

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2010-509656
(P2010-509656A)

(43) 公表日 平成22年3月25日(2010.3.25)

(51) Int. Cl.	F I	テーマコード (参考)
G06K 9/20 (2006.01)	G06K 9/20 340C	5B029
G06F 17/30 (2006.01)	G06F 17/30 220Z	5B075
G06F 17/21 (2006.01)	G06F 17/21 530A	5B109

審査請求 有 予備審査請求 未請求 (全 39 頁)

(21) 出願番号 特願2009-535346 (P2009-535346)
 (86) (22) 出願日 平成19年11月5日 (2007.11.5)
 (85) 翻訳文提出日 平成21年6月29日 (2009.6.29)
 (86) 国際出願番号 PCT/US2007/023233
 (87) 国際公開番号 W02008/057473
 (87) 国際公開日 平成20年5月15日 (2008.5.15)
 (31) 優先権主張番号 11/592,268
 (32) 優先日 平成18年11月3日 (2006.11.3)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 11/644,009
 (32) 優先日 平成18年12月22日 (2006.12.22)
 (33) 優先権主張国 米国 (US)

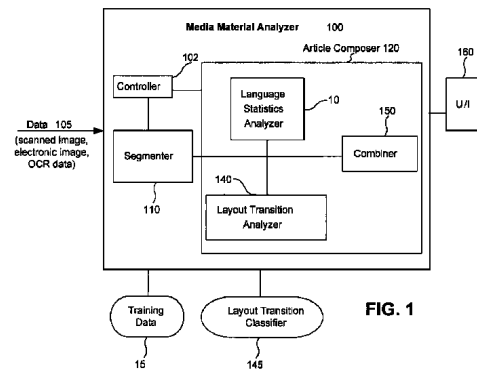
(71) 出願人 502208397
 グーグル インコーポレイテッド
 アメリカ合衆国 カリフォルニア州 94
 043 マウンテン ビュー アンフィシ
 アター パークウェイ 1600
 (74) 代理人 100078282
 弁理士 山本 秀策
 (74) 代理人 100062409
 弁理士 安村 高明
 (74) 代理人 100113413
 弁理士 森下 夏樹
 (72) 発明者 ファーマニーク, ラルフ
 カナダ国 エヌ5エックス-3ティー9
 オンタリオ, ロンドン, ウィンダミア
 ロード 43-683

最終頁に続く

(54) 【発明の名称】 連続する記事部分の媒体資料解析

(57) 【要約】

本発明は、複数のページにわたって連続する記事を有する媒体資料を解析するシステムおよび方法に関する。媒体資料アナライザは、セグメンタと記事コンポーザを含む。セグメンタは、媒体資料内のコラム状の本文テキストに関連するブロックセグメントを識別する。記事コンポーザは、言語統計情報および連続移行情報に基づいて、識別されたブロックセグメントのうちのいずれが、媒体資料内の複数のページにわたって広がる連続する記事に属するかを決定する。



【特許請求の範囲】**【請求項 1】**

レイアウトを有し、かつ複数のページにわたって広がる 1 つ以上の連続する記事を含む媒体資料を表すデータを解析する媒体資料アナライザであって、

(a) 該媒体資料のページ内のコラム状の本文テキストと関連するブロックセグメントを識別するセグメントと、

(b) 言語統計情報および連続移行情報に基づいて、該識別されたブロックセグメントのいずれが該媒体資料内の複数のページにわたって広がる連続する記事に属するかを決定する記事コンポーザと

を備えている、媒体資料アナライザ。

10

【請求項 2】

前記記事コンポーザは、連続レイアウト移行アナライザを含み、該連続レイアウト移行アナライザは、第 1 のページ内の候補となる連続する記事部分に関連する最後のブロックセグメントを識別し、該最後のブロックセグメントの下の 1 つ以上の項目を識別し、該識別された 1 つ以上の項目の少なくとも 1 つの特徴に基づいて、該最後のブロックセグメントを分類し、そして、決定ツリーを適用して、該最後のブロックセグメントが連続する記事内にある確率を示す 1 つ以上の連続移行特徴を選び出す、請求項 1 に記載の媒体資料アナライザ。

【請求項 3】

前記記事コンポーザは、連続言語統計アナライザを含み、該連続言語統計アナライザは、前記連続する記事の前記第 1 のページから連続するページにおける記事部分に対して、言語統計情報を計算し、そして該計算された連続する言語統計情報に基づいて、前記候補となる連続する記事部分内の最後のブロックセグメントが、連続する記事部分を有する確率を決定する、請求項 2 に記載の媒体資料アナライザ。

20

【請求項 4】

前記連続レイアウト移行アナライザは、さらに、連続するページ内の候補となる連続する記事と関連する第 1 のブロックセグメントを識別し、該第 1 のブロックセグメントの上の 1 つ以上の項目を識別し、該 1 つ以上の識別された項目の少なくとも 1 つの特徴に基づいて、該第 1 のブロックセグメントを分類し、そして決定ツリーを適用して、該第 1 のブロックセグメントが連続する記事内にある確率を示す 1 つ以上の連続移行特徴を選び出し、該適用された決定ツリーに基づいて、該第 1 のブロックセグメントが連続する記事である確率を決定する、請求項 3 に記載の媒体資料アナライザ。

30

【請求項 5】

前記連続言語統計アナライザは、さらに、前記第 1 のブロックセグメントを有するページよりも前のページ内の記事部分に対して、言語統計情報を計算し、該計算された、より前のページの言語統計情報に基づいて、前記候補となる連続する記事部分内の該第 1 のブロックセグメントが連続する記事部分を有する確率を決定する、請求項 4 に記載の媒体資料アナライザ。

【請求項 6】

前記連続レイアウト移行アナライザは、さらに、候補となる一対の最後および第 1 のブロックセグメントに対する連続移行特徴を識別し、該一対の最後および第 1 のブロックセグメントに対して一組の連続移行特徴を決定し、そして決定ツリーを適用して、該一組の決定された連続移行特徴に基づいて、該候補となる一対の最後および第 1 のブロックセグメントが、前記媒体資料内の複数のページにわたる同一の連続する記事に属する確率を決定する、請求項 5 に記載の媒体資料アナライザ。

40

【請求項 7】

前記言語統計情報は、単語頻度情報を備え、前記連続言語統計アナライザは、前記最後のブロックセグメント内のテキストおよび連続するページ上の前記記事部分内のテキストにおける単語頻度に基づいて、マッチスコアを計算する、請求項 3 に記載の媒体資料アナライザ。

50

【請求項 8】

前記言語統計情報は、単語頻度情報を備え、前記連続言語統計アナライザは、前記第 1 のブロックセグメント内のテキストおよびより前のページ上の前記記事部分内のテキストにおける単語頻度に基づいて、マッチスコアを計算する、請求項 5 に記載の媒体資料アナライザ。

【請求項 9】

レイアウトを有し、かつ複数のページにわたって広がる 1 つ以上の連続する記事を含む媒体資料を表すデータを解析するコンピュータ実装された方法であって、

(a) 該媒体資料のページ内のコラム状の本文テキストに関連するブロックセグメントを識別することと、

(b) 言語統計情報および連続移行情報に基づいて、該識別されたブロックセグメントのいずれが該媒体資料内の複数のページにわたって広がる連続する記事に属するかを決定することと

を包含する、方法。

10

【請求項 10】

前記記事決定ステップは、

候補となる連続する記事部分内のブロックセグメントに対する連続レイアウト移行情報を解析することと、

該候補となる連続する記事部分内のテキストに対する言語統計を解析することと

を含む、請求項 9 に記載の方法。

20

【請求項 11】

前記連続レイアウト移行情報を解析するステップは、

第 1 のページ内の候補となる連続記事部分に関連する最後のブロックセグメントを識別することと、

該最後のブロックセグメントより下の 1 つ以上の項目を識別することと、

該識別された 1 つ以上の項目の少なくとも 1 つの特徴に基づいて、該最後のブロックセグメントを分類することと、

決定ツリーを適用して、該最後のブロックセグメントが連続する記事内に存在する確率を示す 1 つ以上の連続移行特徴を選ぶことと

を包含する、請求項 10 に記載の方法。

30

【請求項 12】

前記言語統計解析ステップは、

前記連続する記事の第 1 のページから連続するページ上の記事部分に対する、言語統計情報を計算することと、

該計算された連続する言語統計情報に基づいて、前記候補となる連続する記事部分内の前記最後のブロックセグメントが連続する記事部分を有する確率を決定することと

を包含する、請求項 11 に記載の方法。

【請求項 13】

前記連続レイアウト移行情報を解析するステップは、

連続するページ内の候補となる連続する記事に関連する第 1 のブロックセグメントを識別することと、

該第 1 のブロックセグメントより上の 1 つ以上の項目を識別することと、

該 1 つ以上の識別された項目の少なくとも 1 つの特徴に基づいて、該第 1 のブロックセグメントを分類することと、

決定ツリーを適用して、該第 1 のブロックセグメントが連続する記事内に存在する確率を示す 1 つ以上の連続移行特徴を選び出し、そして該適用された決定ツリーに基づいて、該第 1 のブロックセグメントが連続する記事内に存在する確率を決定することと

をさらに包含する、請求項 12 に記載の方法。

40

【請求項 14】

前記言語統計解析ステップは、

50

前記第1のブロックセグメントを有するページよりも前のページ内の記事部分に対して、言語統計情報を計算することと、

該計算されたより前のページの言語統計情報に基づいて、前記候補となる連続する記事部分内の該第1のブロックセグメントが連続する記事部分を有する確率を決定することとをさらに包含する、請求項13に記載の方法。

【請求項15】

前記連続レイアウト移行解析ステップは、

候補となる一对の最後および第1のブロックセグメントを識別することと、

該一对の最後および第1のブロックセグメントに対する一組の連続移行特徴を決定することと、

決定ツリーを適用して、該一組の決定された連続移行特徴に基づいて、該候補となる一对の最後および第1のブロックセグメントが、前記媒体資料内の複数のページにわたる同一の連続する記事に属する確率を決定することと

をさらに含む、請求項12に記載の方法。

【請求項16】

前記言語統計情報は、単語頻度情報を備え、前記連続言語統計解析ステップは、前記最後のブロックセグメント内のテキストおよび連続するページの前記記事部分内のテキストにおける単語頻度に基づいて、マッチスコアを計算することを含む、請求項12に記載の方法。

【請求項17】

前記言語統計情報は、単語頻度情報を備え、前記連続言語統計解析ステップは、前記第1のブロックセグメント内のテキストおよび以前のページの前記記事部分内のテキストにおける単語頻度に基づいて、マッチスコアを計算することを含む、請求項14に記載の方法。

【請求項18】

レイアウトを有する媒体資料内の複数のページにわたって広がる連続する記事を構成する記事コンポーザであって、

連続レイアウト移行アナライザと、

連続言語統計アナライザと

を備え、該連続レイアウト移行アナライザは、異なるページ上の候補となる記事の最後のブロックセグメントおよび第1のブロックセグメントが同じ連続する記事内に存在する確率を示す1つ以上の連続移行特徴を選び出すために、決定ツリーを適用し、

該連続言語統計アナライザは、異なるページ上の異なる記事部分に対する言語統計情報を計算し、該計算された言語統計情報に基づいて、候補となる記事部分の第1および最後のブロックセグメントが連続する記事部分を有する確率を決定し、それにより、該記事コンポーザは、解析された連続レイアウト移行特徴および該計算された言語統計に従って、該第1および最後のブロックセグメントが同じ連続する記事に属する確率に基づいて、複数のページにわたる連続する記事を構成することが可能である、記事コンポーザ。

【発明の詳細な説明】

【技術分野】

【0001】

本願は、2006年11月3日出願の米国出願第11/592,268号(代理人整理番号第2525.0010000)の一部継続出願であり、該米国出願の全内容は本明細書において参照により援用される。

【0002】

(発明の分野)

本発明は、媒体資料のコンピュータ補助による解析に関する。

【背景技術】

【0003】

(発明の背景)

10

20

30

40

50

(関連技術)

文書および印刷された資料の解析を行うか、または補助するために、コンピュータがますます使用されている。レイアウト解析技術およびシステムは、文書中のテキストおよび画像の位置および相対的な配列を解析するために使用されてきた。このような文書レイアウト解析は、多くの文書画像化用途において重要であり得る。例えば、文書レイアウト解析は、レイアウトベースの文書検索、光学文字認識を用いたテキスト抽出、および文書画像の電子形式への変換の一部として、使用され得る。文書レイアウト解析は、概して、単純な文書(例えば、ビジネスレターまたは一列の報告書)において最良に機能し、かつレイアウトが複雑であるか、または可変であるときには、困難であり得るかまたは機能不可能でさえあり得る。例えば、自動の文書レイアウト解析または半自動の文書レイアウト解析は、しばしば、複雑なレイアウトに分類され、かつ再ソートがレイアウトの手動解析に対してなされなければならない。

10

20

30

40

50

【0004】

レイアウトに配列された本文のテキストのコラムを有する媒体資料は、文書レイアウト解析に対する特別な挑戦を生み出す。例えば、新聞のレイアウトは、概して、非常に複雑であり、多くの記事と論理的な要素とがページ上で一緒に接近して組み合わせられる。新聞の構造を理解することは、記事の文脈、パターンマッチングおよび可能性としては新聞のスタイル、すなわちコンピュータよりも人間に対して自然な要素によって、人間によって自然に行われる。自動的な方法は、概して、図形的な特徴または幾何学的な特徴のみにほとんど依存しており、その結果、全ての新聞にわたって機能する一貫した一組の単純な罫線がないので、多くの間違いを生じる。このような限定された自動的な方法は、媒体資料の2つ以上のページにわたって連続する記事を解析するさらなる困難性を有する。

【発明の概要】

【発明が解決しようとする課題】

【0005】

レイアウトを有する媒体資料を解析する向上したシステムおよび方法が必要である。

【課題を解決するための手段】

【0006】

(発明の概要)

本発明は、レイアウトを有する資料媒体を解析するシステムおよび方法に関する。

【0007】

一実施形態において、媒体資料アナライザは、セグメントと記事コンポーザとを含み得る。セグメントは、媒体資料内のコラム状の本文のテキストと関連するブロックセグメントを識別する。一例において、セグメントは、画像データ内の画素データを解析して、類似の画素値変化の複雑性(pixel value change complexity)を有する領域を識別する。画素値変化は、画素から水平方向および垂直方向に沿って識別される。記事コンポーザは、識別されたブロックセグメントのいずれが媒体資料中の1つ以上の記事に属するかを決定する。記事コンポーザは、言語統計情報、レイアウト移行情報、または言語統計情報およびレイアウト移行情報の両方に基づいて、候補となるブロックセグメントが同一の記事に属するかどうかを決定し得る。

【0008】

別の実施形態において、記事コンポーザは、言語統計アナライザを含み得る。言語統計アナライザは、言語統計情報に基づいて、セグメントから出力されたブロックセグメントのいずれが、媒体資料中の1つ以上の記事に属するかを決定する。特に、言語統計アナライザは、セグメントによって出力された候補となるブロックセグメントに対して、言語統計を計算し、そして、言語統計情報におけるオーバーラップに基づいて、候補となるブロックセグメントが同一の記事に属する確率を決定する。

