

(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2014-106961

(P2014-106961A)

(43) 公開日 平成26年6月9日 (2014. 6. 9)

(51) Int.Cl.		F I			テーマコード (参考)
G06K 9/62	9/62	(2006.01)	G06K 9/62	630Z	5B064
G06T 7/00	7/00	(2006.01)	G06T 7/00	350A	5L096

審査請求 未請求 請求項の数 16 O L (全 22 頁)

(21) 出願番号	特願2013-118680 (P2013-118680)	(71) 出願人	511072895
(22) 出願日	平成25年6月5日 (2013. 6. 5)		キング・アブドゥルアジズ・シティ・フォー・サイエンス・アンド・テクノロジー (ケイ・エイ・シー・エス・ティ)
(31) 優先権主張番号	13/685, 088		KING ABDULAZIZ CITY FOR SCIENCE AND TECHNOLOGY (KACST)
(32) 優先日	平成24年11月26日 (2012. 11. 26)		サウジアラビア、11442 リヤド、ビィ・オウ・ボックス・6086、ザ・ナショナル・センター・フォー・テクノロジー・ディベロップメント
(33) 優先権主張国	米国 (US)	(74) 代理人	110001195 特許業務法人深見特許事務所

最終頁に続く

(54) 【発明の名称】 アラビア語テキストを自動的に認識するためのコンピュータによって実行される方法、およびコンピュータプログラム

(57) 【要約】

【課題】アラビア語テキストの認識において、適切にテキスト特徴を抽出する。

【解決手段】アラビア語の文字のラインがデジタル化されることにより、各々が2進数で表現された画素値に関連付けられた二次元の画素の配列が形成される。画素値は2進数で表現される。さらに、アラビア語の文字のラインが複数のライン画像へと分割されて、複数のライン画像の中の1つにおいて複数のセルが規定される。複数のセルの各々は、隣接した画素のグループを有する。さらに、複数のライン画像の中の1つにおいて複数のセルの各々の画素の画素値がシリアル化されることにより、2値セル番号が形成される。また、複数のライン画像の中の1つにおける複数のセルから取得された2値セル番号に従ってテキスト特徴ベクトルが形成される。そして、テキスト特徴ベクトルが隠れマルコフモデルに送られることによりアラビア語の文字のラインが認識される。

【選択図】図1

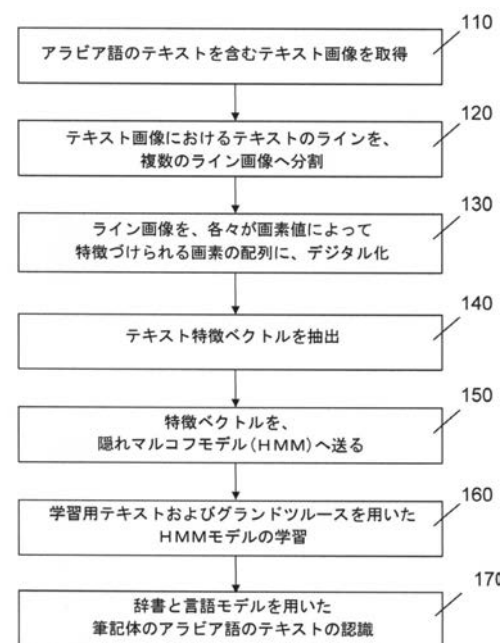


Figure 1

【特許請求の範囲】**【請求項 1】**

アラビア語テキストを自動的に認識するための、コンピュータによって実行される方法であって、

アラビア語の文字のラインを含むテキスト画像を取得することと、

当該アラビア語の文字のラインをデジタル化することにより、各々が 2 進数で表現された画素値に関連付けられた二次元の画素の配列を形成することとを備え、前記二次元の画素の配列は、第 1 の方向における複数の行と第 2 の方向における複数の列とを含み、

前記方法は、さらに、

画素の列における画素のストリング中の同じ画素値の連続する画素の頻度をカウントすることを備え、各々が異なる画素値を有する隣接した画素のストリングは、それらの間での遷移によって規定され、前記カウントすることは、さらに、

列の遷移数が予め定められた足切遷移番号に達したときに、当該列における同じ画素値の連続する画素の頻度のカウントを停止することと、

前記画素の列におけるストリングから取得される頻度カウントを用いてテキスト特徴ベクトルを形成することと、

当該テキスト特徴ベクトルを隠れマルコフモデルに送ることによりアラビア語の文字のラインを認識することとを備える、方法。

【請求項 2】

前記アラビア語の文字のラインは、複数のアラビア語の単語を含む、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 3】

前記テキスト特徴ベクトルは、前記画素の列における連続する画素のストリングから取得された一連の頻度カウントによって形成される、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 4】

前記予め定められた足切遷移番号は、前記アラビア語の文字のラインをデジタル化するステップに先立つ、アラビア語テキストについての統計的解析によって取得される、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 5】

前記予め定められた足切遷移番号は 6 である、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 6】

前記二次元の配列における画素値は、単一のビットの 2 進数で表現される、請求項 1 に記載のコンピュータによって実行される方法。

【請求項 7】

前記頻度をカウントすることは、

列における最初の 1 またはそれ以上の画素の画素値が「0」であるときに、第 1 の頻度カウントの値に「0」を割り当てることを含み、前記第 1 の頻度カウントの次に、当該列の初めに画素値「0」を有する連続した画素の数が続く、請求項 6 に記載のコンピュータによって実行される方法。

【請求項 8】

前記頻度をカウントすることは、

列の頂点の 1 またはそれ以上の画素の画素値が「1」であるときに、第 1 の頻度カウントの値として「0」を割り当てることを含み、前記第 1 の頻度カウントの次に、当該列の初めに画素値「1」を有する連続した画素の数が続く、請求項 6 に記載のコンピュータによって実行される方法。

【請求項 9】

コンピュータに以下のことを実行させるためのコンピュータ読取可能なプログラムであって、プログラムコード関数を含み、前記プログラムコード関数は、コンピュータに、

10

20

30

40

50

アラビア語の文字のラインを含むテキスト画像を取得させ、

アラビア語の文字のラインをデジタル化させることにより、各々が2進数で表現された画素値に関連付けられた二次元の画素の配列を形成させ、前記二次元の画素の配列は、第1の方向における複数の行と第2の方向における複数の列とを含み、

前記プログラムコード関数は、さらに、前記コンピュータに、画素の列における画素のストリング中の同じ画素値の連続する画素の頻度をカウントさせ、各々が異なる画素値を有する隣接した画素のストリングはそれらの間での遷移によって規定され、前記カウントするステップは、さらに、前記列における遷移の数が予め定められた足切遷移番号に到達したときに、同じ画素値の連続する画素の頻度のカウントを停止することを含み、

前記プログラムコード関数は、前記コンピュータに、

前記画素列におけるストリングから取得される頻度カウントを用いてテキスト特徴ベクトルを形成することと、

前記テキスト特徴ベクトルを隠れマルコフモデルに送ることによりアラビア語の文字のラインを認識することとを実行させる、コンピュータプログラム。

【請求項10】

前記アラビア語の文字のラインラインは、複数のアラビア語の単語を含む、請求項9に記載のコンピュータプログラム。

【請求項11】

前記テキスト特徴ベクトルは、前記画素の列における連続する画素のストリングから取得された一連の頻度カウントによって形成される、請求項9に記載のコンピュータプログラム。

【請求項12】

前記予め定められた足切遷移番号は、前記アラビア語の文字のラインをデジタル化するステップに先立つ、アラビア語テキストについての統計的解析によって取得される、請求項9に記載のコンピュータプログラム。

【請求項13】

前記予め定められた足切遷移番号は6である、請求項9に記載のコンピュータプログラム。

【請求項14】

前記二次元の配列における画素値は、単一のビットの2進数で表現される、請求項9に記載のコンピュータプログラム。

【請求項15】

前記頻度をカウントするステップは、

列における最初の1またはそれ以上の画素の画素値が「0」であるときに、第1の頻度カウントの値に「0」を割り当てることを含み、前記第1の頻度カウントの次に、当該列の初めに画素値「0」を有する連続した画素の数が続く、請求項9に記載のコンピュータプログラム。

【請求項16】

前記頻度をカウントするステップは、

列の頂点の1またはそれ以上の画素の画素値が「1」であるときに、第1の頻度カウントの値として「0」を割り当てることを含み、前記第1の頻度カウントの次に、当該列の初めに画素値「1」を有する連続した画素の数が続く、請求項9に記載のコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本特許出願は、同一発明者によって2011年12月14日に出願され、同一譲受人の、係属する米国特許出願第13/325,789号、名称「効果的なアラビア語テキスト特徴の抽出に基づく、アラビア語テキスト認識のためのシステムおよび方法」の継続出願であり、当該出願についての優先権を主張する。米国特許出願第13/325,789号

