

(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2017-105175

(P2017-105175A)

(43) 公開日 平成29年6月15日(2017.6.15)

(51) Int.Cl.			F I			テーマコード (参考)	
B41J	5/30	(2006.01)	B41J	5/30		Z	2C061
G06F	3/12	(2006.01)	G06F	3/12	344		2C187
G06K	9/00	(2006.01)	G06K	9/00		Z	5B064
B41J	29/38	(2006.01)	G06F	3/12	302		
			G06F	3/12	324		

審査請求 未請求 請求項の数 10 O L (全 17 頁) 最終頁に続く

(21) 出願番号 特願2016-225518 (P2016-225518)
 (22) 出願日 平成28年11月18日 (2016.11.18)
 (31) 優先権主張番号 14/960, 986
 (32) 優先日 平成27年12月7日 (2015.12.7)
 (33) 優先権主張国 米国 (US)

(特許庁注：以下のものは登録商標)

1. イーサネット

(71) 出願人 596170170
 ゼロックス コーポレイション
 XEROX CORPORATION
 アメリカ合衆国、コネチカット州 068
 56、ノーウォーク、ピーオーボックス
 4505、グローバー・アヴェニュー 4
 5
 (74) 代理人 100079049
 弁理士 中島 淳
 (74) 代理人 100084995
 弁理士 加藤 和詳
 (72) 発明者 ジェロウム・ポウヤドウ
 フランス共和国 グルノーブル 3800
 0 リュ・デ・ラ・ペ 1

最終頁に続く

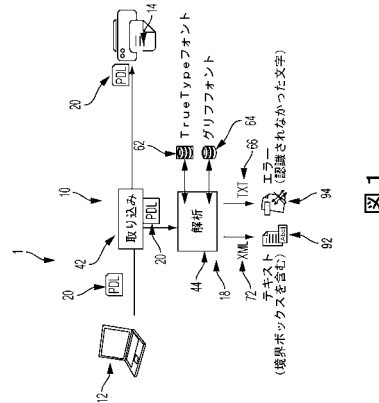
(54) 【発明の名称】 ページ記述言語文書からの直接文字認識

(57) 【要約】 (修正有)

【課題】テキストベースのPDL文書の印刷時に、そのテキストベースのPDL文書から文字を直接認識できるシステムを提供する。

【解決手段】システム10は、印刷文書に関して印刷ドライバによって生成され、一連のグリフから形成されたテキストランを描画するための描画命令セットを含むPDL(ページ記述言語)文書20を取り込む取り込みコンポーネント42と、PDLフォーマットに従ってPDL文書の各グリフのすべてのテキスト関連演算子を捕捉して直接文字認識を行う解析システム44と、を備える。解析システムは、PDL文書を解析して各グリフの描画命令を抽出し、グリフ特徴データベース64と比較する。一致するものが見つかった場合は、グリフに関連付けられたテキスト文字を抽出し、抽出されたテキスト文字の要約を生成する。一致するものが見つからなかった場合は、見つけ損ねたテキスト文字のログを取る。

【選択図】図1



【特許請求の範囲】**【請求項 1】**

ページ記述言語 (P D L) 文書からテキストを抽出する方法であって、
印刷される文書に関して、印刷ドライバによって生成された P D L ファイルを取り込む
ステップであって、前記 P D L ファイルが、前記印刷される文書にある一連のグリフから
形成されたテキストランに関する描画命令セットを含む P D L 文書を含む、ステップと、
前記 P D L 文書を解析して、グリフセットのそれぞれに関して、描画命令を抽出する、
ステップと、
前記グリフセットにある各グリフに関して、前記グリフの前記描画命令を、グリフ特性
のデータベースと比較して、一致するグリフがあるかどうかを判定する、ステップと、
前記描画命令とグリフ特性の前記データベースとの間で一致するものが見つかった場合
、前記グリフに関連付けられたテキスト文字を抽出するステップと、
前記一致するグリフに関連付けられた前記抽出されたテキスト文字の要約を生成するス
テップと、
を含む、方法。

10

【請求項 2】

取り込む、解析する、比較する、抽出する、および生成する、前記ステップのうちの少
なくとも 1 つが、プロセッサで実行される、請求項 1 に記載の方法。

【請求項 3】

前記描画命令とグリフ特性の前記データベースとの間で一致するものが見つからなかつ
た場合、見つけ損ねたテキスト文字のログを取るステップをさらに含む、請求項 1 に記載
の方法。

20

【請求項 4】

前記グリフ特性セットが、グリフセットのそれぞれに関するランレングス特徴ベクトル
を含む、請求項 1 に記載の方法。

【請求項 5】

コンピュータによって実行されると請求項 1 に記載の方法を行う命令を格納する非一時
的記録媒体を備える、コンピュータプログラム製品。

【請求項 6】

請求項 1 に記載の方法を行うための命令を格納するメモリと、前記命令を実行するた
めに、前記メモリと通信可能なプロセッサとを備える、文書のページ記述言語 (P D L) か
らテキストを抽出するためのシステム。

30

【請求項 7】

文書のページ記述言語 (P D L) からテキストを抽出するためのシステムであって、
印刷される文書に関して、印刷ドライバによって生成された P D L ファイルを取り込む
、取り込みコンポーネントであって、前記 P D L ファイルが、前記印刷される文書にある
一連のグリフから形成されたテキストランに関する描画命令セットを含む、取り込みコン
ポーネントと、

前記 P D L ファイルを解析して、各グリフの前記描画命令を捕捉する、解析器と、
各グリフの前記描画命令を、グリフ特性のデータベースと比較する、比較コンポーネン
トと、

40

グリフ特性の前記データベースとの前記描画命令の前記比較に基づいて前記描画命令と
グリフ特性の前記データベースとの間で一致するものが見つかった場合、各グリフに関連
付けられたテキスト文字を抽出する、抽出コンポーネントと、

各グリフに関連付けられた前記抽出されたテキスト文字のテキスト要約を生成する、要
約コンポーネントと、

前記取り込みコンポーネントと、解析器と、比較コンポーネントと、抽出コンポーネン
トと、要約コンポーネントとを実装するプロセッサと、

を備える、システム。

【請求項 8】

50

前記描画命令の、グリフ特性の前記データベースとの比較に基づいて、前記描画命令とグリフ特性の前記データベースとの間で一致するものが見つからなかった場合、見つけ損ねたテキスト文字に関するエラーログを生成するログ取得コンポーネントをさらに備える、請求項7に記載のシステム。