【0009】

さらなる実施形態において、記事コンポーザは、レイアウト移行アナライザを含み得る。レイアウト移行アナライザは、セグメントによって出力された候補となるブロックセグ

メントにおけるレイアウト移行特徴を解析し、そしてレイアウト移行解析に基づいて、候補となるブロックセグメントが媒体資料内の同一の記事に属するかどうかを決定する。一例において、レイアウト移行特徴は、垂直方向および水平方向の移行特徴を含む。

【0010】

本発明の一面に従って、コンピュータ実装された方法は、レイアウトを有する媒体資料を表すデータを解析する。この方法は、媒体資料内のコラム状の本体テキストと関連するブロックセグメントを識別することと、言語統計情報およびレイアウト情報に基づいて、識別されたブロックセグメントのいずれが、媒体資料内の1つ以上の記事に属するかを決定することとを含み得る。

【0011】

さらなる実施形態において、ネットワークを介して（例えば、ウェブを介して）、ブラウザを通して、レイアウトを有する媒体資料を探索するシステムが提供される。ブラウザは、サーチ要求を満たすことにおいて識別された同一の記事内の1つ以上のブロックセグメントからテキストを受信し得る。

【0012】

さらなる実施形態において、レイアウトを有する媒体資料を表し、かつ複数のページにわたって広がる1つ以上の連続する記事を含むデータを解析する媒体資料アナライザが提供される。媒体資料アナライザは、媒体資料ページ内のコラム状本体テキストに関連するブロックセグメントを識別するセグメントと、言語統計情報および連続移行情報に基づいて、識別されたブロックセグメントのいずれが、媒体資料内の複数のページにわたって広がる連続する記事に属するかを決定する記事コンポーザとを含む。

【0013】

なおさらなる実施形態において、レイアウトを有する媒体資料を表し、かつ複数のページにわたって広がる1つ以上の連続する記事を含むデータを解析するコンピュータ実装された方法が提供される。この方法は、媒体資料ページ内のコラム状本体テキストに関連するブロックセグメントを識別することと、言語統計情報および連続移行情報に基づいて、識別されたブロックセグメントのいずれが、媒体資料内の複数のページにわたって広がる連続する記事に属するかを決定することとを含む。

【0014】

さらに、一実施形態において、レイアウトを有する媒体資料内の複数のページにわたって広がる連続する記事を構成する記事コンポーザは、連続レイアウト移行アナライザと、連続統計アナライザとを含む。連続レイアウト移行アナライザは、1つ以上の連続移行特徴を選び出すために決定ツリーを適用し、この1つ以上の連続移行特徴は、異なるページ上の候補となる記事部分の最後および第1のブロックセグメントが、同一の連続記事に存在する確率を示す。連続言語統計アナライザは、異なるページ上の異なる記事部分に対する言語統計情報を計算し、そして、該計算された統計情報に基づいて、候補となる記事部分の第1および最後のブロックセグメントが連続する記事部分を有する確率を決定する。このようにして、記事コンポーザは、解析された連続レイアウト移行特徴および計算された言語統計に従って、第1および最後のブロックセグメントが同一の連続する記事に属する確率に基づき、複数のページにわたる連続する記事を構成し得る。

【0015】

本発明のさらなる実施形態、特徴および利点と、本発明の様々な実施形態の構造および動作とが、添付の図面を参照して以下に詳細に記載される。

【0016】

特許または出願のファイルは、カラーで制作された少なくとも1つの図面を含む。カラーの図面を有する、この特許または特許出願刊行物のコピーは、要求し、必要な料金を支払うと、官庁より提供される。

【0017】

本発明の実施形態は、添付の図面を参照して記載される。図面において、同様の参照番号は同一の要素または機能的に類似の要素を示し得る。ある要素が第1に現れる図面は、

10

20

30

40

50

対応する参照番号のもっとも左の桁によって概して示される。

【図面の簡単な説明】

【0018】

【図1】図1は、本発明の実施形態に従う、媒体資料アナライザの図である。

【図2】図2は、本発明の実施形態に従う、媒体資料を解析する方法の図である。

【図3】図3は、図2の方法における、ブロックセグメント識別ステップを実行する例示的なルーチンを示す図である。

【図4】図4は、図3のルーチンに従って識別されたブロックセグメントを有する媒体資料の画像を示し、カラーを含む、図である。

【図5】図5は、本発明のさらなる実施形態に従う、OCRデータからテキストを抽出し、そしてブロックセグメント領域を調整するステップを示すフローチャート図である。

【図6】図6は、本発明の一実施形態に従う、言語統計を解析して、記事内のブロックセグメントを識別する方法を示す図である。

【図7A】図7Aおよび図7Bは、本発明の一実施形態に従う、訓練モードにおいてレイアウト移行アナライザの動作を示すフローチャート図である。図7Aは、垂直方向の移行特徴を決定するために訓練モードにおいて動作する方法を示す。図7Bは、水平方向の移行特徴に基づいて、訓練モードにおいて動作する方法を示す。

【図7B】図7Aおよび図7Bは、本発明の一実施形態に従う、訓練モードにおけるレイアウト移行アナライザの動作を示すフローチャート図である。図7Aは、垂直方向の移行特徴を決定するために訓練モードにおいて動作する方法を示す。図7Bは、水平方向の移行特徴に基づいて、訓練モードにおいて動作する方法を示す。

【図8】図8は、本発明の一実施形態に従う、レイアウト移行分類子を訓練し、かつ構築するために用いられ得る、候補となるブロックセグメントを示す例示的な媒体資料である。

【図9A】図9Aおよび図9Bは、本発明の一実施形態に従う、実行モードにおけるレイアウト移行アナライザの動作を示すフローチャート図である。図9Aは、媒体資料レイアウト内の垂直方向移行特徴に基づいた実行モード動作を示す。図9Bは、レイアウト内の水平方向移行特徴に基づいた実行モード動作を示す。

【図9B】図9Aおよび図9Bは、本発明の一実施形態に従う、実行モードにおけるレイアウト移行アナライザの動作を示すフローチャート図である。図9Aは、媒体資料レイアウト内の垂直方向移行特徴に基づいた実行モード動作を示す。図9Bは、レイアウト内の水平方向移行特徴に基づいた実行モード動作を示す。

【図10A】図10A～図10Dは、カラーを含み、本発明の例示的な実施形態に従って解析された新聞のページを含む例示的な媒体資料を示す。

【図10B】図10A～図10Dは、カラーを含み、本発明の例示的な実施形態に従って解析された新聞のページを含む例示的な媒体資料を示す。

【図10C】図10A～図10Dは、カラーを含み、本発明の例示的な実施形態に従って解析された新聞のページを含む例示的な媒体資料を示す。

【図10D】図10A～図10Dは、カラーを含み、本発明の例示的な実施形態に従って解析された新聞のページを含む例示的な媒体資料を示す。

【図11】図11は、本発明のさらなる実施形態に従う、World Wide Webを介して、レイアウトを有する媒体資料をサーチするシステムを示す図である。

【図12】図12は、本発明の一実施形態に従う、媒体資料アナライザによって解析されたデータのサーチにおけるサーチ結果の例示的な表示を示す図である。

【図13】図13は、本発明のさらなる実施形態に従う、連続する記事部分を解析し得る媒体資料アナライザの図である。

【図14A】図14A～図14Eは、本発明の一実施形態に従う、図13の媒体資料アナライザにおける記事コンポーザの動作を示すフローチャート図である。

【図14B】図14A～図14Eは、本発明の一実施形態に従う、図13の媒体資料アナライザにおける記事コンポーザの動作を示すフローチャート図である。

10

20

30

40

50

【図 1 4 C】図 1 4 A ~ 図 1 4 E は、本発明の一実施形態に従う、図 1 3 の媒体資料アナライザにおける記事コンポーザの動作を示すフローチャート図である。

【図 1 4 D】図 1 4 A ~ 図 1 4 E は、本発明の一実施形態に従う、図 1 3 の媒体資料アナライザにおける記事コンポーザの動作を示すフローチャート図である。

【図 1 4 E】図 1 4 A ~ 図 1 4 E は、本発明の一実施形態に従う、図 1 3 の媒体資料アナライザにおける記事コンポーザの動作を示すフローチャート図である。

【図 1 5】図 1 5 は、本発明の実施形態を実装するために使用され得る例示的なコンピュータシステムの図である。

【発明を実施するための形態】

【0019】

10

【表 1】

目次

概観

媒体資料アナライザ

媒体資料アナライザの動作

ブロックセグメント化

記事構成

言語統計

レイアウト移行

訓練モード

垂直方向の移行

水平方向の移行

実行モード

表示例

ワールドワイドウェブへの応用

さらなる特徴および利点

さらなる応用一連続する記事

例示的なコンピュータシステム実装

結論

20

30

(実施形態の詳細な説明)

本発明は、特定の用途に対する例示的な実施形態を参照して本明細書に記載されるが、本発明がそれらの実施形態に限定されないことが理解されるべきである。本明細書において提供される教示を利用する機会を有する当業者は、本発明の範囲内のさらなる修正、用途および実施形態、ならびに本発明が大いに有効であるさらなる分野を認識するだろう。

40

【0020】

(概観)

本発明は、レイアウトを有する媒体資料を解析するシステムおよび方法に関する。例として、本文テキストのコラムを伴うレイアウトを有する媒体資料を含むが、それに限定はされない。このような例は、新聞、雑誌、カタログ、小冊子、パンフレットおよび他のタイプの印刷資料を含むがこれらに限定はされない。

【0021】

(媒体資料アナライザ)

50

図1は、本発明の実施形態に従う媒体資料アナライザ100を示す。媒体資料アナライザ100は、コントローラ102と、セグメンタ110と、記事コンポーザ120とを含む。記事コンポーザ120は、純粋な言語統計モード、純粋なレイアウト移行モードまたは二つの組み合わせにおいて動作し得る。

【0022】

図1に示される実施形態において、記事コンポーザ120は、言語統計アナライザ130と、レイアウト移行アナライザ140と、コンバイナ150とを含む。媒体資料アナライザ100は、データ105、訓練データ135およびレイアウト移行分類子(classifier)145を受信し得るか、またはそれらにアクセスし得る。媒体資料アナライザ100はまた、ユーザインターフェース160に連結され得る。

10

【0023】

データ105は、媒体資料の画像データを含み得る。このような画像データは、電子的またはスキャンされた画像データと、画像データから抽出された光学文字認識(OCR)データとを含み得る。データ105は、任意のタイプのファイルフォーマットで提供され得る。

【0024】

訓練データ135は、媒体資料内の記事に属するブロックセグメントのポジティブな例およびネガティブな例を含み得る。レイアウト移行分類子145は、候補となるブロックセグメントが媒体資料内の記事に属するように分類されることを可能にする移行特徴情報を含むデータ構造を含むがこれらに限定されない。このようなデータ構造は、決定ツリー

20

【0025】

セグメンタ110は、データ105内の媒体資料内のコラム状の本文テキストに関連するブロックセグメントを識別する。記事コンポーザ120は、言語統計情報および/またはレイアウト移行情報に基づいて、識別されたブロックセグメントのうちのいずれが媒体資料のうちの1つ以上の記事に属することを決定する。

【0026】

一実施形態において、言語統計アナライザ130は、セグメンタ110によって出力された候補となるブロックセグメントに対する言語統計を計算する。次いで、言語統計アナライザ130は、言語統計情報のオーバーラップに基づいて、候補となるブロックセグメントが同一の記事に属する確率を決定する。

30

【0027】

レイアウト移行アナライザ140は、レイアウト移行特徴およびセグメンタ110によって出力された候補となるブロックセグメントをさらに解析する。次いで、レイアウト移行アナライザ140は、レイアウト移行特徴に基づいて、候補となるブロックセグメントが媒体内の同一の記事に属するかどうかを決定する。

【0028】

コンバイナ150は、言語統計アナライザ130およびレイアウト移行アナライザ140によって解析された候補となるブロックセグメントが同一の記事に属するかどうかを識別する。一例において、コンバイナ150は、言語統計アナライザ130によって決定された確率と、レイアウト移行アナライザ140から出力されたレイアウト移行特徴に基づく、ブロックが同一の記事に属するかどうかという決定との両方の出力に基づいて、候補となるブロックセグメントが同一の記事に属するかどうかを識別する。

40

【0029】

あるいは、コンバイナ150は、言語統計アナライザ130のみによって決定された確率に基づいて、同一の記事に属する候補となるブロックセグメントを識別し得る。コンバイナ150は、また、レイアウト移行アナライザ140のみによって解析されたレイアウト移行特徴に基づいて、候補となるブロックセグメントが同一の記事に属するかどうかを決定し得る。

【0030】

50

コントローラ102は、セグメンタ110と記事コンポーザ120とを制御および管理する。ユーザからのさらなる制御は、ユーザインターフェース160を介して提供され得る。例えば、ユーザは、動作を開始し得るか、またはデータ105、訓練データ135もしくはレイアウト移行分類子145の入力を開始し得る。ユーザは、媒体資料アナライザ100と相互作用して、訓練データ135の生成またはレビューを助け得る。例えば、ユーザは、訓練データ135の質を向上させるために、所与の媒体資料内の記事に属するブロックセグメントのポジティブな例とネガティブな例とを選択し得る。ユーザはまた、レイアウト移行分類子145を構築または修正するために、媒体資料アナライザ100と相互作用し得る。

【0031】

媒体資料アナライザ100はまた、スキャンされたデータ105または媒体資料アナライザ100から出力されたデータの画像を表示し得る。表示のための出力データは、媒体資料アナライザ100の解析に従って構成されたハイライトされたブロックセグメントを示すために解析された媒体資料の表示を含み得る。特定のレイアウトに対して、ユーザはフィードバックを提供し得るか、またはハイライトされたブロックセグメントを選択し得る。他のタイプの情報は、この記載を提供された当業者に対して明らかであるように、表示され得る。

【0032】

媒体資料アナライザ100（その構成要素モジュールを含む）は、ソフトウェア、ファームウェア、ハードウェアまたはそれらの任意の組み合わせにおいて実装され得る。媒体資料アナライザ100は、コンピュータ、ワークステーション、分散コンピューティングシステム、埋め込みシステム、スタンドアロン電子デバイス、ネットワーク化されたデバイス、モバイルデバイス、セットトップボックス、テレビ、あるいは他のタイプのプロセッサまたはコンピュータシステムを含むが、これらに限定されない任意のタイプの処理デバイス上で実行するように実装され得る。

【0033】

媒体資料アナライザ100は、また、種々の用途において使用され得る。データ105において自動的に、または半自動的に動作させることによって、媒体資料アナライザ100は、格納された画像データ（例えば、アーカイブされた媒体資料）を解析し得る。マイクロフィッシュ、フィルムおよび他のストレージ媒体は、入力用に画像データを取得するためにスキャンされ得る。任意のファイルフォーマットの電子ファイルがまた入力され得る。解析は、ユーザからの最小の入力で、またはユーザからの入力なしで、自動的に、または半自動的に実行され得る。このようにして、媒体資料アナライザ100は、種々の媒体に対するブロックセグメントから成り立つ記事を構成するために使用され得る。次いで、媒体資料アナライザ100は、レイアウトを有する媒体資料内の記事から成り立つテキストデータのブロックセグメントを出力し得る。このような出力は、媒体資料のコンテンツをレビューまたはサーチすることを望むローカルユーザおよびリモートユーザに、配信され得るか、または格納され得る。

