

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2006-106972

(P2006-106972A)

(43) 公開日 平成18年4月20日(2006.4.20)

(51) Int. Cl.

G06K 9/00 (2006.01)

F I

G06K 9/00

P

テーマコード (参考)

5B064

審査請求 未請求 請求項の数 8 O L (全 16 頁)

(21) 出願番号 特願2004-290386 (P2004-290386)

(22) 出願日 平成16年10月1日 (2004. 10. 1)

(71) 出願人 000001007

キヤノン株式会社

東京都大田区下丸子3丁目30番2号

(74) 代理人 100076428

弁理士 大塚 康德

(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

(72) 発明者 鶴沢 充

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

Fターム(参考) 5B064 AA01 AB09 BA01 CA03

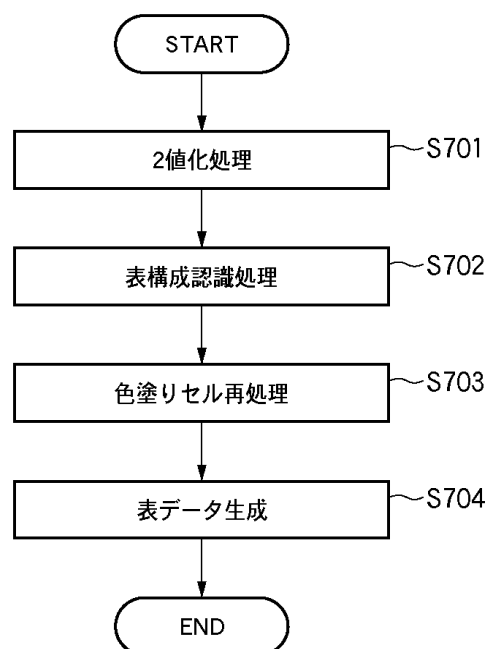
(54) 【発明の名称】 画像処理方法及び画像処理装置

(57) 【要約】

【課題】 本来文字情報が埋め込まれているセルが2値化処理により塗りつぶされることを防止し、表中の全ての文字情報を抽出する。

【解決手段】 原稿上の表領域における表枠の2値データを抽出し、表枠の2値データに基づいて表枠の表構成を認識し、その表枠の中より文字情報を抽出し、表枠を再処理する際に、表枠の表構成の認識結果に応じて文字情報を抽出する領域を選択し、該選択位置における文字情報を再抽出し、その表枠を再処理する。

【選択図】 図7



【特許請求の範囲】**【請求項 1】**

原稿上の表領域における表枠の 2 値データを抽出する工程と、
前記表枠の 2 値データに基づいて表枠の表構成を認識する工程と、
前記表枠の中より文字情報を抽出し、表枠を再処理する工程とを有し、
前記表枠の表構成の認識結果に応じて文字情報を抽出する領域を選択し、該選択位置における文字情報を再抽出し、前記表枠を再処理することを特徴とする画像処理方法。

【請求項 2】

前記表構成の認識結果は、矩形図形の集合として表を表現し、前記文字情報を抽出する領域は該矩形図形単位で選択されることを特徴とする請求項 1 記載の画像処理方法。

10

【請求項 3】

前記再処理は、少なくとも 2 値化処理か像域分離処理の何れかであることを特徴とする請求項 1 記載の画像処理方法。

【請求項 4】

前記再処理された表枠を表現するベクトルデータを生成する工程を有することを特徴とする請求項 1 記載の画像処理方法。

【請求項 5】

原稿上の表領域における表枠の 2 値データを抽出する抽出手段と、
前記表枠の 2 値データに基づいて表枠の表構成を認識する認識手段と、
前記表枠の中より文字情報を抽出し、表枠を再処理する再処理手段とを有し、
前記表枠の表構成の認識結果に応じて文字情報を抽出する領域を選択し、該選択位置における文字情報を再抽出し、前記表枠を再処理することを特徴とする画像処理装置。

20

【請求項 6】

前記再処理手段は、前記表枠でないと想定される領域について、テキスト情報が入っていないと判定した場合は再処理しないことを特徴とする請求項 5 記載の画像処理装置。

【請求項 7】

請求項 1 記載の画像処理方法をコンピュータに実行させるためのプログラム。

【請求項 8】

請求項 7 記載のプログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

30

【技術分野】**【0001】**

本発明は、スキャナなどの入力装置より読み込まれた紙文書を編集可能な電子データへ変換する技術に関し、特に紙文書中の表枠オブジェクトを解析し、罫線に置き換える技術に関する

【背景技術】**【0002】**

近年、情報の電子化が進み、文書を紙ではなく電子化して保存、あるいは送信するシステムが急速に普及している。特に、フルカラーの文書を保存、送信に適した電子データとしては、紙原稿を文字、表、図等のオブジェクトへ像域分離し、各オブジェクトに適した形態でデータ化したベクトルデータが適しており、データ量を削減し、更にオブジェクトを再利用することが可能となる。

40

【0003】

ここで、文書中の表オブジェクトについては、像域分離処理により表中のテキスト情報を抽出する一方、表枠を 2 値画像データとして抽出することが可能である。この 2 値画像データとして抽出された表枠は、表枠の解析処理により、罫線などで表現されるベクトルデータに変換される（例えば、特許文献 1 参照）。そして、ベクトルデータに変換された表枠はベクトルデータとしてデータサイズが小さくなるだけでなく、表枠としての再利用性も高い。

【特許文献 1】特開平 5 - 12489 号公報

50

【発明の開示】

【発明が解決しようとする課題】

【0004】

表枠及び表中のテキストを含めて、文書中のオブジェクトを２値化して抽出する際に、画像領域で２値化、もしくはオブジェクト毎に２値化するために、それぞれ２値化の閾値は画像領域、もしくはオブジェクト領域により設定される。

【0005】

例えば、表領域について２値化を行った場合、図２５に示す（ａ）のように、色の濃いセルがあると、図２５に示す（ｂ）のように、色の濃いセルは塗りつぶされた画像として生成されるという問題が生じる。また、２値化した際に、セル内にノイズが多く発生する

10

【0006】

尚、このようなセルについては、正常に文字情報を抽出することができない。元原稿がカラー原稿である場合は、各セル色の相関関係によって２値化によりセルの塗りつぶしが発生する場合が多く、表枠の正確な抽出は困難であり、非常に深刻な問題となっている。

【0007】

本発明は、上述の課題を解決するためになされたもので、本来文字情報が埋め込まれているセルが２値化処理により塗りつぶされることを防止し、表中の全ての文字情報を抽出することを目的とする。