【請求項9】

各グリフをビットマップとしてレンダリングする、レンダリングコンポーネントと、各グリフの前記ビットマップに基づいて特徴ベクトルを抽出する、特徴抽出部と、前記抽出された特徴ベクトルとグリフ特性の前記データベースとの間の類似度を計算する、類似度コンポーネントと、

をさらに備え、

前記抽出コンポーネントが、前記抽出された特徴ベクトルとグリフ特性の前記データベースとの間の前記類似度に基づいて各グリフに関連付けられたテキスト文字を抽出する、請求項7に記載のシステム。

【請求項10】

文書のページ記述言語(PDL)からテキストを抽出する方法であって、

テキスト文字のそれぞれに関連付けられたグリフ特性セットによって定義される前記テキスト文字の参照データベースを提供するステップと、

印刷ドライバによって生成されたPDL文書を受信するステップと、

プロセッサにより、前記PDL文書を解析して、前記PDL文書にあるテキスト描画プリミティブを識別する、ステップと、

前記テキスト描画プリミティブの、前記グリフ特性セットとの比較に基づいて、前記参照データベースにある前記テキスト描画プリミティブによって表現されるテキスト文字を識別するステップと、

前記テキスト描画プリミティブによって表現される前記テキスト文字を抽出するステップと、

前記抽出されたテキスト文字に基づいて情報を出力するステップと、

を含む、方法。

【発明の詳細な説明】

【技術分野】

【0001】

本例示的な実施形態は、テキスト抽出の分野に関し、詳細には、光学文字認識を必要とすることなく、文書からのテキスト抽出に利用できる。

【背景技術】

【0002】

ページ記述言語(PDL)は、印刷される文書を、印刷機に依存しないフォーマットで記述する。電子PDL文書が送られる印刷機において、文書は紙などの印刷媒体の上にレンダリングされる。ここでは、文書のテキストの内容を、PDLフォーマットで取り込むことが望ましい事例がいくつかある。そのような事例としては、印刷されるべきではない、文書上の機密情報を検出するためのセキュリティチェックの実行、印刷コストを課金するための、文書を印刷させた顧客の検出、私的な印刷ジョブと公的な印刷ジョブとを区別するため、またはカラー印刷機が適切な文書のために使用されているかをチェックするための利用制御、印刷しているユーザおよびその理由を検出するための監査、ならびに後の検索を加速するためにすべての印刷文書をアーカイブして索引付けする状況における索引付けが挙げられる。

【0003】

しかしながら、光学文字認識(OCR)などの、文字認識のための既存の方法は、PDLフォーマットの文書の処理に適用できない。1つの理由としては、印刷される文書のケースにおいて、印刷レンダリングエンジン、すなわちラスライメージプロセッサ(RIP)が、どの文字が印刷されるかを知る必要がないことが挙げられる。RIPは、全体的な一体となった結果が人間の読み手にとって意味をなすテキストに見えるように、印刷ペー

10

20

30

40

50

ジ上のどこにインクのドットを置くかを分かっているだけでよいだけである。大半の P D L および大半の文書で、テキストをレンダリングするための表記法および A P I が存在し得る。しかしながら、表示される文字の実際の「値」は、レンダリングに関係しない。よって、この情報は、印刷機に届くデータに含まれていない。

【 0 0 0 4 】

P D L をテキストにする既存のツールは、既知の識別子（通常、文字列）とそれらが表すグリフとの間に対応構造が存在すると仮定することによって、P D L 文書からテキストを抽出する。このことはよくあることとはいえ、常にというわけではなく、表現される文字セットを拡張するのにこの種の対応を使わない文書は多い。こうしたケースでは、テキスト抽出は、誤った結果を生じる。この種の対応を使わない他のケースでは、テキスト抽出は不可能であると一般的に考えられている。加えて、入手可能なツールは、P o s t S c r i p t コンピュータ言語の上だけで動作する傾向にあるが、P o s t S c r i p t コンピュータ言語は、現在使用されている多くの入手可能な P D L のうちの 1 つに過ぎない。

10

【 0 0 0 5 】

テキスト抽出のための別の手法は、印刷機に P D L 文書が到着した時に、または印刷機自体において、P D L 文書を画像に変換することと、テキストを再構築するために O C R 技法を適用することとを必要とする。しかしながら、レンダリングして O C R にかける手法では、文書全体に対して O C R を行う前にすべてのページをレンダリングすることから、特に、文書に含まれるページ数が多い場合には、時間がかかり得る。この処理時間は、セキュリティ検出の状況において特に問題となり得、そのような状況では、ルールベースのエンジンが文書内で特定の単語を探すのであるが、その単語が出現するのが最初のページだとしても、その単語を検出できるまでには文書全体のレンダリングおよび O C R を待たねばならない。

20

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 6 】

したがって、テキストベースの P D L 文書の印刷時に、そのテキストベースの P D L 文書から文字を直接認識できるシステムおよび方法のニーズがある。

【 課題を解決するための手段 】

30

【 0 0 0 7 】

本例示的な実施形態の一実施態様によれば、ページ記述言語（P D L）文書からテキストを抽出する方法が提供される。本方法は、印刷される文書に関して、印刷ドライバによって生成された P D L ファイルを取り込むステップを含む。P D L ファイルは、印刷される文書において一連のグリフから形成されているテキストランに関する描画命令セットを含む、P D L 文書を含む。P D L 文書を解析して、グリフセットのそれぞれに関する描画命令を抽出する。グリフセットの各グリフに関して、グリフの描画命令を、グリフ特性のデータベースと比較して、データベースに一致するグリフがあるかどうかを判定する。描画命令とグリフ特性のデータベースとの間で一致するものが見つかった場合、一致するデータベースグリフに関連付けられたテキスト文字を抽出する。抽出した、一致するグリフに関連付けられたテキスト文字の要約を生成する。

40

【 0 0 0 8 】

本方法のステップのうちの一つ以上は、プロセッサによって実装され得る。

【 0 0 0 9 】

本例示的な実施形態の別の実施態様によれば、文書のページ記述言語（P D L）からテキストを抽出するためのシステムは、印刷文書に関して印刷ドライバによって生成された P D L ファイルを取り込む、取り込みコンポーネントを含む。P D L ファイルは、印刷文書における一連のグリフから形成されたテキストランを描画するための描画命令セットを含む。解析器は、P D L ファイルを解析して、各グリフの描画命令を捕捉する。比較コンポーネントは、各グリフの描画命令を、グリフ特性のデータベースと比較する。抽出コン