10

20

30

40

50

は、同一発明者によって2009年4月27日に出願された、米国特許出願第12/430,773号、名称「効果的なアラビア語テキスト特徴の抽出に基づく、アラビア語テキスト認識のためのシステムおよび方法」の継続出願であり、その開示は、ここに参照により組み込まれる。

【背景技術】

【0002】

発明の背景

本願は、概して、アラビア語テキストの自動的な認識に関する。

【0003】

テキスト認識、つまり、テキストの自動読取は、パターン認識の一分野である。テキスト認識の目的は、印刷されたテキストを、人間の精度で、かつ、より速く、読取ることである。多くのテキスト認識の方法は、テキストが個々の文字へと分離できることを前提としている。このような技術では、タイプライタで打たれた、または、活字に組まれたラテン語については首尾よくいくが、アラビア語のような筆記体には、信頼できる程度に適用することはできない。これまでのアラビア語の手書きテキストの認識についての研究によれば、アラビア語の単語を個々の文字へとセグメント化する試みにおける困難性が確認されている。




【0004】

【数1】

アラビア語は、テキスト認識のアルゴリズムにいくつかの課題を提供している。アラビア語の手書き文字は、本質的に、筆記体であり、分離された文字をブロック体で記述することは容認できない。しかも、アラビア文字の形は、文脈に依存することがあり、つまり、単語における当該文字の位置に依存し得る。たとえば、次の

「」に示されるような文字は、4つの異なる外形、つまり、「」におけ

る独立形「」、「」における語頭形「」、「」における語中形

「」、「」における語尾形「」を有する。さらに、すべてのアラビア文字が、単語の中で連結されるわけではない。単語内において、空間が特定の文字を分割する場合もあるため、単語間の境界を自動的に決定することは困難な場合がある。

【0005】

アラビア語テキストの認識には、統計モデルのような異なる分類体系が適用されてきた。しかしながら、適切にテキスト特徴を抽出することは、未だ、正確なアラビア語テキストの認識を達成することにおいての主要な障害のままである。

【発明の概要】

【発明が解決しようとする課題】

【0006】

発明の概要

概略的な側面において、本願発明は、アラビア語テキストを自動的に認識するための方法に関する。当該方法は、アラビア語の文字のラインを含むテキスト画像を取得することと、アラビア語の文字のラインをデジタル化することにより、各々が画素値に関連付けられた二次元的の画素の配列を形成することとを含み、画素値は2進数で表現され、上記方法は、さらに、アラビア語の文字のラインを複数のライン画像へと分割することと、複数のライン画像の中の1つにおいて複数のセルを規定することとを含み、複数のセルの各々は、隣接した画素のグループを有し、上記方法は、さらに、複数のライン画像の中の1つ

10

20

30

40

50

において複数のセルの各々の画素の画素値をシリアル化することにより2値セル番号を形成することと、複数のライン画像の中の1つにおける複数のセルから取得された2値セル番号に従ってテキスト特徴ベクトルを形成することと、テキスト特徴ベクトルを隠れマルコフモデル(Hidden Markov Model)に送ることによりアラビア語の文字のラインを認識することを含む。

【0007】

他の概略的な局面において、本願発明は、アラビア語テキストを自動的に認識するための方法に関する。当該方法は、アラビア語の文字のラインを含むテキスト画像を取得すること、アラビア語の文字のラインをデジタル化することにより、各々が2進数で表現された画素値に関連付けられた二次元の画素の配列を形成することを含み、二次元の画素の配列は、第1の方向における複数の行と、第2の方向における複数の列とを含み、上記方法は、さらに、画素の列において同じ画素値を有する連続する画素の頻度をカウントすることと、画素の列から得られた頻度カウントを利用してテキスト特徴ベクトルを形成することと、当該テキスト特徴ベクトルを隠れマルコフモデルへ送ることによりアラビア語の文字のラインを認識することを含む。

10

【0008】

他の概略的な局面において、本願発明は、アラビア語テキストを自動的に認識するための方法に関する。当該方法は、アラビア語の文字のラインを含むテキスト画像を取得することと、アラビア語の文字のラインをデジタル化することにより、各々が画素値に関連付けられた二次元の画素の配列を形成することと、当該アラビア語の文字のラインを複数のライン画像へと分割することと、当該複数のライン画像の少なくとも1つを小型化することにより小型化されたライン画像を生成することと、小型化されたライン画像の各々の列の画素の画素値をシリアル化することにより一連のシリアル化された番号を形成することとを含み、一連のシリアル化された番号はテキスト特徴ベクトルを形成し、上記方法は、さらに、当該テキスト特徴ベクトルを隠れマルコフモデルへ送ることによりアラビア語の文字のラインを認識することを含む。

20

【0009】

他の概略的な局面において、本願発明は、コンピュータ読取可能なプログラムコード関数を含むコンピュータプログラムに関し、当該コード関数は、コンピュータに、アラビア語の文字のラインを含むテキスト画像を取得させ、アラビア語の文字のラインをデジタル化させることにより、各々が画素値に関連付けられた二次元の画素の配列を形成させ、当該画素値は2進数で表現され、上記コード関数は、上記コンピュータに、さらに、アラビア語の文字のラインを複数のライン画像へと分割させ、複数のライン画像の中の1つにおける複数のセルを規定させ、複数のセルの各々は隣接する画像のグループを有し、上記コード関数は、上記コンピュータに、さらに、複数のライン画像の中の1つにおける複数のセルの各々の画素の画素値をシリアル化させ、複数のライン画像の中の1つにおける複数のセルから取得された2進数のセル番号に応じてテキスト特徴ベクトルを形成させ、当該テキスト特徴ベクトルを隠れマルコフモデルに送ることによりアラビア語の文字のラインを認識させる。

30

【0010】

システムの実現は、以下に示されたもののうち1またはそれ以上を含む場合がある。上記方法は、さらに、2進数のセル番号を10進数のセル番号へと変換することと、複数のライン画像の中の1つにおける複数のセルから取得された10進数のセル番号をシリアル化することにより一連の10進数のセル番号を形成することと、複数のライン画像の中の1つにおける複数のセルから取得された一連の10進数のセル番号に従ってテキスト特徴ベクトルを形成することとを含み得る。二次元の画素の配列は、第1の方向における複数の行と、第2の方向における複数の列とを含み得る。アラビア語の文字のラインは、実質的に上記第1の方向に沿って並び得る。複数のライン画像は、上記第1の方向に沿って連続的に並び得る。複数のライン画像の中の少なくとも1つは、第1の方向におけるM個の行によって定義される高さ、第2の方向におけるN個の列によって規定される幅とを有

40

50

し得る。MおよびNは、整数である。二次元の画素の配列は、N行の画素を含み得る。Nは、2とおよそ100との間の範囲にあり得る。Nは、3とおよそ10との間の範囲にあり得る。二次元の画素の配列における画素値は、単一のビットの2進数で表現され得る。二次元の画素の配列における画素値は、マルチビットの2進数で表現され得る。隠れマルコフモデルは、隠れマルコフモデルツールキットとして実装され得る。

【0011】

本願において記述されるシステムおよび方法は、アラビア語テキストにおける特徴の抽出のための、包括的な、定量的な、かつ正確な技術を提供する。開示されたアラビア語の文字の認識は、いくつかの従来技術よりも、より効率的であり、かつ計算時間が短い。開示されたシステムおよび方法は、さらにいくつかの従来技術よりも、より単純かつ私

10

【0012】

発明は複数の実施例を参照することにより具体的に示され記述されているが、形式上の種々の変更や詳細は、発明の精神および範囲を離れることなくなされ得ることが、当業者によって、理解されるであろう。

