に1個ずつ、「 $(\text{識別したい帳票のブロック高さ}) / (\text{マスター帳票のブロック高さ}) < \text{変倍率 } Y + 0.027 + 1.9 / (\text{マスター帳票のブロック高さ})$ 」を満たすかどうかをチェックし、満たさない場合には、ブロックレイアウトは相似形でないとして当該処理を終了し、ステップS 1 0 9にもどる。また、満たす場合には、変倍率 Yのテーブルブロックサイズへの適用可と判断する。

30

【 0 0 4 5 】

ステップS 3 0 5およびステップS 3 0 9にてテーブルブロックサイズへの適用可能と判断された変倍率( X、 Y )について、ステップS 3 1 0で、ブロックの左上角、「マスター帳票のX座標×変倍率( X ) 識別したい帳票のX座標±10、かつマスター帳票のY座標×変倍率( Y ) 識別したい帳票のY座標±10」の条件式を満たしていれば、識別したい帳票とマスター帳票のレイアウトは相似形であると判断し、帳票レイアウト相似形チェック処理のステップS 2 2 3へ進む(ステップS 3 1 1)。また、条件を満たさない場合には、相似形でないとして当該処理を終了し、ステップS 1 0 9にもどる。

40

【 0 0 4 6 】

ステップS 2 2 1またはステップS 3 1 1にて相似形であると判断された場合には、ステップS 2 2 3でページレイアウトの変倍によるペナルティを以下の式で決定する。

【 0 0 4 7 】

【数2】

$$PX = \sin\left(\frac{\pi}{4} - \tan^{-1} \text{変倍率} \delta X\right) \times 60.0$$

$$PY = \sin\left(\frac{\pi}{4} - \tan^{-1} \text{変倍率} \delta Y\right) \times 60.0$$

$$PXY = \sin(\tan^{-1} \text{変倍率} \delta X - \tan^{-1} \text{変倍率} \delta Y) \times 120.0$$

#### 【 0 0 4 8 】

10

P X、P Yは各成分の変倍によるペナルティ、P X Yは、両成分の変形度によるペナルティをあらわす。

#### 【 0 0 4 9 】

変倍率 X、Yが1ならば、P X、P Yは0である。つまり、変倍していないのでペナルティを課せないことを意味する。

#### 【 0 0 5 0 】

P X Yは、X、Y成分が均等に変倍したときには0になるし、X成分が1より大きく、Y成分が1より小さく変倍するように、変倍によるレイアウトの変形が大きくなればペナルティが大きくなるように調整する式である。

#### 【 0 0 5 1 】

20

レイアウトが相似形である場合は、以上の計算式でペナルティを与えて、従来のページ書のマッチングによるペナルティを0とする。逆に、相似形でない場合は、従来通りのページ書のマッチングによるペナルティを与える。

#### 【 0 0 5 2 】

ページ書のマッチングの次に、テーブルブロックの詳細構造、その次にテキストブロックの文字比較を行うが、これらのブロックを検出する際には、相似形チェックで求めた変倍率を使用した計算式を使用する。

#### 【 0 0 5 3 】

例えば、図2の(A)、(B)に示すような識別したい帳票とマスタ帳票について、本実施形態による帳票レイアウト相似形チェック処理で、相似形だと判定され、変倍率(X、Y)が得られたとする。

30

#### 【 0 0 5 4 】

マスタ帳票の(X1、Y1)のブロックに対応する識別したい帳票のブロックは、(X×X1、Y×Y1)で正確な位置を求めることができる。

#### 【 0 0 5 5 】

この計算式で検出したブロックがテーブルブロックの場合は、テーブルの各罫線情報が帳票ページの変倍率と同様に変倍されているので、マスタ帳票の罫線情報(Lx、Ly)を(Lx×X、Ly×Y)に変倍して、識別したい帳票の罫線情報と比較することで、正確なテーブルブロックの詳細構造のマッチングを行うことができる。