【課題を解決するための手段】

20

【0008】

本発明の画像処理方法は、原稿上の表領域における表枠の２値データを抽出する工程と、前記表枠の２値データに基づいて表枠の表構成を認識する工程と、前記表枠の中より文字情報を抽出し、表枠を再処理する工程とを有し、前記表枠の表構成の認識結果に応じて文字情報を抽出する領域を選択し、該選択位置における文字情報を再抽出し、前記表枠を再処理することを特徴とする。

【0009】

また、本発明の画像処理装置は、原稿上の表領域における表枠の２値データを抽出する抽出手段と、前記表枠の２値データに基づいて表枠の表構成を認識する認識手段と、前記表枠の中より文字情報を抽出し、表枠を再処理する再処理手段とを有し、前記表枠の表構成の認識結果に応じて文字情報を抽出する領域を選択し、該選択位置における文字情報を再抽出し、前記表枠を再処理することを特徴とする。

30

【発明の効果】

【0010】

本発明によれば、カラー原稿より表枠中のノイズを除去し、表枠を正確に抽出し、セル内のノイズ、塗りつぶしにより文字情報が抽出されないことを防止できる。また、正確な表枠ベクトル情報を抽出できる。

【発明を実施するための最良の形態】

【0011】

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

40

【実施例１】

【0012】

図１は、実施例１における文書処理装置の外観を示す図である。図１において、１０１はコンピュータ装置であり、後述するフローチャートを参照して説明する処理を実現するためのプログラムを含む、文書の電子化処理プログラムを実行する。また、コンピュータ装置１０１は、ユーザに状況や画像を表示するためのディスプレイ装置１０２と、ユーザの操作を受け付けるキーボードやマウス等のポインティングデバイスを含んで構成される入力装置１０３とを付随する。このディスプレイ装置１０２としては、ＣＲＴやＬＣＤ等が用いられる。１０４はスキャナ装置であり、文書画像を光学的に読み取って電子化し、得られた画像データをコンピュータ装置１０１に送る。尚、スキャナ装置１０４としては

50

、カラスキャナを用いるものとする。

【0013】

図2は、実施例1における文書処理装置の構成の一例を示すブロック図である。図2において、201はCPUであり、後述するROM又はRAMに格納された制御プログラムを実行することにより、後述する電子化処理を含む各種機能を実現する。202はROMであり、CPU201によって実行される各種制御プログラムや制御データが格納されている。203はRAMであり、CPU201によって実行される各種制御プログラムを格納したり、CPU201が各種処理を実行するのに必要な作業領域が定義されている。

【0014】

204は外部記憶装置であり、詳細は後述する実施例1における処理をCPU101によって実現するための制御プログラムや、スキャナ装置104で読み取って得られた文書画像データ等を格納する。そして、205はコンピュータバスであり、上述した各構成を接続するものである。

【0015】

図3は、文書処理装置における文書の電子化処理の概要を示す図である。ここで、電子化処理の流れは、まず入力部301において、電子化の対象であるカラー文書をスキャナ装置104によって読み込み、画像データとして外部記憶装置204に格納する。次に、2値化処理302において、後段の像域分離処理、アウトライン生成処理のために、外部記憶装置204に格納された文書の画像データに対して2値化処理を施す。そして、像域分離処理303では、2値化処理302で得られた2値画像から、文字、図、表、枠、線などの要素を抽出し、各領域に分割する。

【0016】

次に、ベクトル化処理304において、領域分割された画像データに対して、文字部は文字認識部305で文字認識を行い、アウトライン作成部306でアウトラインベクトルデータへ変換する。また、表、枠の要素については表処理部308でアウトライン化し、アウトラインを罫線化する。尚、アウトライン作成部306で変換された画像データは、各オブジェクトの輪郭線が滑らかな曲線により表現される高画質で、解像度に依存しない、かつ編集容易なベクトルデータへ変換される。

【0017】

一方、その他の図、写真画、背景については、例えば背景については、圧縮部309でJPG圧縮など各々に適した形態で保持、圧縮する。

【0018】

次に、電子文書作成処理310は、分割された要素毎の属性に基づいて文字認識データや表構造データを用い、それぞれ変換された画像データに基づき電子化文書を作成する。そして、出力部311は生成された電子化文書を外部記憶装置204に格納する。

【0019】

尚、出力部311の出力形態は外部記憶装置204への格納に限られるものではなく、ディスプレイ装置102へ表示したり、不図示のネットワークインターフェースを介してネットワーク上の他の装置へ出力したり、不図示のプリンタへ送出したりすることも可能である。

【0020】

ここで、図1及び図2に示す文書処理装置において実行される文書の電子化処理（図3参照）における各処理の詳細について、以下順に説明する。

【0021】

[2値化処理]

2値化処理302では、入力された文書画像データより輝度情報を抽出し、その輝度値のヒストグラムを作成する。ヒストグラム上より複数の閾値を設定し、各々の閾値で2値化された2値画像上の黒画素の連結等を解析することで最適な閾値を導出し、その閾値による2値画像を得る。

【0022】

10

20

30

40

50

〔 像域分離処理 〕

像域分離処理 3 0 3 とは、図 4 に示す左側の読み取られた 1 ページのイメージデータをオブジェクト毎の塊（ブロック）として認識し、各々の塊を文字 / 図画 / 写真 / 線 / 表等の属性に判定し、図 4 に示す右側のように、異なる属性（TEXT / PICTURE / PHOTO / LINE / TABLE）を持つ領域に分割する処理である。

【 0 0 2 3 】

像域分離処理 3 0 3 では、2 値化処理 3 0 2 で得られた 2 値画像より、黒画素の輪郭線追跡を行って黒画素輪郭で囲まれる画素の塊を抽出する。また、面積の大きい黒画素の塊については、内部にある白画素に対しても輪郭線追跡を行い、白画素の塊を抽出し、更に一定面積以上の白画素の塊の内部からは再帰的に黒画素の塊を抽出する。

10

【 0 0 2 4 】

このようにして得られた黒画素の塊を、大きさ及び形状で分類し、異なる属性を持つ領域へ分類していく。例えば、縦横比が 1 に近く、大きさが一定の範囲のものを文字相当の画素塊とし、更に近接する文字が整列良くグループ化可能な部分を文字領域、扁平な画素塊を線領域、一定の大きさ以上で、かつ四角系の白画素塊を整列よく内包する黒画素塊の占める範囲を表領域、不定形の画素塊が散在している領域を写真領域、それ以外の任意形状の画素塊を図画領域、などとする。

【 0 0 2 5 】

図 5 は、像域分離処理 3 0 3 で分離された各ブロックに対するブロック情報と入力ファイル情報を示す図である。図 5 に示すように、ブロック情報は、各ブロックの属性、座標（X，Y）、幅（W）、高さ（H）、OCR 情報を含み、属性 1 は文字、属性 2 は図画、属性 3 は表、属性 4 は線、属性 5 は写真である。そして、入力ファイル情報は、ブロック総数 N（図 5 に示す例では、ブロック 1 ～ブロック 6 までの 6 である）を有する。

