

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3950498号
(P3950498)

(45) 発行日 平成19年8月1日(2007.8.1)

(24) 登録日 平成19年4月27日(2007.4.27)

(51) Int. Cl.		F I			
G06T	7/60	(2006.01)	G06T	7/60	200C
G06K	9/20	(2006.01)	G06K	9/20	340J

請求項の数 11 (全 25 頁)

(21) 出願番号	特願平8-221835	(73) 特許権者	000001007
(22) 出願日	平成8年8月6日(1996.8.6)		キヤノン株式会社
(65) 公開番号	特開平9-171557		東京都大田区下丸子3丁目30番2号
(43) 公開日	平成9年6月30日(1997.6.30)	(74) 代理人	100125254
審査請求日	平成15年8月1日(2003.8.1)		弁理士 別役 重尚
(31) 優先権主張番号	08/514, 250	(74) 代理人	100118278
(32) 優先日	平成7年8月11日(1995.8.11)		弁理士 村松 聡
(33) 優先権主張国	米国 (US)	(74) 代理人	100138922
			弁理士 後藤 夏紀
		(74) 代理人	100136858
			弁理士 池田 浩
		(74) 代理人	100135633
			弁理士 二宮 浩康

最終頁に続く

(54) 【発明の名称】 イメージ処理方法及び装置

(57) 【特許請求の範囲】

【請求項1】

入力されたイメージ内に含まれる連結成分に基づいて特定される領域の外接矩形を得る外接矩形取得工程と、

前記入力されたイメージ内に含まれる連結成分に基づいて特定される領域の輪郭情報を得る輪郭情報取得工程と、

前記外接矩形取得工程で得た複数の外接矩形のうち、重複する外接矩形を判別する重複判別工程と、

前記重複判別工程で重複すると判別された外接矩形が、当該外接矩形に対応する所望の連結成分のイメージデータ抽出に影響するかどうか判別する影響判別工程と、

前記重複判別工程で重複すると判別された外接矩形に関して、前記輪郭情報に基づいて複数の矩形に分解する分解工程と、

前記分解工程で複数の矩形に分解された場合には、当該分解された複数の矩形に基づいてイメージデータを抽出し、前記分解工程で複数の矩形に分解されなかった場合には、当該分解されなかった外接矩形に基づいてイメージデータを抽出する抽出工程と、

前記抽出工程で抽出されたイメージデータを処理する処理工程とを有し、

前記分解工程では、前記影響判別工程で影響すると判別された場合、前記重複判別工程で重複すると判別された外接矩形を前記輪郭情報に基づいて複数の矩形に分解し、前記影響判別工程で影響しないと判別された場合、前記重複判別工程で重複すると判別された外接矩形の分解を行わないことを特徴とするイメージ処理方法。

10

20

【請求項 2】

前記外接矩形取得工程では、各外接矩形の左上角と右下角の座標を求める工程を含み、前記重複判別工程では各外接矩形の座標同士を比較することにより、重複する外接矩形を判別することを特徴とする請求項 1 に記載のイメージ処理方法。

【請求項 3】

前記輪郭情報取得工程では、前記入力されたイメージ内に含まれる連結成分に基づいて特定される領域のタイプに応じて、前記輪郭情報を得ることを特徴とする請求項 1 に記載のイメージ処理方法。

【請求項 4】

前記領域のタイプがテキストタイプである場合、テキスト連結成分の頂部、底部、左側及び右側の縁に基づいて、前記輪郭情報を得ることを特徴とする請求項 3 に記載のイメージ処理方法。

10

【請求項 5】

前記領域のタイプが内部成分を持たない非テキストタイプである場合、連結成分の輪郭を連結することにより、前記輪郭情報を得ることを特徴とする請求項 3 に記載のイメージ処理方法。

【請求項 6】

前記領域のタイプが内部成分を有する非テキストタイプである場合、外側の輪郭と、内部の白輪郭とに基づいて、前記輪郭情報を得ることを特徴とする請求項 3 に記載のイメージ処理方法。

20

【請求項 7】

前記影響判別工程で影響すると判別された場合の前記重複すると判別された外接矩形が、テキストタイプ領域の外接矩形と非テキストタイプ領域の外接矩形とである場合に、前記分解工程では、前記非テキストタイプ領域の輪郭情報と前記テキストタイプ領域の連結成分とに基づいて、複数の矩形に分解することを特徴とする請求項 1 に記載のイメージ処理方法。

【請求項 8】

前記影響判別工程で影響すると判別された場合の前記重複すると判別された外接矩形が、テキストタイプ領域の外接矩形と非テキストタイプ領域の外接矩形とである場合に、前記分解工程では、前記非テキストタイプ領域の輪郭情報と前記テキストタイプ領域の輪郭情報とに基づいて、複数の矩形に分解することを特徴とする請求項 1 に記載のイメージ処理方法。

30

【請求項 9】

前記影響判別工程で影響すると判別された場合の前記重複すると判別された外接矩形が内部成分を持たない非テキストタイプ領域の外接矩形同士である場合、前記分解工程では、前記非テキストタイプ領域の輪郭情報に基づいて、複数の矩形に分解することを特徴とする請求項 1 に記載のイメージ処理方法。

【請求項 10】

前記影響判別工程で影響すると判別された場合の前記重複すると判別された外接矩形が内部成分を有する非テキストタイプ領域である場合、前記分解工程では、前記内部成分の輪郭情報と、非テキストタイプ領域内部の白輪郭の輪郭情報とに基づいて、複数の矩形に分解することを特徴とする請求項 1 に記載のイメージ処理方法。

40

【請求項 11】

入力されたイメージ内に含まれる連結成分に基づいて特定される領域の外接矩形を得る外接矩形取得手段と、

前記入力されたイメージ内に含まれる連結成分に基づいて特定される領域の輪郭情報を得る輪郭情報取得手段と、

前記外接矩形取得手段で得た複数の外接矩形のうち、重複する外接矩形を判別する重複判別手段と、

前記重複判別手段で重複すると判別された外接矩形が、当該外接矩形に対応する所望の

50

連結成分のイメージデータ抽出に影響するかどうか判別する影響判別手段と、

前記重複判別手段で重複すると判別された外接矩形に関して、前記輪郭情報に基づいて複数の矩形に分解する分解手段と、

前記分解手段で複数の矩形に分解された場合には、当該分解された複数の矩形に基づいてイメージデータを抽出し、前記分解手段で複数の矩形に分解されなかった場合には、当該分解されなかった外接矩形に基づいてイメージデータを抽出する抽出手段と、

前記抽出手段で抽出されたイメージデータを処理する処理手段とを有し、

前記分解手段は、前記影響判別手段で影響すると判別された場合、前記重複判別手段で重複すると判別された外接矩形を前記輪郭情報に基づいて複数の矩形に分解し、前記影響判別手段で影響しないと判別された場合、前記重複判別手段で重複すると判別された外接矩形の分解を行わないことを特徴とするイメージ処理装置。

10

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、文書頁のテキスト領域と非テキスト領域を区別するために該文書頁のイメージを分析するブロック選択（または特徴抽出）装置であって、テキスト領域の各グループの周囲と非テキスト領域の各グループの周囲とに外接する矩形を定め、その後で重複する矩形を同定し、それらの矩形の重複する部分を分解することによって該重複する部分の除去を可能にする装置に関する。

【0002】

20

【従来の技術】

従来の特徴抽出装置においては、文書頁の異なるタイプのイメージデータは、まず、該イメージデータがテキストタイプであるかまたは非テキストタイプであるかに応じて同定され、次にイメージデータの領域が、種類に応じて一緒にグループ化（または「ブロック化」）される。イメージデータの各ブロックが、その後で、更なる処理のために抽出される。即ち、従来の特徴抽出装置は、ブロック選択ルーチンを実行して、前記イメージデータ内の連結成分を同定し、非テキストタイプの連結された成分からテキストタイプの連結成分を分離し、別々に該テキスト及び非テキスト連結成分を好ましくは矩形のブロック内にグループ化し、そして非テキスト連結成分を更なる分析に付して、表データ、グラフデータ、線画イメージ、ハーフトーンイメージ、フレーム等の非テキスト連結成分の特定の種

30

【0003】

一度イメージデータの特徴が抽出されると階層木を、該木の各ノードがブロック選択中に定められたイメージデータの各ブロック矩形グループに1対1対応するように設けて、設定することができる。階層木の各ノードには、属性情報がイメージデータの各ブロック矩形片毎に記憶される。即ち、該属性情報は、イメージ情報のブロックがテキストタイプであるか非テキストタイプであるかについての情報を含み、テキストタイプ情報は、さらに題名領域、見出し領域、テキスト領域、などに分類され、非テキストタイプ情報は、さらに、表情情報、線画情報、グラフ情報、ハーフトーンイメージ情報などに分類されていることが可能である。なお、階層木中のノードの位置は、文書頁内のイメージ情報の対応するブロック矩形の位置を暗に記憶している。文書のイメージ情報内の矩形ブロックの場所との組合せで、階層木により、光学式文字認識（OCR）、データ及び/またはイメージ圧縮、データルーチング、データ抽出、保存、検索などのようなその後の適切な処理のための情報の各ブロックの抽出が可能になる。例えば、テキストデータとして指定されたイメージデータのブロックは、適切なOCR処理に付しうるが、ピクチャデータと指定されたイメージデータのブロックは、データ圧縮に付し得るであろう。その結果、多様な異なる文書頁のいずれのイメージデータであれ、任意に入力し、オペレータの介入なしに自動的に処理可能となる。

40

【0004】

イメージデータの各ブロックの秩序だった処理のために、多くの従来の特徴抽出装置は、

50

ブロックが重複していないという仮定に依存していた。そのような仮定は、図2に示した文書などの、非常に多くの文書において正しい。該図から分かるように、代表的な文書頁1は、2コラムフォーマットで配列され、題名領域2、水平な線3、各々テキストデータの行を含むテキスト領域4、5及び6、ハーフトーンピクチャ領域7、フレーム領域8、ならびに表9を含む。米国特許出願第07/873、012号「文字識別の方法及び装置」ならびに米国特許出願第08/338、781号「頁分析装置」に記載されたブロック選択技術、によると、文書頁1の各領域が、その中に位置するイメージデータのタイプに従って同定及び指定され、次いでイメージデータがその各々のタイプに基づいて区分される。ブロック選択処理の結果、各イメージ領域の矩形ブロックが図3に示すように作成され、対応する階層木が形成される。即ち、図3に示すように、文書1に対応するイメージデータ11の場合は、ブロック選択により、題名ブロック12、テキストブロック14、15及び16などの多様なテキストタイプ領域ならびに線ブロック13、ハーフトーンイメージブロック17、フレーム領域18及び表領域19等の多様な非テキストブロックが定められる。

10

【0005】

図3に示すように、非テキストブロックのあるものはその中にテキストブロックを含みうるが、矩形ブロックのいずれも他のブロックに重複していない。例えば、フレーム領域18は、非テキスト線画領域18a及びテキスト領域18bを含み、表領域19は、総じて19aで指定されるテキストタイプ表記載事項を含む。

【0006】

20

【発明が解決しようとする課題】

上に述べたように、イメージデータ11内のイメージ領域のブロック矩形指定に基づき、各情報領域のイメージデータをその後の適切な処理のために抽出することができる。さらに、そして同じく前述したように、矩形ブロックのいずれも他と重複していないので、抽出は直接的である。