#### 【 0 0 5 6 】

40

以上、記述した中での数値は、数多くの帳票サンプルを使用した統計値であり、帳票識別の環境によっては、変更してもかまわない。

#### 【 0 0 5 7 】

#### [ 実施形態2 ]

原点ずれが生じると、変倍のみでは正しく認識できない恐れがある。

#### 【 0 0 5 8 】

以下に図面を参照して本発明の実施形態のうち、識別したい帳票とマスタ帳票の原点位置がずれた場合の識別処理について詳細を説明する。

#### 【 0 0 5 9 】

なお、帳票書式識別装置は図1と同様のものを使用し、図3と同様の書式データを作成す

50

る。したがって、図1と図3の内容は実施形態1と重複することから説明は省略する。

【0060】

本実施形態の帳票書式識別装置、特に図1のプロセッサ12が実行する各種制御処理のうち、実施形態1と異なる処理を中心に説明する。

【0061】

本実施形態の帳票レイアウトの相似形チェックは、ステップS108にて行う。図10～図13を使用して、詳細にその処理を説明する。

【0062】

図10のステップS403からS407までは、図5のステップS203からS207までと同じ処理を行う。

10

【0063】

ステップS408\_\_1ではブロック左上角X座標の比較処理を行う。この処理の詳細を図11を参照して説明する。

【0064】

すなわち、ステップS409では、ブロック情報のX成分の相似形チェックをすべく、識別したい帳票の左上角X座標を縦軸に、マスタ帳票の左上角X座標を横軸にして、相関係数を計算する。

【0065】

ここで、相関係数の算出にあたっては、相関係数の分子の値をチェックし（ステップS411）、相関係数の分子が14以下であれば、別プログラムで変倍率Xを求めている。これは、相関係数の誤差が大きくなるからで、詳細は実施形態1と同じであるため、説明は省略する。

20

【0066】

相関係数の分子が14以下の場合は、信頼できないとして（識別したい帳票のページ幅）／（マスタ帳票のページ幅）＝変倍率Xとしている（ステップS411\_\_2）。ただし、帳票をはさみ等で切り取った場合にできるズレによる誤差を考慮して、帳票ページ幅の差分が10ピクセル以内であれば、変倍率Xを1.0とする（ステップS411\_\_1、S411\_\_3）。

【0067】

上記でもとめた変倍率Xを用いて、ステップS411\_\_4にて、原点ずれ量shiftXを「（識別したい帳票の先頭ブロックの左上角X座標）－（マスタ帳票の先頭ブロックの左上角X座標）×変倍率X」より算出する。

30

【0068】

ステップS411にて相関係数の分子が14以上であり、かつ相関係数が0.9996以上あれば、X成分に関しては相似形と判断する（ステップS413）。

【0069】

一方、ステップS411にて相関係数の分子が14以上であっても、相関係数が0.9996より小さい場合には、相似形でないと判断して相似形チェック処理を終了し、ステップS109へもどる。

【0070】

相関係数が1に近いほど、前記の情報が直線上に並んでいるとみなすことができる。つまり、変倍されている可能性はあるが、X成分の並びは相似形であると考える。

40

【0071】

次に、これらのデータが直線上に並んでいるので、その傾きを計算する（回帰直線の傾き）。この傾きがX成分の変倍率Xを示す（ステップS415）。傾き＝1の時は、識別したい帳票とマスタ帳票のX成分は同じであり、傾きが1より小さいと、マスタ帳票のX成分の増加率が識別したい帳票の比べて大きいことになり、マスタ帳票の幅が識別したい帳票の幅に比べて拡大しているといえる。一方、回帰直線の縦軸との切片がX座標の原点ずれ量shiftXになる（ステップS417）。

【0072】

50

ステップS 4 1 9でこの変倍率  $X$  が、テーブルブロックのサイズまで適用できるかをチェックしている。すなわち、ステップS 2 0 3で整列した帳票のブロック情報を順に1個ずつ、「(識別したい帳票のブロック幅)/(マスタ帳票のブロック幅)=変倍率  $X$ 」を満たすかどうかをチェックし、満たさない場合には、ブロックレイアウトは相似でないとして当該処理を終了し、ステップS 1 0 9にもどる。また、満たす場合には、変倍率  $X$  のテーブルブロックサイズへの適用可と判断し、処理を終了する。

【0073】

ステップS 4 0 8 \_\_ 1と同様に、ステップS 4 0 8 \_\_ 2では、ブロック左上角Y座標の比較処理を行う。ステップS 4 0 8 \_\_ 2の詳細フローチャートを図12に示す。

【0074】

すなわち、ステップS 4 1 0では、ブロック情報のX成分の相似形チェックをすべく、識別したい帳票の左上角Y座標を縦軸に、マスタ帳票の左上角Y座標を横軸にして、相関係数を計算する。