20

【 0 0 2 6 】

尚、各ブロックに対して、より鮮明な 2 値画像を得ようとした場合は、ここでブロック毎に上述した 2 値化処理を行っても良い。

【 0 0 2 7 】

〔 文字認識部 〕

文字認識部 3 0 5 では、文字単位で切り出された画像に対して、パターンマッチングの一手法を用いて認識を行い、対応する文字コードを得る。この認識処理は、文字画像から得られる特徴を数十次元の数値列に変換した観測特徴ベクトルと、予め字種毎に求められている辞書特徴ベクトルとを比較し、最も距離の近い字種を認識結果とする処理である。この特徴ベクトルの抽出には種々の公知手法があり、例えば文字をメッシュ状に分割し、各メッシュ内の文字線を方向別に線素としてカウントしたメッシュ数次元ベクトルを特徴とする方法がある。

30

【 0 0 2 8 】

像域分離処理 3 0 3 で抽出された文字領域に対して文字認識を行う場合、まず該当領域に対して横書き、縦書きの判定を行い、それぞれ対応する方向に行を切り出し、その後、文字を切り出して文字画像を得る。この横書き、縦書きの判定は、該当領域内で画素値に対する水平 / 垂直の射影を取り、水平射影の分散が大きい場合には横書き領域と判定し、垂直射影の分散が大きい場合には縦書き領域と判定すれば良い。また、文字列及び文字への分解は、横書きならば水平方向の射影を利用して行を切り出し、更に切り出された行に対する垂直方向の射影から、文字を切り出すことで行う。縦書きの文字領域に対しては、水平と垂直を逆にすれば良い。尚、この時、文字のサイズが検出できる。

40

【 0 0 2 9 】

〔 アウトライン生成部 〕

アウトライン作成部 3 0 6 では、像域分離処理 3 0 3 で図画或いは線、表領域とされた領域を対象に、抽出された画素塊の輪郭をベクトルデータに変換する。具体的には、輪郭をなす画素の点列を角と看做される点で区切って、各区間を部分的な直線或いは曲線で近似する。角とは曲率が極大となる点であり、曲率が極大となる点は、図 6 に示すように、

50

任意点 P_i に対して左右 k 個の離れた点 P_{i-k} , P_{i+k} の間に弦を引いたとき、この弦と P_i の距離が極大となる点として求められる。更に、点 P_{i-k} , P_{i+k} 間の弦の長さ / 弧の長さを R とし、 R の値が閾値以下である点を角とみなすことができる。角によって分割された後の各区間は、直線は点列に対する最小二乗法など、曲線は 3 次スプライン関数などを用いてベクトル化することができる。

【 0 0 3 0 】

また、対象が内輪郭を持つ場合、像域分離処理で抽出した白画素輪郭の点列を用いて、同様に部分的直線或いは曲線で近似する。

【 0 0 3 1 】

[表処理部]

表処理部 308 では、表中のセル及びその構成を認識し、表枠を罫線により表現する等、セル毎に編集可能な電子データへ変換する。尚、表部は、像域分離処理 303 により、表枠として表枠中の文字部と分離して抽出されている。

【 0 0 3 2 】

図 7 は、表処理部 308 における表処理を示すフローチャートである。まず、像域分離処理 303 で分離された表部を 2 値化し、表枠の 2 値データを得る (ステップ S701)。ここでは、像域分離処理 303 より入力された表枠の 2 値データをそのまま出力しても良い。また、像域分離処理 303 で得られた表領域に対して、再度 2 値化処理、表枠抽出処理を行い、正確に表枠を抽出しても良い。

【 0 0 3 3 】

次に、ステップ S701 で得られた表枠について表構成を認識する (ステップ S702)。図 8 は、表枠を認識する処理を示すフローチャートである。まず、表枠について上述したアウトライン処理によりアウトライン化し、滑らかな直線及び曲線で表枠を表現する (ステップ S801)。次に、表枠のアウトラインよりセルを表現しているセルアウトラインを抽出する (ステップ S802)。ここで、元々アウトラインは、外輪郭と内輪郭に分類されているが、まず外輪郭のうち、表全体の外枠を構成しているアウトラインを抽出する。尚、表の内部に表が存在するような場合もあるので、外枠は複数抽出される場合もある。

【 0 0 3 4 】

次に、外枠の内側に存在するセルを構成しているアウトラインを抽出する。尚、ここでの処理はアウトラインよりそのアウトラインがセルを構成するサイズであるか否かを判定し、更にアウトラインを図形認識処理によりアウトラインが矩形図形、もしくは三角図形を構成しているか否かを判定する。尚、矩形図形、三角図形、もしくは矩形図形の集合と判定されたアウトラインをセルアウトラインとする。図 9 に示す (a), (b), (c) は、矩形図形、三角図形、矩形図形の集合と判定されたアウトラインの一例である。

【 0 0 3 5 】

次に、ステップ S702 で外枠、もしくはセルを構成していると判定されたセルアウトラインについてセル図形へ変換する (ステップ S803)。具体的には、まず、ステップ S802 で外枠を構成するアウトラインと内部セルを構成するアウトラインが抽出されているが、外枠を構成するアウトラインの角度が全て 90° で表現される図形であると判定された場合、 90° 角の間を直線で表現した図形へ変換する。次に、外枠の内側のセルを構成すると判定されたセルについて、セル図形へ変換する。ここで、セル図形とは矩形図形である。

【 0 0 3 6 】

例えば、図 9 に示すセルアウトラインをセル図形へ変換した例を図 10 に示す。図 9 に示すセルアウトライン (a), (b), (c) はそれぞれ図 10 に示す (a), (b), (c) のように変換される。ここで、図 9 に示す (a) のセルアウトラインは矩形図形の当てはめ処理によりそのまま図 10 に示す (a) となる。図 9 に示す (b) のような三角図形のセルアウトラインも同様に、矩形図形の当てはめ処理を行う。この三角図形に対する矩形図形の当てはめ処理では、三角を構成するセルアウトラインを囲むようにセル矩形

10

20

30

40

50

を当てはめる。

【0037】

尚、当てはめられた矩形図形は最終的にその位置関係よりセル図形同士統合され一つのセル図形として抽出する。例えば、図10に示す(b)のセル図形は、図11に示すセル図形に統合され、一つのセル図形として表現される。統合されないセル図形も当然あり、それらの図形についてはそのまま三角アウトラインに当てはめられたセル図形をセル図形とする。