【0034】

自動的に、または半自動的に動作することによって、複数の媒体資料アナライザ100は、大量の媒体資料を解析するために、使用され得、かつ縮尺が合わせられ得る。このようにして、媒体資料レイアウト内のコンテンツは、広範囲のユーザに対して、ローカルに、そしてネットワークを介してリモートに利用可能にし得る。媒体資料アナライザ100は、ユーザが、図書館、大学、政府機関、会社および他の場所において、ローカルにまたはリモートにアクセスされる媒体資料内のテキストデータをレビューすることを可能にし得る。媒体資料アナライザ100は、サーチエンジン、ウェブポータルまたは他のウェブサイトと共に使用され得、リモートユーザが、レイアウトを有する媒体資料をレビューおよびサーチすることを可能にする。

【0035】

媒体資料アナライザ100、ならびにセグメンタ110および記事コンポーザ120を

10

20

30

40

50

含む構成要素の動作は、図 2 に示される媒体資料データを解析する方法に関して、以下にさらに詳細に記載される。

【 0 0 3 6 】

(媒体資料アナライザの動作)

さらなる実施形態に従って、媒体資料からのデータを解析する方法 2 0 0 が提供される (図 2)。簡潔さのために、方法 2 0 0 は、媒体資料アナライザ 1 0 0 を参照して記載されるが、必ずしも媒体資料アナライザ 1 0 0 の構造に限定されることを意図されない。

【 0 0 3 7 】

(ブロックセグメント化)

ステップ 2 1 0 において、特徴に従って、媒体資料内のコラム状の本文テキストと関連するブロックセグメントが識別される。図 3 は、ステップ 2 1 0 のブロックセグメント化を実行する例示的なルーチンをさらに詳細に示す (ステップ 3 0 5 ~ ステップ 3 2 0)。明確さのために、ブロックセグメント化ルーチンは、図 4 に示されるように、新聞のページの例示的な画像 4 1 0 に関して記載される。この例において、データ 1 0 5 は、画像 4 1 0 を表す画素データを含む。画素データは、特定のピクチャ要素 (画素) 位置における画像の強度を表す画素値からなる。画素値は、グレイスケール、カラー、バイナリまたは他のタイプの画素データを含むが、これらに限定されない任意のタイプの画素値であり得る。

10

【 0 0 3 8 】

ループ 3 0 5 において、セグメンタ 1 1 0 は、類似の画素値変化複雑性を有する領域を識別するために、画像データ内の画素を解析する。画素の全てまたは画素のサンプルが解析され得る。セグメンタ 1 1 0 は、解析されている各画素から、水平方向および垂直方向に沿った画素値変化を解析する。類似の画素値変化複雑性のこれらの領域は、ブロックセグメントを含み得る。特に、のど空き (gutter) および他の境界を有するレイアウトにおいて配列されたコラム状の本文テキストを覆う領域に対して、本文テキストのブロックセグメントが得られる。

20

【 0 0 3 9 】

一例において、セグメンタ 1 1 0 は、バイナリ (純粋な黒および白の) 画像における一致する複雑性の領域を見出すテクスチャ方法を実行する。画像内の各画素位置に対して、セグメンタ 1 1 0 は、カラーが両方の側で n 回変化するまで、水平方向 (左および右) を見る必要がある距離のログを計算する。テキストの領域は、相対的に一致する小さい値を有するが、のど空きおよび他の単純な領域はより大きい値を有する。 $n = 1$ に設定することは、例示的な実行長のアプローチを与える。図 4 の例において、 $n = 2$ が用いられ、ページにわたって水平方向に広がるブロックセグメントという結果をもたらす (画像 4 2 0 を参照)。 n に対するより大きな値は、よりむらのない領域をもたらすが、同様にカットオフされる境界をもたらす。この計算は、垂直方向において繰り返され、 $n = 2$ のとき、ページにわたって垂直方向に広がるブロックセグメントをもたらす (画像 4 3 0 を参照)。記事内の支配的な本文テキスト (例えば、ページの見出しではなく記事内の本文テキスト) を探すとき、水平方向および垂直方向の計算は、最終的なセグメント化された画像を得るために、共に加算され得、次いで、閾値化され得る (カラー画像 4 4 0 を参照)。さらなる例において、異なるサイズのテキストの領域を区別するために見る場合には、セグメンタ 1 1 0 は、最小値から開始し得、許容値で `floor-fill` を行い得る。

30

40

【 0 0 4 0 】

データ 1 0 5 は、また、媒体資料レイアウトの画像データに関連するテキストデータを含み得る。例えば、媒体資料内のテキストを表す光学文字認識 (O C R) データが提供され得る。あるいは、媒体資料アナライザ 1 0 0 は O C R モジュール (図示せず) を含み得、この O C R モジュールは、スキャンされた画像データまたは電子画像データに関連する O C R データを生成する。

【 0 0 4 1 】

ループ 3 1 5 において、画像データから抽出された O C R テキストデータに対して、ス

50

ステップ 310 において、セグメンタ 110 は、テキストデータを、類似の画素値変化複雑性を有するように識別された対応する画像領域と関連付ける。セグメンタ 110 はまた、テキストデータ内のテキストサイズを識別し得、特に、コラム状の本文テキストに関連する本文テキストサイズを識別し得る。このようにして、本文テキストサイズを有する記事に属する候補であるテキストデータのブロックセグメントが識別され得る。

【0042】

図 5 に示されるさらなる実施形態において、ループ 315' (ステップ 510 ~ ステップ 550) は、ループ 315 の代わりに用いられる。ステップ 510 において、セグメンタ 110 は、単語を見出すために、画像データから抽出された OCR テキストデータを解析する。セグメンタ 110 は、テキストデータ内で見出された単語を、類似の画素値変化複雑性を有するとして、ステップ 310 において識別された初期の組の領域にマッピングする。セグメンタ 110 は、領域のテキストデータにおけるテキストサイズを決定し、特に、コラム状の本文テキストと関連する本文テキストサイズを識別し得る (ステップ 530)。セグメンタ 110 は、どの領域がコラム状の本文テキストと関連するかを決定する (ステップ 540)。

10

【0043】

次いで、セグメンタ 110 は、コラム状の本文テキスト (本文セグメントとも呼ばれる) を有する初期の組の領域を調整し、マッピングされた単語の分布に基づいて、最終の組の画像領域を得る。OCR セグメント化が良好であるとき、この調整は、画素値変化および閾値解析をただ用いることなく、OCR によって見出された領域における値の分布を見ることによって、問題のあるレイアウトエリアを固定することに役立ち得る。

20

【0044】

特徴に従って、セグメンタ 110 は、コラム状の本文テキストを対応する最終の組の画像領域とさらに関連付ける。例えば、どのテキストが本文テキストであるかを決定するために、テキストサイズを、所与の許容値を有する支配的なテキストサイズと比較する比較がなされ得る。

【0045】

画素値変化に基づくセグメント化は、例示的に上記され、本発明を必ずしも限定することを意図されない。現在公知の、または将来的に開発される他のセグメント化技術がデータ 105 をセグメント化するために使用され得る。

30

【0046】

(記事構成)

ステップ 220 において、特徴に従って、記事コンポーザ 120 は、言語統計情報およびレイアウト移行情報に基づいて、いずれの候補となるブロックセグメントが同一の記事に属するかを決定する。この言語統計およびレイアウト移行の組み合わせは、片方のみで達成し得る精度を超えて、精度を向上させ得る。

【0047】

(言語統計)

一実施形態において、ステップ 220 はルーチン 600 (ステップ 610 ~ ステップ 630) を含む。言語統計アナライザ 130 は、ルーチン 600 を実行する。ルーチン 600 は、セグメンタ 110 によって出力された複数の対の候補となるブロックセグメントに対するマッチスコアを計算する。マッチスコアは、スコア関数に従って計算される。特定のブロックのテキスト (好ましくは 30 を超える単語を有する) に対して、各単語が、全体の言語資料 (corpus) に対してブロック内にどの程度あるかを計算する。単語が記事内の X パーセントの単語と、言語資料内の Y パーセントの単語とを形成する場合、用いられる正しい式は $\log((X/Y) + 1)$ である。各ブロックに対して、値のベクトルが得られる。複数の対のブロック間のコサイン距離を得るために、これらのベクトルを用いることは、1 (正しいマッチ) ~ 0 (単語にオーバーラップしない) の範囲のスコアを与える。

40

【0048】

50

ステップ620において、言語統計アナライザ130は、計算されたマッチスコアおよび訓練データ135に基づいて、複数の対の候補となるブロックセグメントが同一の記事に属する一組の言語統計確率を計算する。この訓練データ135は、訓練の組および/またはユーザ入力から得られた確率データを含む。このような訓練は、利用可能な場合には同一の媒体資料によって実行され得る(例えば、図4の例に対して、他の発行者の新聞の1000ページに対する画像データ)。他の場合には、言語統計目的のための訓練が異なる媒体資料によってなされ得る。

【0049】

例えば、ステップ620において、複数の対のブロック間のスコア関数(ステップ610において計算される)を仮定すると、言語統計アナライザ130は、2つの任意のブロックが同一の記事からである確率を計算する。言語のために、大きな集合の記事へのアクセスがある場合には、様々なサイズのブロックに記事を分割する。同一の記事からの複数の対を、ポジティブな例として使用し、別の記事からの複数の対を、ネガティブな例として使用する。このような例がない場合には、明確に同一の記事であるブロック(同一のセグメント化された領域)およびほぼ明確に異なるブロック(異なるページ/発行者または離れている)を選び出すためにOCR化された文書自体をその代りに使用し得る。特有の単語長およびコサイン距離を有する一对のブロックを仮定する場合には、類似の例を見て、どの程度の割合がポジティブな例であるかを見出す。データ点の数に依存して、これは、kernel smootherまたはlocal regressionによって向上する。

10

20

【0050】

最終的に、ステップ630において、言語統計アナライザ130は、決定された確率に基づいて同一の記事に属するブロックを識別する。例えば、確率が50%を超える場合に、ブロックセグメントは、同一の記事に属するように識別され得る。精度が相対的に重要である一例においては、確率が90%を超える場合に、ブロックセグメントが同一の記事に属するように識別される。これらは例示である。他の確率閾値が使用され得る。

【0051】

これは、自動的または半自動的(半管理された習得タスク)であり得、このことはテキストの一部のブロックを仮定すると、一对のブロックが同一の記事に由来する確率を出力する。

30

【0052】

このマッチスコア関数ならびにコサイン距離および単語頻度の使用は、例示的であり、本発明を限定することを意図していない。現在公知の、または将来開発される他の関連のある技術が、2つのテキストのブロックの関連性を決定するか、またはスコア付けするために使用され得る。

【0053】

(レイアウト移行)

特徴に従って、レイアウト移行解析が、訓練モードにおいて、または実行モードにおいて実行され得る。一実施形態において、ステップ220は、訓練モジュール700および実行モジュール900を含む。レイアウト移行アナライザ140は、訓練モードまたは実行モードにおいて動作する。訓練モードにおいて、レイアウト移行アナライザ140は、媒体資料の複数のサンプルからの収集されたデータに対して動作し、レイアウト移行分類子145を構築する。実行モードにおいて、レイアウト移行アナライザ140は、レイアウト移行分類子145を、解析される媒体資料レイアウト内のデータに適用する。

40

【0054】

(訓練モード)

レイアウト移行アナライザ140は、訓練モジュール700を実行する。レイアウト移行アナライザ140は、垂直方向の移行(図7A、ステップ710~760)および水平方向の移行(図7B、ステップ770~796)について媒体資料のレイアウトを解

50

析する。

【0055】

(垂直方向の移行)

ステップ710において、レイアウト移行アナライザ140は、1つのブロックがもう1つのブロックよりも上にあり、垂直方向に整列されたブロックの間に本文テキストのブロックがないように垂直方向に整列された、複数の対の本文テキストのブロックセグメントを発見する(図7A)。これらの垂直方向に整列された複数の対のブロックセグメントは、セグメント110から出力されたブロックセグメント内に発見され得る。例えば、図8に示される新聞のページにおいて、一对のブロックセグメント810、830は、垂直方向に整列されているものとして識別され得る。

10

【0056】

複数の対の垂直方向に整列されたブロックセグメントの間におけるレイアウトに配置された項目が識別される(ステップ720)。そして、ブロックセグメントの対は、(少なくとも1つの)介在項目(intervening item)のうちのいずれか1つ以上に基づいて分類される(ステップ730)。例えば、垂直方向に整列されたブロックセグメント810、830の場合、介在項目820は、水平方向の罫線、テキストのライン、下線、水平方向の罫線、テキストのライン、水平方向の罫線である。そして、1つの分類は、これらの項目の特定の移行特徴、例えば、水平方向の罫線と、16ptのテキストのラインと、下線と、水平方向の罫線と、24ptのテキストのラインと、水平方向の罫線とによって分離されたブロックであり得る。

20

【0057】

次に、レイアウト移行アナライザ140は、分類されたブロックセグメントに対して、一組の移行特徴を計算する(ステップ740)。例えば、介在項目820を有するブロックセグメント810、830の場合、計算された移行特徴の組は、ブロックの垂直方向の分離の全て、どの程度良好にブロックが整列するか、ブロックの幅に対する罫線の幅、テキストのフォントサイズ、ブロックの幅に対するテキストのラインの幅、等であり得る。新聞のレイアウトのデータの一例においては、ブロックセグメントの対に対して用いられ得る垂直方向の移行特徴(例えば、ほぼ同一の平均テキストサイズの本文テキストの領域であって、それぞれの上に配置され、頂部ブロックおよび底部ブロックと称される)のリストは、(1)頂部ブロックおよび底部ブロックの平均の幅、(2)頂部ブロックと底部ブロックとの間の垂直方向の距離、(3)頂部ブロックの幅と底部ブロックの幅との間に本文でないテキストのブロックが存在するときの、頂部ブロックの幅と底部ブロックの幅との間の平均の幅の小部分としての差、(4)頂部ブロックおよび底部ブロックの左拡張部(left extent)、(5)頂部ブロックおよび底部ブロックの右拡張部(right extent)、(6)頂部ブロックおよび底部ブロックの頂部、(7)頂部ブロックおよび底部ブロックの底部、(8)頂部ブロックおよび底部ブロックの頂部と底部との間の距離、(9)頂部ブロックおよび底部ブロックの左と右との間の距離、(10)これらの頂部ブロックおよび底部ブロックにおける平均フォントサイズ、ならびに(11)これらの頂部ブロックおよび底部ブロックにおける最大フォントサイズ、を含む。

30

【0058】

これらの例の垂直方向の移行特徴は、例示的なものに過ぎず、本発明を限定することを意図されていない。計算されるべき正確な組の移行特徴は、訓練される分類子145に必要なとされる所望の精度と、媒体資料の複雑性とに依存する。より単純なレイアウトは、より少ない計算されるべき移行特徴を必要とし得る。精度が比較的重要であるより複雑なレイアウトまたはアプリケーションは、計算されるべきより多くの組の移行特徴を見込み得る。計算されるべき正確な組の特徴は、手動または自動で変更され得る。手動の変更は、ユーザインターフェース160を介するユーザからの入力に基づいて実行され得る。