50

ポーネントは、描画命令の、グリフ特性のデータベースとの比較に基づいて、描画命令とグリフ特性のデータベースとの間で一致するものが見つかった場合、各グリフに関連付けられたテキスト文字を抽出する。要約コンポーネントは、抽出した、各グリフに関連付けられたテキスト文字のテキスト要約を生成する。プロセッサは、取り込みコンポーネントと、解析器と、比較コンポーネントと、抽出コンポーネントと、要約コンポーネントとを実装する。

【0010】

本例示的な実施形態の別の実施態様によれば、文書のページ記述言語（PDL）からテキストを抽出する方法は、テキスト文字のそれぞれに関連付けられたグリフ特性セットによって定義される、テキスト文字の参照データベースを提供するステップと、印刷ドライバによって生成されたPDL文書を受信するステップとを含む。プロセッサによって、PDL文書を解析して、PDL文書におけるテキスト描画プリミティブを識別する。本方法は、テキスト描画プリミティブの、グリフ特性セットとの比較に基づいて、参照データベースにおいて、テキスト描画プリミティブによって表現されるテキスト文字を識別するステップと、テキスト描画プリミティブによって表現されるテキスト文字を抽出するステップと、抽出されたテキスト文字に基づいて情報を出力するステップとをさらに含む。

10

【0011】

本方法のステップのうちの一つ以上は、プロセッサによって実装され得る。

【図面の簡単な説明】

【0012】

20

【図1】図1は、ページ記述言語（PDL）文書からの直接文字認識のためのシステムおよび方法の概略図である。

【図2】図2は、本例示的な実施形態の一実施態様による、ページ記述言語（PDL）文書からの直接文字認識のためのシステムの機能ブロック図である。

【図3】図3は、図2のシステムの追加的な実施態様を示すブロック図である。

【図4】図4は、グリフ特性によって定義されるテキスト文字の参照データベースの記号的表現の図である。

【図5】図5は、本例示的な実施形態の別の実施態様による、PDL文書からの直接文字認識のための例示的な方法を示す流れ図である。

【発明を実施するための形態】

30

【0013】

本例示的な実施形態の実施態様は、一般的な事務文書の印刷時または印刷時付近に、一般的な事務文書にローカル文字認識を適用するための、方法、装置、およびコンピュータ可読媒体に関する。本方法は、同様のフォント技法が使用される状況およびグリフを文字にする情報が失われている場合の他の状況において使用可能であるほど十分に汎用的である。

【0014】

本例示的な方法は、ページ記述言語（PDL）文書からの直接文字認識を提供する。本方法は、ユーザが印刷ジョブを開始した時、例えばアプリケーションから「印刷」をクリックすることによって作成されるPDLファイルを解析するステップを含み得る。解析されるPDL文書は、描画命令から構成され得る。処理される記号（グリフ）の各視覚的表現に関して、グリフ特性を、参照データベースと照合することによって、それが表現する文字を見つける。参照データベースは、フォントのデータベースにおけるグリフルックアップに基づくものであり得る。フォントは、グリフの2次スプライン制御点によって直接認識が行われる、TrueType/OpenType技術を使用できる。他のケースでは、半厳密な、視覚的類似度によって認識が行われる。

40

【0015】

一実施形態では、本方法は、PDLファイルからテキスト関連演算子セットを捕捉するステップを含み得る。

【0016】

50

図1は、ページ記述言語(PDL)文書からの直接文字認識のためのシステム10が動作する、印刷ネットワーク1における文書処理の概略図である。一般的に、システム10では、開示される実施形態と整合するコンピュータソフトウェアが使用され得る。ネットワークは、コンピュータ12などの、1つ以上のネットワーク化されたPDL生成装置と、従来の通信プロトコルおよび/またはデータポートインターフェースを用いて情報を交換できるようにする通信リンク16(図2)を通して接続される、1つ以上の印刷機14とを含み得る。

【0017】

ネットワーク1内の各計算装置12は、コンピュータワークステーションもしくはデスクトップコンピュータ、ラップトップもしくはポータブルコンピュータ、携帯機器、またはネットワーク化環境で使用可能な任意の他の計算装置であり得る。

10

【0018】

図2に示すように、PDL文書20は、ユーザが計算装置12において印刷ジョブを要求した時に、印刷ドライバ22によって作成される。元の文書24を印刷するには、計算装置12のユーザは、例えば、印刷機14などの送信先印刷機と、元の文書を印刷するためのいくつかの仕上げオプション(例えば、両面モード、ホチキス留め、カラーもしくは単色、およびページ選択)とをユーザが選択できる、通常「印刷ダイアログ」の形態である、特定のアプリケーションコマンドを利用する。印刷オプションの選択がなされた時(例えば、ユーザが、アイコンまたはキーボードのキーなどの「印刷」セクタをクリックした時)、アプリケーションは、印刷ドライバ22によって提供される仮想ディスプレイの形態を用いて、文書の各ページをデジタル形式でレンダリングする。特に、印刷ドライバ22は、アプリケーションによって提供された描画命令を、印刷機固有の描画命令に変換する。これらの印刷機固有の描画命令は、それら独自の構文を有し、その構文は、元のアプリケーションの描画構文とは異なる。このような印刷機固有の描画命令セットはPDL文書20を構成し、PDL文書20は、PDLファイル26で印刷機に送信され得る。

20

【0019】

印刷ドライバ22は、コンピュータ12の一部であってもよいし、計算装置12と印刷機14との中間にある別個の計算装置上に配置されてもよい。

【0020】

従来のネットワーク印刷システムのように、PDLファイル26は、特定のアプリケーションにおいてオープンな、固有のフォーマットの元の文書24から、印刷ドライバによって、PSまたはPCL6などの特定のPDLフォーマットで作成される。例えば、事務文書は、その種類に応じて、ワードプロセッサ、スプレッドシートハンドラ、またはスライドプレゼンテーションアプリケーションで開くことができる。この元の文書24のフォーマットは、公知であってもよいし、そうでなくてもよいが、いずれにせよ、元の文書を作成、修正、および/または閲覧するアプリケーションに固有のものである。大半の印刷機は、そのような固有のフォーマットを直接理解できない。