【0075】

ここで、相関係数の算出にあたっては、相関係数の分子の値をチェックし(ステップS 4 1 2)、相関係数の分子が14以下であれば、別プログラムで変倍率  $Y$  を求めている。これは、相関係数の誤差が大きくなるからで、詳細は実施形態1と同じであるため、説明は省略する。

【0076】

相関係数の分子が14以下の場合は、信頼できないとして(識別したい帳票のページ高さ)/(マスタ帳票のページ高さ)=変倍率  $Y$  としている(ステップS 4 1 2 \_\_ 2)。ただし、帳票をはさみ等で切り取った場合にできるズレによる誤差を考慮して、帳票ページ幅の差分が10ピクセル以内であれば、変倍率  $Y$  を1とする(ステップS 4 1 2 \_\_ 1、S 4 1 2 \_\_ 3)。

【0077】

上記でもとめた変倍率  $Y$  を用いて、ステップS 4 1 2 \_\_ 4にて、原点ずれ量  $shift Y$  を「(識別したい帳票の先頭ブロックの左上角Y座標)-(マスタ帳票の先頭ブロックの左上角Y座標)×変倍率  $Y$ 」より算出する。

【0078】

ステップS 4 1 2にて相関係数の分子が14以上であり、かつ相関係数が0.9996以上あれば、Y成分に関しては相似形と判断する(ステップS 4 1 4)。

【0079】

一方、ステップS 4 1 2にて相関係数の分子が14以上であっても、相関係数が0.9996より小さい場合には、相似形でないと判断して相似形チェック処理を終了し、ステップS 1 0 9へもどる。

【0080】

相関係数が1に近いほど、前記の情報が直線上に並んでいるとみなすことができる。つまり、変倍されている可能性はあるが、Y成分の並びは相似形であると考ええる。

【0081】

次に、これらのデータが直線上に並んでいるので、その傾きを計算する(回帰直線の傾き)。この傾きがY成分の変倍率  $Y$  を示す(ステップS 4 1 6)。傾き=1の時は、識別したい帳票とマスタ帳票のY成分は同じであり、傾きが1より小さいと、マスタ帳票のY成分の増加率が識別したい帳票の比べて大きいことになり、マスタ帳票の高さが識別したい帳票の高さに比べて拡大しているといえる。一方、回帰直線の縦軸との切片がY座標の原点ずれ量  $shift Y$  になる(ステップS 4 1 8)。

【0082】

ステップS 4 2 0でこの変倍率  $Y$  が、テーブルブロックのサイズまで適用できるかをチェックしている。すなわち、ステップS 2 0 3で整列した帳票のブロック情報を順に1個ずつ、「(識別したい帳票のブロック高さ)/(マスタ帳票のブロック高さ)=変倍率  $Y$ 」を満たすかどうかをチェックし、満たさない場合には、ブロックレイアウトは相似で

10

20

30

40

50

ないとして当該処理を終了し、ステップS 1 0 9にもどる。また、満たす場合には、変倍率 Yのテーブルブロックサイズへの適用可と判断し、処理を終了する。

【0083】

ステップS 4 2 1では、X成分、Y成分両方ともに相似形であり、変倍率( X、 Y)を取得できた帳票のみ、レイアウトが相似形であると判断し、それ以外の場合には相似でないとして、当該処理を終了し、ステップS 1 0 9にもどる。

【0084】

さて、ステップS 4 0 7でブロックの個数が2個以下の場合は、図13のブロック個数が2個以下の場合の処理を行う(ステップS 5 0 1)。

【0085】

すなわち、ステップS 5 0 1でブロック個数が2個以下の場合の処理が開始され、ステップS 5 0 2にて識別したい帳票のページ幅とマスター帳票のページ幅を比較し、差分が10ピクセルより大きい場合には、変倍率  $X = (\text{識別したい帳票のページ幅}) / (\text{マスター帳票のページ幅})$  とし、差分が10ピクセル以内の場合には変倍率  $X = 1.0$  とする(ステップS 5 0 2、S 5 0 3、S 5 0 4)。

【0086】

さらにステップS 5 0 5にてこの変倍率 Xが、テーブルブロックのサイズまで適用できるかをチェックしている。すなわち、ステップS 4 0 3で整列した帳票のブロック情報を順に1個ずつ、「(識別したい帳票のブロック幅) / (マスター帳票のブロック幅) = 変倍率 X」を満たすかどうかをチェックし、満たさない場合には、ブロックレイアウトは相似でないとして当該処理を終了し、ステップS 1 0 9にもどる。また、満たす場合には、変倍率 Xのテーブルブロックサイズへの適用可と判断する。