【0007】

ブロック矩形領域が互いに他と重複するときは困難が生じる。そのような重複は、例えば、文書の全体の外見に深く影響する編集書式の単純な変更とともに生じ得る。かくして、図4に示すように、フレーム領域8を文書の2つのコラムにまたがる8aで示される位置に移動する編集書式の単純な変更が行われている。従来のブロック選択及び特徴抽出技術では、図5に示すように、ブロック18aがブロック14及びブロック15に重なるブロック矩形フォーマットが得られる結果となる。そのように定められた矩形ブロックに基づいてイメージデータを抽出するときになると、これらのブロックの重複は困難を引き起こす。即ち、領域4のテキストイメージデータにのみ対応すると想定されているブロック14についてイメージデータを抽出する時になると、非テキストフレーム領域18aの不要な部分も抽出されることが分かる。同様に、テキストタイプイメージデータのみを含むと推定されている領域15のイメージデータを抽出するときになると、非テキストフレーム領域18aの不要なもう一つの部分が抽出されることが分かる。かくして、特徴抽出またはブロック選択技術により設定されるブロック矩形領域における重複は、現在までのところ困難を生じている。

30

40

【0008】

本発明は、上述の問題点に鑑み、特徴抽出またはブロック選択技術により設定されるブロック矩形領域における重複の問題を克服できるイメージ処理方法及び装置を提供することを目的とする。

【0009】

【課題を解決するための手段】

本発明のイメージ処理方法は、上記目的を達成するため、入力されたイメージ内に含まれる連結成分に基づいて特定される領域の外接矩形を得る外接矩形取得工程と、前記入力されたイメージ内に含まれる連結成分に基づいて特定される領域の輪郭情報を得る輪郭情報取得工程と、前記外接矩形取得工程で得た複数の外接矩形のうち、重複する外接矩形を

50

判別する重複判別工程と、前記重複判別工程で重複すると判別された外接矩形が、当該外接矩形に対応する所望の連結成分のイメージデータ抽出に影響するかどうか判別する影響判別工程と、前記重複判別工程で重複すると判別された外接矩形に関して、前記輪郭情報に基づいて複数の矩形に分解する分解工程と、前記分解工程で複数の矩形に分解された場合には、当該分解された複数の矩形に基づいてイメージデータを抽出し、前記分解工程で複数の矩形に分解されなかった場合には、当該分解されなかった外接矩形に基づいてイメージデータを抽出する抽出工程と、前記抽出工程で抽出されたイメージデータを処理する処理工程とを有し、前記分解工程では、前記影響判別工程で影響すると判別された場合、前記重複判別工程で重複すると判別された外接矩形を前記輪郭情報に基づいて複数の矩形に分解し、前記影響判別工程で影響しないと判別された場合、前記重複判別工程で重複すると判別された外接矩形の分解を行わないことを特徴とする。

10

また、本発明のイメージ処理装置は、入力されたイメージ内に含まれる連結成分に基づいて特定される領域の外接矩形を得る外接矩形取得手段と、前記入力されたイメージ内に含まれる連結成分に基づいて特定される領域の輪郭情報を得る輪郭情報取得手段と、前記外接矩形取得手段で得た複数の外接矩形のうち、重複する外接矩形を判別する重複判別手段と、前記重複判別手段で重複すると判別された外接矩形が、当該外接矩形に対応する所望の連結成分のイメージデータ抽出に影響するかどうか判別する影響判別手段と、前記重複判別手段で重複すると判別された外接矩形に関して、前記輪郭情報に基づいて複数の矩形に分解する分解手段と、前記分解手段で複数の矩形に分解された場合には、当該分解された複数の矩形に基づいてイメージデータを抽出し、前記分解手段で複数の矩形に分解されなかった場合には、当該分解されなかった外接矩形に基づいてイメージデータを抽出する抽出手段と、前記抽出手段で抽出されたイメージデータを処理する処理手段とを有し、前記分解手段は、前記影響判別手段で影響すると判別された場合、前記重複判別手段で重複すると判別された外接矩形を前記輪郭情報に基づいて複数の矩形に分解し、前記影響判別手段で影響しないと判別された場合、前記重複判別手段で重複すると判別された外接矩形の分解を行わないことを特徴とする。

20

本発明は、上述の困難に対してブロック選択及び特徴抽出処理の結果得られる重複矩形を同定し、該重複矩形をより小さい非重複矩形に分解することにより対処する。好ましくは、重複する矩形が全て分解されるのではなく、むしろ該重複矩形がまず分解が必要かどうかを判別するために分析されて、第1の矩形についてイメージデータを抽出するときに第2の矩形中の不要なイメージデータを抽出するのを回避する。

30

【0010】

好ましい形態においては、重複矩形の分解は、各矩形内のイメージデータの輪郭ペアに基づく。即ち、多くの水平の行にわたるイメージデータの場合、隙間のない輪郭が定められ、そこから各イメージ行について輪郭ペアを得ることが可能であるが、該輪郭ペアはその線についてイメージデータの最左端及び最右端を定める。輪郭ペアは外接矩形内のイメージデータの3つの互いに排他的な異なる種類のイメージデータの各々について別様に得られる。即ち、輪郭ペアは、(1)テキストデータ、(2)内部対象物を全く含まない非テキストデータ及び(3)内部対象物を含む非テキストデータについて、別様に得られる。

【0011】

40

イメージデータの重複する矩形ブロックは、その重複がしかるべきイメージ抽出に干渉するかどうかを判別するために分析される。もし重複がしかるべきイメージ抽出に干渉しない場合には、分解は行われず、輪郭ペアは用いられない。他方、矩形ブロックにおける重複がしかるべきイメージ抽出に干渉すると判別された場合は、輪郭ペアを用いて重複矩形を分解する。即ち、分解は、関係のある重複矩形の種類に応じて別様に行われる。即ち、分解は、(1)他のタイプ(即ちテキストまたは非テキスト)ブロックに重複するテキストブロック、(2)内部コンテンツを含まず、他の種類のブロックに重複する非テキストブロック、及び(3)内部コンテンツを含み他の種類のブロックに重複する非テキストブロックについて、別様に行われる。

【0012】

50

かくして、本発明の１つの形態においては、テキスト領域と非テキスト領域の両方を含むイメージデータを処理する処理方法が提供される。本方法は、ブロック選択を行って、イメージデータ内のテキストタイプ領域の各ブロックの周囲及びイメージデータ内の非テキスト領域の各ブロックの周囲に外接する矩形を得、各テキスト及び非テキストブロックについて輪郭ペアを得るステップを含む。その次に、外接する矩形が重複するかどうかが判別される。輪郭ペアに基づいて重複矩形を分解し、非重複矩形の場合、外接矩形に基づいて、重複矩形の場合、分解された矩形に基づいて、イメージデータが抽出される。次に抽出されたイメージデータが処理される。

【 0 0 1 3 】

この簡単な概要は、本発明の本質が速やかに理解されるように提供したものである。本発明のより完全な理解は、添付の図面に基づく好ましい形態の下記の説明を参照することにより得ることができる。

【 0 0 1 4 】

【発明の実施の形態】

図 6 は、本発明に係る重複矩形分析及び分解を含む本発明の代表的な実施の形態の外観を示す図である。図 6 には、プログラム化された汎用コンピュータが示されているが、本発明は、他のイメージ処理装置に組み込み可能な、専用、ROM ベースまたは据え付け型の装置などの他の装置において具体化可能であることは理解されねばならない。

【 0 0 1 5 】

図 6 には、マイクロソフト Windows OS 等のウインドウ型オペレーションシステムを有する IBM PC または PC 互換機等のコンピュータ機器 4 1 0 が示してある。コンピュータ機器 4 1 0 には、モノクロまたはカラーディスプレイモニタ 4 1 2 が設けられ、ユーザに対しイメージを表示する。コンピュータ機器 4 1 0 には更に、文書イメージファイル等のデータファイルならびにブロック選択及び重複矩形分析分解プログラム等のアプリケーションプログラムファイルを記憶する固定ディスクドライブ 4 1 1 が設けられている。さらにコンピュータ機器 4 1 0 には、テキストデータを入力し、ディスプレイ 4 1 2 の表示画面に表示される対象の操作を可能にするキーボード 4 1 3 と、ディスプレイ 4 1 2 に表示される対象を指示したり操作したりするための、マウス等のポインティング・デバイス 4 1 4 とが設けられる。

【 0 0 1 6 】

複数頁を有する文書はスキャナ 4 1 6 によって入力される。スキャナ 4 1 6 は文書の各頁または他のイメージをスキャンして、これらの頁のビットマップイメージデータをコンピュータ機器 4 1 0 に供給する。該イメージデータは、圧縮または非圧縮フォーマットでディスク 4 1 1 に記憶される。

【 0 0 1 7 】

コンピュータ機器 4 1 0 が処理した文書を出力するために従来のカラープリンタ 4 1 8 が設けられる。

【 0 0 1 8 】

更に、ローカル・エリア・ネットワークとインターフェースするためのネットワーク・インターフェース 4 2 4、及びファクシミリ／データモデムを介してファクシミリ・メッセージ及び他のデータファイルを送受信するためのファクシミリ／モデム・インターフェース 4 2 6 が設けられる。そのようなインターフェースは、ドキュメントイメージデータを入力するためのスキャナ 4 1 6 に加えて、またはその代わりに用いることができる。

【 0 0 1 9 】

オペレータの指示に従って且つウインドウ型オペレーションシステムの制御のもとで、デスクトップ・パブリッシング・プログラム、ドローイング・アプリケーション・プログラム、ブロック選択アプリケーション等の記憶されたアプリケーション・プログラムが選択的に起動され、データを処理したり操作したりする。また、オペレータの指示に従って且つこれらの記憶されたアプリケーション・プログラムに基づいて、イメージをモニタ 4 1 2 に表示したり、モニタ 4 1 2 に表示されているイメージをプリンタ 4 1 8 で印刷するよ

10

20

30

40

50

うにコマンドが発せられる。

【0020】

最も典型的には、本発明を具体化するブロック選択プログラムや重複矩形分析・分解プログラムを含むコンピュータディスク411に記憶されたアプリケーション・プログラムは、フロッピーディスク、CD-ROM、またはコンピュータ掲示板などのコンピュータ読み出し可能媒体から各アプリケーションをダウンロードしてディスク411に記憶したものである。

【0021】

図7は、コンピュータ機器410の内部構成を示す詳細ブロック図である。図7に示すように、コンピュータ機器410は、コンピュータバス521に接続された、プログラム式マイクロプロセッサ等から成る中央演算処理装置(CPU)520を有する。また、スキャナ・インターフェース522、プリンタ・インターフェース523、ネットワーク・インターフェース524、ファクシミリ・モデム・インターフェース526、ディスプレイ・インターフェース527、キーボード・インターフェース528、及びマウス・インターフェース529が、コンピュータバス521に接続されている。

10

【0022】

CPU520がアクセスできるように、ランダム・アクセス・メモリ(RAM)等の主メモリ530がコンピュータバス521に接続されている。特に、ディスク411に記憶されたアプリケーション・プログラムに関連する命令等、記憶されたアプリケーション・プログラムの命令列を実行するときに、CPU520は、これらの命令列を、ディスク411(またはネットワーク・インターフェース524を介してまたはフロッピー・ディスク・ドライブ(図示せず)を介してアクセスされる媒体等の他のコンピュータ読み出し可能な媒体)から主メモリ530にロードして、主メモリ530からこれらの記憶されたアプリケーション・プログラムの命令列を読み出して実行する。