30

【0021】

元の文書をPDLフォーマットに変換した後、PDLデータ20は、選択された印刷機14に送信され、必ずしもそうではないが、各ページの画像表現を含み得る。PDLフォーマットどうしは構文および表現が異なるが、すべてのPDLフォーマットは、コンピュータグラフィクスプリミティブに基づいている。これらのプリミティブは、その最終的なゴールが、この場合紙である「表示面」上にグラフィックエレメントを描画するというものである、プログラム命令である。このようなプリミティブは、様々なフォーマットのビットマップを入力画像とし、これらを紙の上にレンダリングさせる、画像描画プリミティブを含む。他のプリミティブは、例えば、線分、円、四角形、およびベジェ曲線などの幾何形状を表現するために使用されるグラフィックプリミティブを含む。PDLの高度化レベルに応じて、これらの命令はまた、なんらかの形態のテクスチャ表現(すなわち、幾何学的形状を埋めるために使用されるパターン)ならびに様々な線の属性(例えば、色、幅、および線種)も含み得る。また、テキスト文字列(「ラン」ともいう)のグリフを描画

40

50

するために使用される、テキスト描画プリミティブも含む。一部の言語は、描画プリミティブの挙動を修正するために使用される、非描画プリミティブを含むこともある。一般的に、非描画プリミティブは、縦向きから横向きへと印刷の向きを切り替えるために使用される、行列演算やページ回転演算子などであり、座標系を修正するために使用される。非描画プリミティブはまた、様々な装置制御演算子、例えば、出力トレイ選択または仕上げオプションも修正できる。

【0022】

大半のPDLフォーマットにはテキスト描画演算が存在するが、だからといって、以下に詳述するように、所与の印刷物シート上のすべてのテキストがテキスト命令によって描画されているとは言えない。しかしながら、テキスト命令が使用されるならば、レンダリング装置14は、いかなる他のグラフィクスでもなくテキストを印刷していることを一般的に把握している。

10

【0023】

多くのケースでは、テキスト演算子は、テキストラン(すなわち、単語の境界にかかり得る、あるいはかかり得ない、文字シーケンス)を構成する実際の文字ではなく、紙などの物理的媒体28上にレンダリングされるグリフの順番に割り当てられた識別子を引数とすることによって作用する。例えば、印刷機14は、次の疑似コード: `operator DrawText(「hello」)` といった命令を受信するのではなく、`operator DrawText(12, 1, 15, 15, 21)` などの、文字のグリフコードを識別する命令を受信し得る。この例では、12は、その時使用中のフォントで文字「h」を表現するグリフのインデックスであり、1は「e」のインデックスであり、15は「l」のインデックスであり、21は「o」のインデックスである。

20

【0024】

対応する文字をレンダリングするためのインデックスおよび命令を格納する、グリフィンデックス30を、例えばPDLファイル26で印刷機14に送信してもよい。このインデックスは、フォントに固有でないことが多く、したがって、簡素な、インデックスをグリフにする対応表を使用できない。大半の印刷ドライバは、PDL文書20を生成する際に、グリフィンデックス30を徐々に構築する。一般的に、インデックスは、ジョブの際、ある文字を初めて所与のフォントで描画しなければならない時はいつでも割り当てられる。例えば、上記の例では、「e」は、ジョブの際、選択されたフォントに関して、表示される最初の文字になるため、インデックス1を得、以下、hは12番目であるためインデックス12を得るなどとなる。このことから明らかのように、テキスト演算子をフックすることができる場合であっても、グリフィンデックスだけからでは、描画される文字に戻る実質的な方法はない。

30

【0025】

単一のテキスト描画演算子によって処理されたとしても、グリフは、多くのフォーマットで格納できることに留意されたい。典型的なフォーマットとしては、純粋なビットマップ、TrueType輪郭(これは、基本的には、グリフ曲線を描画するために使用されるベジエ点セットである)、および/またはPostScript PDLとPDFなどのそのパリエーションとで使用されるPostScript命令が挙げられる。このケースでは、グリフ記述に組み込まれたPostScript命令サブセットを用いてグリフ曲線を描画する。他のパリエーションとしては、例えば、Adobe Type 2フォント、CFFフォント、およびChameleonフォントが挙げられ、これらのフォーマットは、独自仕様であってもよいし、そうでなくてもよい。

40

【0026】

図示の直接文字認識システム10は、本例示的な方法を行うための命令18を格納するメモリ31と、命令を実行するためにメモリと通信可能なプロセッサ32(または複数のプロセッサ)とを含む。システムは、サーバコンピュータ34などの1つ以上の計算装置にわたって分散させてもよい。他の実施形態では、システム10は、計算装置12もしくは印刷機14でホストされてもよいし、ネットワーク1上に分散させてもよい。1つ以上

50

の入出力（「I/O」）装置36、38によって、システムがコンピュータ12および/または印刷機14などの外部装置と通信できる。システムのハードウェアコンポーネント31、32は、データ/制御バス40によって通信可能に接続される。コンピュータ12、印刷機14、および/またはサーバコンピュータ34との間で送受信される情報には、使用されるプロトコルに従って、データ、コマンド、言語データファイルのロケーションおよびフォーマットに関する情報、ケーブルリテリ要求、ステータス要求、応答、および/または肯定確認を含むことができる。

【0027】

例示的な命令18は、取り込みコンポーネント42と、解析システム44と、テキスト要約コンポーネント46と、ログ取得コンポーネント48と、キャッシュコンポーネント50と、検証コンポーネント52と、テキスト処理コンポーネント54と、出力コンポーネント56とを含むが、提供されるコンポーネントは、より少なくてもよいし、より多くてもよいし、または異なってもよい。

10

【0028】

取り込みコンポーネント42は、印刷ドライバ22と印刷機14との中間の、ネットワーク1上、例えば、サーバコンピュータ34上、またはサーバコンピュータと通信可能に接続される別個の計算装置上に配置される。取り込みコンポーネントは、PDLファイル26を、これが印刷機に向かう途中で捕捉し、PDL文書20のコピーとグリフィンデックス30とを解析システムに転送する。解析システム44は、PDLフォーマットに従って取り込んだPDL文書20を解析する。

20

【0029】

解析システム44は、PDLフォーマットに従ってPDL文書20の各グリフのすべてのテキスト関連演算子を捕捉し、直接文字認識を行う。図示の解析システムは、TrueTypeフォントデータベース62およびグリフ特徴データベース64のうちの少なくとも1つを含む、参照データベース60にアクセスできる。