【0087】

同様に、ステップS 5 0 7にて識別したい帳票のページ高さとのマスター帳票のページ高さを比較し、差分が10ピクセルより大きい場合には、変倍率  $Y = (\text{識別したい帳票のページ高さ}) / (\text{マスター帳票のページ高さ})$  とし、差分が10ピクセル以内の場合には変倍率  $Y = 1.0$  とする(ステップS 5 0 7、S 5 0 8、S 5 0 9)。

【0088】

さらにステップS 5 1 0にてこの変倍率 Yが、テーブルブロックのサイズまで適用できるかをチェックしている。すなわち、ステップS 4 0 3で整列した帳票のブロック情報を順に1個ずつ、「(識別したい帳票のブロック高さ) / (マスター帳票のブロック高さ) = 変倍率 Y」を満たすかどうかをチェックし、満たさない場合には、ブロックレイアウトは相似形でないとして当該処理を終了し、ステップS 1 0 9にもどる。また、満たす場合には、変倍率 Yのテーブルブロックサイズへの適用可と判断する。

【0089】

ステップS 5 0 6およびステップS 5 1 1にてテーブルブロックサイズへの適用可能と判断された変倍率( X、 Y)について、ステップS 5 1 2で、ブロックの左上角が、「マスター帳票のX座標×変倍率( X) + 原点ずれ量 shift X 識別したい帳票のX座標 ± 10、かつマスター帳票のY座標×変倍率( Y) + 原点ずれ量 shift Y 識別したい帳票のY座標 ± 10」の条件式を満たしていれば、識別したい帳票とマスター帳票のレイアウトは相似形であると判断し、帳票レイアウト相似形チェック処理のステップS 4 2 2へ進む(ステップS 5 1 3)。また、条件を満たさない場合には、相似形でないとして当該処理を終了し、ステップS 1 0 9にもどる。

【0090】

ステップS 4 2 1またはステップS 5 1 3にて相似形であると判断された場合には、ステップS 4 2 2でページレイアウトの変倍によるペナルティを以下の式で決定する。

【0091】

【数3】

10

20

30

40

$$PX = \sin\left(\frac{\pi}{4} - \tan^{-1} \text{変倍率} \delta X\right) \times 60.0$$

$$PY = \sin\left(\frac{\pi}{4} - \tan^{-1} \text{変倍率} \delta Y\right) \times 60.0$$

$$PXY = \sin(\tan^{-1} \text{変倍率} \delta X - \tan^{-1} \text{変倍率} \delta Y) \times 120.0$$

PX、PYは各成分の変倍によるペナルティ、PXYは、両成分の変形度によるペナルティをあらわす。 10

【0092】

変倍率 X、Yが1ならば、PX、PYは0である。つまり、変倍していないのでペナルティを課せないことを意味する。

【0093】

PXYは、X、Y成分が均等に変倍したときには0になるし、X成分が1より大きく、Y成分が1より小さく変倍するように、変倍によるレイアウトの変形が大きくなればペナルティが大きくなるように調整する式である。

【0094】

また、原点ずれ量によるペナルティは、 $PX1 = \text{原点ずれ量}(\text{shift} X) \times 0.22$ 、 $PY1 = \text{原点ずれ量}(\text{shift} Y) \times 0.22$ よりもとめる。 20

【0095】

レイアウトが相似形である場合は、以上の計算式でペナルティを与えて、従来のページ書式のマッチングによるペナルティを0とする。逆に、相似形でない場合は、従来通りのページ書式のマッチングによるペナルティを与える。

【0096】

ページ書式のマッチングの次に、テーブルブロックの詳細構造、その次にテキスト・ブロックの文字比較を行うが、これらのブロックを検出する際には、相似形チェックで求めた変倍率を使用した計算式を使用する。

【0097】

例えば、図9のような帳票AとBがあり、帳票Aが識別したい帳票、帳票Bがマスタ帳票と仮定する。 30

【0098】

本実施形態による帳票レイアウト相似形チェック処理で、相似形だと判定され、変倍率(X、Y)、原点ずれ量(shift X、shift Y)が得られたとする。

【0099】

帳票Bの(X1、Y1)のブロックに対応する帳票Aのブロックは、 $(X \times X1 + \text{shift} X, Y \times Y1 + \text{shift} Y)$ で正確な位置を求めることができる。