20

【0023】

図1は、コンピュータで実現される方法を示すフローチャートであって、文書頁を表す入力イメージが入力され、該入力イメージに対してブロック選択が実行されて文書イメージ内のテキスト及び非テキストブロックの周囲の外接矩形を定め、輪郭ペアを各ブロックについて得、外接矩形を分析してそれらが重複しているか否かを判別し、重複していた場合は、該重複がイメージ抽出に影響するかどうかを判別し、次に、それらがイメージ抽出に影響すると判別された場合は、重複矩形を分解する。即ち、図1に示すように、分解は、重複矩形が、テキストブロック、内部コンテンツを含まない非テキストブロック、または内部コンテンツを含む非テキストブロックに外接するかどうかに応じて別様に行われる。図1に示された処理ステップは、上述のようにコンピュータ読み出し可能な媒体に記憶され、それらの処理ステップを主メモリ530にロードしてそこから実行するCPU520によって実行される。

30

【0024】

即ち、ステップS501で、文書頁を表すイメージデータが入力される。前述したように、イメージデータはスキャナ416を介して入力してよく、またはディスク411上や、ネットワーク・インターフェース424もしくはファクシミリ/モデム・インターフェース426を介して記憶されたイメージデータファイルから得てもよい。

40

【0025】

ステップ502において、入力イメージに対してブロック選択が実行されて、テキスト及び非テキストブロックの周囲の外接矩形を得、また各ブロックについて属性情報と位置情報を記憶する階層木を得る。好適なブロック選択技術は、前述の米国出願第07/873,012号及び第08/388,781号に記載されており、その内容は、完全に記述したのと同じようにここに引用により加入する。

【0026】

ステップS503においては、ステップS502のブロック選択中に既に輪郭ペアが算出されていない程度だけ、輪郭ペアが、ステップS502でブロック選択により定められた

50

各ブロック矩形領域について得られる。即ち、各ブロックは、通常多数の水平の線にわたる。ブロックの中に入るそのような行の各々について、該ブロック中の最左端及び最右端の領域を定める輪郭ペアがステップS503で得られる。輪郭ペアは、矩形ブロック内のデータの種類に応じて別様に得られる。即ち、輪郭ペアは、該ブロックが(1)テキストブロック、(2)内部コンテンツを含まない非テキストブロックまたは(3)内部コンテンツを含む非テキストブロックであるかどうかに応じて、別様に定められる。ステップS503にしたがって輪郭ペアを得る方法については、図8及び図9を参照して下記において詳述する。

【0027】

ステップS504では、ステップS502で得た外接矩形が重複するかどうかを判別する。外接矩形が重複するかどうかの判別は、図13を参照して、下記に詳述する。もし重複する外接矩形がなければ、本発明によりさらに処理を行う必要はなく、フローは、ステップS509までスキップする。

10

【0028】

他方、ステップS504で外接矩形が重複すると判別されたときは、ステップS505で、重複がイメージ抽出に影響するかどうかを判別する。一般的に、ほとんどの重複は、イメージ抽出に影響するが、中にはイメージ抽出に影響しないものも存在しうる。例えば、重複するけれども、第1の外接矩形が第2の外接矩形のイメージデータを全く含まないかまたはその逆であるという点で、互いに排他的である2つの外接矩形がブロック選択により定められることはあり得る。そのような外接矩形は、重複しているが、イメージ抽出には影響を与えない。重複がイメージ抽出に影響するかどうかを判別する方法についての詳細は、図14(A)~(C)を参照して詳述する。

20

【0029】

ステップS505で重複がイメージ抽出に影響を与えないと判別された場合は、本発明によりさらに処理を行う必要はなく、フローは、ステップS509までスキップして進む。他方、重複がイメージ抽出に影響する場合は、フローは、ステップS506、S507及びS508に進み、各ステップで重複ブロックをそれぞれ分解してイメージ抽出に影響しないようにする。

【0030】

ステップS506、S507及びS508では、重複するブロックの種類に応じて重複ブロックが別様に分解される。かくして、ステップS506では、他の種類のブロック、即ちテキストまたは非テキストブロックに重複するテキストブロックが分解される。ステップS507では、内部コンテンツを含まない重複非テキストブロックが分解される。そしてステップS508では、内部コンテンツを含む非テキストブロックが分解される。階層木は、ステップS506、S507またはS508においてどのような分解が行われたかに基づいて更新される。

30

【0031】

重複テキストブロック及び非テキストブロックをステップS506またはS507またはS508で分解した後で、フローはステップS509に進み、イメージデータがイメージデータの境界を定める外接矩形に応じて抽出される。かくして、例えば、重複矩形が存在しない状況では、イメージデータは、ステップS502でブロック選択により定められた外接矩形に応じて抽出される。他方、重複矩形が存在し、ステップS506、S507またはS508で分解された場合は、各ブロックのイメージ抽出がそれらのステップで決定された分解ブロックに応じて実行される。

40

【0032】

フローは、次に抽出イメージデータが適切に処理されるステップS510に進む。例えば、上述のように、テキストタイプイメージデータを抽出する時は、適切な処理は、テキストデータ内の文字イメージの同一性を判別するための光学式文字認識であってよい。同様に、非テキスト表領域用の適切な処理には、該表内に含まれるテキスト用のOCR処理が含まれてよい。さらに別の例として、非テキストハーフトーンイメージデータの適切な処

50

理には、より小さな記憶領域にハーフトーンピクチャを保存することを可能にする簡単なイメージ圧縮が含まれてよい。

【0033】

図1の処理は、必要に応じて、ブロック選択及びイメージ抽出が望まれる各文書頁のイメージデータについて繰り返される。

【0034】

図8は、図1の上記ステップS503で簡単に説明したように、イメージ情報の各矩形ブロックについて輪郭ペアを得る方法を示すフローチャートであり、図9は、「輪郭ペア」が何を意味するかを説明する図である。

【0035】

簡単に図9を参照すると、ブロック選択ステップS502で定められたイメージデータの任意のブロック40は、垂直方向に多数の走査線を含み、該走査線の各々はブロックを水平に横切って延びている。任意に形作られた連結成分41は、例えば、文字のイメージまたはその他のイメージであってよいが、やはり多数の走査線を横切って垂直方向に延びている。「輪郭ペア」は、対象物41を含む各走査線について定められる。走査線上の各輪郭ペア（対象物が違う位置で走査線を横切る場合は各走査線には2個以上の輪郭ペアが存在しうる）は、正確に2つのポイント：対象物が走査線上で始まる第1（または左側）の点と対象物が走査線上で終わる第2（または右側）の点とを含む。かくして、例えば、図9を参照すると、走査線iには2つの輪郭ペアが含まれ、第1のものは（a, b）からなり、第2のものは（c, d）からなる。輪郭ペア（a, b）には、対象物41が走査線i上で始まる第1（左側）の点aと対象物41が輪郭線i上で終了する第2（右側）の点bが含まれる。同様に、輪郭ペア（c, d）には、対象物41が走査線i上で始まる第1（左側）の点cと対象物が輪郭線i上で終了する第2（右側）の点dが含まれる。走査線「j」の場合は、3つの輪郭ペア、即ち、（e, f）、（g, h）及び（k, l）があり、走査線「k」の場合は、ただ1つの輪郭ペア（m, n）がある。輪郭ペアの内側にあり輪郭ペアを含むイメージは、対象物41の連結成分に属するイメージである。したがって、対象物のイメージは、各走査線に沿う輪郭ペアにのみもついで抽出することができる。

【0036】

図9からよく理解されるように、イメージ中の各対象物が輪郭ペアで表わされれば、イメージのどの部分が次のイメージ処理のために抽出されなければならないかについてもはや混乱はないであろう。他方、かなりのメモリスペースがイメージ内の各輪郭ペアを記憶するのに必要とされ、特にイメージが何千もの輪郭ペアが必要とされるであろうテキストから構成されている時にはそうである。さらに、輪郭ペアの使用は、ブロック矩形フォーマットでより自然に提示されるイメージの非直観的な表現になるために、ユーザにとっては不便である。勿論、ブロック矩形フォーマットは、各矩形に対して単に左上の角と右下の角のみが設定に必要であるので、ずっと少ないメモリ記憶装置条件しか要求しない。重複矩形が発生するときに生じる上述の欠点は、下記においてより詳細に説明されるように、部分的には、輪郭ペアのしかるべき使用法により対処される。

【0037】

図8を参照すると、3つの異なる種類のブロック：テキストブロック、内部成分を含まない非テキストブロック及び内部成分を含むテキストブロックの各々について、輪郭ペアを導出するフローチャートが示されている。よく理解されるように、イメージ内で出会ういかなるブロックもこれらの互いに排他的な3つのカテゴリの一つに属する。図8からさらによく理解されるように、各異なる種類のブロックの輪郭ペアは、それぞれが他の種類のブロックについてから導出される方法とは、別様に導出される。

【0038】

かくして、ステップS601では、ステップS502（図1）で導出された矩形ブロックがテキストブロックであるかどうか判別するために調べられる。テキストブロックである場合は、次に輪郭ペアがステップS602、S603及びS604により得られるが、これらのステップにおいては、該ブロックの連結成分の頂部、底部、左側及び右側の縁がま

10

20

30

40

50

ず得られ、それらの4つの縁が組み合わされて、ブロック内の全ての連結成分の隙間のない輪郭にされ、輪郭ペアがその隙間のない輪郭から作成される。この処理は、図10(A)~(F)に示されている。図10(A)を参照すると、テキスト連結成分を含む矩形ブロックが描かれている。12のテキスト連結成分が示されているが、この数字は、通常生じるよりもずっと少ないものであり、簡潔さのためにのみ示してある。図10(B)~(E)では、図8のステップS602により、全ての連結成分の頂部の縁、底部の縁、左側の縁及び右側の縁がそれぞれ得られる。次に、ステップS603により図10(F)に示されるように、頂部、底部、左側及び右側の縁が組み合わされて全ての連結成分が隙間のない輪郭にされる。4つの縁を組み合わせることで、隙間のない輪郭が、テキストブロック内のテキストの全てを隙間なく囲む閉じたループに形成されることはよく理解されるであろう。最後に、輪郭ペアが図10(F)の隙間のない輪郭の輪郭ペアを得ることにより、得られる(ステップS604)。

10

【0039】

ステップS602ないしS604においては、テキストブロックを囲む隙間のない輪郭を形成するため連結成分の矩形境界を用いると、各連結成分に基づいて隙間のない輪郭を算出する場合に多すぎる時間を費やすことなく良好な結果が得られるであろうと感じられる。しかし、この人為的に作られた輪郭は、他の対象物の輪郭と重複する結果となった場合に、該重複に関わる連結成分の各々の輪郭ペアを用いて重複する部分を修正することができる。

【0040】

20

図8に戻って、ステップS601で矩形ブロックがテキストブロックでないと判別された場合は、フローは、ステップS605に進み、矩形ブロックが内部成分を含まない非テキストブロックであるかどうか判別される。矩形ブロックが内部成分を含まない非テキストブロックである場合は、フローはステップS606及びS607に切り替わり、そこで、輪郭ペアがブロック内の各連結成分について得られ、そのようにして得られた輪郭ペアが連結される。この処理は、図11に示されるが、該図においては、3つの任意な連結成分46、47及び48が示されている。輪郭ペアは、該連結成分の各々について導出され、輪郭ペアを結び組み合わせて、全体に49で示されるように、全非テキストブロック囲む単一の輪郭を形成すべきかどうか考慮される。もし49で示されるような連結が望まれる場合は、連結のための空の道筋、即ち、重複を引き起こさない道筋が見出されねばならない。その後で、連結された輪郭線が輪郭ペアとして出力される。他方、連結が望まれていない場合は、各分離した連結成分の輪郭ペアが出力される。