【0013】

図面の簡単な説明

以下の図面は、出願書類に組込まれかつその一部を形成し、本願発明の実施例を説明し、かつ、明細書とともに、発明の本質を説明するために供される。

【図面の簡単な説明】

20

【0014】

【図1】本開示におけるアラビア語のテキスト認識の工程を説明するためのフロー図である。

【図2】アラビア語テキストを含むテキスト画像を説明する図である。

【図3A】テキスト画像を、各々が複数の画素を含む複数のライン画像へと分割することを説明する図である。

【図3B】図3Aに示されたライン画像の一部分における、画素および画素値を説明する図である。

【図3C】図3Aに示されたライン画像の一部分における、画素および画素値を説明する図である。

30

【図4】本願に従ったテキスト特徴抽出の方法を説明する図である。

【図5】図4に示されたテキスト特徴抽出の工程を説明するフロー図である。

【図6】本願に従ったテキスト特徴抽出の他の方法を説明する図である。

【図7A】本開示に従った他のテキスト特徴抽出方法を説明する図である。

【図7B】本開示に従った他のテキスト特徴抽出方法を説明する図である。

【図7C】本開示に従った他のテキスト特徴抽出方法を説明する図である。

【図7D】本開示に従った他のテキスト特徴抽出方法を説明する図である。

【図8】図7A～図7Dに示されたテキスト特徴抽出の工程を説明するフロー図である。

【発明を実施するための形態】

【0015】

40

発明の詳細な説明

図1は、本発明に従ったアラビア語のテキスト認識の概略的な流れを説明する。図1～図3Cを参照して、アラビア語のテキスト文書から、テキスト画像200が取得される(図1のステップ110)。テキスト画像200におけるアラビア語テキストは、複数のテキストライン211～214に配置され得、その各々は、筆記体のアラビア語の文字のストリングを含む。テキストライン211～214は、複数のライン画像311～313へと分割される(図1のステップ120)。ライン画像311, 312, または313は、それから、各々が画素値を割り当てられた画素321～323へと分割される(図1のステップ130)。ライン画像311, 312, または313の幅は、2画素と100画素との間の範囲にあり得、または、3画素と10画素との間の範囲にあり得る。ライン画像

50

3 1 1 , 3 1 2 , または 3 1 3 は、完全な文字、部分的な文字、または結合した文字を含み得る。

【 0 0 1 6 】

画素値は、特定の画素の位置でのテキスト画像 2 0 0 の明度値を表わす。ある実装では、明度値が高いことは、白色背景に位置し得る画素における明るい画像の色（または、低密度）を表す。明度値が低いことは、一筆のアラビア語の文字（a stroke of an Arabic character）内に位置し得る暗い画像の色（または、高密度）を表わす。画素値は、2 進数、1 0 進数、および 1 6 進数のような、異なる計数法で表現されてもよい。

【 0 0 1 7 】

図 3 A ~ 図 3 C を参照して、ライン画像 3 1 1 は、複数の画素 3 2 1 - 3 2 3 を含む画像部分 3 2 0 を含む。画素 3 2 1 - 3 2 3 の各々は、2 進数の画素値「 0 」または「 1 」を割り当てられている。画素値「 1 」は、白色の背景を表わす。画素値「 0 」は、一筆のアラビア語の文字内にある、暗画像色（つまり、低い明度値）を表わす。開示されたシステムおよび方法は、2 進数で表わされたマルチビットの画素値にも適合可能であり、当該 2 進数で表わされたマルチビットの画素値は、多段階のトーンレベル（たとえば、グレースケール）で、画像濃度を表わし得る。

【 0 0 1 8 】

本開示に従うと、テキスト特徴ベクトルは、テキストライン 2 1 1 またはライン画像 3 1 1 - 3 1 3 から抽出され得る（図 1 のステップ 1 4 0 ）。テキスト特徴抽出のさまざまな実装の詳細については、以下に、図 4 ~ 図 8 に関連付けられて、議論される。テキスト特徴ベクトルの厳密な形態は、以下に記載されるように、抽出方法によって変化し得る。

【 0 0 1 9 】

ステップ 1 4 0 において取得された特徴ベクトルは、次に、隠れマルコフモデル（HMM）に送られる（図 1 のステップ 1 5 0 ）。本開示では、HMM は、隠れマルコフモデルツールキット（HTK）によって実装される場合があり、それは、隠れマルコフモデルを構築し操作するための移植可能なツールキットである。HTK は、語彙集がなく、学習用サンプル文字からのモデルおよび文法に依存する。HMM は、確率解釈を提供し、特徴ベクトルにおいて見い出されたパターンにおける変化を許容し得る。HTK の機能性の大部分は、C ソースコードで利用可能なライブラリモジュールに組込まれ得る。これらのモジュールは、従来のコマンドライン形式のインターフェイスで動作するように設計されているため、HTK ツールの実行を制御するためのスクリプトの記述がシンプルになる。

【 0 0 2 0 】

HMM は、既知のアラビア語の単語を含むテキスト画像から取得された特徴ベクトルを用いることによって、学習させることができる（データ転記）（図 1 のステップ 1 6 0 ）。HTK は、学習用サンプルのための文字モデルとグランドツルース（ground truth）とともに提供される。文字のモデル化のためのコンポーネントは、特徴ベクトルとそれに対応するグランドツルースとを利用し、文字モデルを評価する。学習用サンプルによって生成された観察結果は、モデルパラメータを調整するのに用いられるが、テスト用のサンプルによって生成された観察結果は、システムの性能を調査するのに利用される。モデルの各状態は、アルファベットの組における字を表わし、各特徴ベクトルは、1 つの観察結果に相当する。HTK 学習ツールは、準備された学習用データを利用して文字モデルパラメータを調整し、既知のデータ転記を予測することができる。

【 0 0 2 1 】

HMM パラメータは、学習用画像セグメントのためのグランドツルースから推定された。このセグメント化は輪郭にも適用されて、セグメント化のポイントを発見し、これらのセグメントから特徴を抽出し、そして、特徴ベクトルを観察シーケンスに伝達し得る。セグメント化を基礎とした技術は、単語の画像と文字列とを一致させるためのダイナミックプログラミングに利用される。学習段階では、テキスト画像に相当するテキストであるグランドツルースと一体となった、走査されたテキストのラインが、入力として取得される。そして、各ラインは、狭い縦割りの窓へと分割され、そこから特徴ベクトルが抽出され

10

20

30

40

50

る。

【 0 0 2 2 】

学習した H M M は、辞書および言語モデルを利用して、特徴ベクトルにおけるアラビア語テキストを認識するために用いられる（図 1 のステップ 1 7 0）。認識段階は、最も高い尤度の文字シーケンスを見つけるための学習段階において推定された異なる知識源とともに用いられる特徴ベクトルを抽出するのと同じ工程に引き続く。認識ツールは、あるモデルから他のモデルへの遷移確率を記述するために、ネットワークを必要とする。辞書および言語モデルが当該ツールに入力され、認識装置が正しい状態シーケンスを出力するのに役立つことができる。

【 0 0 2 3 】

いくつかの実施形態では、図 3 A ~ 図 5 を参照して、ライン画像 3 1 1 - 3 1 3 は、各々が画素値によって特徴付けられる画素 3 2 1 - 3 2 3 の配列へとデジタル化される（図 5 のステップ 5 1 0）。ライン画像 3 1 1 は、図 4 に示されるように、複数のセル 4 1 0 - 4 6 0 へと分割される（図 5 のステップ 5 2 0）。セル 4 1 0 - 4 6 0 の各々は、3 × 3 画素の配列のような、隣接する画素のグループを含む。たとえば、セル 4 2 0 は、画素 4 2 2, 4 2 3 および他の画素を含む。

【 0 0 2 4 】

次に、各々のセルの画素値が、2 進数のセル番号で表わされる（図 5 のステップ 5 3 0）。各セルにおける画素値は、まず、シリアル化される。たとえば、セル 4 2 0 における 9 つの画素 3 2 2 - 3 2 3 は、連続する 3 行の順に、次のようにシリアル化される：1, 1, 1, 1, 0, 0, 1, 0, 0。一連の 2 進数の画素値は、その後、9 ビットの 2 進数のセル番号へとマップされる。画素 3 2 2 の画素値は、最上位ビットにマップされ、画素 3 2 3 の画素値は、最下位ビットにマップされる。結果として、セル 4 2 0 における画素値は、2 進数で表わされる 9 ビットのセル番号 1 1 1 1 0 0 1 0 0 で表わされる。同様に、セル 4 1 0 - 4 6 0 における画素値が、それぞれが 0 と 5 1 1 との間の範囲にある、2 進数で表わされるセル番号 4 8 0 へと変換される。

【 0 0 2 5 】

ライン画像 3 1 1 のセルにおける、2 進数のセル番号は、次に、1 0 進数のセル番号 4 9 0 へと変換される（図 5 のステップ 5 4 0）。1 0 進数のセル番号 4 9 0 は、その後、ライン画像 3 1 1 のための特徴ベクトルを形成するためにシリアル化される（図 5 のステップ 5 5 0）。ステップ 5 2 0 - 5 5 0 は、別のライン画像のために繰返される。別のライン画像 3 1 1 - 3 1 3 からの特徴ベクトルは、その後、隠れマルコフモデルへと送られ、テキストラインにおけるアラビア語の文字を認識する（図 5 のステップ 5 6 0）。