40

【0059】

ステップ750において、レイアウト移行アナライザ140は、ブロックセグメントが同一の記事にある確率を決定する。ステップ730における各分類と、一組の垂直方向の

50

特徴とに対し、一連の数字によって概略される多数の垂直方向の特徴が存在する。また、言語統計アナライザ 130 の出力から、移行が記事の一部である確率が知られる。

【0060】

ここで、セグメントが同一の記事内にないときに、ブロックセグメントの区分の尤度を最大化するために、レイアウト移行分類子 145、例えば決定ツリーが、自動的に形成され得る（ステップ 760）。この分類子 145 は、各垂直方向の移行に対して、2 つの垂直方向に整列されたブロックが融合されるかどうかを決定するために用いられ得る。このようにして、垂直方向に整列されたブロックセグメントは、ここで、最大のコラムの集合であり、水平方向の移行の解析のための準備が来ている。垂直方向の移行特徴に基づいてレイアウト移行分類子 145 を形成するこの訓練は、利用可能なデータの集合体（例えば、1 つ以上の画像からの複数のブロックセグメント）にわたって実行され得る。一例では、本発明を限定することを意図しないが、訓練は、画像の大きな集合、例えば新聞の異なる記事からの 100 ページ以上にわたって実行され得、レイアウト移行分類子 145 を形成し得る。

10

【0061】

（水平方向の移行）

ステップ 770 において、レイアウト移行アナライザ 140 は、複数の対の水平方向に整列された本文テキストのブロックセグメントを発見し、1 つのブロックはその他のブロックの近くにあり、本文テキストのブロックは、水平方向に整列されたブロックの間には存在しない（図 7B）。これらの複数の対の水平方向に整列されたブロックセグメントは、セグメント 110 によって出力されたブロックセグメント内に発見される。例えば、図 8 に示されている新聞のページにおいて、一对のブロックセグメント 840、850 は、水平方向に整列されているとして識別され得る。

20

【0062】

複数の対の水平方向に整列されたブロックセグメントの間のレイアウトに配置された介在項目が識別される（ステップ 780）。そして、複数の対のブロックセグメントは、（少なくとも 1 つの）介在項目のうちの任意の 1 つ以上に基づいて、分類される（ステップ 790）。例えば、水平方向に整列されたブロックセグメント 840、850 の場合、介在項目はのど空きである。そして、1 つの分類は、これらの（少なくとも 1 つの）項目の特定の移行特徴（例えば、のど空きおよびその幅）によって分離されたブロックであり得る。

30

【0063】

次に、レイアウト移行アナライザ 140 は、分類されたブロックセグメントに対して、一組の移行特徴を計算する（ステップ 792）。例えば、ブロックセグメント 840、850 ならびにそれらの介在項目の場合、計算される一組の移行特徴は、のど空きおよびその幅、ブロックの水平方向の水平方向の分離の全体、ブロックがどの程度良好に整列されているか、等であり得る。新聞のレイアウトデータの一例において、ほぼ同一の平均テキストサイズであって、互いに近くに配置される複数の対のブロックセグメントまたは本文テキストの領域（左ブロックおよび右ブロックとも称される）に対して用いられ得る水平方向の移行特徴のリストは、（1）右ブロックの右縁と左ブロックの左縁との間の距離、（2）左ブロックおよび右ブロックの頂部の垂直方向の整列、（3）左ブロックと右ブロックとの間の水平方向の距離、（4）2 つの左ブロックと右ブロックとの幅の間の差、（5）本文ではないテキストの近傍に対する関係、を含む。本文ではないテキストの近傍に対するそのような関係は、例えば、左ブロックの頂部の最も近くの本文ではないテキストのブロックと、右ブロックの頂部の最も近くのブロックとを発見し、近傍の本文ではないテキストのブロックの各々に対して、2 つの左ブロックおよび右ブロックの最も遠くからの本文ではないテキストのブロックの垂直方向の距離、2 つの左ブロックおよび右ブロックの最も近くからの本文ではないテキストのブロックの垂直方向の距離、左ブロックを越えた本文ではないテキストのブロックの左の範囲、右ブロックを越えた本文ではないテキストのブロックの右の範囲、2 つの左ブロックおよび右ブロックの頂部の平均からの本文

40

50

ではないテキストのブロックの距離、2つの左ブロックおよび右ブロックの底部の平均からの本文ではないテキストのブロックの距離、本文ではないテキストのブロックの幅、本文ではないテキストのブロックの高さ、本文ではないテキスト内のフォントサイズおよび本文ではないテキストのブロックの単語数、を含む。

【0064】

これらの例の水平方向の移行特徴は、例示的なものであり、本発明を限定することを意図されていない。計算される完全な組の移行特徴は、訓練される分類子145に必要なとされる所望の精度と、媒体資料の複雑性とに依存する。より単純なレイアウトは、より少ない計算されるべき移行特徴を必要とし得る。精度が比較的重要であるより複雑なレイアウトまたはアプリケーションは、計算されるべきより大きな組の移行特徴を見込み得る。計算されるべき正確な組の特徴は、手動または自動で変更され得る。手動の変更は、ユーザインターフェース160を介するユーザからの入力に基づいて実行され得る。

10

【0065】

ステップ794において、レイアウト移行アナライザ140は、ブロックセグメントが同一の記事にある確率を決定する。ステップ790における各分類と、一組の水平方向の特徴とに対し、一連の数字によって概略される多数の移行特徴が存在する。また、言語統計アナライザ130の出力から、移行が記事の一部である確率が知られる。ここで、セグメントが同一の記事内にないときに、ブロックセグメントの区分の尤度を最大化するために、レイアウト移行アナライザ145、例えば決定ツリーが、自動的に形成され得る(ステップ796)。この分類子145は、水平方向の移行の各々を決定し、2つの水平方向に整列されたブロックのグループが同一の記事内に存在するかどうかを決定するために用いられ得る。このようにして、ブロックセグメントの最大のコラムは、それらが同一の記事に属するとき、水平方向にさらにグルーピングされる。水平方向の移行特徴に基づいてレイアウト移行分類子145を形成するこの訓練は、利用可能なデータの集合体(例えば、1つ以上の画像からの複数のブロックセグメント)にわたって実行され得る。一例では、本発明を限定することを意図しないが、訓練は、画像の大きな集合、例えば新聞の異なる記事からの100ページ以上にわたって実行され得、レイアウト移行分類子145を形成し得る。

20

【0066】

1つの利点は、このレイアウト移行解析が、構成される記事の精度を向上させるための言語統計解析を補足し得るということである。レイアウト手段において垂直方向および水平方向の移行特徴に基づいてブロックセグメントを分類する分類子145を用いると、一対のブロックセグメントが言語統計解析のみに基づいたときには無関係に見えるがいくつかの強く関連する対のパターンにフィットするという場合でさえも、コンバイナ150は、その対を同一の記事の一部として配置し得る。このようにして、言語統計とレイアウト移行とのこの組み合わせは、言語統計とレイアウト移行との一方のみで達成し得る精度を超えて、精度を向上させ得る。

30

【0067】

上述のように、発明者がテキストの2つのブロックがどのように関連しているかを評価するために単語の頻度を用いて実行した一部の例において、媒体資料解析ルーチン200は、2つのブロックが同一の記事からのものだったかどうかを示すことに関して、約90%の精度であった。媒体資料解析ルーチン200は、一般に正しい規則を発見するために、ページの大きな集合にわたってこれらの予測を組み合わせる。例えば、2つのブロックが16~20ptのヘルベティカ(Helvetica)テキストによって分離されるときに、これら2つのブロックが通常関連のないテキストを有している場合に、これは、記事を分離するものである可能性がある。

40

【0068】

この方法は、言語統計確率の評価を形成し、レイアウト統計特徴に基づいて、承認される一組の確率を形成し、このプロセスは、新しいデータが解析されるときに繰り返され、各実行可能性が、幾分か多くの情報を追加し、分類子の精度を向上させる。

50

【0069】

訓練分類子145がデータの集合体にわたって訓練され、構築されると、訓練分類子145は、実行モードにおいて動作するレイアウト移行アナライザ140によって用いられ得る。

【0070】

(実行モード)

実行モードは、図7に関連して上述された訓練モードに類似している。レイアウト移行アナライザ140は、実行モートルーチン900を実行する。レイアウト移行アナライザ140は、垂直方向の移行(図9A、ステップ710~740および910)と水平方向の移行(図9B、ステップ770~792および920)とについて、媒体資料のレイアウトを解析する。

10

【0071】

実行モードにおいて、レイアウト移行アナライザ140は、上述のように垂直方向に整列された複数の対のブロックセグメントに対して、ステップ710~740を実行する。分類子145を構築する代わりに、レイアウト移行アナライザ140は、分類子145(例えば、決定ツリー)を適用して、垂直方向に整列されたブロックセグメントが同一の記事に属するかどうかを決定する(ステップ910)。

【0072】

同様に、実行モードにおいて、レイアウト移行アナライザ140は、上述のように水平方向に整列された複数の対のブロックセグメントに対して、ステップ770~792を実行する。レイアウト移行アナライザ140は、分類子145(例えば、決定ツリー)を適用して、水平方向に整列されたブロックセグメントが同一の記事に属するかどうかを決定する(ステップ920)。

20

【0073】

訓練モードと実行モードとの間のこの分割は、例示的なものであり、本発明の実施形態を限定することを意図されていない。別の実施形態において、実行モード中に媒体資料アナライザ100の実行の間に出力された結果は、新しいデータ105が解析されるときに分類子が定期的に更新され得るように、分類子145を改変するために用いられ得る。

【0074】

プロセスの流れを垂直方向の流れと垂直方向の流れとの2つのタイプに分離することによってレイアウトを解析し、その後、垂直方向の移行と水平方向の移行とに何が似ているかと、それらの間に何が存在しているかに基づいて、垂直方向の移行と水平方向の移行とを合算する方法は、固有の特徴であるが、本発明を限定することは意図されていない。代替的に、レイアウト移行解析は、プロセスの流れを2つの部分に分離することなしに、レイアウト移行特徴に基づいてなされ得る。また、本文テキストのブロックを見る代わりに、ページ上の全ての要素に対してツリー構造を形成し、任意の2つのタイプの要素の間の移行の規則を形成するように試みることができる。

30

【0075】

分類子145を構築する際に用いられ得る複数の機械学習アプローチが存在する。本記載を与えられた当業者には明白なように、決定ツリーの実装の他に、`linear separator after a basis expansion`、`k-means clustering`、`kernel smoothing methods`等を用いることができる。別のアプローチは、単純に、特徴を離散化し、これらのバケットに分類し、見てきたケースが十分な例であったことを望むことである。

40

【0076】

(表示例)

図10A~10Dは、本発明の実施形態にしたがって解析される新聞のページを含む、例示的な媒体資料を示す。

【0077】

図10Aは、言語統計解析とレイアウト移行解析とに基づいて、媒体資料アナライザ1

50

00によって解析された新聞の第一面(f r o n t p a g e)の表示である。本文テキストを含むブロックセグメントは、媒体資料アナライザ100によって解析されるときに、本文テキストが属する対応する記事内でハイライトされる。同一の記事内のコラム状の本文テキストに対応するブロックセグメントは、アナライザ100がどのようにしてデータを細分化し、適切なセグメントを用いてどのように記事を構成したのかを示すために、同一のカラーによってカラー表示されたり、または陰影付けされたりする。図10Bは、レイアウト移行分類子とレイアウト移行アナライザとを有する媒体資料アナライザ100によって解析される、比較的トリッキーなレイアウトの新聞の内部ページの例の表示である。

【0078】

図10Cおよび図10Dは、本発明の一実施形態にしたがって言語統計解析(純粋な言語統計モード)に基づいて解析される媒体材料のハイライトされた例を示している。図10Cは、新聞の第一面の例を示しており、ブロックセグメントがハイライトされ、第一面に記事がある。図10Dは、新聞の内部のページを示しており、記事内でブロックセグメントがハイライトされている。この例では、同一の記事内のブロックセグメントが同一の色でハイライトされているが、本発明は、そのように限定されるわけではない。その他のタイプのハイライト表示(例えば、グレイスケールの陰影付け、境界、テクスチャあるいはその他のマークまたは印)が、カラーの代わりに、またはカラーに加えて用いられ得る。また、ハイライト表示は、必要に応じて用いられなくてもよく、記事のセグメントまたはその一部分のみが表示されることがあり得る。

【0079】

(ワールドワイドウェブへの応用)

本発明のさらなる実施形態にしたがうと、ワールドワイドウェブを介して、レイアウトを有する媒体資料をサーチするシステムが提供される(図11)。図11に示されているように、媒体資料をサーチするシステム1100は、クライアント1110と、ウェブサーバ1130と、サーバ1140と、データベース1145とを含む。クライアント1110は、ネットワーク1120を介して、ウェブサーバ1130に結合されている。ネットワーク1120は、ローカルエリアネットワーク、中規模エリアネットワーク、またはワイドエリアネットワークを含むがそれらには限定されない任意のタイプの1つ以上の任意のネットワーク、例えばインターネットであり得る。一例において、クライアント1110は、ネットワーク1120を介して通信するブラウザを含み得る。任意のタイプのブラウザが用いられ得る。ウェブサーバ1130は、サーバ1140に結合されている。

【0080】

サーバ140は、上述のように、媒体資料アナライザ100を含むか、それに結合されている。サーバ140はまた、データベース1145に結合されている。データベース1145は、媒体資料アナライザ100をサポートするためにデータを格納する任意のタイプのデータベースまたはメモリである。データベース1145は、上述のように、例えば、訓練データ135と、レイアウト移行分類子145と、データ105とを格納し得る。データベース1145はまた、画像データそのものを表すデータを含む媒体資料アナライザ100からの任意の出力を、そして媒体資料アナライザ100によって識別された記事に属するブロックセグメントと共に、格納し得る。勿論、特定の用途に依存して、インデックス付けおよびその他の操作が実行され得、それにより、出力されたデータは、サーチ要求またはその他のタイプのデータ要求を満足するように、容易に検索される。

【0081】

操作中、ユーザは、クライアント1110においてサーチクエリを入力し得る。そして、クライアント1110におけるブラウザは、ネットワーク1120を介して、ウェブサーバ1130にサーチクエリを転送する。ウェブサーバ1130は、サーバ1140と通信し、オプションとして、媒体資料アナライザ100と直接的に通信する。一実施形態において、媒体資料アナライザ100は、画像データそのものを表すデータと共に、記事に属するブロックセグメントを識別するメタデータを出力する。この出力は、データベース

10

20

30

40

50

1145に格納される。サーバ1140は、サーチ要求を満たすように、キーワードまたはサーチ用語について、データベース1145をサーチする。そして、サーバ1140は、ウェブサーバ1130に、サーチ要求を満足する結果を返す。そして、ウェブサーバ1130は、表示のために、満足されたサーチ結果を、クライアント1110におけるブラウザに転送する。このようにして、ウェブサーバ1130およびサーバ1140は、連携して動作し、任意のサーチエンジン、ポータルまたはウェブサイトの一部分となり得る。