30

【0041】

図8に戻り、ステップS605で、矩形ブロックが内部成分を含まない非テキストブロックであると判別されなかった場合は、カテゴリは互いに排他的であるので、該矩形ブロックは必然的に内部成分を有する非テキストブロックでなければならない。したがって、フローは、ステップS609及びS610(ステップS608は単に完全さのために示してあり、実際には実行を要しない)に進む。ステップS609とS610においては、非テキスト対象物の外側の輪郭の輪郭ペアのみならず内側白輪郭の輪郭ペアも得られ、ブロックの内部成分を抽出するのを補助する。

40

【0042】

即ち、図12(A)を参照して、フレーム対象物、表対象物または線画対象物等の内部成分を有する非テキスト成分は、対象物の最も外側の輪郭の内側に含まれる白輪郭を有してよい。内側白輪郭は、内部対象物を抽出する補助のために用いられる。かくして、図12(A)に示されるように、そしてステップS609で説明されたように、テキスト成分などの内部成分52を含む任意の非テキストブロック50について、まず輪郭ペアが、該非テキスト対象物の最も外側の輪郭を定める連結成分について得られる。したがって、図12(A)に示してある状況では、輪郭ペアが連結成分51について得られる。その後、ステップS610にしたがって最も外側の輪郭の内側白輪郭について輪郭ペアが得られる。そのような内側白輪郭が53で示されている。(内側白輪郭を得る方法についての詳

50

細な説明は米国特許出願第 07 / 873 , 012 号でなされており、その内容は上述のように引用によりここに加入される。) 最も外側の輪郭を定める連結成分と内側白輪郭の両方についての輪郭ペアが、次に、必要な輪郭ペア情報として出力される。

【0043】

米国特許出願第 07 / 873 , 012 号は、図 12 (C) に示す 4 方向パターンでの内側白輪郭の導出を記述しているが、図 12 (D) に示される 8 方向パターンで内側白輪郭を導出することも可能である。8 方向導出は、非テキストの最も外側の輪郭が単に垂直に配された縁だけではなく斜めの縁も有するような図 12 (B) に示されるような状況の場合に有利である。勿論、該 8 方向パターンは、図 12 (A) に示されるような状況パターンでも使用可能である。8 方向導出は、53a で示されるように斜めの縁が存在する場合でも、内側白輪郭を良好に定めることを可能にする。しかし、8 方向検索パターンが用いられると、各ステップにおいて白輪郭が外側の黒い境界で完全に包まれているかどうかを判別するために白輪郭トレースを検査しなければならない。

10

【0044】

要約すると、図 8 及び図 9 は、ステップ S503 で簡単に触れた処理の詳細な処理を示すが、それによりブロック選択により同定された各矩形ブロックについて輪郭ペアが得られ、またそれにより、輪郭ペアは、該ブロックがテキストブロックであるか、内部成分を含まない非テキストブロックであるか、内部成分を含む非テキストブロックであるかに応じて別様に得られる。

【0045】

20

図 13 及び図 14 (A) ~ (C) は、重複する外接矩形があるかどうかを判別し、重複矩形が存在する場合は、重複がイメージ抽出に影響するかどうかを判別する図 1 の処理ステップ S504 及び S505 を説明する図である。

【0046】

即ち、ステップ S505 では、ブロック選択により同定されたブロックに重複するものがあるかどうかを判別される。そのような重複は、各ブロックを定める 2 つの座標 (即ち、左上の角と右下の角) を他のブロックの対応する座標と比較することにより判別可能である。かくして、図 13 は、4 つのブロック、即ちテキストブロック 54、テキストブロック 55、テキストブロック 56、及び非テキストブロック 57 を含む任意の文書のイメージ 53 を示す。各ブロックについての左上の角及び右下の角と他のブロックについての対応する座標との比較により、文書 53 には重複するブロックがないことが示される。したがって、イメージ 53 中のブロックについては分解は必要でなく、文書 53 の処理は、直接イメージデータの抽出 (ステップ S509) に進むことが可能である。

30

【0047】

図 14 (A) ~ (C) は、矩形ブロックにおける重複がイメージ抽出に影響するかどうかを説明する図である。かくして、図 14 (A) においては、任意の文書のイメージ 60 は 3 つのブロック、即ち、ピクチャブロック 61、ピクチャブロック 62、及びテキストブロック 63 を含む。テキストブロック 63 の左上の角及び右下の角の座標 (即ち、座標 (X1, Y1) 及び (X2, Y2)) の比較により、テキストブロック 63 が非テキストブロック 61 及び非テキストブロック 62 に重複していると判別される。しかし、図 14 (A) の状況下では、テキストブロック 63 がブロック 61 及び 62 に重複しているといっても、テキストブロックがそれらに重複している領域には、ブロック 61 及び 62 のイメージデータは存在しない。かくして、重複にも関わらず、ブロック 63 のイメージ抽出は影響を受けず、テキストブロック 63 の分解は必要とされない。他方、非テキストブロック 61 及び 62 のイメージ抽出は、両方とも、重複により影響される。即ち、ブロック 61 の矩形座標に基づいてイメージデータを抽出する場合は、必要なピクチャデータが抽出されるのみならずブロック 63 から不要なテキストデータ部分も抽出されるであろう。したがって、非テキストブロック 61 及び 62 の両方について、ブロック分解がステップ S506、S507 または S508 にしたがって適切な処理として、必要となる。(ここで、非テキストブロック 61 及び 62 は内部成分を含まないと仮定すると、その場合は、分

40

50

解処理はステップ S 5 0 7 にしたがって行われるであろう)。

【 0 0 4 8 】

図 1 4 (B) は、テキストブロック 6 5 及び非テキストブロック 6 6 を含む任意の文書 6 4 について、重複がブロック 6 5 と 6 6 間に存在する状況を示す。この状況では、図 1 4 (A) の状況とは違って、重複のために、ブロック 6 5 または 6 6 のいずれについてもイメージ情報を抽出するのは、同じく他のブロックについての不要なイメージデータを得ることなしには不可能である。したがって、ブロック 6 5 及び 6 6 の両方がステップ S 5 0 6 ないし S 5 0 8 にしたがって分解に付される。

【 0 0 4 9 】

図 1 4 (A) 及び (B) に示される重複は、編集スタイルによるものであり、グラフィックが単一のページ上でテキストと混合され、インデントされ位置決めされたために生じたものである。しかし、重複は他の原因によるものであり得、したがって、編集スタイルによる重複に厳格に限定されるべきものではない。例えば、重複は、イメージデータ内のスキュー (斜行) によっても発生可能であり、スキューは、意図的なもの (やはり編集スタイルによる) またはある角度で文書を走査したことによる非意図的なもののいずれかであり得る。この状況は、任意の文書 6 7 が第 1 のテキストブロック 6 8 及び第 2 のテキストブロック 6 9 を含む図 1 1 C に示されている。両方のテキストブロックが斜めになっており、そのスキューのため、ブロック 6 8 と 6 9 の間に重複が形成される。その重複のために、他のブロックからの不要な情報も抽出することなしには一つのブロックからのイメージ情報を抽出できない。したがって、ブロック 6 8 及び 6 9 の両方の分解が必要である。

【 0 0 5 0 】

[重複するテキストブロックの分解]

他のテキストブロックまたは他のイメージタイプのブロックに重複するテキストをしかるべく抽出するために、抽出されるべきテキストを含むテキストブロックがより小さい非重複の矩形に分解される。ステップ S 5 0 6 について上で論じた重複するテキストブロックを分解する工程は、図 1 5 (A) ~ (I)、1 6 (A) 及び (B)、図 1 7 ~ 図 1 9、ならびに図 2 0 (A) 及び (B) を参照してより詳細に下記で論じる。

【 0 0 5 1 】

一度重複テキストブロック領域が存在すると判別され、該重複テキストブロック領域がイメージ抽出に影響すると判別された場合は、重複状態がどのように存在しているかが、どの領域が直ぐに分解されるべきか (非重複領域) そしてどの領域がさらに分解を要するか (重複領域) を判別するために吟味される。図 1 5 (A) ~ (I) に示されるように、2つのブロックに重複が発生する仕方には 9 通りの例がある。(これらの 9 つの例の鏡像及び回転が存在してもよい)。例えば、2つの矩形が何らかの重複を有する場合、図 1 5 (D) に示す重複状態等のように、その水平及び垂直の縁が完全にまたは一部重複しているかも知れず、一つの矩形のいずれかの縁が、他の矩形の縁内に完全に入っているかも知れない。

【 0 0 5 2 】

重複領域に関わる領域は、図 1 5 (A) ~ (C) 及び 1 5 (E) ~ (I) の各重複状態に示されるように、少なくとも 2 つのそして多くとも 4 つの矩形に直ぐに分解可能である。勿論、図 1 5 (D) に示すように、2つのブロックが完全に互いに重複する場合がある。分解の第 1 のステップは、重複していない領域を分解して 1 以上の非重複矩形に分解することにより刈り込んで取り去ることである。図 1 5 (A) ~ (C) 及び 1 5 (E) ~ (I) に示すように、非重複領域は非ボールド体の輪郭で示されている。例えば、図 1 5 (A) に示されるように、テキストブロック 1 2 0 は、重複しない第 1 及び第 2 の領域に分解可能であり、テキスト領域 1 2 1 は、テキストブロック 1 2 0 と重複しない第 1 及び第 2 の領域に分解可能である。その結果、直ぐに分解可能な 4 つの非重複領域が得られる。しかし、重複領域 1 2 2 の場合は、更なる分解が必要となる。この工程は下記においてより詳細に論じられるであろう。

【 0 0 5 3 】

10

20

30

40

50

まず、明瞭さのために、分解されるべきテキストブロックを「テキストブロック」と呼び、対象物または該テキストブロックの矩形領域が重複する非テキスト領域を「重複対象物」と呼ぶ。この場合に重複対象物は、テキストブロック、非テキストブロック、または既に分解された矩形でありうる。重複対象物が分解された矩形である場合は、最初の分解処理が、テキストブロックを分解するのに十分であるべきである。重複ブロックがテキストまたは非テキストブロックであれば、分解の第2のステップが、テキストまたは非テキストブロックの「輪郭ペア」を用いて必要となるであろう。

【0054】

さて、図16(A)を参照すると、非テキストブロック131に重複するテキストブロック130の例が示されている。図15(A)で注意した重複状態のように、イメージ抽出されるべきテキストが重複領域132内に存在し、一度非重複ブロックが分解されるとその結果図16(B)に示す重複領域132が得られる。重複領域132が次に、より小さい非重複矩形を作るために更なる分解に付され、テキストブロックが重複領域からすぐに抽出することが可能となる。

10

【0055】

かくして、重複領域132をさらに分解する工程を図17～図19に示すフローチャートを参照して論じる。ステップS1400においては、領域132内のテキストブロックの全ての連結成分が集められる。一度、それらが集められてしまうと、各成分が図20(A)に示すような重複領域に入るように刈り込まれる。例えば、図20(A)に示すように、ブロック151が刈り込まれ、重複領域内のブロックの部分のみが残る。即ち、ブロック151について図20(A)に示されたものは、ブロック全体の一部にすぎない。同じことが、ブロック152及び非テキストイメージ153の残余の部分についても当てはまる。