【 0 0 2 6 】

図 4 ~ 図 5 と併せて記述された上記の抽出方法は、図 1 において説明された処理のためのテキスト特徴抽出の実装を表す。上記のテキスト特徴抽出方法は、データストリングにおけるマルチビットの画素値および他の数値表現に適合することが理解されるべきである。たとえば、画素値は、テキスト画像におけるグレースケール情報（または、マルチトーン）を取り込むことのできる、3 ビットまたは 5 ビットの 2 進数によって表わされ得る。マルチビットの画素値は、ストロークのエッジに沿ったテキスト特徴の記述の精度を改善し得る。

【 0 0 2 7 】

さらに、2 進数の代わりに、画素値は、最小値と最大値との間のいかなる数値範囲によっても表わされ得る。いくつかの実装においては、画素値は、[0, 1] または [- 1, 1] のような、所定の範囲に比例した（または、正規化された）値となり得る。そして、画素値は、量子化され得る。特徴ベクトルは、ステップ 5 3 0 - 5 5 0 と同様に取得され得る。

【 0 0 2 8 】

いくつかの実施形態では、図 6 を参照して、ライン画像 6 1 0 は、分解能において縮小され（つまり、小型化され）、これにより、小型化されたライン画像 6 2 0 が形成される

10

20

30

40

50

。たとえば、ライン画像 6 1 0 は、6 0 画素の高さを有し得る。小型化されたライン画像 6 2 0 は、1 / 3 倍の寸法で、2 0 画素の高さを有し得る。小型化されたライン画像 6 2 0 は、各々が画素値によって表わされる画素の配列 6 3 0 を形成するために、デジタル化される。配列 6 3 0 における各列の画素値は、2 進数を形成するために、シリアル化される。異なる列からの 2 進数は特徴ベクトルを形成するデータストリング 6 4 0 を形成する。テキストラインのライン画像から取得された特徴ベクトルは、隠れマルコフモデルへ送られ、これにより当該テキストラインにおけるアラビア語の文字を認識することができる (図 5 のステップ 5 6 0)。

【 0 0 2 9 】

図 7 A , 図 7 B , および図 8 を参照して、ライン画像 7 0 0 は、ステップ 5 1 0 (図 5) と同様に、画素の配列へとデジタル化される (図 8 のステップ 8 1 0)。画素は、複数の列に配置される。画素値は、値「1」または値「0」を有する、単一のビットの 2 進数によって表わされる。各列の画素値がシリアル化されることにより、単一のビットの 2 進数の列が形成される (図 8 のステップ 8 3 0)。

【 0 0 3 0 】

次に、図 7 C および図 7 D に示されるように、値「0」および値「1」の、同じ 2 進数の画素値を有する連続した画素の頻度が、計算される (図 8 のステップ 8 4 0) 。当該頻度は、足切遷移番号 (cut off transition number) まで、カウントされる。当該頻度を表形式化して、頻度カウント 7 5 0 および 7 6 0 を形成する (図 8 のステップ 8 5 0) 。

【 0 0 3 1 】

【 数 2 】

0	1
0	1
0	1
1	0
1	0
0	1
0	1
1	0

【 0 0 3 2 】

以外同じ数の遷移を有する 2 つの画素の列を区別するために、列の最上部の画素から値「1」の数のカウントを開始することによって、頻度カウントが実行される。左側の列では、初めは、画素値「1」のカウントは「0」であり、「3」カウントの画素値「0」が続く。当該 2 つの列におけるコンプリメンタリ画素値は、結果として、次の頻度カウントのようになる：

【 0 0 3 3 】

【 数 3 】

0	3
3	2
2	2
2	1
1	0
0	0

【 0 0 3 4 】

各列の初めにおける、当初の画素カウントが、本発明の精神から逸脱することなく、画素値「0」について行なうこともできることが、理解されるべきである。

【0035】

表形式の頻度カウント750, 760(図7C, 図7D)における各行は、白色の背景(画素値「1」を有する)から暗テキスト領域(画素値「0」を有する)への、またはその逆の、画素値における遷移を表わしている。データを圧縮するために、頻度カウントが、最大遷移番号で切り捨てられている。

【0036】

表形式の頻度カウント750, 760の各列における頻度カウントは、特徴ベクトルを形成している(図8のステップ860)。したがって、本実施の形態では、各列は、ベクトルと称することもできる。ライン画像におけるさまざまな列からの特徴ベクトルが、隠れマルコフモデルへ送られる(図8のステップ870)。

【0037】

最大遷移番号は、アラビア語テキストの大標本についての統計的解析によって決定される。表1に示されるように、およそ99.31%の列が、6以下の遷移を有している。換言すれば、テキスト画像の大多数が、足切遷移番号として6を選択することにより適切に特徴付けられ得る。

【0038】

【表1】

表1:コーパスにおける遷移カウント

列内の遷移番号	列のカウント	パーセント
0	3003663	18.44%
1	95418	0.59%
2	7694625	47.24%
3	74196	0.46%
4	4231776	25.98%
5	45013	0.28%
6	1028765	6.32%
<= 6		99.31%
7	7403	0.04%
8	94771	0.57%
9	900	0.01%
10	9543	0.05%
12	1367	0.01%
12を超える遷移		0.01%

【0039】

HMMをベースとしたシステムを構築するときには、このシステムの学習および検査において用いられる特徴ベクトルのタイプが最初に規定される。特徴ベクトルは、継続タイプと分離タイプとに分類されることができる。継続タイプの特徴ベクトルを利用するシステムでは、上記モデルに送られる係数の配列が、またある場合はマトリクスが、利用される。分離タイプの特徴ベクトルが利用されるシステムでは、単一の係数が、上記モデルに送られる。ベクトル量子化手段が、継続タイプのベクトルを分離タイプのベクトルに変換し、これは、HTKに伴うHQuantツールとHCoppyツールとが用いられることによってなされる。HQuantは、後に分離タイプのベクトルを生成するHCoppyとともに用いられる学習用データからコードブックを構築するために用いられる。コードブックの構築は、システムのサイズに応じて当該システムの性能に影響を及ぼし、また、その構築に利用されたデータの量に影響を受ける。HQuantは、コードブックの構築に、線形ベクトル量子化アルゴリ

ズム (Liner Vector Quantization Algorithm) を利用し、これは、計算するのには計算コストが高いアルゴリズムである。本開示では、ユニークベクトル量子化 (Unique Vector Quantization (UVQ)) という名前の新しい方法が導入され、これにより、演算時間が削減され、そして、システムの性能が改善される。この方法は、特徴ベクトルの繰返しを削除することによって、線形ベクトル量子化アルゴリズム (Liner Vector Quantization Algorithm) を利用するコードブックの構築に利用される特徴ベクトルの数を減らすことおよび、各特徴ベクトルのたった一つのコピーを保持するために用いられる特徴ベクトルの数を減らすことに焦点を当てている。表 2 に示されるように、コーパス内の特徴ベクトルの数は、大幅に削減されている。

【0040】

【表 2】

表 2: コーパス内のユニークベクトルカウント

コーパス ラインカウント	ベクトルカウント	ユニークベクトル カウント	縮小率
10,000 ライン	12,285,426	413,410	96.64%
15,000 ライン	16,288,252	591,673	96.37%

【0041】

2000 個の異なるライン画像の特徴ベクトルのすべてを用いてコードブックを構築しようとしたとき、このコードブックについて構築できる最大のサイズが 728 であることを発見した。ユニーク特徴ベクトルのみから 1024 サイズのコードブックを構築するのに 1 時間 30 分を要したのに対し、このコードブックの構築にはおよそ 9 時間を要した。モノラルモデル (mono models) を用いたこれらの実験からの認識速度は、表 3 に示される。ユニークな特徴ベクトルが線形ベクトル量子化アルゴリズムとともに用いられると、コードブックのサイズは増大する。計算速度は 6 倍に上昇し、認識速度は上昇した。

【0042】

【表 3】

表 3: ユニークベクトルカウントの認識率

コードブックタイプ	コードブック サイズ	構築時間	認証率
UVQ 利用無し	728	9 時間	83.59%
UVQ 利用	1024	1 時間 30 分	85.22%

【0043】

上述の方法は、言及された特定の例に限定されるものではないことが、理解されるべきである。設定は、発明の精神から逸脱することなく変更され得る。たとえば、足切遷移番号は、6 以外にも選択され得る。ライン画像の高さおよび幅は、当該ライン画像内のセルのサイズと同様に、上述の例とは異なるものにされ得る。テキスト特徴ベクトルの形態は、抽出方法に応じて変更され得る。たとえば、特徴ベクトルは、2 進数、10 進数、または他の記数法で記述された数値の形態を取り得る。

【0044】

今回開示された実施の形態およびその変形例はすべての点で例示であって制限的なものではないと考えられるべきである。本発明の範囲は上記した説明ではなくて特許請求の範囲によって示され、特許請求の範囲と均等の意味および範囲内でのすべての変更が含まれることが意図される。実施の形態およびその変形例において開示された技術は、可能な限り単独でも組み合わせても実施され得ることが意図される。