【0030】

一部の実施形態では、解析システム44は、図3に示すように構成され得る。この実施形態では、解析システムは、解析器78と、比較コンポーネント80と、抽出コンポーネント82と、レンダリングコンポーネント84と、特徴抽出部86と、類似度コンポーネント88とを含む。ただし、含まれているコンポーネントは、より少なくてもよいし、より多くてもよいし、または異なってもよい。

30

【0031】

TrueTypeフォントデータベース62は、既定の文字セットのそれぞれに関して、TrueType輪郭のリストを含み得る。例えば、TrueTypeデータベース62は、多数の、例えば、少なくとも10の、または少なくとも100、または少なくとも1000の入手可能なTrueTypeフォントから生成されたグリフ点座標の集まりであり得る。TrueTypeデータベース62に格納されたTrueTypeフォントは、多くのビジネス環境で使用されるフォントの大半を表現する、一般的な事務文書において一般的に利用されているフォントである。TrueTypeフォントにおけるグリフ90のアウトラインは、図4に示すように、直線線分と2次ベジェ曲線点とから作られる。

40

【0032】

グリフ特徴データベース64は、既定の文字セットのそれぞれに関して、ストックフォントから生成されたグリフ特徴セットを含み得る。よって、グリフ特徴データベース64は、予め計算したベクトルの集まりを含み得る。予め計算した特徴は、標準的なストックフォントから生成できる、または顧客固有のフォントに関して生成できる。グリフ特徴データベース64を作成するためには、図3に示すように、入力グリフがそれぞれ、各入力グリフの最外輪郭を正確に取り囲むビットマップ92としてレンダリングされる。ビットマップは、バイナリ値を有してもよいし、非バイナリ値を有してもよい。文字「a」のビットマップ92は、図4に描かれ、アレイの形態を取っている。各要素の値は、画像のその部分の色に対応する。文字「a」は、12×14のマトリクスで表現でき、ここで、マ

50

トリクス内の値は、ピクセル（画素）の輝度を示している。より大きな値は、より明るい領域に対応し、より低い値はより暗い領域に対応する。連続する同じ色のピクセルの数の計数、すなわちランレングス特徴は、各グリフに関してビットマップから抽出され、特徴ベクトルとしてグリフ特徴データベース64に格納され得る。

【0033】

一部の実施形態では、解析器78は、個々のテキスト文字を表現するグリフなどの個々のオブジェクトを識別し得る。グリフおよび関連しているテキスト文字は、一般的に、PDL文書20内においてテキスト描画プリミティブで符号化される。テキスト描画プリミティブは、テキストランにおけるグリフを描画するために使用される描画命令セットである。解析器78は、例えばPostScript（「PS」）またはPrinter Command Language（「PCL6」）などの、使用されるPDLフォーマットに従ってPDL文書20を解析でき、すべてのテキスト関連演算子またはグリフ描画命令を捕捉できる。他のPDLフォーマットとしては、限定するものではないが、PCL5、BBJL、Portable Document Format（「PDF」）、およびXML Paper Specification（「XPS」）が挙げられる。

10

【0034】

直接文字認識のための一実施形態では、テキスト演算子文字列引数における各グリフに関して、解析器78は、グリフインデックス30に基づいてグリフ形状を取得する。多くの文書がTrueTypeフォントを利用していることから、グリフインデックス30は、TrueTypeフォントを使用し得る。本実施形態では、グリフ識別器78が、PDL文書20におけるTrueTypeグリフに関する各描画命令を識別する。取り込んだPDL文書20から取得したグリフ形状がTrueType輪郭リストである場合、比較コンポーネント80は、描画命令を、グリフ特性のTrueTypeデータベース68と比較して、TrueTypeフォントデータベース62に一致するリストがあるかどうかを判定する。特に、比較コンポーネント80は、PDL文書20の描画命令にある入力グリフの直線線分および2次ベジェ曲線点を、TrueTypeデータベース62に格納された各文字90の直線線分および2次ベジェ曲線点と比較する。点のリスト間で一致するものが見つかった場合、テキスト文字が見つかり、抽出コンポーネント82によって抽出できる。

20

【0035】

直接文字認識のための別の実施形態では、例えば、取得したグリフ形状がTrueTypeフォントではない（例えば、解析器78が、TrueTypeデータベース68にあるTrueType輪郭ではない、グリフ形状に関する描画命令を識別する）場合、または一致するものが見つからなかった場合、グリフ特徴データベース64にアクセスできる。本実施形態では、レンダリングコンポーネント84は、ストックフォントに関しては、グリフを、その輪郭を正確に取り囲むビットマップ92（図3）としてレンダリングする。例えば、ビットマップの各ピクセルは、「on」または「off」のバイナリ値を有し、「on」は、グリフの輪郭内に実質的に含まれるピクセルを表し、「off」は背景ピクセルを表す。特徴抽出部86は、レンダリングコンポーネントによって生成されたビットマップを受信し、データベース64のストックフォントと同様の仕方で、生成されたビットマップ92から特徴セットを、例えば、1つ以上の特徴ベクトルの形態で抽出する。特徴ベクトルは、ビットマップからベクトルとして収集された、連続する同じ色のピクセルの数の計数、すなわちランレングス特徴であり得る。特徴ベクトルは、図4に示すアレイと同様に、アレイとして格納され得る。類似度コンポーネント88は、グリフ特徴データベース64の文字セットのそれぞれに関して、抽出された特徴ベクトルとグリフ特徴データベース64にある特徴ベクトルとの間の類似度を計算する。類似度の計算は、抽出された特徴ベクトルとグリフ特徴データベース64に格納された、予め計算されたベクトルの集まりとの間の距離の計算であり得る。例えば、ユークリッド距離またはバッチャリヤ距離が計算され得る。ただし、他の標準的な特徴抽出または距離の方法を代替的に使用できる。