20

【0056】

ステップS1401では、重複領域132内にある重複非テキスト対象物の全ての輪郭ペアが集められ、上述したようにして非重複部分が重複部分から刈り込まれる。ステップS1402では、連結成分が重複非テキストイメージ内に入らない矩形の組にグループ化される。この点については、ステップS1402での連結成分を非重複矩形の組にグループ化する処理は二つの異なる方法、即ち方法A及び方法Bによって行うことが可能である。方法Aにおいては、領域132内のテキストブロックの刈り込まれた連結成分がステップS1404で用いられ、方法Bにおいては、重複テキストブロックの刈り込まれた連結成分から導出される輪郭ペアが用いられる。方法Aまたは方法Bのいずれかより、テキストブロックが非テキストイメージに重複しない矩形にさらに分解される。

30

【0057】

かくして、方法Aの分解を用いて、ステップS1405において、テキストブロックが水平なテキストブロックであるかどうか判別される。ステップS1405において、水平なテキストブロックであると判別されたときは、フローは、ステップS1407に進み、そこで、水平方向に沿う全ての連結成分が一緒にグループ化される。一度水平方向の全ての連結成分がグループ化されると、フローは、ステップS1408に進み、そこで、垂直方向に沿う全ての連結成分がグループ化される。

40

【0058】

ステップS1409において、重複非テキスト対象物の輪郭ペアのいずれとも重複しない重複領域内の連結成分が残存しているかどうか判別される。もし輪郭ペアと重複しない更なる連結成分が存在しないならば、ステップS1412において、グループ化処理が終了する。グループ化された水平方向連結成分及びグループ化された垂直方向連結成分が次いで、組み合われて4つの非重複矩形154、155、156及び157を形成する。これらの矩形の座標は、非重複領域内のテキストと同様にして重複領域からテキストを抽出するのに利用される。しかし、もっと多くの連結成分が存在する場合は、フローはステップS1407に戻る。

【0059】

50

ステップS 1 4 0 5で、テキストブロックが水平方向のテキストブロックではない場合は、フローは、ステップS 1 4 1 4に進み、そこで連結成分がまず垂直方向に沿って連結される。一度、それらが垂直方向に沿って連結されると、フローは、ステップS 1 4 1 5に進み、その時点で、全ての連結成分が水平方向に沿って一緒にグループ化される。ステップS 1 4 0 9におけると同様に、ステップS 1 4 1 6において、重複非テキストイメージの輪郭ペアのいずれとも重複しない連結成分が残存しているかどうか判別される。もし重複対象物の輪郭ペアのいずれとも重複しない成分が残存していないと判別された場合は、グループ化処理がステップS 1 4 1 2で終了する。グループ化された水平方向連結成分及びグループ化された垂直方向連結成分は、次に、組み合わせられて非重複矩形を形成する。これらの矩形の座標は、記憶され、非重複領域のテキストと同様に重複領域からテキストを抽出するのに使用される。重複しない連結成分が存在する場合は、フローはステップS 1 4 1 4に戻る。

10

【0060】

図18の連結成分のグループ化の結果、今や更なる処理のために抽出可能となったブロック化された領域が得られる。更なる分解により今度は重複領域内のテキストデータが直ぐにしかるべき後処理のために抽出可能となる。

【0061】

一方、方法Bが用いられた(ステップS 1 4 0 6)場合は、ステップS 1 4 2 0で、輪郭ペアが集められた連結成分から導出されるか、即ちテキストブロックの輪郭ペアが図20(B)に示すように重複領域に入るように刈り込まれる。なお、該グループ化は、テキストブロックの連結成分ではなくて、輪郭ペアにより生じるため、図20(A)に示されているのとは少し異なっている。一度輪郭ペアが刈り込まれた連結成分から導出されると、フローは、ステップS 1 4 2 1に進む。ステップS 1 4 2 1では、全ての垂直方向の連結された輪郭ペアと一緒にグループ化される。ステップS 1 4 2 2では、一度連結された垂直方向の輪郭ペアの全てが垂直方向に沿ってグループ化されると、次に水平方向の輪郭ペアと一緒にグループ化される。ステップS 1 4 2 3では、垂直方向及び水平方向グループ輪郭ペアから作られた矩形が重複対象物の輪郭ペアのいずれとも重複しないような垂直方向又は水平方向の連結輪郭ペアが残存するかが判別される。もし、ステップS 1 4 2 3において、残存する連結輪郭ペアが存在しないと判別されると、グループ化処理はステップS 1 4 2 4で終了する。垂直方向及び水平方向両方のグループ化された連結輪郭ペアは重複対象物の輪郭ペアのいずれともまた前に分解された矩形の非グループ化輪郭ペアのいずれとも重複しない矩形に分解される。

20

30

【0062】

上述の処理の結果として、重複テキストブロック内のテキストは、さらにテキストブロックの重複領域をより小さな非重複矩形に分解することにより抽出することができる。

【0063】

[内部コンテンツを含まない非テキストブロックの分解]

上記ステップS 5 0 7で論じた内部コンテンツを含まない重複非テキストブロックを分解する処理を図21(A)~(C)及び図22を参照してさらに詳細に論じる。

【0064】

非テキストブロックの矩形領域がもう一つの非テキストブロックと重複する場合は、非テキストブロックの一方をより小さな非重複矩形に分解することが可能である。より小さな非重複矩形の組によって提供される情報に基づいて該ブロック内の重複非テキストイメージを直ぐに抽出することが可能である。

40

【0065】

明瞭さのために、分解されるべき非テキストブロックを「非テキストブロック」と呼び、該非テキストブロックにより重複されている領域を「重複対象物」と呼ぶ。

【0066】

さて、図21(A)を参照すると、非テキストイメージブロック160及び161が領域162で重複している。上記において重複テキストブロックについて論じたように、非重

50

重複領域は直ぐに最大4つの非重複矩形に分解される。非テキストブロック160及び161の「輪郭ペア」に基づいて、図22のフローチャートに示される処理を用いて、重複領域162、非テキストブロック160及び161内のより小さな非重複矩形の組を作ることが可能である。

【0067】

かくして、ステップS1700において、非テキストブロック160及び重複対象物161の輪郭ペアを用いて、非重複領域を刈り込んで領域162内に入るイメージのみが残存するようにする。即ち、図15(A)を参照して前述したように、矩形120及び121が重複領域122から刈り込まれる。同様に、ブロック123及び124が刈り込まれ、重複領域のみが残される。

10

【0068】

図21(C)に示されるように、重複非テキストブロックが刈り込まれ、分解される。ステップS1701においては、重複領域162において輪郭ペア情報が存在するか否かが判別される。輪郭ペア情報が重複領域に存在しないとステップS1702で判別されると、さらに分解を行う必要はなく、処理が終了する。しかし、重複領域162に輪郭ペア情報が存在するとフローはステップS1703に進む。

【0069】

ステップS1703においては、矩形重複領域(重複領域の境界領域)に接触する全ての輪郭ペアが非テキストブロック及び重複対象物の両方について集められる。それらの集められた輪郭ペアは図21(B)に示す重複領域内に入るように刈り込まれる。

20

【0070】

次に、ステップS1704においては、非テキストブロック160の刈り込まれた輪郭ペアが集められる。ステップS1705では、非テキストブロック160及び重複ブロック161の重複領域162が、連結された垂直方向の輪郭ペアをグループ化し、次いで、垂直方向に沿うグループ化が終了した後で水平方向の輪郭ペアがグループ化されることにより、分解される。一度垂直方向の輪郭ペアがグループ化され、水平方向の輪郭ペアがグループ化されると、重複対象物161の輪郭ペアのいずれとも重複しないように矩形が形成される。非テキストブロック160の重複領域162を分解した結果を図21(C)に示す。

【0071】

図21(C)に示すように、非テキストブロック160の重複領域162は2つのより小さい非重複矩形163及び164に分解されている。

30

【0072】

[内部コンテンツを含む非テキストブロックの分解]

ステップS508で説明したように、テキストなどの内部コンテンツを含む非テキストブロックを分解する処理について図23(A)~(D)及び図24ならびに図25(A)及び(B)を参照してより詳細に論じる。

【0073】

初めは、フレーム、表、ピクチャ等のブロックには、その矩形領域が、該フレーム、表またはピクチャ内の白輪郭によって包まれる異なる種類のイメージデータを包んでいるものがあってもよい。例えば、フローチャートの場合は、その中にテキスト(アクション)を含む非テキストイメージ(工程ボックス)があり、したがって、イメージは、イメージまたはテキストが、他を抽出しないでしかるべく抽出され得るように、テキストとは別に定められねばならない。

40

【0074】

コンテンツを有するまたは有しない非テキストイメージをしかるべく抽出するために、コンテンツを有する非テキストブロックがコンテンツと重複しない最小数のより小さな外接矩形に分解されねばならない。それらのより小さな外接矩形から、内容イメージというよりむしろブロック化されたイメージがしかるべく抽出可能になる。

【0075】

50

さて図 2 3 (A) を参照すると、テキストコンテンツを含む非テキストイメージの例が示されている。図 2 3 (A) は、非テキストイメージ 1 8 0 ~ 1 8 2 を含むフローチャートの例である。非テキスト領域 1 8 0、1 8 1 及び 1 8 2 の各々の内側に、それぞれ 3 つの白輪郭領域 1 8 6、1 8 7 及び 1 8 8 があり、各白輪郭内に、テキストブロック 1 8 3、1 8 4 及び 1 8 5 がある。

【 0 0 7 6 】

図 2 3 (B) に示される非テキストイメージのみを得るために、非テキストイメージ 1 8 0、1 8 1 及び 1 8 2 をより小さな外接矩形によって分解できるように、非テキストイメージ、白輪郭及びテキストブロックの輪郭ペアが集められる。

【 0 0 7 7 】

かくして、図 2 3 (B) に示すイメージをしかるべく得るために、下記のステップが図 2 4 のフローチャートに示されたように実行される。

【 0 0 7 8 】

図 2 3 (B) に示されるイメージの分解を開始するために、非テキストイメージの各行の輪郭ペアが、該輪郭ペアの一つを白輪郭の輪郭ペアで置き換えることにより修正される。例えば、図 2 3 (C) に示すように、行 i の輪郭ペアは、白輪郭の対応する行で各行の輪郭ペアを置き換えることにより修正される。即ち、該修正に先立って、非テキストイメージ 1 8 0 及び 1 8 1 の輪郭ペアは、それぞれ (a , b) 及び (c , d) であり、白輪郭 1 8 6 及び 1 8 7 の輪郭ペアはそれぞれ (e , f) 及び (g , h) である。一度輪郭ペアが、各行の輪郭ペアを対応する白輪郭の行に置き換えることにより修正されると、非テキストイメージの輪郭ペアは、(a , e)、(f , b)、(c , g) 及び (h , d) となる。