【図 1】

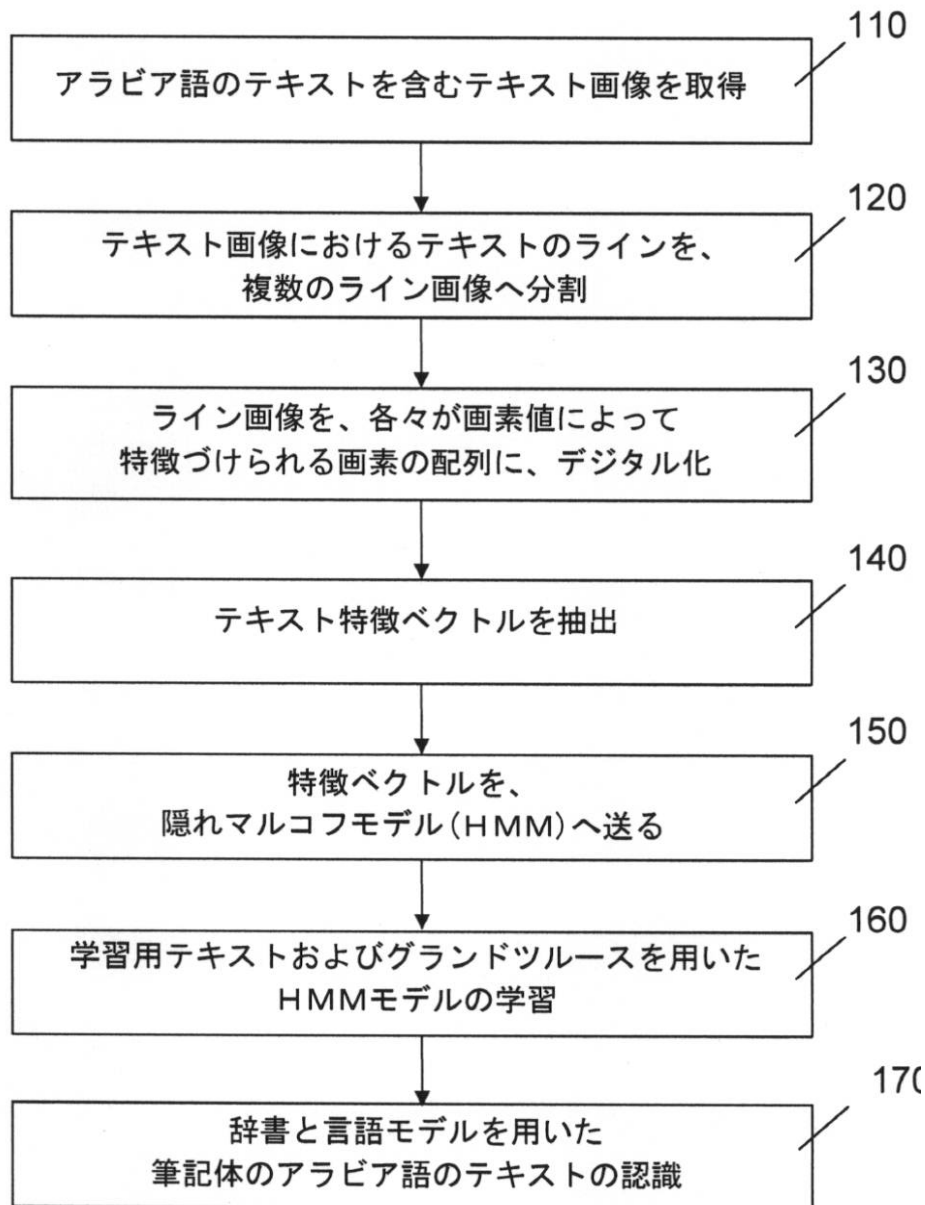


Figure 1

【 図 2 】

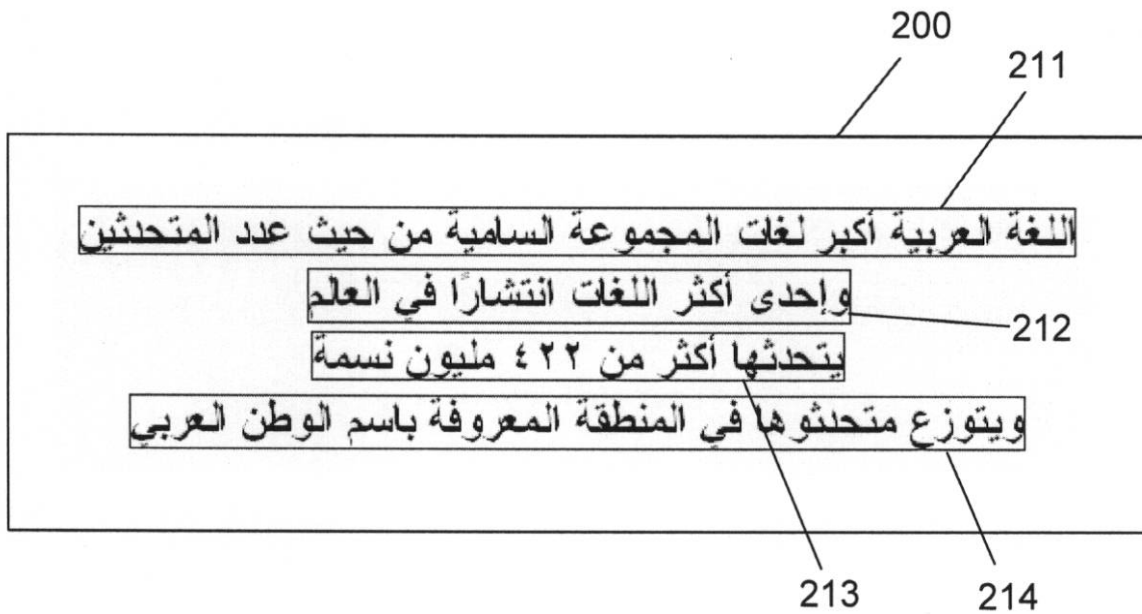


Figure 2

【 図 3 A 】

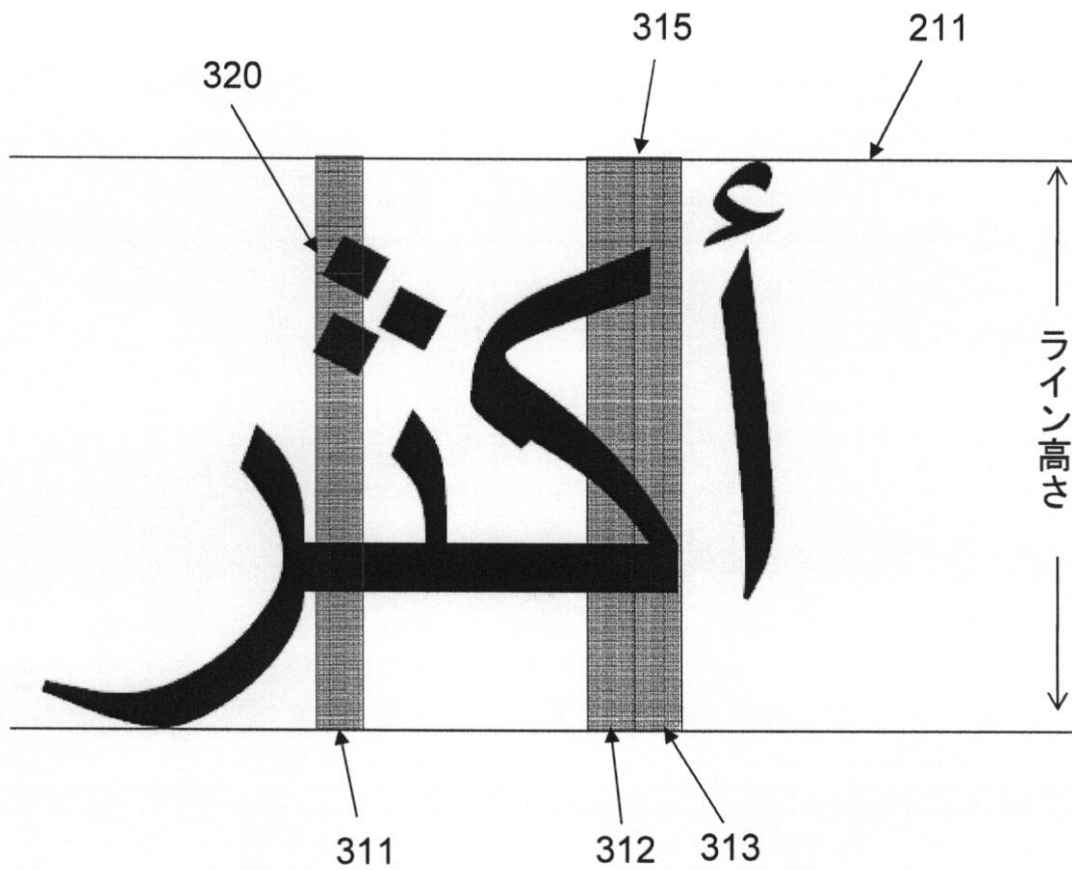
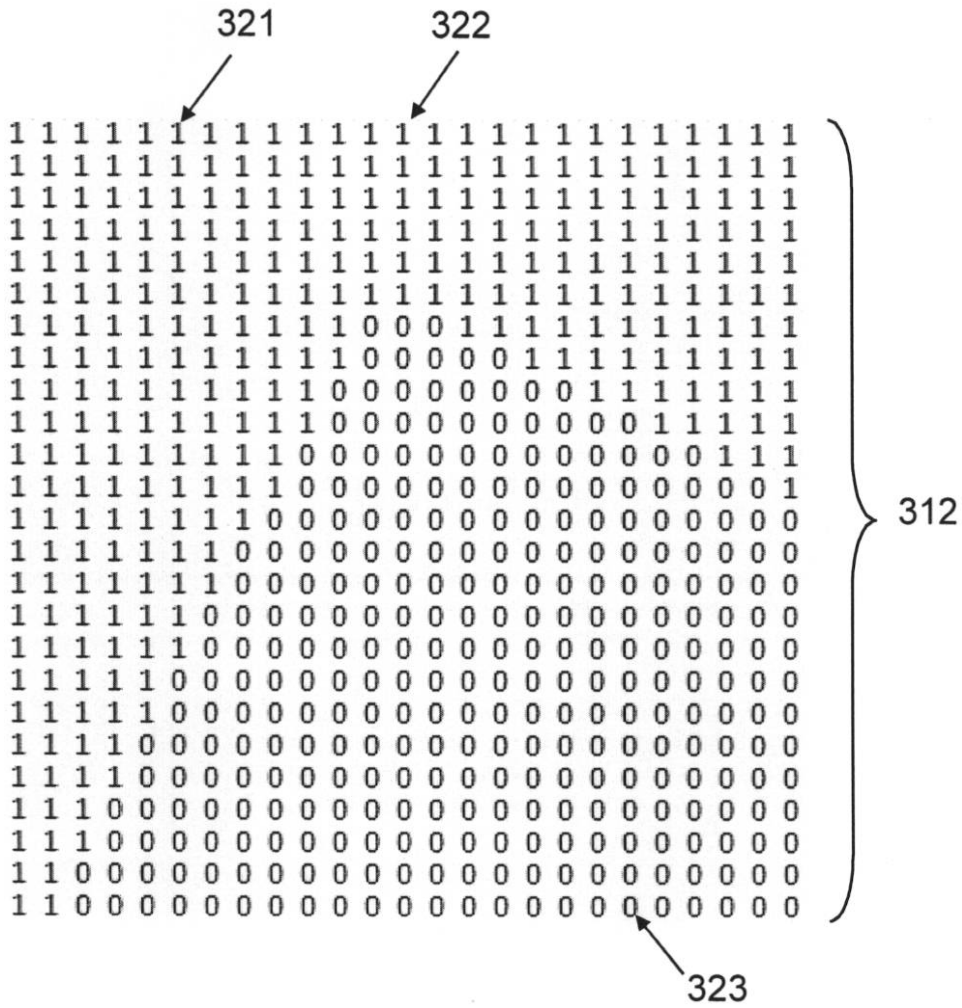