【0100】

この計算式で検出したブロックがテーブルブロックの場合は、テーブルの各罫線情報が帳票ページの変倍率と同様に変倍されているので、帳票Bの罫線情報(Lx、Ly)を $(Lx \times X, Ly \times Y)$ に変倍して、帳票Aの罫線情報と比較することで、正確なテーブルブロックの詳細構造のマッチングを行うことができる。罫線情報は、テーブルブロックの左上角を原点にしているので、帳票ページ原点ずれ量は、テーブルブロックの詳細構造には影響を与えない。 40

【0101】

なお、上述した中での数値は、数多くの帳票サンプルを使用した統計値であり、帳票識別の環境によっては、変更してもかまわない。

【0102】

[実施形態3]

図14に示すように、レイアウト構造が変倍されている場合には、その帳票内のテーブルの罫線情報も同じ率で変倍されている。従って、テーブルブロックの詳細構造である罫線情報の比較に、この変倍率を使用することで、より正確な詳細構造の比較を行うことができる。

#### 【0103】

罫線情報を $L_i$ とすれば、識別したい帳票のテーブルブロックの詳細構造 $L_i = (\text{マスタ帳票のテーブルの詳細構造 } L_i) \times \text{変倍率} (X, Y)$ の関係が成り立つ。

#### 【0104】

テキストブロックの位置の検出する際に、変倍率 $(X, Y)$ を使用することで、比較すべき文字列を正確に知ることができる。しかし、文字の比較は、単なる文字コードの照合だから、変倍率 $(X, Y)$ は不要ではあるが、ペナルティ要素として、文字の大きさを取り入れている場合には、マスタ帳票の文字の大きさに変倍率をかけることで、より正確なマッチングをおこなえる。

#### 【0105】

##### [実施形態4]

識別したい帳票の一部分だけ文字認識を行うために、帳票認識を利用する方法がある。図14(A)、(B)に示すように、マスタ帳票に文字認識を行う領域をあらかじめ設定しているとすると、図14(A)、(B)では、網掛け部分の銀行口座に登録している氏名欄が文字認識する領域である。

#### 【0106】

まず、識別したい帳票を帳票認識することで、マスタ帳票のIDを取得できる。そのIDには、文字認識する領域が対応づけられている。本実施形態では、帳票の $X$ 、 $Y$ 方向の変倍率と帳票ページ原点ずれ量をIDとともに出力することができるので、識別したい帳票の文字認識する領域は、下記の式から修正することができる。

#### 【0107】

IDに対応つけてマスタに登録している文字認識領域を左上角座標 $(X, Y)$ 、幅 $W$ 、高さ $H$ とする。

#### 【0108】

受け取った変倍率が $(X, Y)$ 、原点ずれ量 $(\text{shift } X, \text{shift } Y)$ であれば、識別したい帳票の文字認識領域は、左上角座標 $(X \times X + \text{shift } X, Y \times Y + \text{shift } Y)$ 、幅 $(W \times X)$ 、高さ $(H \times Y)$ となる。

#### 【0109】

##### 【他の実施形態】

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

#### 【0110】

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

#### 【0111】

プログラムコードを供給するための記憶媒体としては、例えば、フロッピディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、磁気テープ、不揮発性のメモリカード、ROMなどを用いることができる。