【 0 0 7 9 】

各行の輪郭ペアを修正した後で、ステップ S 1 9 0 1 において、修正が非テキスト対象物の各輪郭ペアについて繰り返される。一度輪郭ペアの全ての行の修正が完了すると、ステップ S 1 9 0 3 において、修正された輪郭ペアは、ピクチャそのものがしかるべく抽出可能となるより小さな外接矩形の組にグループ化される。即ち、ステップ S 1 9 0 3 では、グループ化処理が、まず垂直方向に連結された輪郭ペアのグループ化が行われ、次に、垂直方向に沿うグループ化が終了した後で、水平方向に連結された輪郭ペアの全てのグループ化が行われる。全体の外接矩形が、テキストブロック 1 8 3、1 8 4 及び 1 8 5、他のグループ化されない輪郭ペアまたは他の前に形成された矩形等の内部ブロック内容物の輪郭ペアのいずれとも重複しないように、さらに修正された輪郭ペアを選択することが不可能になると、グループ化処理は終了する。一度垂直方向に連結された輪郭ペア及び水平方向に連結された輪郭ペアが全て連結されてしまうと、より小さな外接矩形が、非テキストイメージ 1 8 0、1 8 1 及び 1 8 2 の周囲に作られる。得られた分解イメージは図 2 3 (D) に示される。その後、矩形の座標が更なる処理のために記憶され、分解が終了する。

【 0 0 8 0 】

図 2 3 (D) に示すように、図 2 3 (A) のイメージは、最小 1 1 個の外接矩形に分解されている。記憶されたこれらの 1 1 個の外接矩形の座標を利用して、該 1 1 個の矩形、1 8 0、1 8 1 及び 1 8 2 の非テキストイメージがその中のコンテンツを抽出することなくしかるべく抽出可能になる。

【 0 0 8 1 】

図 2 5 (A) は、テキストやピクチャデータ等の、内部に含まれるブロックからイメージデータを抽出するためにどのようにフレームイメージが分解できるかを示す 1 例である。図 2 5 (B) は、上記の方法を用いて、テーブルイメージを、テーブルの内容物を抽出することなく抽出できるようにするために、テーブルを分解することができる方法を示す 1 例である。

【 0 0 8 2 】

図 2 6 (A) は、不規則な形状の輪郭 1 9 2 内に含まれるテキストブロック 1 9 1 を示す。そのような状況では、テキストブロックは、非テキスト輪郭用のブロックに重複しがちである。上述の分解により、テキストブロック 1 9 1 を非重複領域 1 9 1 a、1 9 1 b 及

10

20

30

40

50

び 191c に分解可能であり、その全てが図 26 (B) に示されるような輪郭 192 の白輪郭内にある。そのような分解により、テキストブロックのイメージデータを抽出するときに、輪郭の不要なイメージデータも間違っって抽出されることが、確実になくなる。

【0083】

以上、本発明を特定の態様に関して説明したが、本発明は上記記載に限定されることはなく、抽出されるべき全てのタイプのイメージデータに適用されることができると理解されるべきものである。更に、発明の精神及び範囲から逸脱することなく、当業者によって種々の変更や修正が可能である。

【図面の簡単な説明】

【図 1】重複矩形を分析し分解する方法を示すフローチャートである。

10

【図 2】文書頁の代表的な図である。

【図 3】図 2 にブロック選択処理を行って得られた矩形ブロックの図である。

【図 4】文書頁の代表的な図である。

【図 5】図 4 にブロック選択処理を行って得られた矩形ブロックの図である。

【図 6】本発明を具体化する装置の概観を示す斜視図である。

【図 7】図 6 の装置のブロック図である。

【図 8】イメージデータ内の輪郭ペアを得る方法を説明するフローチャートである。

【図 9】輪郭ペアを説明する図である。

【図 10】テキストタイプイメージデータの場合に、輪郭ペアを得る方法を説明する図である。

20

【図 11】内部成分を含まない非テキストイメージデータの輪郭ペアを得る方法を説明する図であり (A) は元のピクチャを示し、(B) は輪郭を結合した後のピクチャを示す。

【図 12】(A) 及び (B) は、内部成分を含む非テキストイメージデータの場合に輪郭ペアを得る方法を説明する図であり、(C) 及び (D) は、それぞれ 4 方向及び 8 方向輪郭トレース間の相違を説明する図である。

【図 13】重複しない外接矩形を示す図である。

【図 14】重複輪郭がイメージ抽出に影響を与えるかどうかを判別するための分析を説明する図である。

【図 15】2 つのブロックがどのようにして重複可能であるかを説明する図である。

【図 16】(A) は、非テキストブロックに重複するテキストブロックの例を示し、(B) は、(A) に示す重複領域の刈り込まれたものを示す図である。

30

【図 17】重複テキストブロックをより小さな矩形に分解する方法を示すフローチャートである。

【図 18】図 17 のフローチャートの続きである。

【図 19】図 17 のフローチャートの続きである。

【図 20】それぞれ重複するテキストブロック領域を分解する二つの方法を示す図である。

【図 21】二つの重複非テキストブロックを説明する図である。

【図 22】重複非テキストイメージを分解する方法を示すフローチャートである。

【図 23】内部コンテンツを含む非テキストブロックを分解する方法を説明する図である。

40

【図 24】内部コンテンツを含む非テキストブロックをより小さな外接矩形に分解する方法を示すフローチャートである。

【図 25】(A) 分解されたフレーム及び (B) 分解された表の例をそれぞれ示す図である。

【図 26】(A) 不規則な形状の輪郭に含まれるテキストブロック及び (B) その分解後を説明する図である。

【符号の説明】

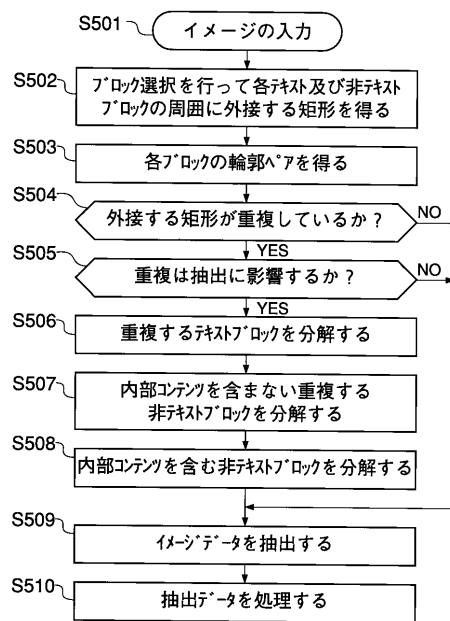
410 コンピュータ機器

411 ディスク

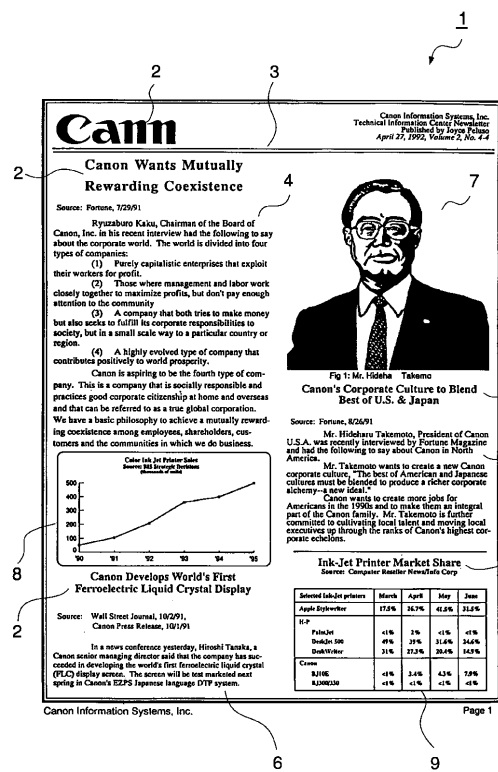
50

5 2 0 C P U
5 2 2 スキャナ・インターフェース

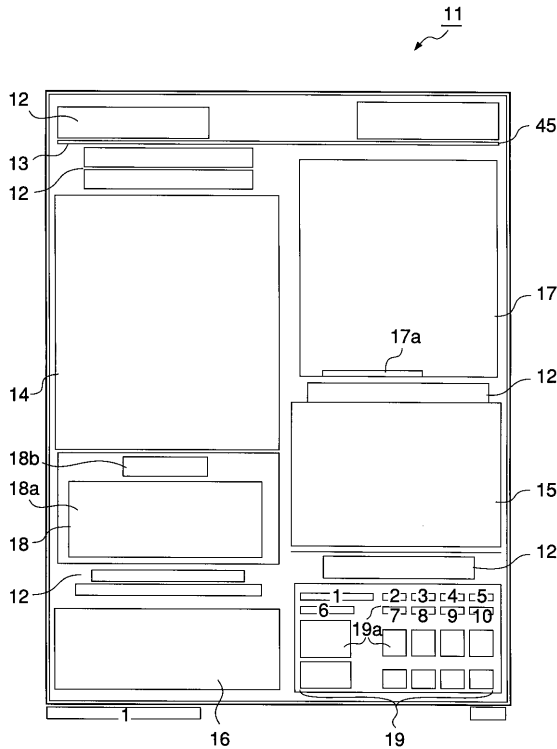
【図 1】



【図 2】



【図 3】



【図 4】

Cam
Canon Information Systems, Inc.
Technical Information Center Newsletter
Published by Joyce Palano
April 27, 1992, Volume 2, No. 4

Canon Wants Mutually Rewarding Coexistence

Source: Fortune, 7/29/91

Ryuzaburo Kaku, Chairman of the Board of Canon, Inc. in his recent interview had the following to say about the corporate world. The world is divided into four types of companies:

- (1) Purely capitalistic enterprises that exploit their workers for profit.
- (2) Those where management and labor work closely together to maximize profits, but don't pay enough attention to the community.
- (3) A company that both tries to make money but also seeks to fulfill its corporate responsibilities to society, but in a small scale way to a particular country or region.
- (4) A highly evolved type of company that contributes positively to world prosperity.

Canon is aspiring to be the fourth type of company. This is a company that is socially responsible and practices good corporate citizenship at home and overseas and that can be referred to as a true global corporation. We have a basic philosophy to achieve a mutually rewarding coexistence among employees, shareholders, customers and the communities in which we do business.

Ryuzaburo Kaku, Chairman of the Board of Canon, Inc. continued to say about the four types of companies. Canon Company - This is a company practicing good corporate and civic policies, a company concerned about its employees and shareholders.

Fig 1: Mr. Hideharu Takemoto
Canon's Corporate Culture to Blend Best of U.S. & Japan

Source: Fortune, 8/26/91

Mr. Hideharu Takemoto, President of Canon U.S.A., was recently interviewed by Fortune Magazine and had the following to say about Canon in North America.

Mr. Takemoto wants to create a new Canon corporate culture. "The best of American and Japanese cultures must be blended to produce a richer corporate alchemy—a new ideal."

Canon wants to create more jobs for Americans in the 1990s and to make them an integral part of the Canon family. Mr. Takemoto is further committed to cultivating local talent and moving local executive positions up through the ranks of Canon's highest corporate echelons.

Canon Develops World's First Ferroelectric Liquid Crystal Display

Source: Wall Street Journal, 10/2/91
Canon Press Release, 10/1/91

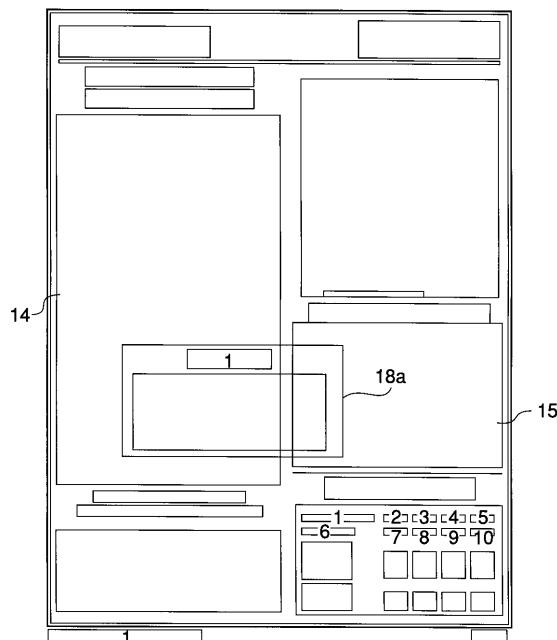
In a news conference yesterday, Hiroshi Tanaka, a Canon senior managing director said that the company has succeeded in developing the world's first ferroelectric liquid crystal (FLC) display screen. The screen will be test marketed next spring in Canon's EZPS Japanese language DTP system.