【0038】

また、図9に示す(c)のような矩形図形の集合として抽出されるセルアウトラインは、図10に示す(c)のように各矩形図形へ分離する。この矩形図形への分離処理では、アウトラインの中の直角をなすであろう角を見つけ出し、その角点の構成から矩形図形へ分解する。尚、更に抽出されたセル全てを用いて表の水平方向及び垂直方向を求め、セル図形全てを求められた水平方向及び垂直方向の成分からなるセル図形へ変換する。このように外枠図形とセル図形は抽出されるが、各セル図形はその構成されるセルアウトラインにより属性情報がつけられる。図12は、セル図形と、それを構成するセルアウトライン及び各属性情報の一例を示す図である。

10

【0039】

次に、ステップS803で変換したセル図形をマッピングするためのマッピング領域を作成する(ステップS804)。ここで、マッピング領域とはセル図形をマッピングするための領域であり、外枠内部の領域がそのままマッピング領域となる。また、外枠の交点を抽出しておく。交点とは、表の罫線と罫線が交差する点のことであり、外枠においては、外枠の角点が多角形となる。図13にマッピング領域と交点の一例を示す。

20

【0040】

次に、セル図形マッピング領域にセル図形をマッピングし(ステップS805)、セル図形の構成を認識する。具体的には、ステップS804で作成されたマッピング領域内にセル図形をマッピングしていき、マッピングされるセル図形より表の罫線と罫線が交差する交点を抽出していくことで、表構成を認識する。即ち、交点の隣接関係を調べていくことで、表構成を認識する。

【0041】

ここで、ステップS805におけるセル図形をマッピングするセル図形マッピング処理の詳細について説明する。

30

【0042】

図14は、セル図形マッピング処理を示すフローチャートである。まず、上述した外枠の交点より注目点を抽出する(ステップS1401)。ここで、注目点とは、右側と下側に隣接して繋がっている交点を持つ交点で、それら3点を含む矩形領域がマッピング領域であり、かつ該矩形領域に対し、まだ何もマッピングされていない交点である。図15にマッピング領域と交点とマッピングセルと注目点の関係を示した例を示す。

【0043】

次に、抽出した注目点に対し、該注目点を左上の交点とするセル図形が存在するか否かを判定する(ステップS1402)。具体的には、未だマッピングされていない全てのセル図形の左上の角点と注目点との距離を調べ、セル図形の左上の角点と注目点との距離が一定値以内であり、最も注目点に近いセル図形を注目点と左上の角点が一一致するセル図形と判定する。ここで、注目点と左上の角点が一一致するセル図形が存在すれば、該セル図形をマッピング領域上にマッピングする(ステップS1403)。

40

【0044】

また、セル図形全ての左上の角点と注目点との距離が一定値以内にあるセル図形が存在しなければ、色塗りセルを作成し、マッピングする(ステップS1404)。色塗りセルは、矩形図形である。まず、注目点とその隣接する右側と下側の交点より少し広げた矩形領域(以後矩形領域Aと呼ぶ)内に、まだマッピングされていないセル図形の角点がないか判定し、角点が存在すれば、その角点を通る水平方向及び垂直方向の直線によって領域

50

を区切る。この区切り作業を矩形領域 A 内に存在する角点全てに対して行い、水平線及び垂直線によって区切られた領域の最も左上にある区切られた矩形図形を色塗りセルとし、マッピングする。

【 0 0 4 5 】

図 1 6 は、矩形領域 A 内に色塗りセルを作成する例を示す図である。図 1 6 に示す例では、矩形領域 A 内に 2 つのセル図形の角点が存在するので、それらの角点を通る水平方向及び垂直方向の直線によって矩形領域 A を区切り、区切られた領域の最も左上にある区切られた矩形図形を塗りつぶしセルとしてマッピングする。

【 0 0 4 6 】

次に、ステップ S 1 4 0 3、S 1 4 0 4 でマッピングされたセル図形よりセル図形上の交点を作成する（ステップ S 1 4 0 5）。交点はこのマッピング図形の角点があるまま交点となるが、もしマッピング図形の角点が、既に存在する交点との距離がある閾値以内であれば、その角点により作成される交点は既に存在すると判断できるため、その角点より新たな交点は作成しない。ここで、マッピング図形の左上の交点は注目点と一致と判断されているため、左上の角点より新たな交点は作成されない。また、マッピング図形の右上の角点より作成される交点は注目点より水平線上にあるとして作成し、左下の角点は注目点より垂直線上にあるとして作成する。

【 0 0 4 7 】

次に、現在抽出されている交点の中で注目点が存在するか否かを判定する（ステップ S 1 4 0 6）。注目点とは、上述したように、右側と下側に隣接し繋がっている交点を持つ交点で、それら 3 点を含む矩形領域がマッピング領域であり、かつ該矩形領域に対しまだ何もマッピングされていない交点である。交点が囲む領域内にセル図形及び塗りつぶしセルがマッピングされていない領域が存在すれば注目点は存在する。尚、注目点が存在しない場合には表構成認識処理を終了とする。また、まだ注目点が存在する場合は、ステップ S 1 4 0 1 に戻り、再度注目点を抽出して一連の処理を繰り返す。

【 0 0 4 8 】

以上の繰り返し処理により、交点の隣接関係が作成され、罫線を表現することが可能となる。図 1 7 に作成された交点の隣接関係と、マッピングされたセル図形及び色塗りセル、またセル図形の場合はその属性情報を記述した例を示す。尚、図 1 7 に示す属性情報により得られる表罫線は図 1 8 に示すようになる。

【 0 0 4 9 】

ここで、図 1 4 に示したセル図形マッピング処理が終了すると、図 8 に示すステップ S 8 0 6 へ進み、罫線の太さ及び位置関係を調節し、交点の正確な位置を抽出する。罫線の太さは、ステップ S 8 0 5 でマッピングされたセル図形のうち、隣接するセル図形の距離から求められる。また、罫線の位置は隣接するセル図形の間となるように調節する。

【 0 0 5 0 】

尚、交点の位置はステップ S 8 0 6 で得られる罫線同士の交わる交点として正確な位置を求める。

【 0 0 5 1 】

[色塗りセル再処理]

図 7 に示すステップ S 7 0 3 では、表構成認識処理（ステップ S 7 0 2）で抽出された色塗りセルについて再処理を行う。

【 0 0 5 2 】

図 1 9 は、色塗りセル再処理を示すフローチャートである。まず、セルの判別を行う（ステップ S 1 9 0 1）。セル判別では、ステップ S 7 0 2 で作成した色塗りセルの区切りについて原稿上でもセルの区切りが存在するか否かを判別し、場合によっては色塗りセルを統合する。具体的には、原画像について、隣り合う色塗りセルの区切り線に対応する近隣画素のエッジ成分を調べ、原画像において区切り線上の罫線があるか否かを判別する。もし、原画像上で罫線が無いと判別された場合は、2 つのセルを統合し一つのセルとする。以上の統合処理を隣接する色塗りセル全てに対して行い、セルを抽出する。