30

40

50

【 0 0 3 6 】

ビットマップから抽出されたグリフ特徴セットとの類似度が閾値を超える、少なくとも1つの特徴セットがデータベース64において見つかった場合、グリフ特徴セットに対して最も類似の特徴セットを有する文字が識別される。これで、文字が見つかり、抽出できる。したがって、PDL文書20の描画命令において符号化されたグリフ特徴とグリフ特性のデータベースとの間で一致するものが見つかった場合、類似度コンポーネントによって計算された距離の計算に基づいて、抽出コンポーネント82は、最も類似のグリフに関連付けられたテキスト文字を抽出する。最初未知の入力グリフ形状とグリフ特徴データベース64との間におけるこの対応は、キャッシュコンポーネント50によってキャッシュ96に格納され得る。この次に解析器44が、同じ入力グリフ形状に関して、PDL文書20において同じ描画命令を識別した時には、システム10は、特徴抽出および距離比較に頼らずとも、キャッシュを使用して、グリフ特徴データベース64において直接一致するものを判定し得る。

10

【 0 0 3 7 】

文字が識別された時または一度解析器システム44がPDL文書20の解析を完了しすべてのテキスト文字を抽出した時、各識別された文字に関して、要約コンポーネント46は、各グリフに関連付けられた抽出されたテキスト文字のテキスト要約92を生成する。所与のテキスト関連演算子に関して解析器によって識別されたそれぞれの一致するものに関して、テキスト要約コンポーネント46は、抽出された文字と文書ページ内の場所とをテキスト要約に加える。テキスト要約92は、計算装置12などの計算装置、および/または装置12を操作するユーザによって読み出し可能とするフォーマットであり得る。テキスト要約92は、拡張マークアップ言語(「XML」)ファイルとして格納され得る。

20

【 0 0 3 8 】

一致するものが見つからなかった場合、所与の場所の候補の文字が識別できなかったという事実が解析器システムによって記録され得る。特に、なんらテキスト文字が認識されない場合、ログ取得コンポーネント48は、見つけ損ねたテキスト文字の記録をつける。ログ取得コンポーネントは、エラーログ94に見つけ損ねたテキスト文字を格納できる。エラーログ94は、例えばテキスト(「TXT」)ファイルとして、メモリ31に格納できる。

【 0 0 3 9 】

検証コンポーネント52は、識別された文字シーケンスにおける候補の単語を識別し、辞書98にアクセスして、候補の単語が存在するかを判定する。存在する場合、候補の単語が検証され、検証された単語シーケンスが出力され得る。

30

【 0 0 4 0 】

処理コンポーネント54は、検証された単語シーケンスの少なくとも一部を処理し、それに基づいて決定をレンダリングし得る。例えば、処理コンポーネントは、1つ以上のキーワードを検索し、文書が印刷許可されているかどうか、および/またはその印刷に関して顧客に課金するかどうかなどを決定する。

【 0 0 4 1 】

出力コンポーネント56は、検証された単語シーケンス、および/またはあらゆるエラー/検証されていない単語、および/または処理コンポーネント54によって出力された決定などの情報をシステム10から出力する。文書を印刷すべきである/すべきではないという決定の場合、情報が、印刷機に送信されて、印刷を許可/阻止し得る。決定が、課金すべき顧客である場合、情報は、課金システム(図示せず)に送信され得る。他の実施形態では、情報は、処理のために別のコンピュータ装置に送信され得る識別された単語シーケンスであり得る。

40

【 0 0 4 2 】

ネットワークリンク16は、サブネット、ローカルエリアネットワーク(LAN)、および/またはインターネットを含み得る。

【 0 0 4 3 】

50

入出力コンポーネント 36、38 はそれぞれ、変復調器（モデム）、ルータ、ケーブル、イーサネットポート、および/またはネットワーク 1 に接続された周辺機器が、例えばネットワーク管理者によって設定されたポリシーに従って有線または無線接続を介して他の装置と通信できるようにする他の通信装置（図示せず）を含み得る。計算装置 12 は、1 つ以上のネットワーク連携システム 16 を介して接続される、複数の PC または複数のワークステーションなどの複数の装置を含み得ることに留意されたい。

【0044】

印刷機または複数の印刷機 14 は、レーザ印刷機、インクジェット印刷機、LED 印刷機、プロッタ、および/またはインクまたはトナーなどのレンダリング媒体を用いて紙などの物理的媒体上に画像をレンダリングできる、任意の他の装置であり得る。印刷機 14 は、コンピュータ印刷機、ファクシミリ装置、デジタル複写機、多機能装置、および/または文書を印刷できる他の装置の形態を取り得る。

10

【0045】

接続部 16 は、計算装置 12 と印刷機 14 とをネットワーク 1 に接続する。接続部 16 は、適切な従来の通信プロトコルおよび/またはデータポートインターフェースを用いる有線または無線接続部として実装され得る。一般的に、接続部 16 は、装置間のデータ伝送を可能にする任意の通信チャネルであり得る。一実施形態では、例えば、装置には、適切な接続部 16 を通してデータを伝送するための、USB（商標）、SCSI、FIREWIRE（登録商標）、および/または BNC ポートなどのデータポート 36 を設けることができる。通信リンクは、無線リンクもしくは有線リンクまたは計算装置 12 と印刷機 14 との間の通信を可能にする任意の組み合わせであり得る。

20

【0046】

印刷機 14 は、ハードウェア、ファームウェア、もしくはソフトウェア、またはこれらの組み合わせによって制御され得る。PDL 文書 20 からの直接文字認識のためのシステム 10 は、例示的なコンピュータ 12 および/または印刷機 14、あるいは図 2 に示すような別個の装置のうちの 1 つ以上に配備され得る。例えば、印刷機 14 は、開示される実施形態と整合する仕方で印刷機 14 が直接文字認識の処理を最適化できるようにするソフトウェアまたはファームウェアを実行し得る。別の実施形態では、システム 10 は、コンピュータ 12 に常駐し、印刷機 14 のための PDL データに対して動作するのであってもよい。一般的に、システムコンポーネントは、システム内の 1 つ以上のコンピュータ 12、34 および/または印刷機 14 上において全体的または部分的に実施され得る。