Ink-Jet Printer Market Share
Source: Computer Reseller News/Info Corp.

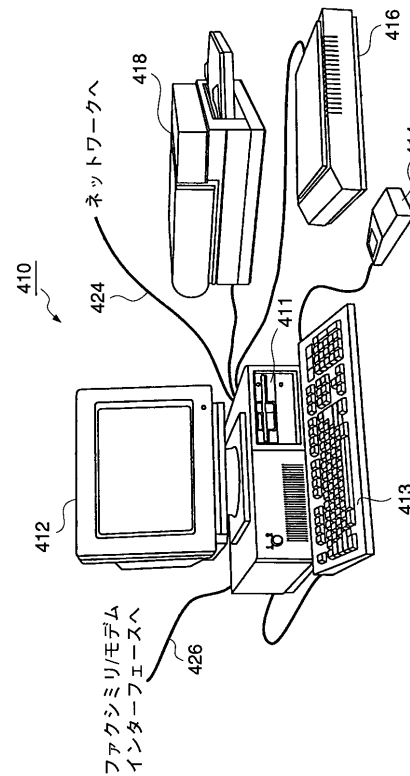
Selected Ink-Jet printers	March	April	May	June
Apple StyleWriter	17.5%	26.7%	42.2%	31.5%
HP	<1%	1%	<1%	<1%
Printronix	4%	3%	31.6%	24.6%
DeskJet 500	31%	37.3%	30.4%	14.5%
Canon	8.11%	<1%	3.4%	4.7%
8.23%/39	<1%	<1%	<1%	<1%