10

20

30

40

50

【 0 0 5 3 】

次に、ステップ S 1 9 0 1 で抽出されたセルについて、一つのセルを一枚の画像とみなし、セル毎に 2 値化処理を行う（ステップ S 1 9 0 2）。セル毎に 2 値化処理を行うことで、塗りつぶされていた色塗りセルより、セル内のテキスト、模様等の 2 値オブジェクトを抽出することが可能である。そして、セル毎に像域分離処理を行う（ステップ S 1 9 0 3）。

【 0 0 5 4 】

尚、像域分離処理は、特にテキストデータを抽出するものである。抽出されたテキストデータが抽出された場合は O C R をかけてテキストコードを抽出し、更にアウトライン化してフォント化し、文字部のデータへ追加する。

10

【 0 0 5 5 】

〔表データ生成処理〕

図 7 に示すステップ S 7 0 4 では、ステップ S 7 0 2、S 7 0 3 で作成されたベクトル情報を使用し、表データを作成する。例えば、図 1 8 に示す表は図 2 4 に示すようになる。図 2 4 において、2 4 0 1 ~ 2 4 0 3 の部位からは、テキスト情報が抽出されている。また、もしテキスト情報がないセルについては、そのまま色塗りセルのままで良いとする。

【 0 0 5 6 】

以上の処理により、表枠は塗りつぶされるようなセルが作成されることを回避し、表中よりテキスト情報を正確に抽出し、表データが作成される。

20

【 0 0 5 7 】

尚、上述の処理では、色塗りセルについて、セル判別を行い、得られる各セルについて 2 値化、像域分離処理を行ったが、より単純に抽出された複数の色塗りセル全ての領域を一枚の原稿と見立て、それに対して 2 値化、像域分離処理を行ってもよく、色塗りセルの領域より文字、線等の情報を得ることが可能である。

【 0 0 5 8 】

〔アプリデータへの変換処理〕

以上の通り、1 頁分のイメージデータを像域分離処理 3 0 3 し、ベクトル化処理 3 0 4 した結果は図 2 0 に示すような中間データ形式のファイルとして変換される。このようなデータ形式は、ドキュメント・アナリシス・アウトプット・フォーマット（D A O F）と

30

【 0 0 5 9 】

図 2 0 は、D A O F のデータ構造を示す図である。図 2 0 において、2 0 0 1 は Header であり、処理対象の文書画像データに関する情報が保持される。2 0 0 2 はレイアウト記述データ部であり、文書画像データ中の文字（TEXT）、タイトル（TITLE）、キャプション（CAPTION）、線画（LINEART）、自然画（PICTURE）、枠（FRAME）、表（TABLE）等の属性毎に認識された各ブロックの属性情報とその矩形アドレス情報を保持する。2 0 0 3 は文字認識記述データ部であり、TEXT、TITLE、CAPTION 等の TEXT ブロックを文字認識して得られる文字認識結果を保持する。2 0 0 4 は表記述データ部であり、TABLE ブロックの構造の詳細を格納する。2 0 0 5 は画像記述データ部であり、PICTURE や LINEART 等のブロックのイメージデータを文書画像データから切り出して保持する。

40

【 0 0 6 0 】

このような D A O F は中間データとしてのみならず、それ自体がファイル化されて保存される場合もあるが、このファイルの状態では、所謂一般の文書作成アプリケーションで個々のオブジェクトを再利用することはできない。

【 0 0 6 1 】

そこで、この D A O F からアプリケーションデータに変換する電子文書作成処理 3 0 9 について説明する。

【 0 0 6 2 】

図 2 1 は、電子文書作成処理の全体の概略を示すフローチャートである。まずステップ

50

S 2 1 0 1において、D A O Fデータの入力を行う。次に、ステップS 2 1 0 2において、アプリデータの元となる文書構造ツリー生成を行う。そして、ステップS 2 1 0 3で、文書構造ツリーに基づいてD A O F内の実データを流し込み、実際のアプリデータを生成する。

【 0 0 6 3 】

図 2 2 は、文書構造ツリー生成処理の詳細を示すフローチャートである。また、図 2 3 は文書構造ツリーを説明するための図である。尚、全体制御の基本ルールとして、処理の流れはマイクロブロック（単一ブロック）からマクロブロック（ブロックの集合体）へ移行する。尚、以下の説明で、「ブロック」はマイクロブロック及びマクロブロック全体を指すものとする。

10

【 0 0 6 4 】

まず、ステップS 2 2 0 1では、ブロック単位に縦方向の関連性に基づいて再グループ化する。スタート直後はマイクロブロック単位での判定となる。ここで、関連性とは、距離が近い、ブロック幅（横方向の場合は高さ）がほぼ同一であることなどで定義することができる。また、距離、幅、高さなどの情報はD A O Fを参照し、抽出する。

【 0 0 6 5 】

図 2 3 は、ページの構成とその文書構造のツリーを示す図である。図 2 3 に示す（ a ）は実際のページ構成、図 2 3 に示す（ b ）はその文書構造ツリーである。

【 0 0 6 6 】

ステップS 2 2 0 1での結果、図 2 3 に示すT 3、T 4、T 5が1つのグループV 1として生成され、T 6、T 7が1つのグループV 2として生成され、図 2 3 に示す（ b ）のように、グループV 1とグループV 2が同じ階層のグループとして生成される。そして、ステップS 2 2 0 2において、縦方向のセパレータの有無をチェックする。セパレータは、例えば物理的にはD A O F中でライン属性を持つオブジェクトである。また、論理的な意味としては、アプリ中で明示的にブロックを分割する要素である。ここでセパレータを検出した場合は、同じ階層で再分割する。

20

【 0 0 6 7 】

次に、ステップS 2 2 0 3において、分割がこれ以上存在し得ないか否かをグループ長を利用して判定する。ここで、縦方向のグループ長がページ高さとなっている場合、文書構造ツリー生成を終了する。また、図 2 3 に示す例の場合、セパレータもなく、グループ高さはページ高さではないのでステップS 2 2 0 4へ進み、ブロック単位で横方向の関連性に基づいて再グループ化する。ここもスタート直後の第一回目はマイクロブロック単位で判定を行うことになる。尚、関連性、及びその判定情報の定義は、縦方向の場合と同じである。

30

【 0 0 6 8 】

図 2 3 に示す例の場合、T 1、T 2でH 1が、V 1、V 2でH 2がV 1、V 2の1つ上の同じ階層のグループとして生成される。そして、ステップS 2 2 0 5において、横方向セパレータの有無をチェックする。図 2 3 に示す例では、S 1があるので、これをツリーに登録し、H 1、S 1、H 2という階層を生成する。