【0082】

図12は、新聞の実施形態と共に用いられえる例示的なディスプレイ1200を示している。ディスプレイ1200に示されているように、フィールド1210は、サーチ結果を入力するために用いられ得る。そして、ボタン1215は、サーチを開始するためにユーザによって選択され得る。フィールド1220は、サーチからの出力結果を表示するように用いられ得る。一実施形態においては、サーチ結果を示すために多数のウィンドウが表示される。例えば、記事内でサーチ結果がヒットした場合、同一の記事からの2つの断片が、2つのウィンドウ1222、1224内に表示され得る。これらの断片は、サーチ用語と、サーチ用語の周辺の情報とを含み得る。これは、単なる例であり、本発明を限定するように意図されていない。1つ以上の記事からの1つ以上の断片が表示され得る。さらに、記事のテキスト全体、または記事内のサーチ用語のみ、または断片、例えばサーチ用語を包囲する領域が、表示され得る。任意の数のヒット、断片またはサーチを満足する所望のテキストが、表示され得る。図10に示されているようなハイライト表示されたブロックセグメントを有する新聞のページの全体の画像（またはその一部分）もまた、表示され得る。

10

20

【0083】

さらなる例にしたがうと、その他のタイプの情報が、ディスプレイ1200内に表示され得る。図12に示されているように、フィールド1230は、関連情報を示すために表示され得る。フィールド1240は、新聞に関する書誌情報、例えば、発行元、新聞が発行された日時、リポーターの署名欄および他の情報等の情報を表示するように用いられ得る。記事、タイトル、および新聞の名前をエリア1254に表示するために、別のフィールド1250が提供される。サーチを実行するユーザが記事をオーダーすること（ボタン1262）および新聞に契約すること（ボタン1264）を可能にするように、追加的な制御フィールド、例えばフィールド1260が提供され得る。ナビゲーション制御もまた提供され得る。例えば、ナビゲーションエリア1270は、ページ番号を表示するジャンプフィールド1272を含み得、該ジャンプフィールドは、ユーザが新聞の異なるページにジャンプすること、または、異なるサーチ結果にジャンプすることを可能にする。ユーザが、ユーザに提示された媒体情報の表示を、スクロールすること、ズームインすること、ズームアウトすること、または、変更することを可能にする、その他のナビゲーション制御（図示されず）が提供され得る。

30

【0084】

（さらなる特徴および利点）

特に、スキャンされた（または電子的に生成された）新聞ならびに関連する資料（例えば、雑誌、カタログ等）のレイアウトをセグメント化および解析するための新規なアプローチが提供される。高度な形態学関連アルゴリズム（*morphology-related algorithm*）は、ページを物理的なブロックに分解する。テキスト情報（ページ内に存在する、または、OCRから抽出される）は、テキストブロックをどのようにして記事に構成すべきか、および、どのようにしてテキストが流れるかを、決定するために用いられ得る。加えて、多数のページのテキスト解析を通して収集された情報が、レイアウト解析のために、集合に特有（*collection-specific*）の幾何学的規則を推察するために用いられ得る。

40

【0085】

レイアウトのセグメント化は、1つの単位として記事または結合した実体（*cohesive entity*）が何であるかを理解し、そしてインデックス付けすることを可能

50

にする。レイアウトのセグメント化はまた、便利にも記事にズームインすること、文脈内の記事を抽出すること、実際にテキストを提示することなしに、テキスト情報をリフローすることを可能にする。言い換えると、レイアウトのセグメント化は、スキャンされた新聞および雑誌をナビゲートするための強力かつ便利なユーザ経験を可能にし、実際、レイアウトのセグメント化はまた、電子ソース（例えば、PDF）に適用される。

【0086】

さらなる特徴にしたがうと、媒体資料アナライザ100の実施形態は、言語統計を用いることにより、様々な幾何学的要素の規則を学習し、例えば新聞のような媒体資料に特有の規則を計算し得る。そのようなアナライザおよび方法は、訓練データ内にどのような例示的なセグメント化も必要とせず、画像とOCR出力とから決定された言語統計から純粋に機能する。

10

【0087】

さらに、その他の上述の制限された幾何学ベースのレイアウト解析とは異なり、本発明の発明者によって、本明細書中に実施形態が提供され、この実施形態は、記事レベルのセグメント化を推察するために用いられるべきテキストデータから言語統計が引き出されることと、新聞/雑誌のページの特定の集合をセグメント化するために用いられ得る幾何学的規則を推察することとを可能にする実施形態が提供される。

【0088】

単純に最終的な推測を提供する代わりに、媒体資料アナライザは、記事に対する最良の推測のリストを保持し、ユーザがユーザインターフェースにおいて、不都合なものがある場合に、一部分の代替物を見ることを可能にする。ユーザは、概して、最も容易に読み取ることが可能なオプションを発見するまで、これを行うことを望み得る。ユーザが選択するものを観察することによって、記事コンポーザは、確率を変更することにより、リアルタイムで選択を更新し、ユーザ選択からトリッキーな領域を学習することを可能にする。

20

【0089】

記事のセグメント化の計算後、記事を表示するための多くの方法が存在する。1つの方法は、ユーザが記事を選択することを可能にし、その時点で、ユーザは、ページのズームイン画面を入手し、ユーザがスクロールホイールを用いて記事の複数の部分にわたってナビゲートすることを可能にする。バウンディングボックスと共に機能して、個々の単語の画像が抽出され、リフローされ得る。これは、新聞に似ているがより読みやすい別個のページを介して行われ得るか、あるいは、ユーザがテキストまたはコラムのサイズを変更し、その結果、新聞がスタイルおよびパラメータにフィットするように「再生成される」ことを可能にすることにより、行われ得る。

30

【0090】

（さらなる応用 - 連続する記事）

コラムの規則を発見するために多くのページにわたる集合体内で言語統計を用いる技術が、いくつかのその他の問題を解決するために用いられ得る。一実施形態にしたがうと、新聞のページ間で連続する記事内のブロックセグメントを決定するために、さらなる解析が実行され得る。異なるページ上の記事部分が同一の連続する記事内に属するかどうかの解析は、言語統計と、連続レイアウト移行情報とを用いて行われる。

40

【0091】

図13に示されているように、媒体資料アナライザ1300は、上述の媒体資料アナライザ100を含み、連続言語統計アナライザ1330と、連続レイアウトアナライザ1340とを有する記事コンポーザ1320をさらに含む。連続言語統計アナライザ1330と連続レイアウトアナライザ1340とは、図13に示されているように、記事コンポーザ1320内に含まれる。これは例示であり、本発明を限定することは意図されていない。例えば、連続言語統計アナライザ1330と連続レイアウトアナライザ1340とは、別個または組み合わせで提供されるか、または、言語統計アナライザ130とレイアウト移行アナライザ140との一部分として、それぞれ追加され得る。連続言語統計アナライザ1330と連続レイアウトアナライザ1340とは、ソフトウェア、ファームウェア

50

、ハードウェアまたはそれらの組み合わせで実行され得る。連続言語統計アナライザ 1 3 3 0 と連続レイアウトアナライザ 1 3 4 0 との機能は、明確化のために別個に記載されるが、1 つのモジュールまたはデバイス内で組み合わせられ得るか、または、より多くのモジュールまたはデバイスにわたって分散され得る。

【 0 0 9 2 】

媒体資料アナライザ 1 3 0 0 は、複数のページにわたって広がっている 1 つ以上の連続する記事を含むレイアウトを有する媒体資料を表すデータを解析する。媒体資料アナライザ 1 3 0 0 は、コントローラ 1 0 5 と、セグメンタ 1 1 0 と、記事コンポーザ 1 3 2 0 とを含む。セグメンタ 1 1 0 は、上述のように、媒体資料のページ内のコラム状の本文テキストに関連するブロックセグメントを識別する。記事コンポーザ 1 3 2 0 は、言語統計情報と連続移行情報とに基づいて、識別されたブロックセグメントのどれが、媒体資料内の複数のページにわたって広がっている連続する記事に属しているかを決定する。

10

【 0 0 9 3 】

コントローラ 1 0 2 は、セグメンタ 1 1 0 と記事コンポーザ 1 3 2 0 とを制御および管理する。ユーザからのさらなる制御が、ユーザインターフェース 1 6 0 を介して提供される。例えば、ユーザは、動作を開始すること、すなわち、データ 1 0 5、訓練データ 1 3 3 5 またはレイアウト移行分類子 1 3 4 5 の入力を開始することができる。ユーザは、媒体資料アナライザ 1 3 0 0 と相互作用することにより、訓練データ 1 3 3 5 を形成またはレビューすることを助け得る。例えば、ユーザは、所与の媒体資料内の複数のページにわたって広がっている連続する記事に属するブロックセグメントのポジティブおよびネガティブな例を用いることにより、訓練データ 1 3 3 5 の品質を向上させ得る。ユーザはまた、媒体資料アナライザ 1 3 0 0 と相互作用することにより、レイアウト移行分類子 1 3 4 5 を構築または修正し得る。

20

【 0 0 9 4 】

訓練データ 1 3 3 5 は、媒体資料内の連続する記事に属するブロックセグメントのポジティブおよびネガティブな例を含み得る。レイアウト移行分類子 1 3 4 5 は、候補となるブロックセグメントが媒体資料内の記事に属するとして分類されることを可能にする連続移行特徴を含むデータ構造を含み得るが、それには限定されない。そのようなデータ構造は、決定ツリーを含み得るが、それには限定されない。図 1 において上述された訓練データ 1 3 5 と分類子 1 4 5 とはまた、媒体資料アナライザ 1 3 0 0 と共に用いられ得、特に、別個に、または、訓練データ 1 3 3 5 と分類子 1 3 4 5 のそれぞれの一部として用いられ得る。

30

【 0 0 9 5 】

一実施形態において、記事コンポーザ 1 3 2 0 は、上述のような言語統計アナライザ 1 3 0、レイアウト移行アナライザ 1 4 0、コンバイナ 1 5 0 を含み、連続レイアウト移行アナライザ 1 3 4 0 と、連続言語統計アナライザ 1 3 3 0 とをさらに含む。連続レイアウト移行アナライザ 1 3 4 0 は、決定ツリー 1 3 4 5 を適用し、異なるページ上の候補となる記事部分の最後のブロックセグメントと第 1 のブロックセグメントとが同一の連続する記事内にある確率を示している 1 つ以上の連続移行特徴をピックアップする。連続言語統計アナライザ 1 3 3 0 は、計算された言語統計情報に基づいて、異なるページ上の異なる記事部分に対する言語統計情報を計算し、候補となる記事部分内の第 1 および最後のブロックセグメントが、連続する記事部分を有する確率を決定する。このようにして、記事コンポーザ 1 3 2 0 は、解析された連続レイアウト移行特徴と、計算された言語統計とにしたがって、第 1 および最後のブロックセグメントが同一の連続する記事に属する確率に基づいて、複数のページにわたって連続する記事を構成し得る。本明細書中の計算される確率は、yes / no またはブール値表示、確率または信頼データを表す数値、あるいは、確率または信頼データを表す値の数値範囲を含み得るが、それらには限定されない。

40

【 0 0 9 6 】

連続レイアウト移行アナライザ 1 3 4 0 と連続言語統計アナライザ 1 3 3 0 との動作は、図 1 4 A ~ E に示されているルーチン 1 4 0 0 に関連して (ステップ 1 4 0 2 ~ 1 4 3

50

6)、以下でさらに詳細に記載される。ルーチン1400は、セグメンタ110がデータ105内のブロックセグメントを識別した後に開始し、言語統計アナライザ130、レイアウト移行アナライザ140およびコンパイン150は、図1~10に関連して上述されたように、媒体材料の複数のページ上の記事内のブロックセグメントを組み合わせている。

【0097】

ステップ1402において、連続レイアウト移行アナライザ1340は、候補となる連続する記事部分に関連する最後のブロックセグメントを識別する。例えば、連続レイアウト移行アナライザ1340は、全てのコラム状のブロックを見て、記事内の最後のブロックセグメントである最後のブロックセグメントを識別する。そして、連続レイアウト移行アナライザ1340は、ブロックセグメントの下にある1つ以上の項目を識別する(ステップ1404)。そのような項目は、最後のブロックの終わりに現れる、単語(例えば、「連続している(continued)」、「連続(cont.)」、「~を参照(see)」または「~に行く(go to)」)あるいはレイアウト項目(例えば、矢印または直線、ドロ잉ボックス)であり得る。

10

【0098】

連続レイアウト移行アナライザ1340は、識別された1つ以上の項目の少なくとも1つの特徴に基づいて、最後のブロックセグメントを分類する(ステップ1406)。項目の特徴は、項目そのもの、または、項目およびレイアウトに関する特性を含み得る。例えば、単語「連続している」の特徴は、用語「連続している」のフォントサイズ、フォントスタイル、間隔または配置、および/または、単純に用語「連続している」そのものの存在であり得る。これらの特徴を分類することは、共通する特徴を有するブロックをそれぞれのグループに配置することを助け得る。例えば、ドロ잉ボックスを下に有するブロックは、1つのグループに分類され得るが、用語「連続している」を下に有するブロックは、別のグループに配置され得る。

20

【0099】

そして、連続レイアウト移行アナライザ1340は、決定ツリー1345を適用することにより、最後のブロックセグメントが連続する記事内にある確率を示す1つ以上の連続移行特徴をピックアップする(ステップ1408)。連続移行特徴は、特定のレイアウトに依存して、最後のブロックが別のページ上に連続する記事部分を有するより高い確率に関連する特徴であり得る。例えば、例示的なレイアウトにおける連続移行特徴は、12ptのイタリック体の用語「連続する(continued on)」と12ptの矢印とに関連する特徴であり得る。これらの特徴および例は、例示的なものであり、本発明を限定することは意図されていない。当業者には明白であり得るように、様々な組み合わせにおけるその他の特徴が、解析されるレイアウトと訓練データとに依存して用いられ得る。

30

【0100】

連続レイアウト移行アナライザ1340が最後のブロックを学習または識別すると(ステップ1402~1408)、連続言語統計アナライザ1330は、以後のページ上の記事部分を、最後のブロックがあるページからサーチする(図14B、ステップ1410)。連続言語統計アナライザ1330は、連続する記事の第1のページから、一連のページ上の記事部分に対する言語統計情報を計算し(ステップ1412)、計算された一連の言語統計情報に基づいて、候補となる連続記事部分内の最後のブロックセグメントが連続する記事部分を有している確率を決定する(ステップ1414)。例えば、言語統計情報は、単語の頻度情報であり得、連続言語統計アナライザ1330は、最後のブロックセグメント内のテキストおよび一連のページ上の記事部分内のテキストにおける単語の頻度に基づいて、マッチスコアを計算し得る。そして、最後のブロックセグメントが連続する記事部分を有する確率が、マッチスコアに基づいて決定され得る。訓練データ1335もまた、マッチスコアから確率を決定するために用いられ得る。

40

【0101】

ブロックが第1のブロックであるかどうかと、連続する記事部分内に存在する可能性が

50

あるかどうかとを決定するために、同様の解析が実行される。図14Cに示されているように、連続レイアウト移行アナライザ1340は、一連のページ内の候補となる連続する記事に関連する第1のブロックセグメントを識別する(ステップ1416)。例えば、連続レイアウト移行アナライザ1340は、全てのコラム状のブロックを見て、記事内の第1のブロックセグメントであるブロックセグメントを識別する。連続レイアウト移行アナライザ1340は、第1のブロックセグメントの上の1つ以上の項目を識別する(ステップ1418)。そのような項目は、第1のブロックの上に見える単語(例えば、「から連続している(continued from)」または「からの連続(cont.from)」)であり得る。