Canon Information Systems, Inc. Page 1

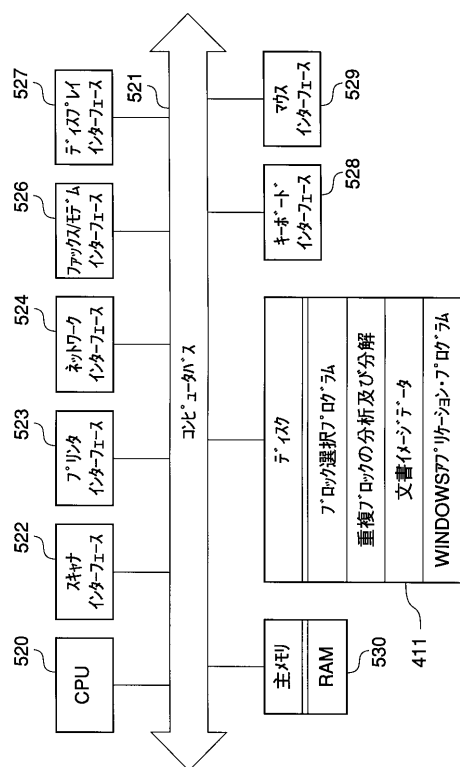
【図 5】



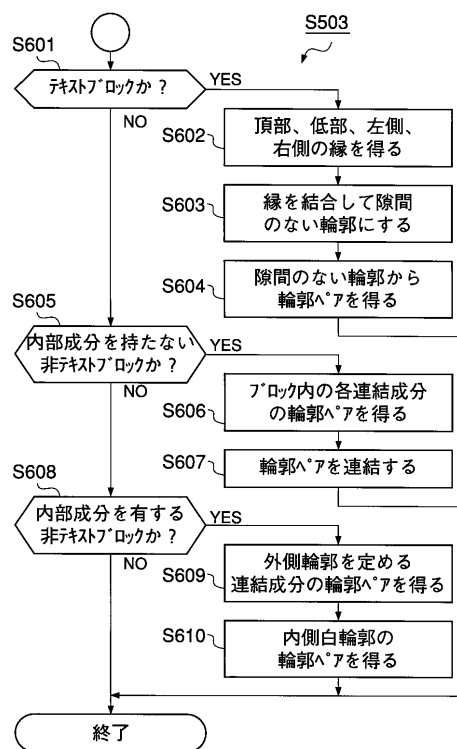
【図 6】



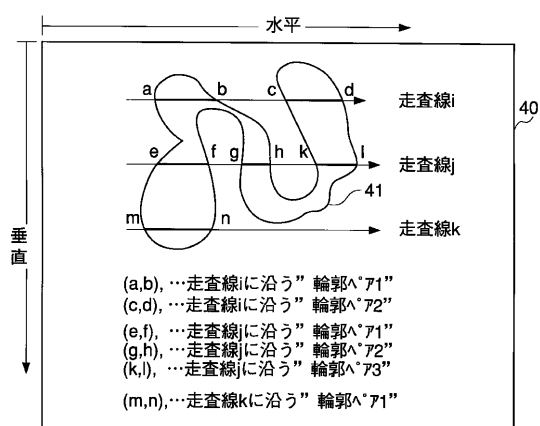
【圖 7】



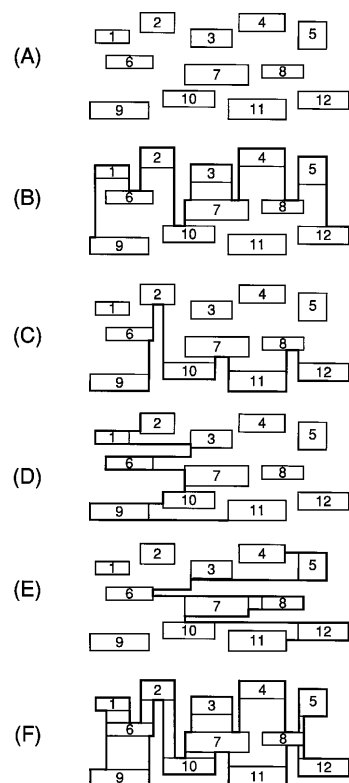
【圖 8】



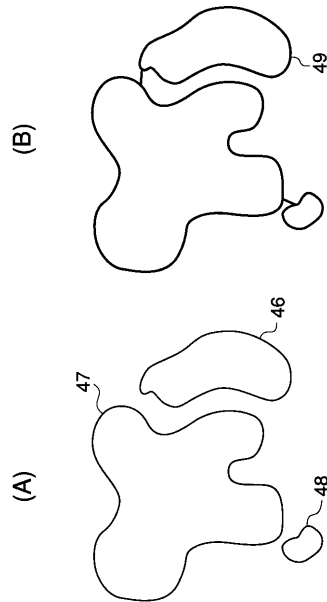
【 図 9 】



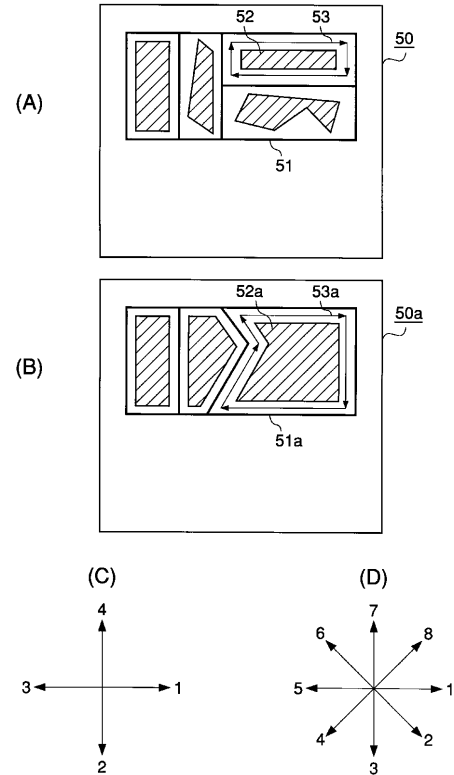
【 図 1 0 】



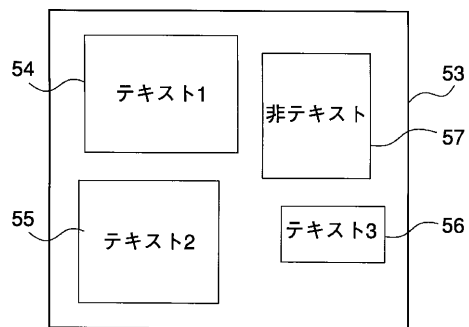
【図 1 1】



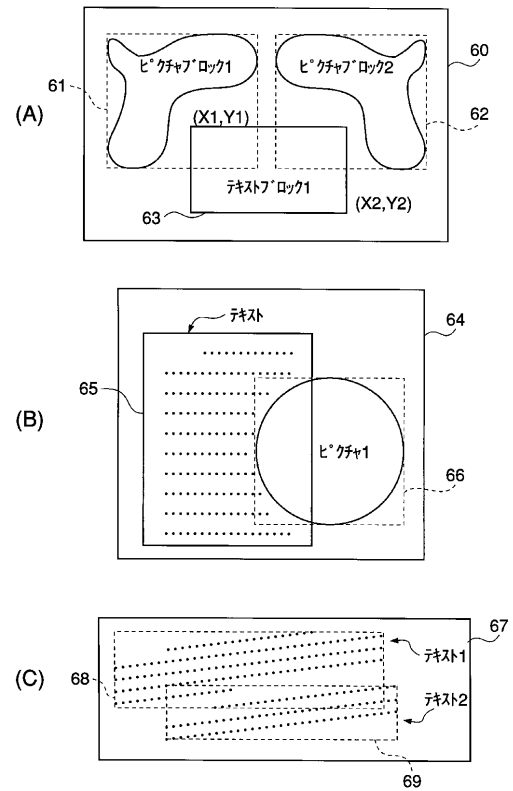
【図 1 2】



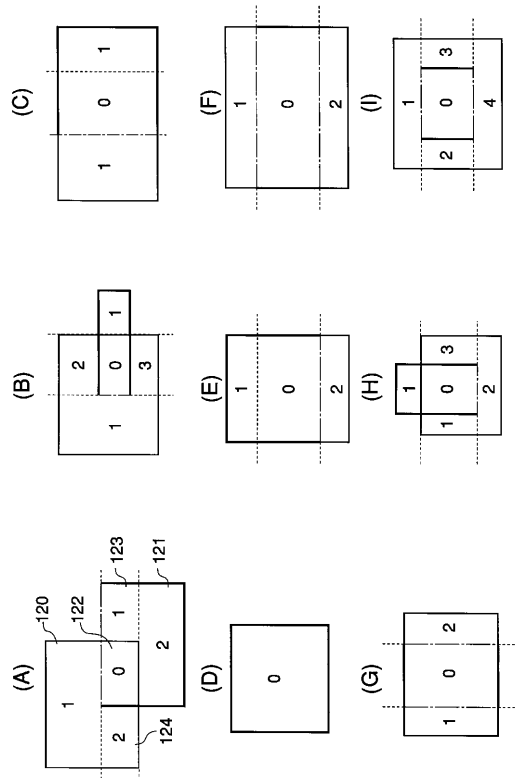
【図 1 3】



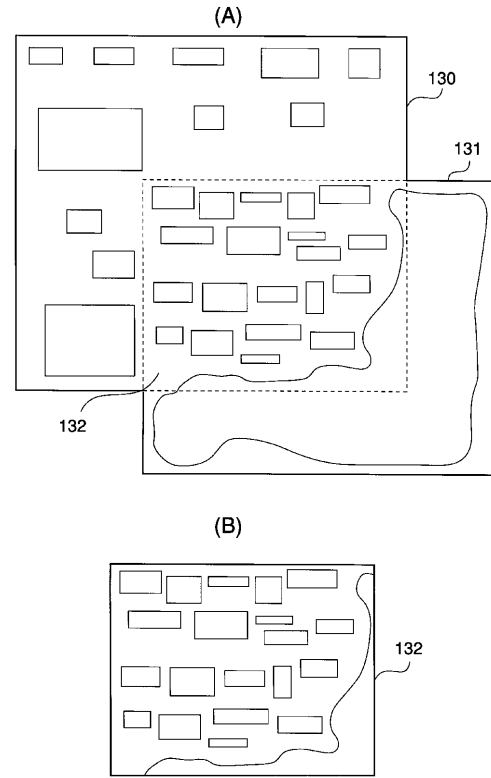
【図 1 4】



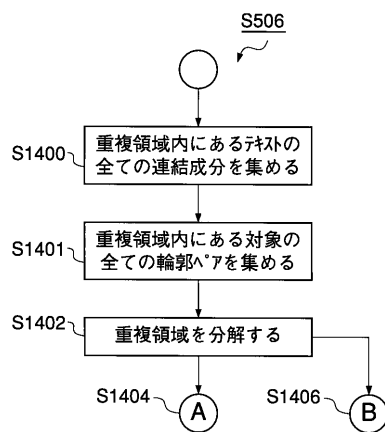
【図 15】



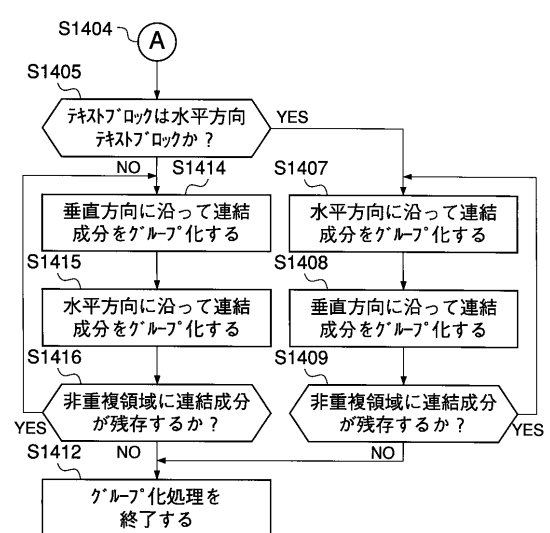
【図 16】



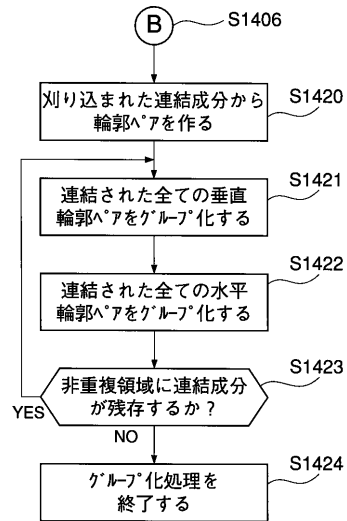
【図 17】



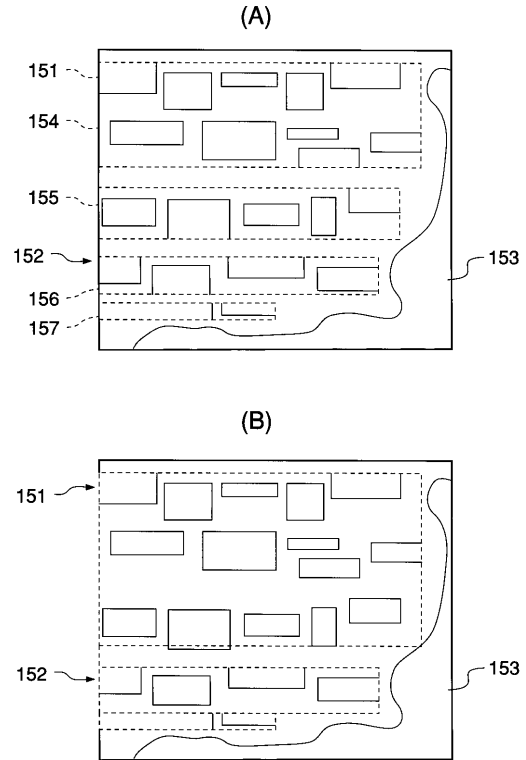
【図 18】



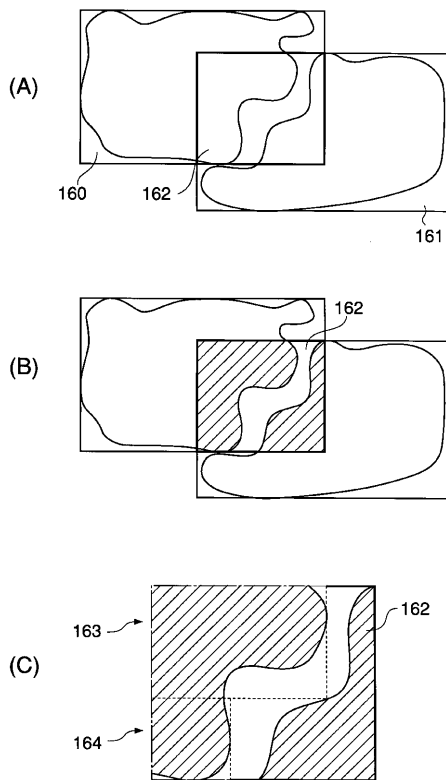
【図 19】



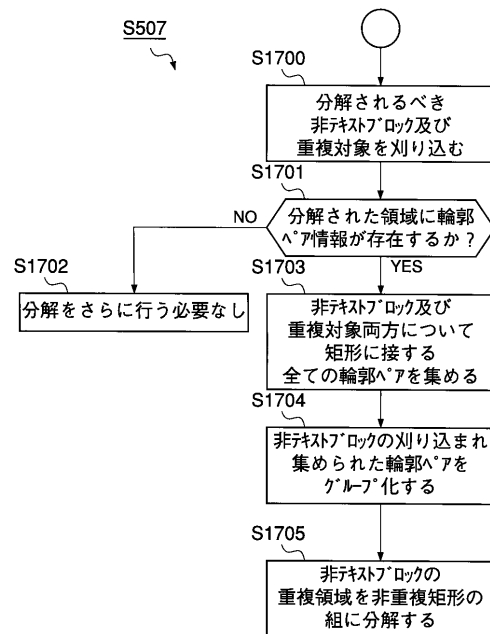
【図 20】



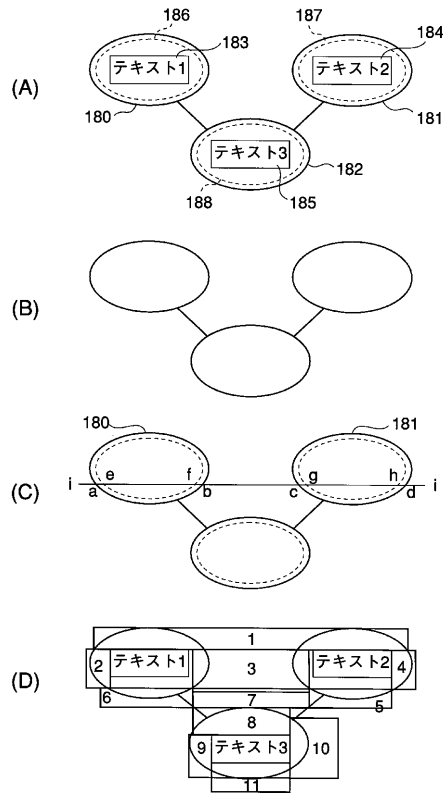
【図 21】



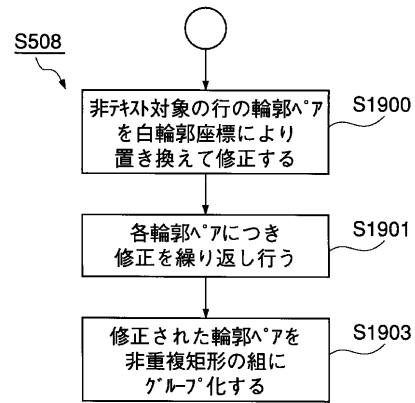
【図 22】



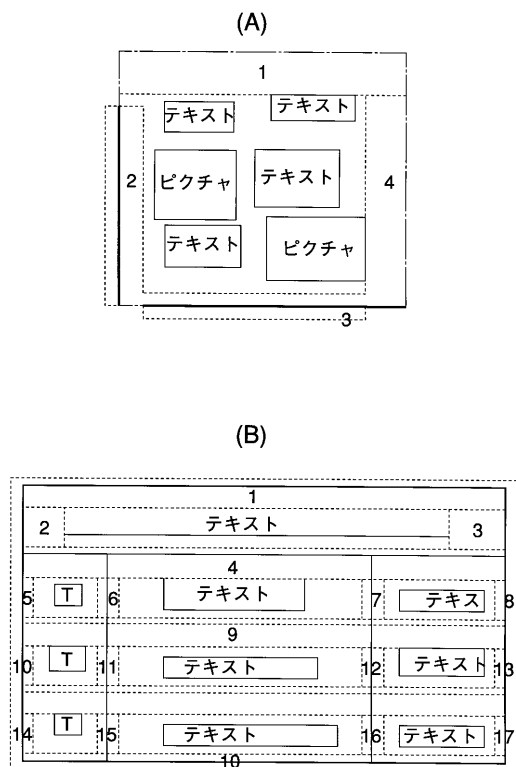
【図 2 3】



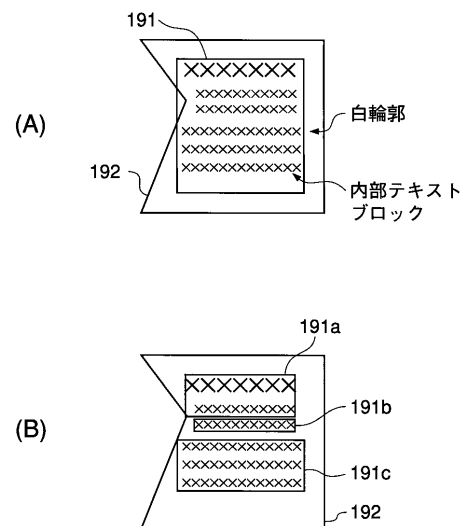
【図 2 4】



【図 2 5】



【図 2 6】



フロントページの続き

(72)発明者 シン ヤン ワン

アメリカ合衆国 カリフォルニア 92680 タスティン マクチャールズ ドライブ 222
1

(72)発明者 矢ヶ崎 敏明

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 飯田 清司

(56)参考文献 特開平08-115380(JP,A)

特開平07-107281(JP,A)

特開平07-234918(JP,A)

特開平06-068301(JP,A)

特開昭62-203280(JP,A)

特開昭62-204380(JP,A)

特開2001-259532(JP,A)

特開2000-348140(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06T 7/00-7/60

G06K 9/20