Figure 3A

【図 3 B】

**Figure 3B**

【図 3 C】

**Figure 3C**

【 図 4 】

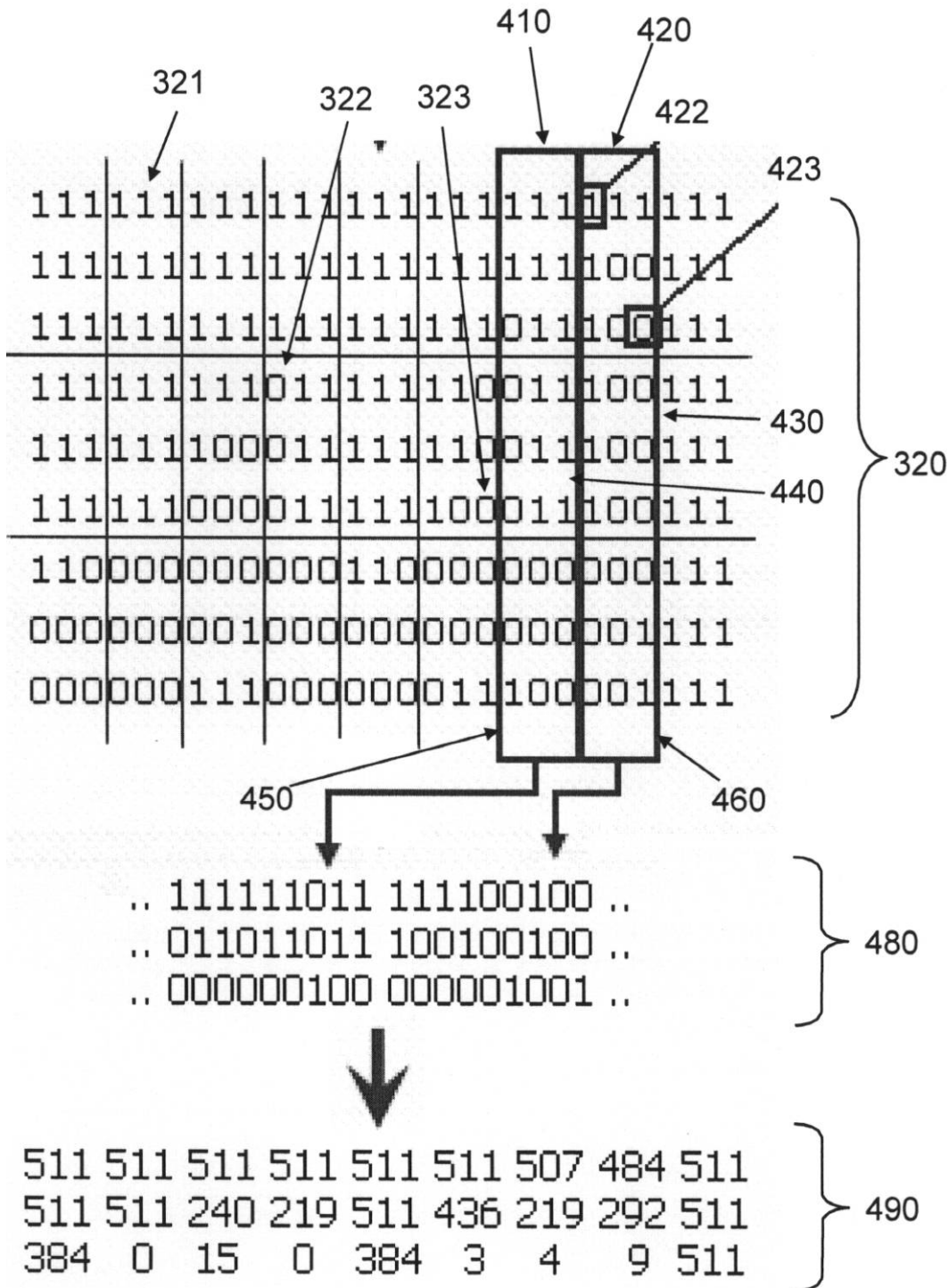


Figure 4

【図 5】

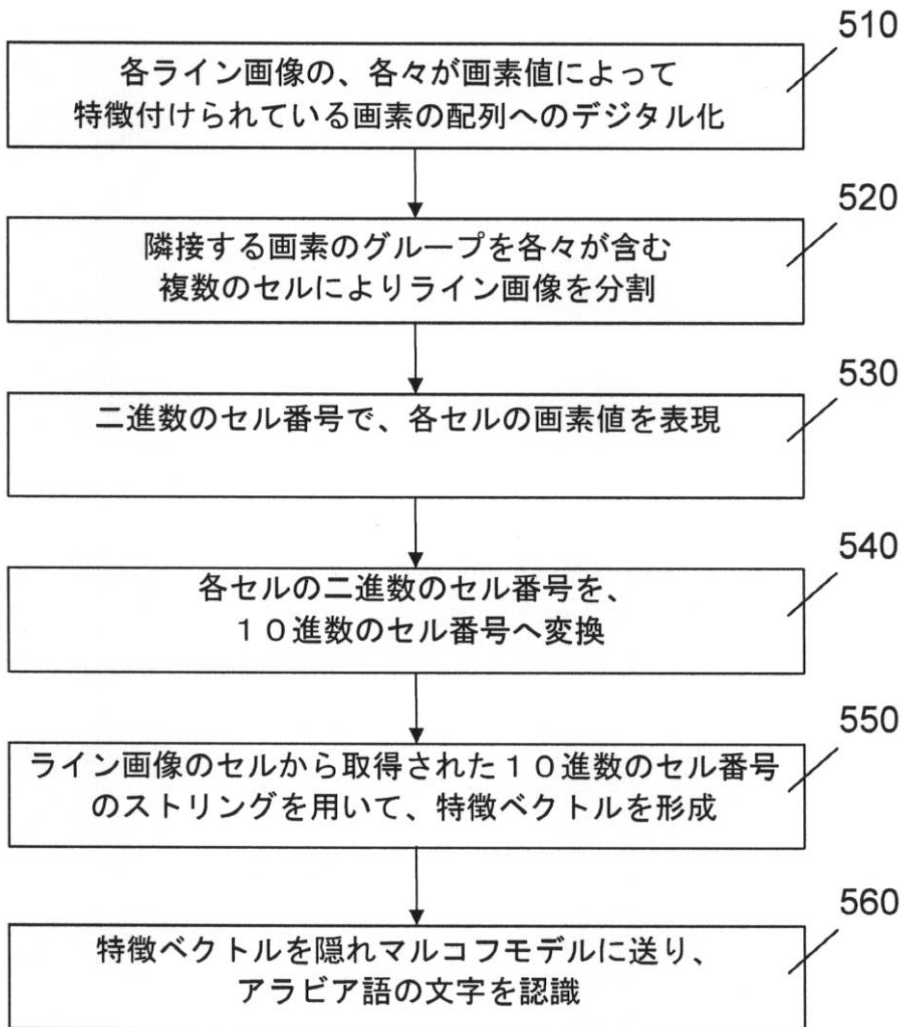


Figure 5

【 図 6 】

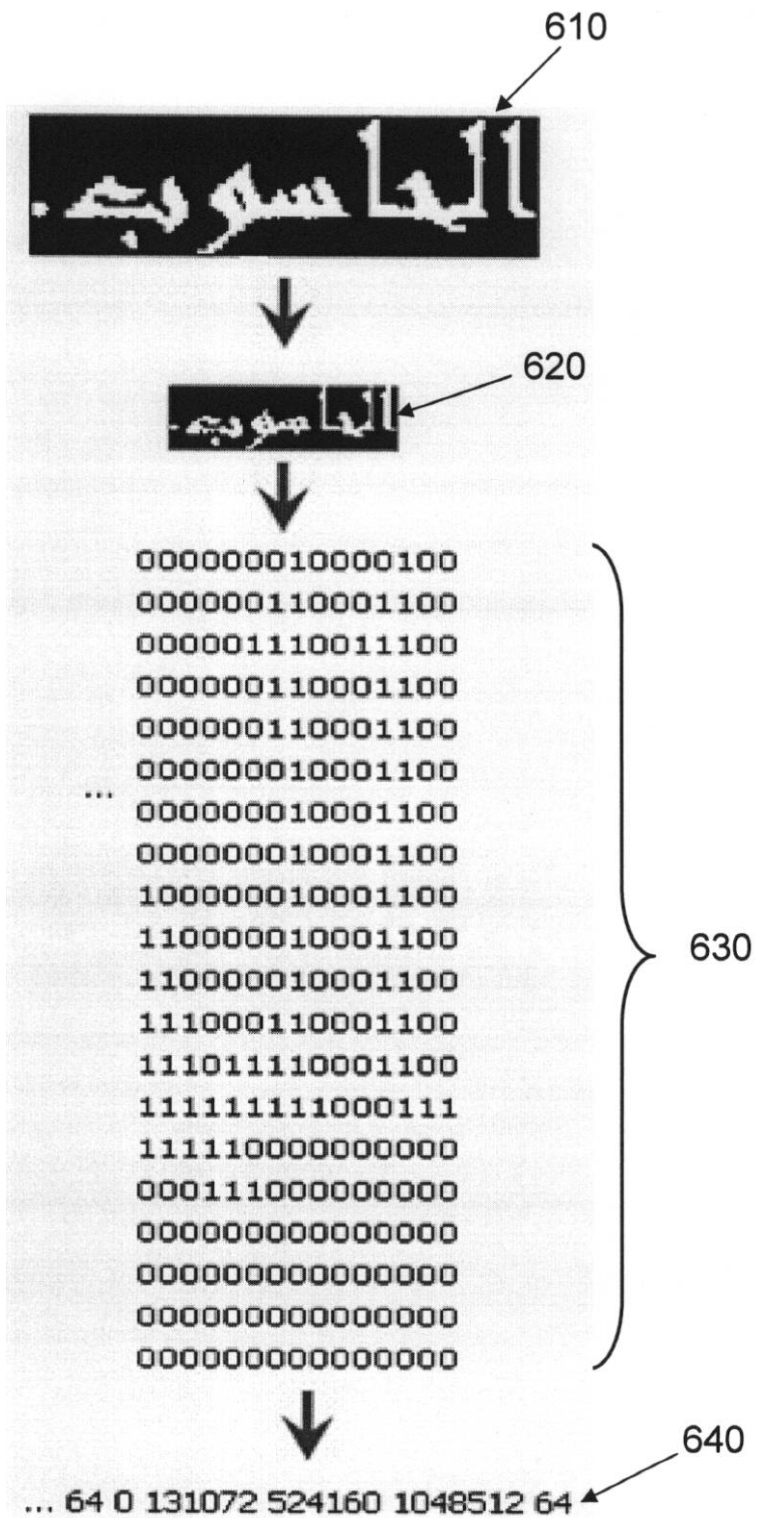


Figure 6

【図 7 A】



Figure 7A