【0102】

そして、連続レイアウト移行アナライザ1340は、1つ以上の識別された項目の少なくとも1つの特徴に基づいて、第1のブロックセグメントを分類する(ステップ1420)。項目の特徴は、項目そのものまたは項目およびレイアウトに関する特性を含み得る。例えば、単語「連続している(continued)」に対する特徴は、用語「連続している(continued)」のフォントサイズ、フォントスタイル、間隔または配置、および/または、単純に用語「連続している(continued)」そのものの存在であり得る。これらの特徴を分類することは、共通する特徴を有するブロックをそれぞれのグループに配置することを助け得る。例えば、ドロ잉ボックスを下に有するブロックは、1つのグループに分類され得るが、用語「連続している(continued)」を下に有するブロックは、別のグループに配置され得る。

【0103】

連続レイアウト移行アナライザ1340は、決定ツリー1345を適用することにより、第1のブロックセグメントが連続する記事内にある確率を示す1つ以上の連続移行特徴をピックアップする(ステップ1422)。そして、連続レイアウト移行アナライザ1340は、適用される決定ツリー1345に基づいて、第1のブロックセグメントが連続する記事内にある確率を決定し得る(ステップ1424)。そのような連続移行特徴は、特定のレイアウトに依存して、第1のブロックが以前のページ上に連続する記事部分を有するより高い確率に関連する特徴であり得る。例えば、例示的なレイアウトにおける連続移行特徴は、12ptのイタリック体の用語「~から連続する(continued from)」に関連する特徴であり得る。これらの特徴および例は、例示的なものであり、本発明を限定することは意図されていない。当業者には明白であり得るように、様々な組み合わせにおけるその他の特徴が、解析されるレイアウトと訓練データとに依存して用いられ得る。

【0104】

連続レイアウト移行アナライザ1340が第1のブロックを学習または識別すると(ステップ1416~1424)、連続言語統計アナライザ1330は、以前のページ上の記事部分を、第1のブロックがあるページからサーチする(図14D、ステップ1426)。連続言語統計アナライザ1330は、第1のブロックのページから、以前のページ上の記事部分に対する言語統計情報を計算し(ステップ1428)、計算された以前のページの言語統計情報に基づいて、候補となる連続記事部分内の第1のブロックセグメントが連続する記事部分である確率を決定する(ステップ1430)。例えば、言語統計情報は、単語の頻度情報であり得、連続言語統計アナライザ1330は、第1のブロックセグメント内のテキストおよび一連のページ上の記事部分内のテキストにおける単語の頻度に基づいて、マッチスコアを計算し得る。そして、第1のブロックセグメントが連続する記事部分である確率が、マッチスコアに基づいて決定され得る。訓練データもまた、マッチスコアから確率を決定するために用いられ得る。

【0105】

最後に、最後および第1のブロックが候補となる連続する記事に対して識別されると、別個のページ上の記事部分の間で連続しているかどうかについて学習され得る。図14Eに示されているように、連続レイアウト移行アナライザ1340は、一対の候補となる最

10

20

30

40

50

後および第1のブロックセグメントを識別し(ステップ1432)、一対の候補となる最後および第1のブロックセグメントに対する一組の連続移行特徴を決定し(ステップ1434)、決定された連続移行特徴の組に基づいて、決定ツリー1345を適用し、一対の候補となる最後および第1のブロックセグメントが媒体資料内の複数のページにわたって同一の連続する記事に属する確率を決定する(ステップ1436)。一組の連続移行特徴は、1つの記事部分から別の記事部分への進行(progression)に関係している特性によって、拡張され得る。例えば、最後のブロックセグメント内の最後の単語、最後のブロックの下の単語または図面、第1のブロック内の第1の単語、第1のブロックの上の単語または図面を考え、様々なセクションの間の単語の重複を見る。決定ツリー1345が形成されたときに、「連続している」を含み、第2のセクション内の第1の3つの単語のうち1つにマッチする単語を有している第1の区分の下に太字が存在するかどうか分かり、そして、これが特定のレイアウトに対して右の連続であり得る可能性が分かる。

10

【0106】

本記載を与えられた当業者には明白であり得るように、上述の実施形態の媒体資料アナライザ1300は、訓練データ1335の使用と共に、または、訓練データ1335の使用を伴わずに、動作され得る。訓練データ1335の使用を伴わないと、連続の決定は、見られている特定のブロックのみに基づき得る。例えば、媒体資料アナライザ1300は、言語スコアが最大になるように、一部の記事内の最後のブロックを取り得、異なるページ上にある一部の記事内の最初のブロックを発見し得る。これは、時々機能するが、多くの可能性のある答えが存在し得、正しいものは、これらの段落内で最良のテキストマッチを有していない可能性があるというだけの理由で、ミスをする比較的高い可能性がある。

20

【0107】

精度を向上させるために、媒体資料アナライザ1300は、所定の訓練データ1335の利点を用いて動作し得るか、または、訓練モードで動作して、訓練データ1335を入手または補足し得る。本質的に、訓練は、図14における実行モードに関連して記載された上述のアプローチを用い、より正確な決定ツリーまたは分類子1345を形成するために、多くのページにわたって結果を組み合わせる。そのためものとして、上述では、異なるページ上の複数の領域の対、1つの領域からその他の領域への移行の記載(述べられている特徴)、および、該対の言語スコアの例があった。これらのスコアを(非常に弱いスコアであったとしても)組み合わせることにより、傾向が抽出され、決定ツリーまたは分類子1345を構築するために用いられ得る。例えば、第1のブロックの下に小さい三角形があったり、または、第1のブロック内に単語「cont'd」があったりし、ブロック間で特定の量の単語が重複している場合には、これらは有効な連続であり得る。

30

【0108】

効率化のために、1つの例では、これらは3つのタスクに分割される。なぜならば、対応のために全ての対のブロックをテストすることは、過度に複雑であり得るからである。

【0109】

代わりに、3つの部分またはテストが存在する：

- (1) 所定のブロックは連続している(continued)か？
- (2) 所定のブロックは連続(continuation)か？
- (3) これらの2つのタイプのブロックの対が与えられたとき、それらは同一の記事か？

40

これを同時に行うことを試みることは、分類子1345が、(3)を理解する必要があるのみならず、(1)および(2)を2つの部分に組み込む必要があり得ることを意味する。単一の分類子のアプローチが可能であるが、上記のものは、より安全であり、テストがより容易である。

【0110】

(例示的なコンピュータシステムの実装)

本発明の様々な局面は、ソフトウェア、ファームウェア、またはそれらの組み合わせに

50

よって実装され得る。図15は、例示的なコンピュータシステム1500を示しており、このコンピュータシステムにおいては、本発明またはその一部分が、コンピュータ読み取り可能なコードとして実行され得る。例えば、図2の方法200を実行する媒体資料アナライザ100、図14の方法1400を実行する媒体資料アナライザ1300が、システム1300内に実装され得る。本発明の様々な実施形態が、この例のコンピュータシステム1500の観点で記載される。この記載を読んだ後、当業者は、その他のコンピュータシステムおよび/またはコンピュータアーキテクチャを用いて、本発明をどのように実装するかを明白に理解し得る。

【0111】

コンピュータシステム1500は、1つ以上のプロセッサ、例えばプロセッサ1504を含む。プロセッサ1504は、特定用途向けプロセッサまたは汎用のプロセッサであり得る。プロセッサ1504は、通信インフラストラクチャ1506（例えば、バスまたはネットワーク）に連結され得る。

10

【0112】

コンピュータシステム1500はまた、メインメモリ1508（好適には、ランダムアクセスメモリ（RAM））を含み得、さらには、2次メモリ1510を含み得る。2次メモリ1510は、例えば、ハードディスクドライブ1512および/またはリムーバブル可能格納ドライブ1514を含み得る。リムーバブル可能格納ドライブ1514は、フロッピー（登録商標）ディスクドライブ、磁気テープドライブ、光ディスクドライブ、フラッシュメモリ等を含み得る。リムーバブル格納ドライブ1514は、周知の方法で、リムーバブル格納ユニット1518からの読み取り、および/または、リムーバブル格納ユニット1518への書き込みを行い得る。リムーバブル格納ユニット1518は、フロッピー（登録商標）ディスク、磁気テープ、光ディスク等を含み得、リムーバブル格納ユニット1518は、リムーバブル格納ドライブ1514によって読み取りおよび書き込みがなされる。当業者は、リムーバブル格納ユニット1518が、コンピュータソフトウェアおよび/またはデータが格納されたコンピュータ使用可能な媒体を含むことを理解し得る。

20

【0113】

代替的な実装において、2次メモリ1510は、コンピュータプログラムまたはその他の命令がコンピュータシステム1500にロードされることを可能にするその他の類似の手段を含み得る。そのような手段は、例えば、リムーバブル格納ユニット1522およびインターフェース1520を含み得る。そのような手段の例は、プログラムカートリッジおよびカートリッジインターフェース（例えば、ビデオゲームデバイスに見られるようなもの）、リムーバブルメモリチップ（例えば、EPROMまたはPROM）および関連ソケット、ならびに、ソフトウェアおよびデータがリムーバブル格納ユニット1522からコンピュータシステム1500に転送されることを可能にするその他のリムーバブル格納ユニット1522およびコンピュータシステム1500を含み得る。

30

【0114】

コンピュータシステム1500はまた、通信インターフェース1524を含み得る。通信インターフェース1524は、ソフトウェアおよびデータが、コンピュータシステムと外部デバイスとの間で転送されることを可能にする。通信インターフェース1524は、モデム、ネットワークインターフェース（例えば、イーサネット（登録商標）カード）、通信ポート、PCMCIAスロットおよびカード等を含み得る。通信インターフェース1524を介して転送されるソフトウェアおよびデータは、信号の形態であり得、この信号は、通信インターフェース1524によって受信されることが可能な電子的、電磁的、光学的、またはその他の信号であり得る。これらの信号は、通信バス1526を介して通信インターフェース1524に提供される。通信バス1526は、信号を搬送し、ワイヤまたはケーブル、光ファイバー、電話回線、携帯電話リンク、RFリンクまたはその他の通信チャネルを用いて実装され得る。

40

【0115】

本書面において、用語「コンピュータプログラム媒体」および「コンピュータ使用可能

50

媒体」は、例えばリムーバブル格納ユニット 1518、リムーバブル格納ユニット 1522、ハードディスクドライブ 1512 にインストールされたハードディスク、ならびに、通信パス 1526 を介して搬送される信号等の媒体を広く意味するように用いられる。また、コンピュータプログラム媒体およびコンピュータ使用可能媒体は、例えばメインメモリ 1508 および 2 次メモリ 1510 等のメモリを意味し得、このメモリは、メモリ半導体（例えば、DRAM 等）であり得る。これらのコンピュータプログラム製品は、コンピュータシステム 1500 にソフトウェアを提供する手段である。

【0116】

コンピュータプログラム（コンピュータ制御論理とも称される）は、メインメモリ 1508 および / または 2 次メモリ 1510 に格納される。コンピュータプログラムはまた、通信インターフェース 1524 を介して受信され得る。そのようなコンピュータプログラムは、実行されたときに、コンピュータシステム 1500 が、本明細書中で議論された本発明を実行することを可能にする。特に、これらのコンピュータプログラムは、実行されたときに、プロセッサ 1504 が、例えば上記で議論された図 2 のフローチャート 200 によって示された方法のステップのような、本発明のプロセスを実行することを可能にする。したがって、そのようなコンピュータプログラムは、コンピュータシステム 1500 のコントローラを代表する。ソフトウェアを用いて本発明が実行される場合、そのソフトウェアは、リムーバブル格納ドライブ 1514、インターフェース 1520、ハードドライブ 1512 または通信インターフェース 1524 を用いることにより、コンピュータプログラム製品に格納されたり、コンピュータシステム 1500 にロードされたりし得る。

10

20

【0117】

本発明の実施形態はまた、任意のコンピュータ使用可能媒体上に格納されたソフトウェアを含むコンピュータ製品に関係し得る。そのようなソフトウェアは、1 つ以上のデータ処理デバイス内で実行されたときに、（少なくとも 1 つの）データ処理デバイスに、本明細書中で議論されたような動作を行わせる。任意のコンピュータ使用可能または読み取り可能な媒体を利用する本発明の実施形態は、現在公知であるか、将来公知になる。コンピュータ使用可能媒体の例は、プライマリ格納デバイス（例えば、任意のタイプのランダムアクセスメモリ）、2 次格納デバイス（例えば、ハードドライブ、フロッピー（登録商標）ディスク、CD ROM、ZIP ディスク、テープ、磁気格納デバイス、光学格納デバイス、MEMS、ナノテクノロジー格納デバイス等）、通信媒体（例えば、有線および無線の通信ネットワーク、ローカルエリアネットワーク、ワイドエリアイントラネット等）を含むが、これらには限定されない。

30

【0118】

（結論）

本発明の例示的な実施形態が提示されてきた。本発明は、これらの例に限定されるものではない。これらの例は、本明細書中では、例示を目的として示されており、限定目的として示されてはいない。当業者は、本明細書中に含まれる教示に基づくことにより、代替案（本明細書中に記載されているものの均等、拡張、変形、逸脱（deviation）等）を明白に理解し得る。そのような代替案は、本発明の範囲および精神に含まれる。