【 0 0 6 9 】

次に、ステップS 2 2 0 6において、分割がこれ以上存在し得ないか否かをグループ長を利用して判定する。ここで、横方向のグループ長がページ幅となっている場合、文書構造ツリー生成を終了する。また、そうでない場合はステップS 2 2 0 1に戻り、再びもう一段上の階層で、縦方向の関連性チェックから繰り返す。図 2 3 に示す例の場合、分割幅がページ幅になっているので、ここで終了し、最後にページ全体を表す最上位階層のV 0が文書構造ツリーに付加される。

40

【 0 0 7 0 】

文書構造ツリーが完成した後、その情報に基づいてアプリデータを生成する（ステップS 2 1 0 3）。図 2 3 に示す例の場合、具体的には、以下のようになる。

【 0 0 7 1 】

50

即ち、H 1 は横方向に 2 つのブロック T 1 及び T 2 があるので、2 カラムとし、T 1 の内部情報 (D A O F を参照、文字認識結果の文章、画像など) を出力後、カラムを変え、T 2 の内部情報出力、その後 S 1 を出力する。次に、H 2 は横方向に 2 つのブロック V 1 及び V 2 があるので、2 カラムとして出力、V 1 は T 3、T 4、T 5 の順にその内部情報を出力、その後カラムを変え、V 2 の T 6、T 7 の内部情報を出力する。

【0072】

以上の処理によりアプリデータへの変換処理を行うことができる。

【0073】

尚、本発明は複数の機器 (例えば、ホストコンピュータ、インターフェース機器、リーダー、プリンタなど) から構成されるシステムに適用しても、1 つの機器からなる装置 (例えば、複写機、ファクシミリ装置など) に適用しても良い。具体的には、複合機や、複写機や、ファクシミリ装置で、高品位に変倍するために、スキャンした画像データを入力し (公衆回線やネットワークから画像データを入力しても良い)、画像データから輪郭ベクトルを抽出し、抽出した輪郭ベクトルを変倍し、変倍された輪郭ベクトルから画像データを生成し、生成した画像データをプリントする際の輪郭ベクトル抽出時に適用できる。

【0074】

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録媒体を、システム或いは装置に供給し、そのシステム或いは装置のコンピュータ (CPU 若しくは MPU) が記録媒体に格納されたプログラムコードを読み出し実行することによっても、達成されることは言うまでもない。

【0075】

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

【0076】

このプログラムコードを供給するための記録媒体としては、例えばフロッピー (登録商標) 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。

【0077】

また、コンピュータが読み出したプログラムコードを実行することにより、前述した実施形態の機能が実現されるだけでなく、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又は全部を行い、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。

【0078】

更に、記録媒体から読み出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書込まれた後、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又は全部を行い、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。