30

【0047】

デジタルプロセッサ 32 は、汎用プロセッサ、専用プロセッサ、または組み込みプロセッサであり得る。プロセッサ 32 は、制御情報と命令とを含むデータをメモリ 31 と交換できる。メモリ 31 は、SDRAM または DRAM などの、任意のタイプのダイナミックランダムアクセスメモリ（「DRAM」）および/または読み出し専用メモリ（ROM）であり得る。命令 18 は、限定されるものではないが、ブートアップシーケンス、PDL の解析、プログラミング言語のためのコンパイラ、自動コード生成ルーチン、解釈されたページ記述言語における機能ベースのオブジェクト操作の処理のための最適化ルーチンなどの PDL を用いて記述された文書処理するルーチン、受信する要求およびメッセージを処理するためのルーチン、発信する応答およびメッセージを作成するためのルーチンを含む、1 つ以上の既定のルーチン、ならびに構成管理のためのルーチン、文書処理のためのルーチン、および他のコードを含む命令を保持し得る。一部の実施形態では、命令 18 のコードは、プロセッサ 32 によって作動させられる前にメモリ 31 にコピーされ得る。任意の PDL 処理ルーチンおよび最適化ルーチンを含む命令 18 は、コンピュータ 12 とネットワークリンク 16 とのうちの 1 つ以上を用いてアップグレード可能であり得る。

40

【0048】

一部の実施形態では、計算装置 12 は、印刷ドライバ 22 を介して印刷機 14 に、PDL を用いて指定した文書 20 の印刷可能データを送信し得る。印刷機 14 は、ルーチン呼び出して、PDL 文書 20 を解析器 44 で解析する。取り込みコンポーネント 42 は、

50

印刷ドライバ 2 2 によって P D L 文書 2 0 が生成された後であるが、印刷のために印刷機 1 4 に送信される前に、P D L 文書 2 0 を捕捉する。解析器 4 4 は、P D L データ 2 0 にある様々なオブジェクト、演算子、および構造を識別し、認識されたオブジェクト、演算子、および / または構造に関連付けられた動作を行うまたはアクションを開始することができる。

【 0 0 4 9 】

図 5 は、図 2 のコンピュータネットワークにおいて行われ得る、P D L 文書からの直接文字認識のための例示的な方法を示している。本方法は、S 1 0 0 で開始する。S 1 0 2 において、テキスト文字のそれぞれに関連付けられたグリフ特性セットによって定義されるテキスト文字の参照データベース 6 0 へアクセスできるようにする。参照データベースは、メモリに格納され得る。テキスト文字の参照データベースは、上述の T r u e T y p e フォントデータベースおよび / またはグリフ特徴データベースを含むことができる。

10

【 0 0 5 0 】

S 1 0 4 において、印刷ドライバによって、P D L フォーマットにおいて指定された文書の印刷ジョブが開始される。

【 0 0 5 1 】

S 1 0 6 において、P D L 文書 2 0 が図 1 および図 2 に示すシステム 1 0 によって受信または取り込まれる。P D L 文書 2 0 は、印刷文書において一連のグリフから形成されたテキストランを描画するための描画命令セットを含み得る。

【 0 0 5 2 】

S 1 0 8 において、P D L 文書 2 0 が解析されて、複数のグリフからなるテキストに関する任意の描画命令を捕捉または識別する。描画命令は、「ラン」としても知られるテキスト文字列のグリフを描画するために使用されるテキスト描画プリミティブの形態であり得る。P D L 解析は、印刷ドライバの出力と印刷機での文書の印刷との間にあるパイプライン上の任意の場所で実行され得る。P D L 解析は、印刷ドライバによって P D L 文書が生成された後、例えば、印刷ドライバ変換の最後の段階であるが、コンピュータまたは印刷サーバのポートモニタで文書が印刷される前に、実行され得る。必要に応じて、P D L 文書は、例えばスイッチまたはルータを介してネットワークから直接取り込まれてもよい。このことは、ネットワークが暗号化されていない場合または暗号化システムが分かっている場合に可能である。

20

30

【 0 0 5 3 】

S 1 1 0 において、S 1 0 8 において識別した描画命令を、グリフ特性のデータベースと比較して、描画命令によって表現されるテキスト文字を見つける。解析するステップが、T r u e T y p e グリフの描画命令を識別する場合、描画命令を、グリフ特性の T r u e T y p e データベースと比較する。上述のように、T r u e T y p e データベースは、単純に、数千の入手可能な T r u e T y p e フォントから容易に生成され得るグリフ点座標の集まりである。入力グリフ点は、一致するものが見つかるまで参照データベースのレコードと比較される。有利なことに、多くのビジネスでは、一般的な事務文書には T r u e T y p e フォントを使用している。必要に応じて、T r u e T y p e データベースは、グリフ検索をさらに加速するために、顧客固有のフォントと補足されまたは置き換えられ得る。

40

【 0 0 5 4 】

S 1 1 2 において、描画命令とグリフ特性のデータベースにあるテキスト文字との間で一致するものが見つかった場合、次いで、S 1 1 4 において、テキスト文字が抽出される。そうでない場合、本方法は、S 1 1 6 に進んでよい。

【 0 0 5 5 】

T r u e T y p e グリフ以外の他の技法が使用される場合、またはグリフ技法が T r u e T y p e であるが、データベースの所与の文書グリフに関して、一致する T r u e T y p e 点がない場合、グリフ形状が、さらなる処理のために、ビットマップとしてレンダリングされる (S 1 1 6) 。

50

【 0 0 5 6 】

S 1 1 8において、ビットマップに基づいて特徴ベクトルが抽出される。例えば、ビットマップから、連続する同じ色のピクセルの数を計数し、ランレングス特徴ベクトルとして格納する。

【 0 0 5 7 】

S 1 2 0において、抽出された特徴ベクトルとグリフ特徴データベースとの間の類似度が計算される。上述のグリフ特徴データベースは、標準的なストックフォントから容易に生成され得る、または必要に応じて顧客固有のニーズのために生成され得る、予め計算されたベクトルの集まりである。S 1 2 0の類似度計算は、自動画像分類ソフトウェアにおいて使用されるものなどの従来の特徴距離計算を用いる、単純な画像類似度予測であり得る。例えば、類似度の計算は、抽出された特徴ベクトルとグリフ特徴データベースに格納された、予め計算されたベクトルの集まりとの間の、ユークリッド距離またはバッタチャリヤ距離などの距離の計算であり得る。テキスト文字の単純な単色の形状に他の標準的な特徴抽出または距離の方法を使用してもよいが、ランレングスベクトルおよびユークリッドもしくはバッタチャリヤ距離は、非常に良好な結果をもたらすのに十分であることが分かっている。