【 図 1 】

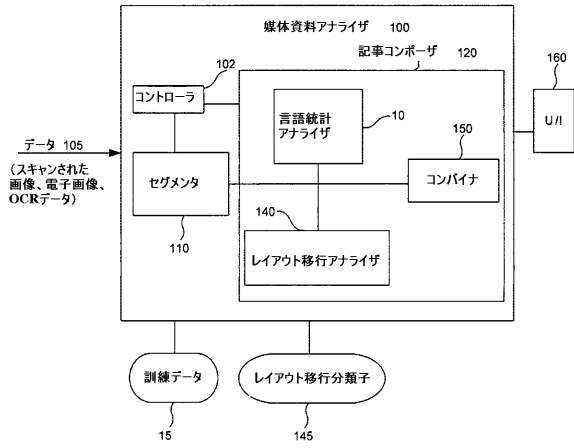


FIG. 1

【 図 2 】

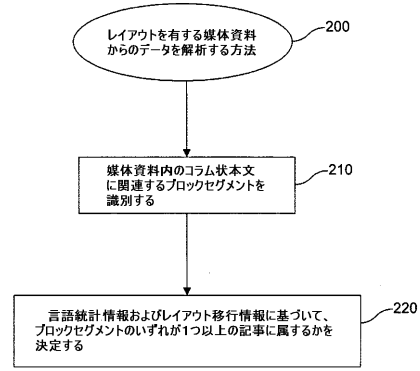


FIG. 2

【 図 3 】

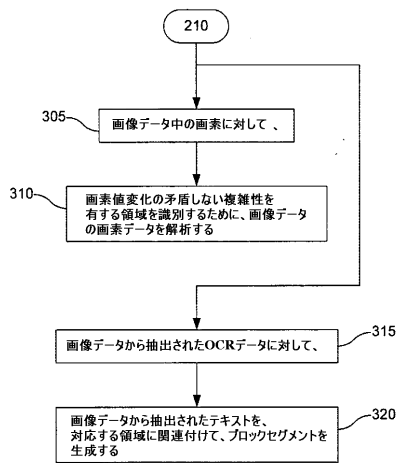


FIG. 3

【 図 4 】

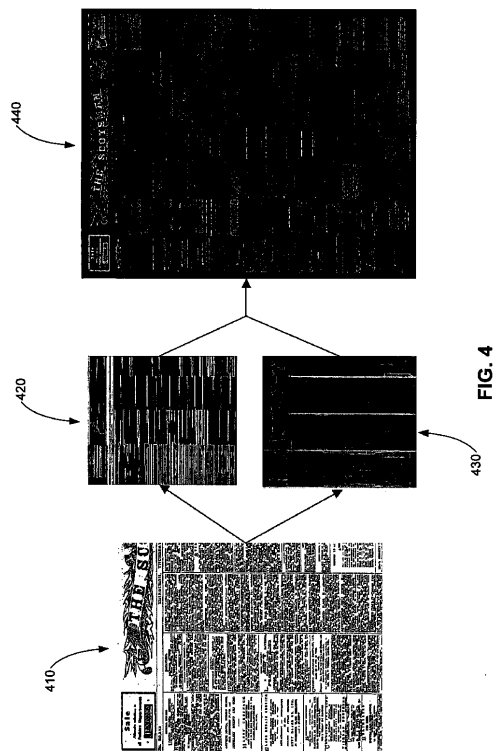


FIG. 4

【 図 5 】

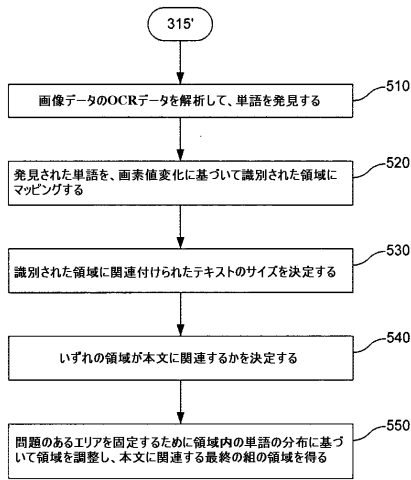


FIG. 5

【 図 6 】

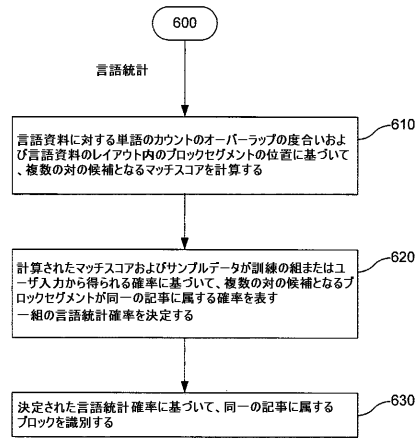


FIG. 6

【 図 7 A 】

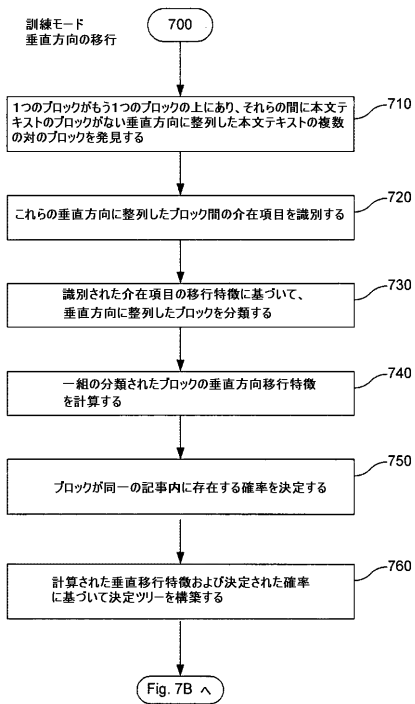


FIG. 7A

【 図 7 B 】

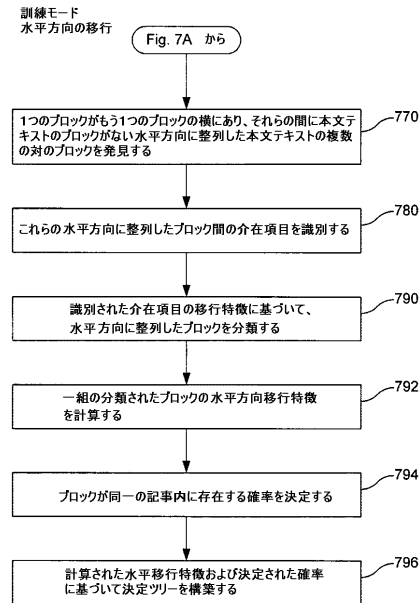


FIG. 7B

【 図 8 】



FIG. 8

【 図 9 A 】

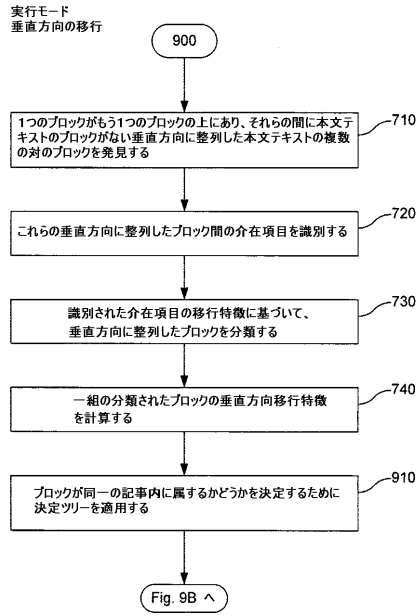


FIG. 9A

【 図 9 B 】

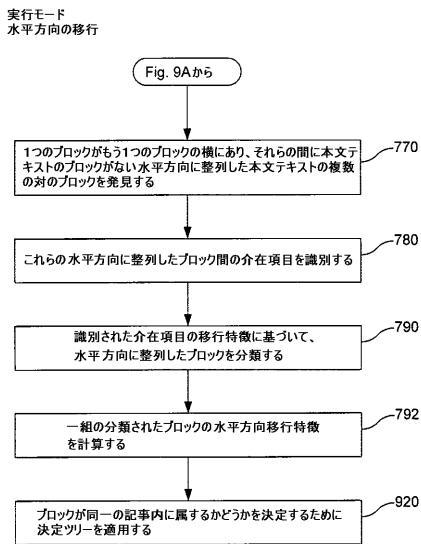


FIG. 9B

【 図 10 A 】



ページ1—最終的な決定ツリー例

FIG. 10A

【図10B】



ページ14 一最終的な決定ツリー例
FIG. 10B

【図10C】



ページ1 一純粋な言語統計例
FIG. 10C

【図10D】



ページ14 一純粋な言語統計例
FIG. 10D

【図11】

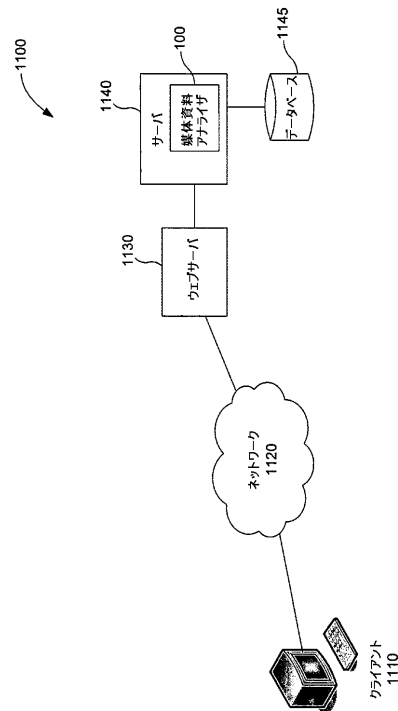


FIG. 11

【 図 1 2 】

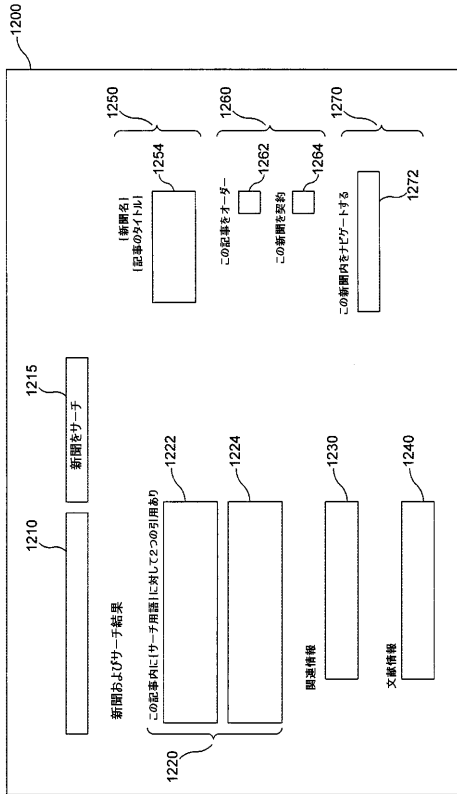


FIG. 12

【 図 1 3 】

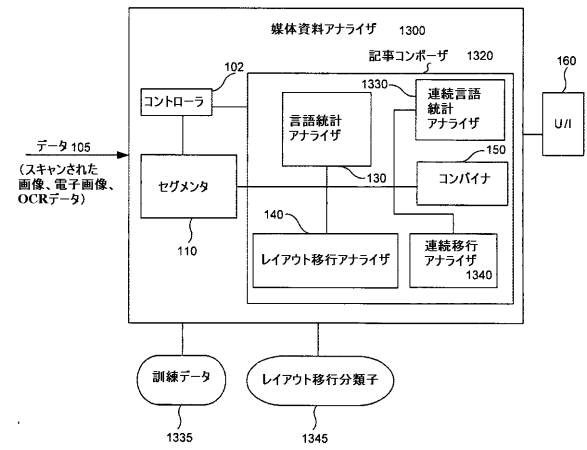
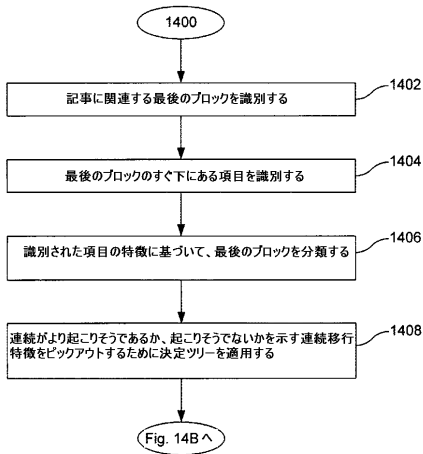


FIG. 13

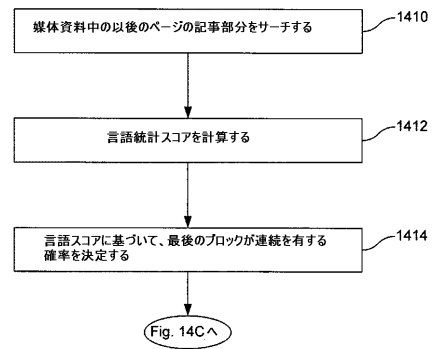
【 図 1 4 A 】



連続移行実行モード

FIG. 14A

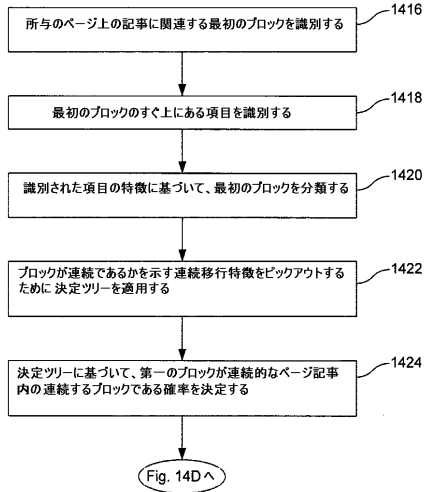
【 図 1 4 B 】



連続移行実行モード

FIG. 14B

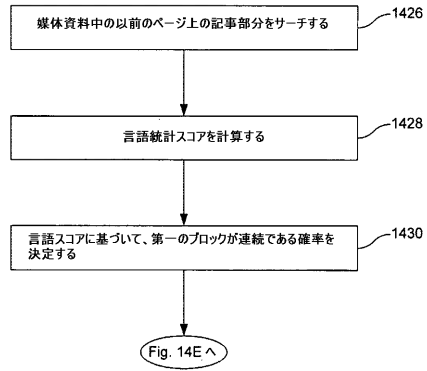
【 図 1 4 C 】



連続移行実行モード

FIG. 14C

【 図 1 4 D 】



連続移行実行モード

FIG. 14D

【 図 1 4 E 】

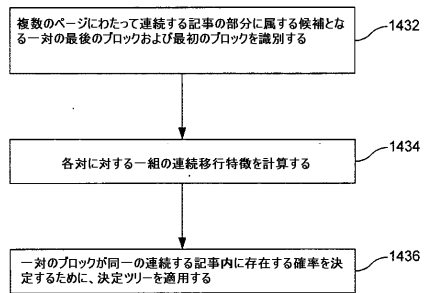


FIG. 14E

【 図 1 5 】

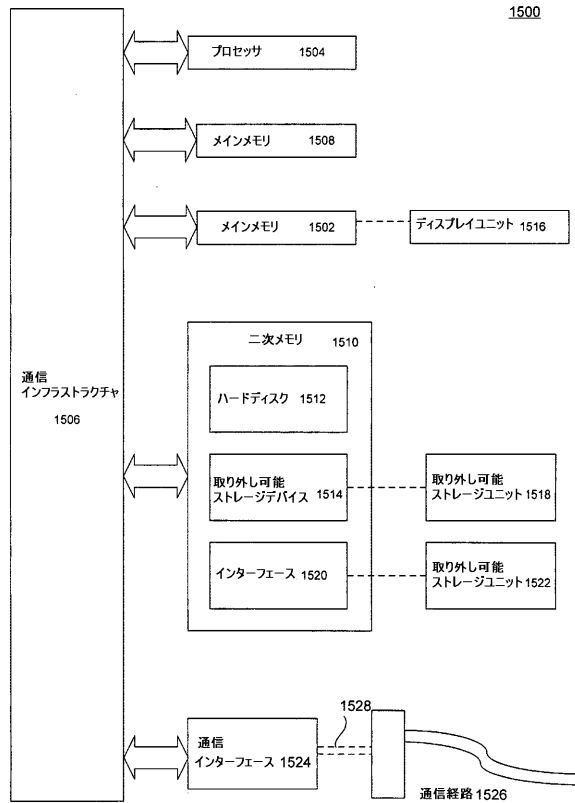


FIG. 15

【手続補正書】

【提出日】平成21年6月29日(2009.6.29)

【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】0017

【補正方法】変更

【補正の内容】

【0017】

本発明の実施形態は、添付の図面を参照して記載される。図面において、同様の参照番号は同一の要素または機能的に類似の要素を示し得る。ある要素が第1に現れる図面は、対応する参照番号のもっとも左の桁によって概して示される。