【図面の簡単な説明】

【0079】

【図 1】実施例 1 における文書処理装置の外観を示す図である。

【図 2】実施例 1 における文書処理装置の構成の一例を示すブロック図である。

【図 3】文書処理装置における文書の電子化処理の概要を示す図である。

【図 4】実施例 1 における像域分離処理を説明するための図である。

【図 5】像域分離処理 303 で分離された各ブロックに対するブロック情報と入力ファイル情報を示す図である。

【図 6】アウトラインベクトルデータへの変換を説明するための図である。

【図 7】表処理部 308 における表処理を示すフローチャートである。

10

20

30

40

50

【図 8】表枠を認識する処理を示すフローチャートである。

【図 9】矩形図形、三角図形、矩形図形の集合と判定されるアウトラインの一例を示す図である。

【図 10】図 9 に示すセルアウトラインをセル図形へ変換した例を示す図である。

【図 11】隣接するセル図形の統合例を示す図である。

【図 12】セル図形を構成するセルアウトラインと各属性情報についての一例を示す図である。

【図 13】マッピング領域と交点の一例を示す図である。

【図 14】セル図形マッピング処理を示すフローチャートである。

【図 15】マッピング領域と交点とマッピングセルと注目点の関係を示す図である。

【図 16】矩形領域 A 内に色塗りセルを作成する例を示す図である。

【図 17】作成された交点の隣接関係と、マッピングされたセル図形及び色塗りセル、またセル図形の場合はその属性情報を記述した例を示す図である。

【図 18】図 17 に示す交点の隣接関係から得られた処理結果を示す図である。

【図 19】色塗りセル再処理を示すフローチャートである。

【図 20】ドキュメント・アナリシス・アウトプット・フォーマット (D A O F) のデータ構造を示す図である。

【図 21】電子文書作成処理の全体の概略を示すフローチャートである。

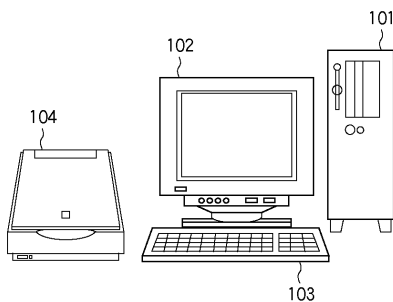
【図 22】文書構造ツリー生成処理の詳細を示すフローチャートである。

【図 23】ページの構成とその文書構造のツリーを示す図である。

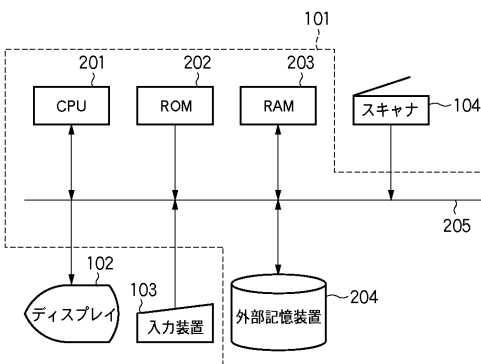
【図 24】本実施例において作成される表枠を示す図である。

【図 25】従来の表処理における問題を説明するための図である。

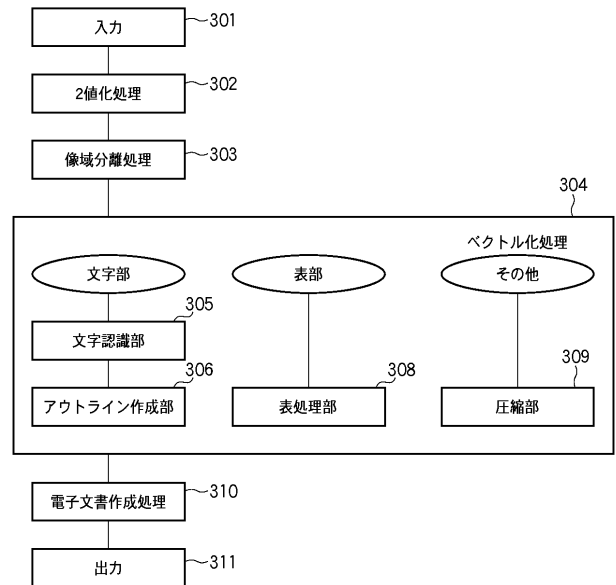
【図 1】



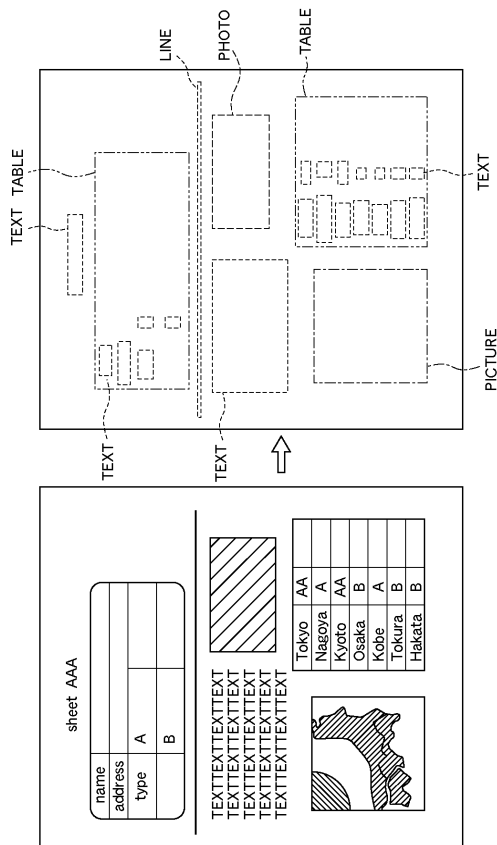
【図 2】



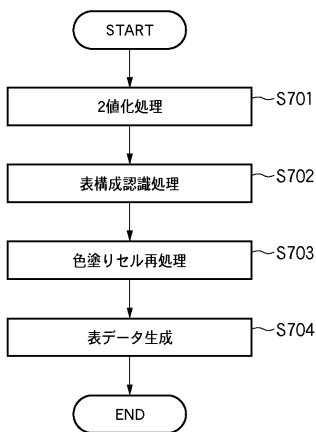
【図 3】



【図 4】



【図 7】



【図 5】

ブロック情報

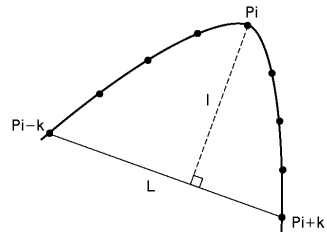
	属性	座標X	座標Y	幅W	高さH	OCR情報
ブロック1	1	X1	Y1	W1	H1	有
ブロック2	3	X2	Y2	W2	H2	有
ブロック3	2	X3	Y3	W3	H3	無
ブロック4	1	X4	Y4	W4	H4	有
ブロック5	3	X5	Y5	W5	H5	有
ブロック6	5	X6	Y6	W6	H6	無

* 属性 1: text 2: picture 3: table 4: line 5: photo

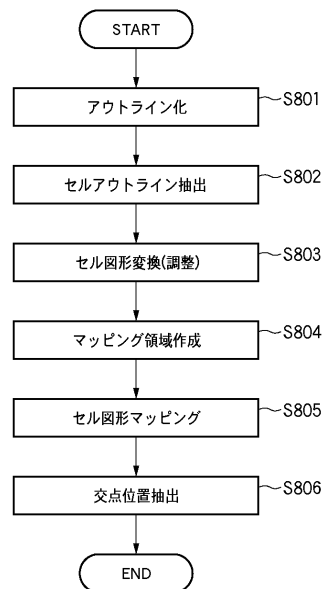
入力ファイル情報

ブロック総数 N (=6)

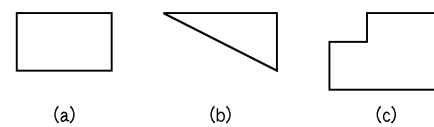
【図 6】



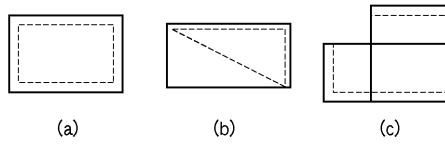
【図 8】



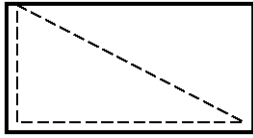
【図 9】



【図 10】



【図 11】

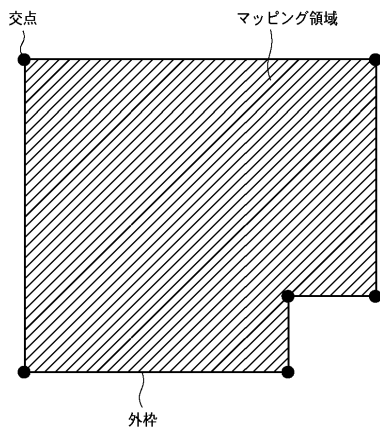


【図 12】

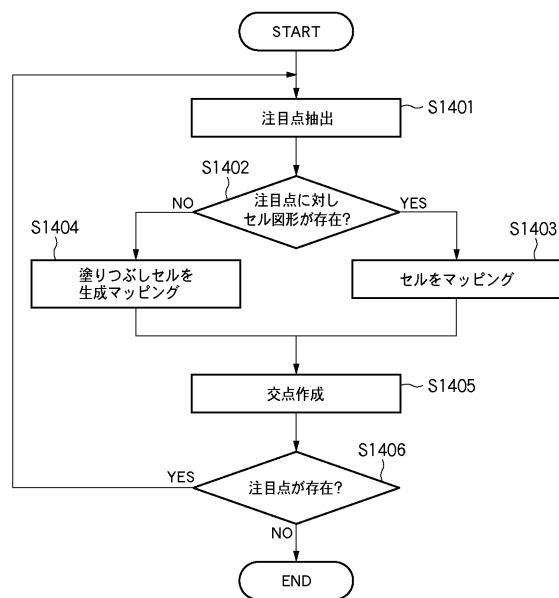
属性名	セルアウトライン
属性1	
属性2	
属性3	
属性4	
属性5	
属性6	
属性7	
属性8	
属性9	