#### 【0112】

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

10

20

30

40

50



うまでもない。

【 0 1 1 3 】

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

【 0 1 1 4 】

【 発明の効果 】

以上説明したように、本発明によれば、異なる変倍率で拡大または縮小された複数の文書が混在する環境でも文書書式を正しく識別することができる。

10

【 図面の簡単な説明 】

【 図 1 】 本発明の実施の一形態に係わる帳票識別装置の概略構成を示すブロック図である。

【 図 2 】 本発明の相似形帳票の一例を示す図である。

【 図 3 】 本発明のマッチング対象となる帳票の一例を示す図である。

【 図 4 】 本発明の処理の概要を示すフローチャートである。

【 図 5 】 本発明の相似形チェック処理を示すフローチャートである。

【 図 6 】 本発明の相似形チェック処理で、ブロック左上角 X 座標の比較処理を示すフローチャートである。

20

【 図 7 】 本発明の相似形チェック処理で、ブロック左上角 Y 座標の比較処理を示すフローチャートである。

【 図 8 】 本発明の相似形チェック処理で、ブロック個数が 2 個以下の場合の処理を示すフローチャートである。

【 図 9 】 本発明の相似形帳票の一例を示す図である。

【 図 10 】 本発明の相似形チェック処理を示すフローチャートである。

【 図 11 】 本発明の相似形チェック処理で、ブロック左上角 X 座標の比較処理を示すフローチャートである。

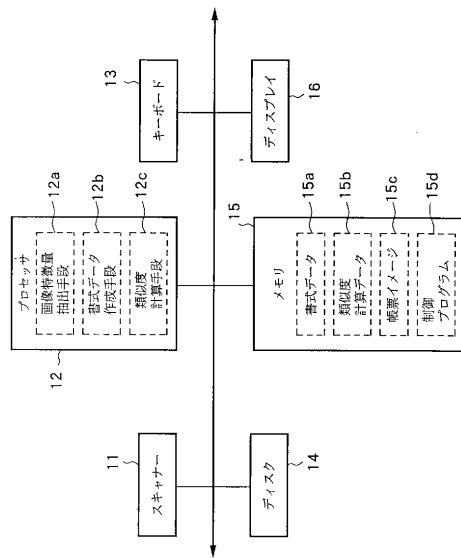
【 図 12 】 本発明の相似形チェック処理で、ブロック左上角 Y 座標の比較処理を示すフローチャートである。

30

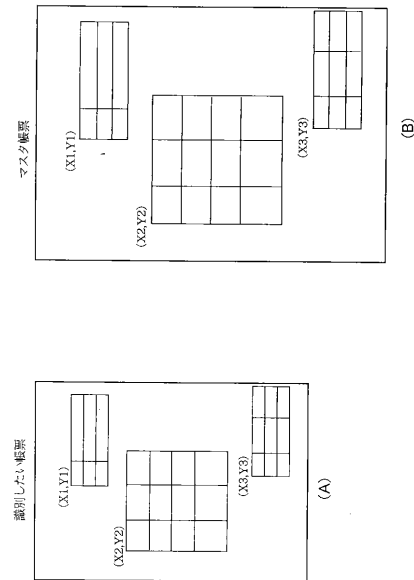
【 図 13 】 本発明の相似形チェック処理で、ブロック個数が 2 個以下の場合の処理を示すフローチャートである。