10

【 0 0 5 8 】

S 1 2 0において、計算された距離の計算結果に基づいて、入力グリフ形状に十分に類似した少なくとも1つの特徴セットが見つかった場合、本方法は、S 1 1 4に進んで、類似度に基づいて入力グリフ形状に関連付けられたテキスト文字を抽出する。場合により、S 1 2 2において、最初未知の入力グリフ形状とグリフ特徴データベースとの間におけるこの対応は、キャッシュに格納され得る。この次にS 1 0 8における解析が、同じ入力グリフ形状に関して、P D L文書2 0において同じ描画命令を識別した時には、システム1 0は、S 1 1 8の特徴抽出およびS 1 2 0の類似度計算に頼らずとも、グリフ特徴データベースにおいて直接一致するものを判定し得る。

20

【 0 0 5 9 】

S 1 2 0において、一致するテキスト文字が見つからない場合、本方法はS 1 2 4に進み、ここで、エラーのログが取られる。見つけ損ねたテキスト文字は、例えばT X Tファイルとしてメモリにログを取られ、格納され得る。

【 0 0 6 0 】

本方法は、S 1 1 4からS 1 2 6に進み、ここで、S 1 1 4において抽出されたテキスト文字とS 1 2 4においてログを取ったエラーとに基づいてテキスト要約が生成される。生成されるテキスト要約は、システム1 0を操作する機械およびユーザの両方に可読であり得、X M Lフォーマットで格納され得る。

30

【 0 0 6 1 】

場合により、S 1 2 8において、抽出されたテキスト文字を含む、S 1 2 6で生成されたテキスト要約は、辞書を用いてクロスチェックされて、抽出されたテキスト文字から形成された単語を検証し得る。例えば、P D L文書は、大文字の「i」の描画命令と同じものになる、小文字の「L」の描画命令を含むことがある。これが当てはまるのは、両方のテキスト文字がまっすぐな垂直の線「I」で描画されるフォントタイプの場合である。この場合、辞書を用いてクロスチェックをすることにより、無効な文字認識を抑制し、S 1 2 6において生成されるテキスト要約において正しいテキスト文字が表現されることを確実にすることによって両義性を低下させることができる。必要に応じて、一般的な辞書または専門辞書を用いてクロスチェックを行って、抽出されたテキスト文字から形成された特定の単語を認識することができる。専門辞書は、専門的なテキスト文字および単語の認識を要求している顧客から提供され得る。S 1 2 4においてログを取ったエラーは、候補の単語を、辞書と比較する際には「任意の文字」とみなされ得る。

40

【 0 0 6 2 】

S 1 3 0において、場合によって検証されるテキスト要約などの情報が出力される。

【 0 0 6 3 】

50

本方法は、S 1 3 2で終了する。

【0064】

本方法は、説明される機能を行うためのソフトウェアで実装され得る。例示的な疑似コードを以下のアルゴリズム1に示す。

【0065】

【表1】

アルゴリズム1	
テキスト演算子「文字列」引数における各グリフに関して	
{引数として与えられるグリフインデックスに基づいてグリフ形状を取得する	10
If (グリフ形状がTrueType輪郭リストである)	
{TrueTypeデータベースにおいて一致するものを探す	
If (一致するものが見つかった)	
{文字が発見された}}	
If (グリフ形状がTrueType輪郭リストではないまたは一致するものがない)	
{グリフをその輪郭を正確に取り囲むビットマップとしてレンダリングする与えられたビットマップの特徴を抽出する	
参照特徴データベースを用いてこれらの特徴の類似度を計算する	20
If (候補に十分に近い少なくとも1つの特徴セットが見つかった)	
{文字が発見された。これは、いかなる参照特徴よりも最も近い}}	
If (一致するものがない)	
{所与の場所で文字を見つけ損ねたという事実の痕跡を保持する}}	

【0066】

図5に示す本方法は、コンピュータ上で実行され得るコンピュータプログラム製品に実装され得る。コンピュータプログラム製品は、ディスク、ハードドライブなどの、制御プログラムが記録される(格納される)非一時的コンピュータ可読記録媒体を含み得る。一般的な形態の非一時的コンピュータ可読媒体は、例えば、フロッピーディスク、フレキシブルディスク、ハードディスク、磁気テープまたは任意の他の磁気記憶媒体、CD-ROM、DVD、または任意の他の光学式媒体、RAM、PROM、EPROM、FLASH-EPROM、または他のメモリチップもしくはカートリッジ、あるいはコンピュータが読み出しおよび使用し得る任意の他の有形の媒体を含む。

【0067】

あるいは、本方法は、音波もしくは光波など、ならびに無線および赤外線データ通信中に生成される媒体などの伝送媒体を使用して、制御プログラムがデータ信号として埋め込まれている送信可能な搬送波などの一時的媒体に実装され得る。

【0068】

本例示的な方法は、1つ以上の汎用コンピュータ、専用コンピュータ(複数可)、プログラムされたマイクロプロセッサもしくはマイクロコントローラおよび周辺集積回路素子、ASICもしくは他の集積回路、デジタル信号プロセッサ、離散要素回路などのハードウェアにより実現されている電子回路または論理回路、PLD、PLA、FPGA、グラフィカルカードCPU(GPU)、またはPALなどのプログラマブルロジックデバイスなどの上に実装され得る。一般的に、図5に示す流れ図を実施できる任意の装置を用いて、本方法を実施できる。

40

【 図 5 】

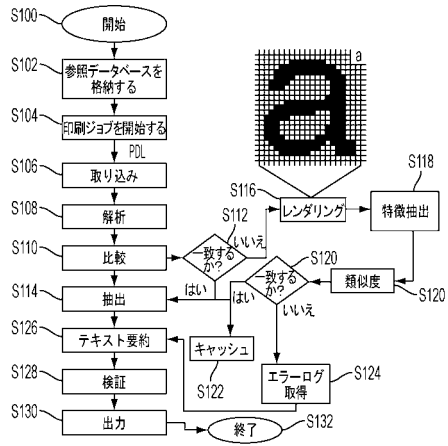


図 5

フロントページの続き

(51)Int.Cl.	F I	テーマコード(参考)
	G 0 6 F 3/12	3 4 5
	G 0 6 F 3/12	3 7 3
	B 4 1 J 29/38	Z

(72)発明者 イヴ・オプノ

フランス共和国 ノートル - ダム - ド - メザージュ 3 8 2 2 0 ル・シャンボール 1 5 2

Fターム(参考) 2C061 AP01 AQ05 AQ06 AR01 AR03 HK03
2C187 AC06 AC08 AE07 BG03 BG19 BG49 CD04 JA07
5B064 AA10 DC26 EA08 EA19