【図 7 B】

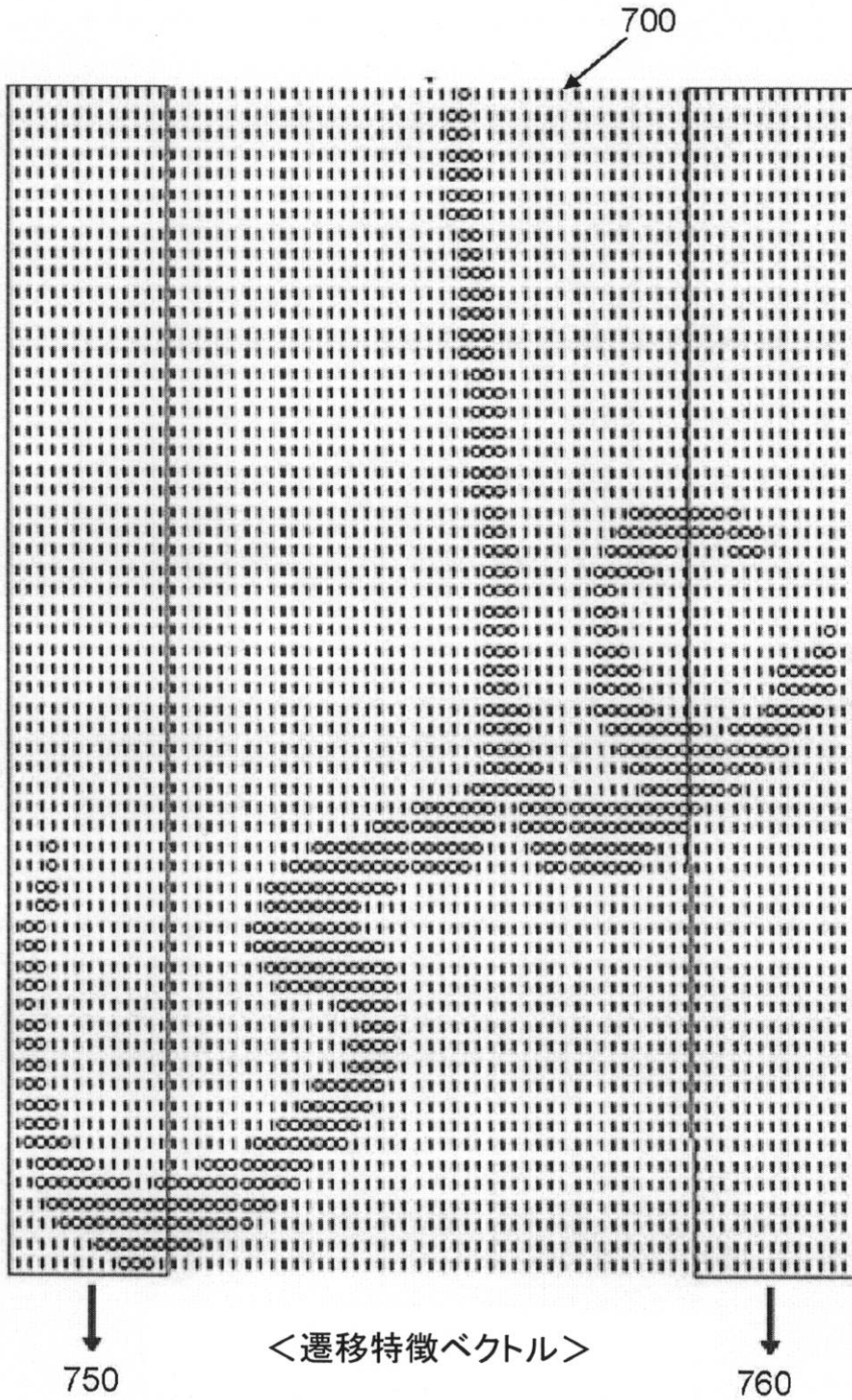


Figure 7B

【図 7 C】

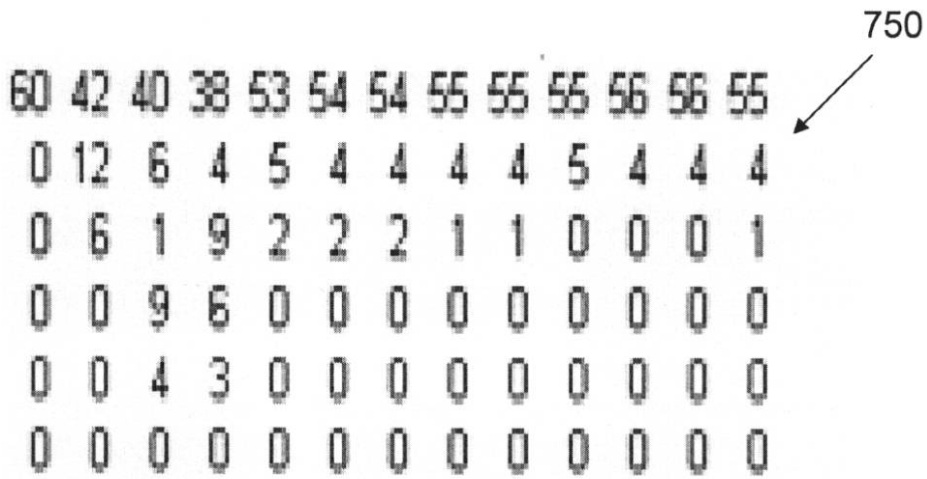


Figure 7C

【図 7 D】

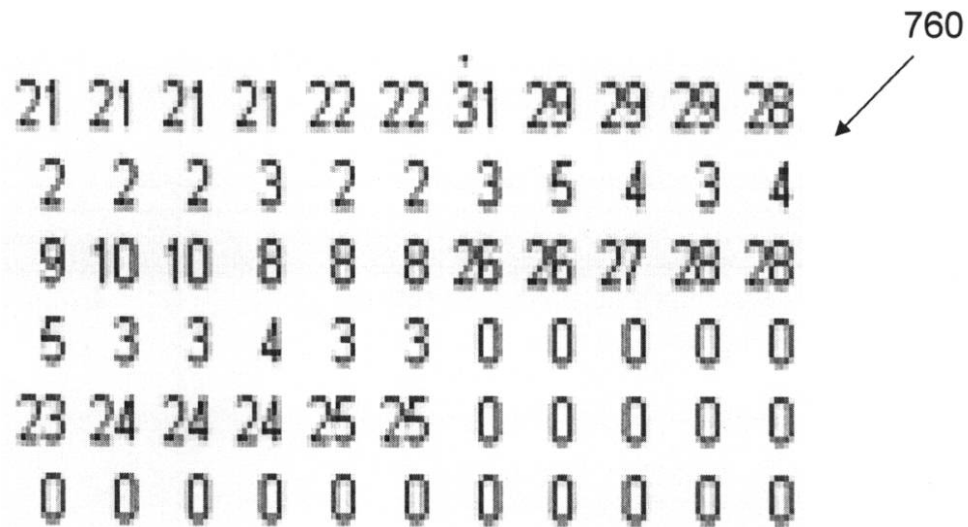


Figure 7D

【図 8】

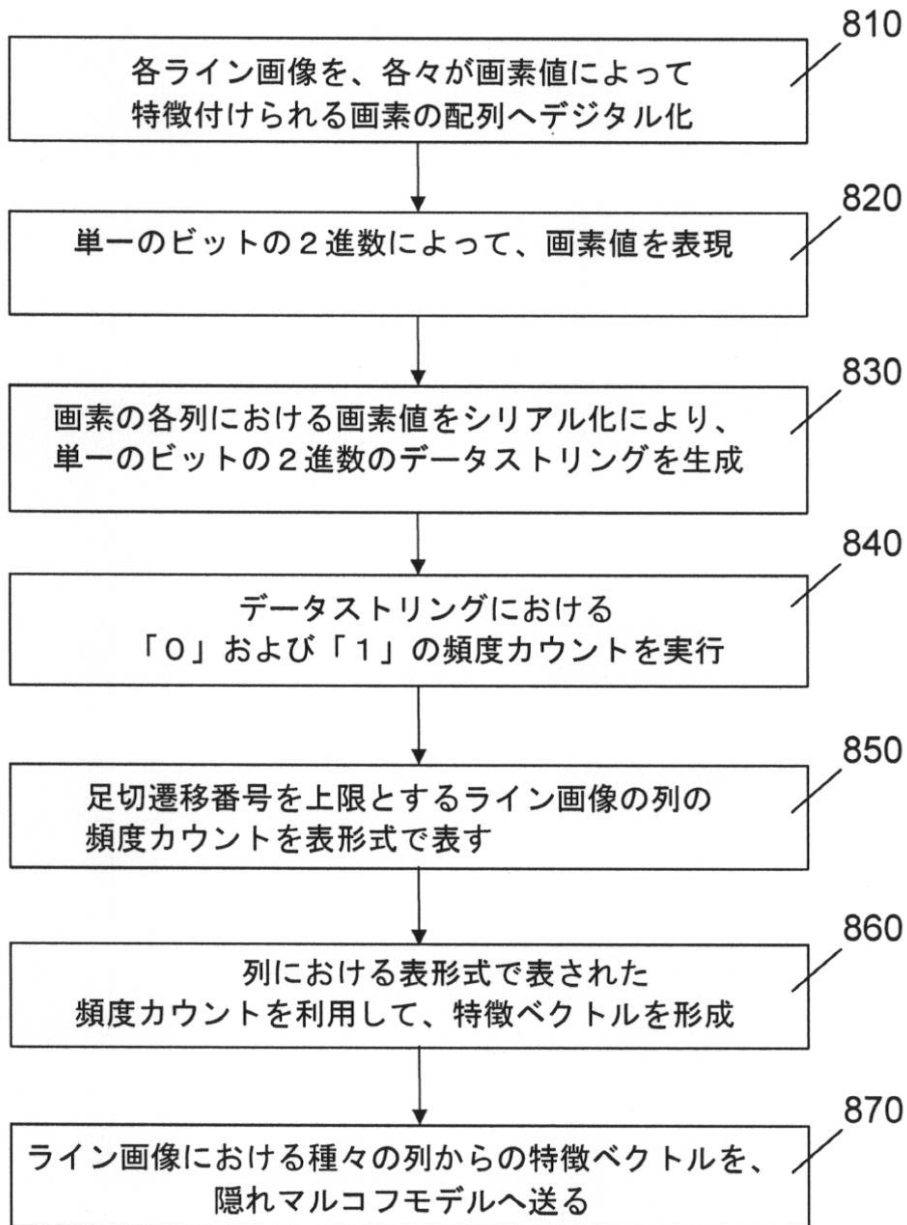


Figure 8

フロントページの続き

(72)発明者 モハメド・エス・ホルシード

サウジアラビア、11442 リヤド、ピィ・オウ・ボックス・6086、キング・アブドゥルア
ジズ・シティ・フォー・サイエンス・アンド・テクノロジー

(72)発明者 フセイン・ケィ・アル・オマリ

サウジアラビア、11442 リヤド、ピィ・オウ・ボックス・6086、キング・アブドゥルア
ジズ・シティ・フォー・サイエンス・アンド・テクノロジー

Fターム(参考) 5B064 AB03 AB19 DA27 DC24 EA07

5L096 BA17 DA02 JA28 KA04