【 図 14 】 本発明の相似形帳票の一例を示す図である。

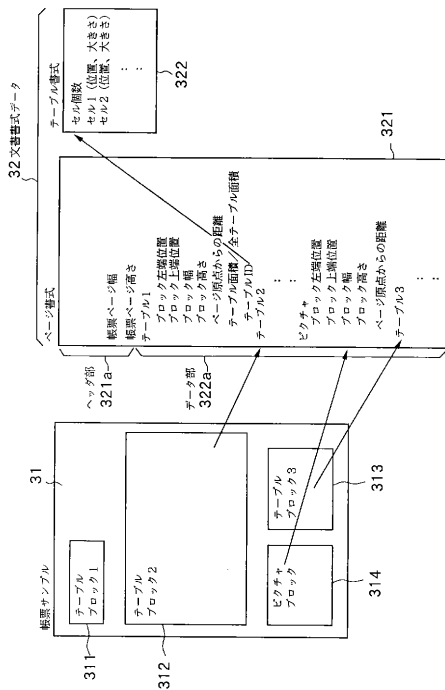
【 図 1 】



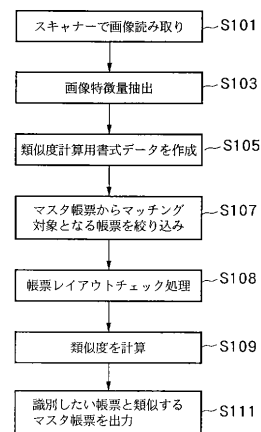
【 図 2 】



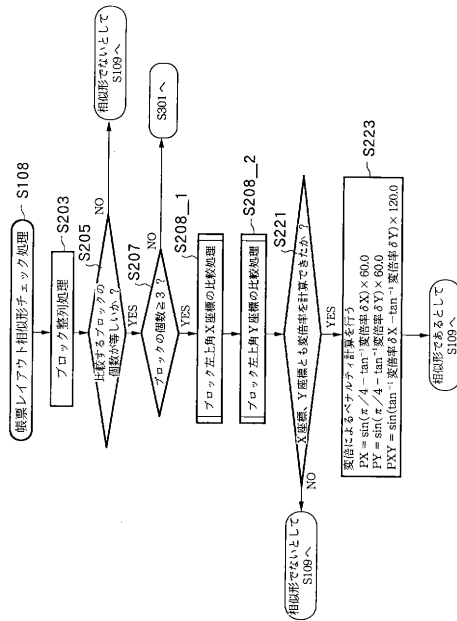
【 図 3 】



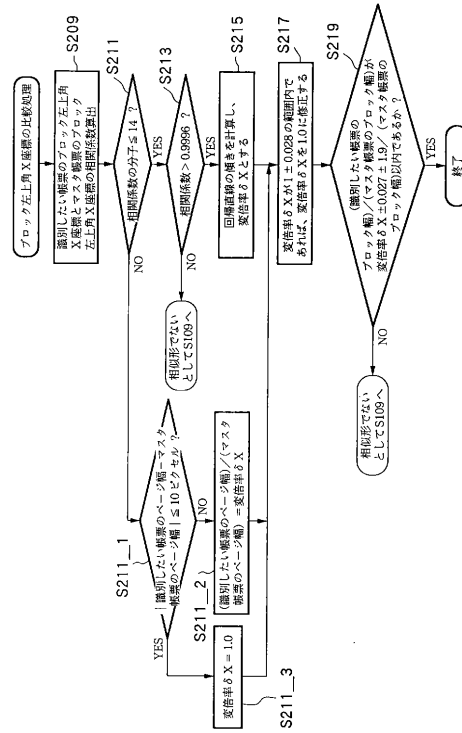
【圖 4】



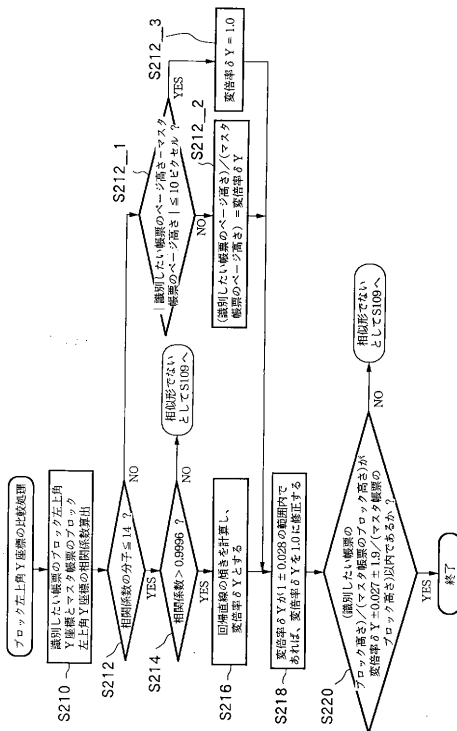
【図 5】



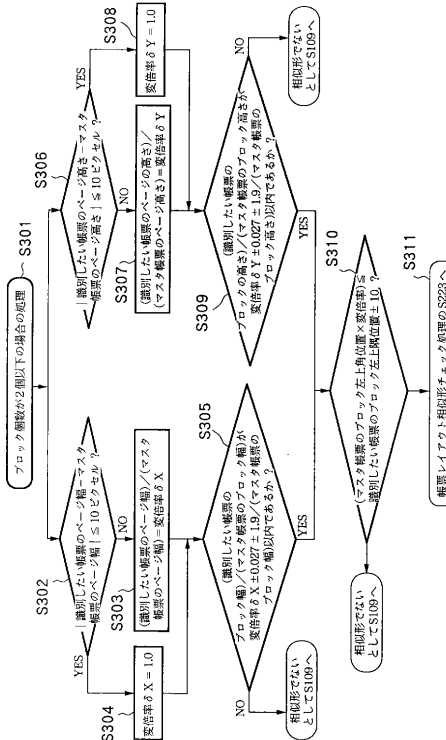
【図 6】



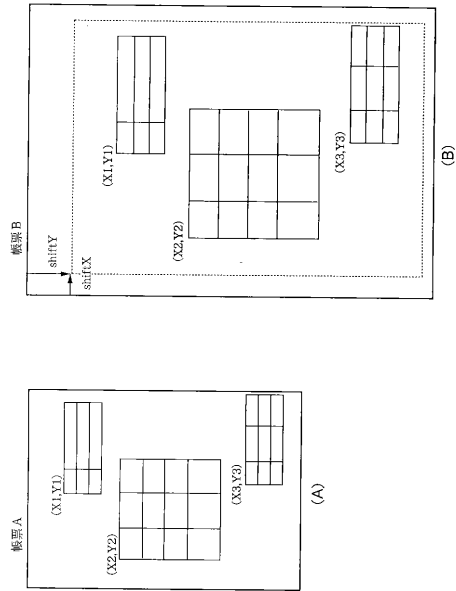