⋮

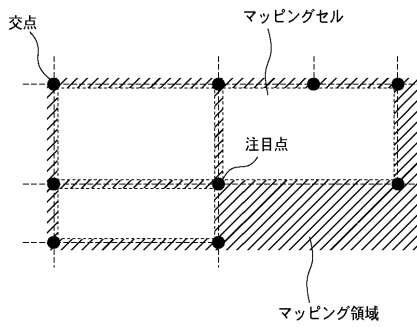
【図 13】



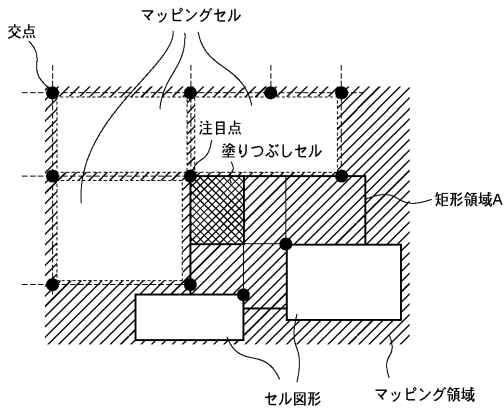
【図 14】



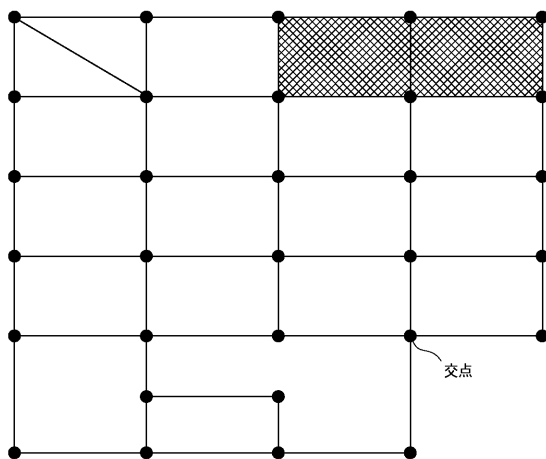
【図 15】



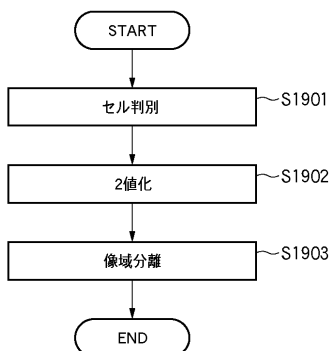
【図 16】



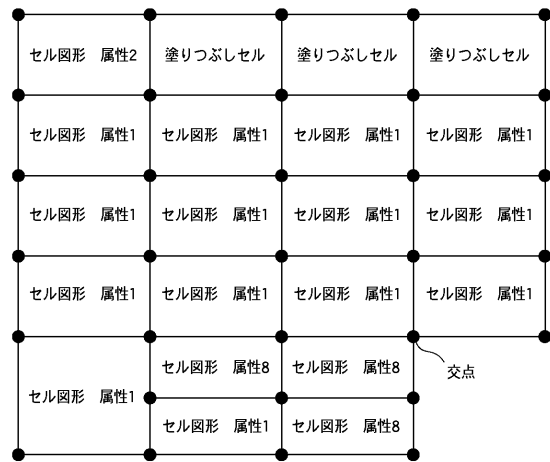
【図 18】



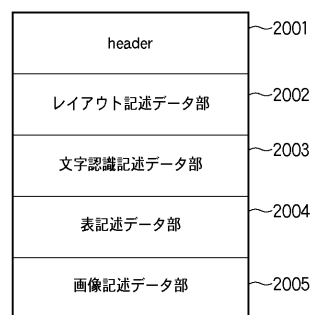
【図 19】



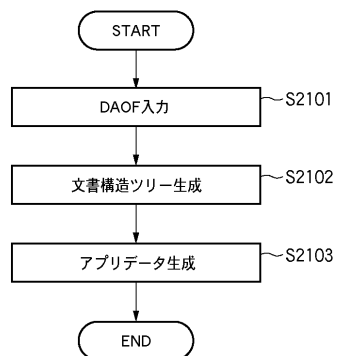
【図 17】



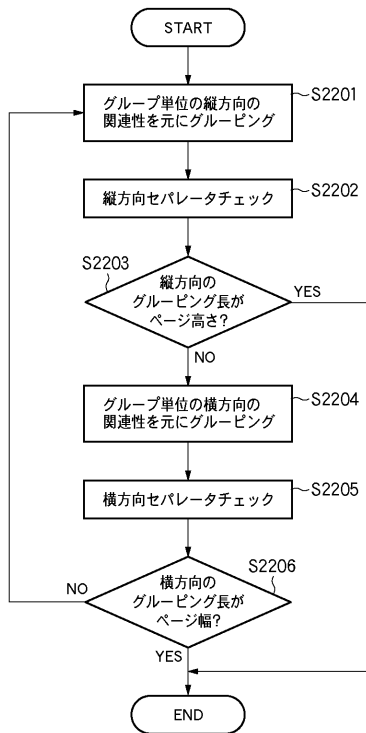
【図 20】



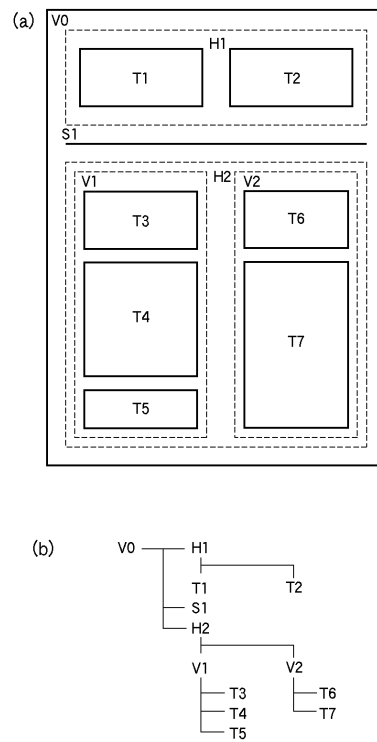
【図 21】



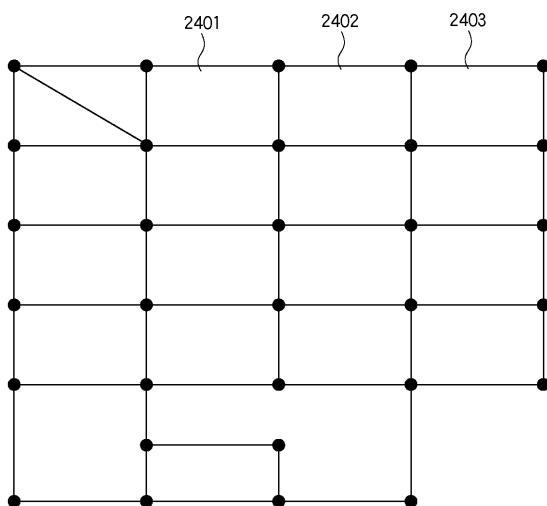
【図 2 2】



【図 2 3】



【図 2 4】



【図 2 5】