【図 7】



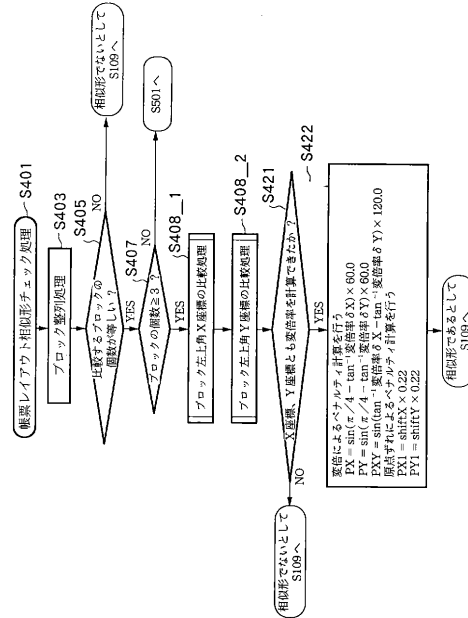
【図 8】



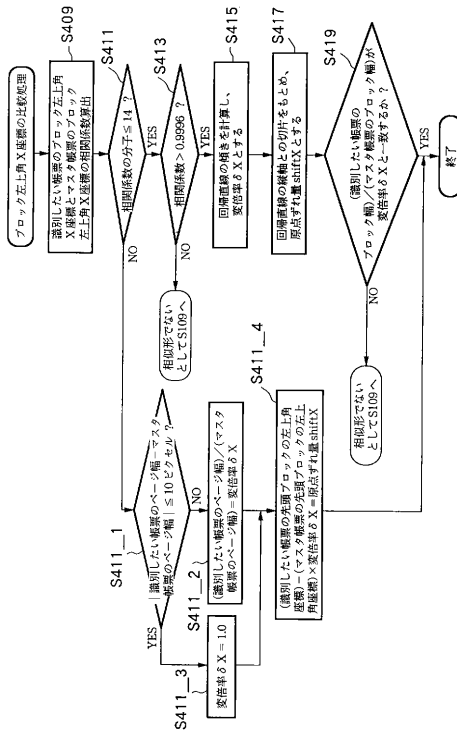
【図 9】



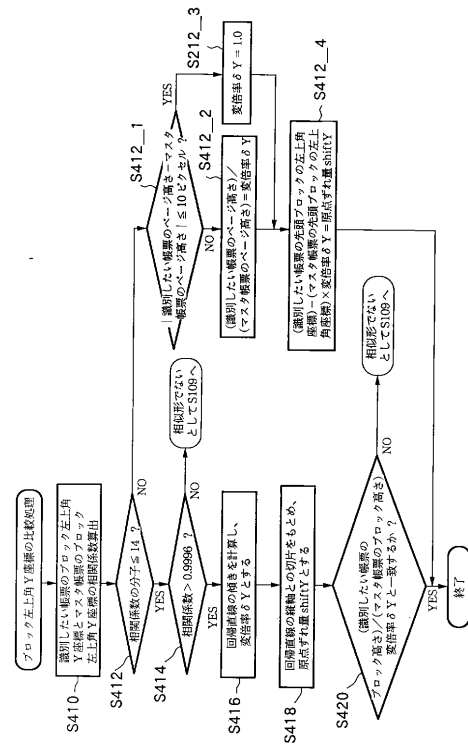
【図 10】



【図 11】



【図 12】





---

フロントページの続き

(72)発明者 数見 健一

東京都大田区下丸子3丁目30番2号 キヤノン株式会社内

審査官 佐藤 実

(56)参考文献 特開平10-027208(JP,A)

特開昭64-005141(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06K 9/00~9/76