例えば、本発明は以下の項目を提供する。

(項目1)

レイアウトを有し、かつ複数のページにわたって広がる1つ以上の連続する記事を含む媒体資料を表すデータを解析する媒体資料アナライザであって、

(a) 該媒体資料のページ内のコラム状の本文テキストと関連するブロックセグメントを識別するセグメントと、

(b) 言語統計情報および連続移行情報に基づいて、該識別されたブロックセグメントのいずれが該媒体資料内の複数のページにわたって広がる連続する記事に属するかを決定する記事コンポーザと

を備えている、媒体資料アナライザ。

(項目2)

上記記事コンポーザは、連続レイアウト移行アナライザを含み、該連続レイアウト移行アナライザは、第1のページ内の候補となる連続する記事部分に関連する最後のブロックセグメントを識別し、該最後のブロックセグメントの下の1つ以上の項目を識別し、該識別された1つ以上の項目の少なくとも1つの特徴に基づいて、該最後のブロックセグメントを分類し、そして、決定ツリーを適用して、該最後のブロックセグメントが連続する記事内にある確率を示す1つ以上の連続移行特徴を選び出す、項目1に記載の媒体資料アナライザ。

(項目3)

上記記事コンポーザは、連続言語統計アナライザを含み、該連続言語統計アナライザは、上記連続する記事の上記第1のページから連続するページにおける記事部分に対して、言語統計情報を計算し、そして該計算された連続する言語統計情報に基づいて、上記候補となる連続する記事部分内の最後のブロックセグメントが、連続する記事部分を有する確率を決定する、項目2に記載の媒体資料アナライザ。

(項目4)

上記連続レイアウト移行アナライザは、さらに、連続するページ内の候補となる連続する記事と関連する第1のブロックセグメントを識別し、該第1のブロックセグメントの上の1つ以上の項目を識別し、該1つ以上の識別された項目の少なくとも1つの特徴に基づいて、該第1のブロックセグメントを分類し、そして決定ツリーを適用して、該第1のブロックセグメントが連続する記事内にある確率を示す1つ以上の連続移行特徴を選び出し、該適用された決定ツリーに基づいて、該第1のブロックセグメントが連続する記事である確率を決定する、項目3に記載の媒体資料アナライザ。

(項目5)

上記連続言語統計アナライザは、さらに、上記第1のブロックセグメントを有するページよりも前のページ内の記事部分に対して、言語統計情報を計算し、該計算された、より前のページの言語統計情報に基づいて、上記候補となる連続する記事部分内の該第1のブロックセグメントが連続する記事部分を有する確率を決定する、項目4に記載の媒体資料アナライザ。

(項目6)

上記連続レイアウト移行アナライザは、さらに、候補となる一对の最後および第1のブロックセグメントに対する連続移行特徴を識別し、該一对の最後および第1のブロックセグメントに対して一組の連続移行特徴を決定し、そして決定ツリーを適用して、該一組の決定された連続移行特徴に基づいて、該候補となる一对の最後および第1のブロックセグメントが、上記媒体資料内の複数のページにわたる同一の連続する記事に属する確率を決定する、項目5に記載の媒体資料アナライザ。

(項目7)

上記言語統計情報は、単語頻度情報を備え、上記連続言語統計アナライザは、上記最後のブロックセグメント内のテキストおよび連続するページ上の上記記事部分内のテキストにおける単語頻度に基づいて、マッチスコアを計算する、項目3に記載の媒体資料アナライザ。

(項目8)

上記言語統計情報は、単語頻度情報を備え、上記連続言語統計アナライザは、上記第1のブロックセグメント内のテキストおよびより前のページ上の上記記事部分内のテキストにおける単語頻度に基づいて、マッチスコアを計算する、項目5に記載の媒体資料アナライザ。

(項目9)

レイアウトを有し、かつ複数のページにわたって広がる1つ以上の連続する記事を含む媒体資料を表すデータを解析するコンピュータ実装された方法であって、

(a) 該媒体資料のページ内のコラム状の本文テキストに関連するブロックセグメントを識別することと、

(b) 言語統計情報および連続移行情報に基づいて、該識別されたブロックセグメントのいずれが該媒体資料内の複数のページにわたって広がる連続する記事に属するかを決定することと

を包含する、方法。

(項目10)

上記記事決定ステップは、

候補となる連続する記事部分内のブロックセグメントに対する連続レイアウト移行情報を解析することと、

該候補となる連続する記事部分内のテキストに対する言語統計を解析することと

を含む、項目9に記載の方法。

(項目11)

上記連続レイアウト移行情報を解析するステップは、

第1のページ内の候補となる連続記事部分に関連する最後のブロックセグメントを識別することと、

該最後のブロックセグメントより下の1つ以上の項目を識別することと、

該識別された1つ以上の項目の少なくとも1つの特徴に基づいて、該最後のブロックセグメントを分類することと、

決定ツリーを適用して、該最後のブロックセグメントが連続する記事内に存在する確率を示す1つ以上の連続移行特徴を選ぶことと

を包含する、項目10に記載の方法。

(項目12)

上記言語統計解析ステップは、

上記連続する記事の第1のページから連続するページ上の記事部分に対する、言語統計情報を計算することと、

該計算された連続する言語統計情報に基づいて、上記候補となる連続する記事部分内の上記最後のブロックセグメントが連続する記事部分を有する確率を決定することと

を包含する、項目11に記載の方法。

(項目13)

上記連続レイアウト移行情報を解析するステップは、

連続するページ内の候補となる連続する記事に関連する第1のブロックセグメントを識別することと、

該第1のブロックセグメントより上の1つ以上の項目を識別することと、

該1つ以上の識別された項目の少なくとも1つの特徴に基づいて、該第1のブロックセグメントを分類することと、

決定ツリーを適用して、該第1のブロックセグメントが連続する記事内に存在する確率を示す1つ以上の連続移行特徴を選び出し、そして該適用された決定ツリーに基づいて、該第1のブロックセグメントが連続する記事内に存在する確率を決定することと

をさらに包含する、項目12に記載の方法。

(項目14)

上記言語統計解析ステップは、

上記第1のブロックセグメントを有するページよりも前のページ内の記事部分に対して、言語統計情報を計算することと、

該計算されたより前のページの言語統計情報に基づいて、上記候補となる連続する記事部分内の該第1のブロックセグメントが連続する記事部分を有する確率を決定することと

をさらに包含する、項目13に記載の方法。

(項目15)

上記連続レイアウト移行解析ステップは、

候補となる一対の最後および第1のブロックセグメントを識別することと、

該一対の最後および第1のブロックセグメントに対する一組の連続移行特徴を決定することと、

決定ツリーを適用して、該一組の決定された連続移行特徴に基づいて、該候補となる一対の最後および第1のブロックセグメントが、上記媒体資料内の複数のページにわたる同一の連続する記事に属する確率を決定することと

をさらに含む、項目12に記載の方法。

(項目16)

上記言語統計情報は、単語頻度情報を備え、上記連続言語統計解析ステップは、上記最後のブロックセグメント内のテキストおよび連続するページの上記記事部分内のテキストにおける単語頻度に基づいて、マッチスコアを計算することを含む、項目12に記載の方法。

(項目17)

上記言語統計情報は、単語頻度情報を備え、上記連続言語統計解析ステップは、上記第1のブロックセグメント内のテキストおよび以前のページの上記記事部分内のテキストにおける単語頻度に基づいて、マッチスコアを計算することを含む、項目14に記載の方法。

(項目18)

レイアウトを有する媒体資料内の複数のページにわたって広がる連続する記事を構成する記事コンポーザであって、

連続レイアウト移行アナライザと、

連続言語統計アナライザと

を備え、該連続レイアウト移行アナライザは、異なるページ上の候補となる記事の最後のブロックセグメントおよび第1のブロックセグメントが同じ連続する記事内に存在する確率を示す1つ以上の連続移行特徴を選び出すために、決定ツリーを適用し、

該連続言語統計アナライザは、異なるページ上の異なる記事部分に対する言語統計情報を計算し、該計算された言語統計情報に基づいて、候補となる記事部分の第1および最後のブロックセグメントが連続する記事部分を有する確率を決定し、それにより、該記事コンポーザは、解析された連続レイアウト移行特徴および該計算された言語統計に従って、該第1および最後のブロックセグメントが同じ連続する記事に属する確率に基づいて、複数のページにわたる連続する記事を構成することが可能である、記事コンポーザ。

【手続補正 2】【補正対象書類名】特許請求の範囲【補正対象項目名】全文【補正方法】変更【補正の内容】【特許請求の範囲】【請求項 1】

レイアウトを有する媒体資料を表すデータを解析する媒体資料アナライザであって、該媒体資料内のコラム状の本文テキストと関連するブロックセグメントを識別するセグメントと、

該セグメントによって出力された候補となるブロックセグメントに対する言語統計を計算し、言語統計情報内のオーバーラップに基づいて、候補となるブロックセグメントが同一の記事に属する確率を決定する言語統計アナライザと

を備えている、媒体資料アナライザ。

【請求項 2】

言語統計情報およびレイアウト移行情報に基づいて、前記識別されたブロックセグメントのうちのいずれが前記媒体資料内の 1 つ以上の記事に属するかを決定する記事コンポーザをさらに備えている、請求項 1 に記載の媒体資料アナライザ。

【請求項 3】

前記記事コンポーザは、連続レイアウト移行アナライザを含み、該連続レイアウト移行アナライザは、第 1 のページ内の候補となる連続する記事部分に関連する前記識別されたブロックセグメントのうちの最後のブロックセグメントを識別し、該最後のブロックセグメントの下に 1 つ以上の項目を識別し、該識別された 1 つ以上の項目の少なくとも 1 つの特徴に基づいて、該最後のブロックセグメントを分類し、そして、決定ツリーを適用して、該最後のブロックセグメントが連続する記事にある確率を示す 1 つ以上のレイアウト移行特徴を選び出し、該識別されたブロックセグメントは、1 つ以上の記事に属し、該 1 つ以上の記事は、該媒体資料内の複数のページにわたって連続し、かつ広がる、請求項 2 に記載の媒体資料アナライザ。

【請求項 4】

前記言語統計情報は、単語頻度情報を備え、前記言語統計アナライザは、言語資料全体に対する各ブロックセグメントにおける単語頻度と、一対の候補となるブロックセグメント間のコサイン距離類似性とに基づいて、該一対の候補となるブロックセグメントに対するマッチスコアを計算する、請求項 2 に記載の媒体資料アナライザ。

【請求項 5】

前記言語統計アナライザは、計算されたマッチスコアと、同一の記事に属するブロックセグメントの所定のポジティブな例と同一の記事に属さないブロックセグメントの所定のネガティブな例とを有するサンプルデータと、に基づいて、一対の候補となるブロックセグメントが前記資料媒体内の同一の記事に属する確率を決定する、請求項 2 に記載の媒体資料アナライザ。

【請求項 6】

前記言語統計アナライザは、訓練データセット内の記事の集合から、前記ポジティブなデータ例とネガティブなデータ例とを自動的に選択する、請求項 5 に記載の媒体資料アナライザ。

【請求項 7】

前記所定のポジティブなデータ例とネガティブなデータ例とは、前記媒体資料の画像からの光学文字認識を介して抽出されたテキストデータの表示から、ユーザによって、ユーザインタフェースにおいて選択される、請求項 5 に記載の媒体資料アナライザ。

【請求項 8】

レイアウトを有する媒体資料を表すデータを解析するコンピュータ実装された方法であって、

該媒体資料内のコラム状の本文テキストに関連するブロックセグメントを識別することと、

言語統計情報およびレイアウト情報に基づいて、該識別されたブロックセグメントのいずれが該媒体資料内の1つ以上の記事に属するかを決定することと

を包含し、

該決定することは、

候補となるブロックセグメントに対する言語統計を計算することと、

言語統計情報におけるオーバーラップに基づいて、比較されるブロックセグメントが同一の記事に属することを決定することと

を含む、方法。

【請求項 9】

前記識別されたブロックセグメントは、1つ以上の記事に属し、該1つ以上の記事は前記媒体資料内の複数のページにわたって連続し、かつ広がる、請求項 8 に記載の方法。

【請求項 10】

前記同一の記事内に存在することが決定された1つ以上のブロックセグメントからテキストを表示することをさらに含む、請求項 8 に記載の方法。

【請求項 11】

前記言語統計情報は、単語頻度情報を備え、前記計算することは、言語資料全体に対する各ブロックセグメントにおける単語頻度と、一对の候補となるブロックセグメント間のコサイン距離類似性とに基づいて、該一对の候補となるブロックセグメントに対するマッチスコアを計算することを含む、請求項 8 に記載の方法。

【請求項 12】

前記確率を決定するステップは、前記計算されたマッチスコアと、同一の記事に属するブロックセグメントの所定のポジティブな例と同一の記事に属さないブロックセグメントの所定のネガティブな例とを有するサンプルデータと、に基づいて、前記一对の候補となるブロックセグメントが前記資料媒体内の同一の記事に属する確率を決定することを含む、請求項 11 に記載の方法。

【請求項 13】

ユーザが、前記ポジティブなデータ例とネガティブなデータ例とを、前記媒体資料の画像からの光学文字認識を介して抽出されたテキストデータの表示から選択することを可能にすることをさらに含む、請求項 12 に記載の方法。

【請求項 14】

前記決定することは、言語統計情報における前記オーバーラップに基づいて決定された確率に基づいて、前記候補となるブロックセグメントが、前記媒体資料内の同一の記事に属するかどうかを識別することを含む、請求項 8 に記載の方法。

【請求項 15】

レイアウトを有する媒体資料を表すデータを解析する媒体資料アナライザであって、

該媒体資料内のコラム状の本文テキストと関連するブロックセグメントを識別するセグメントと、

言語統計情報およびレイアウト移行情報に基づいて、該識別されたブロックセグメントのいずれが該媒体資料内の1つ以上の記事に属するかを決定する記事コンポーザと

を備え、

該記事コンポーザは、レイアウト移行アナライザを含み、該レイアウト移行アナライザは、該セグメントによって出力された候補となるブロックセグメント内のレイアウト移行特徴を解析し、該候補となるブロックセグメントが該媒体資料内の同一の記事に属するか動かを決定し、

該レイアウト移行アナライザは、該レイアウト移行特徴を該候補となるブロックセグメントから計算し、該計算されたレイアウト移行特徴に基づいて、該候補となるブロックセグメントが該媒体資料内の同一の記事内に属するかどうかを決定するために、所定のレイアウト移行分類子を適用する、媒体資料アナライザ。

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/US 07/23233															
A. CLASSIFICATION OF SUBJECT MATTER IPC(8) - G06F 17/00 (2008.04) USPC - 715/255 According to International Patent Classification (IPC) or to both national classification and IPC																	
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) USPC: 715/255 Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched USPC: 715/200, 204, 243, 255, 264, 272 Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) Electronic Databases Searched: USPTO WEST (PGPUB,USPAT,USOCR,EPAB,JPAB); GOOGLE SCHOLAR; DIALOG PRO Search Terms Used: text/printed/media/digital articles/document extraction/metadata/segmentation comparison, multiple/successive column/page spanning, textual keyword/metadata/language/syntax/linguistic analysis, etc.																	
C. DOCUMENTS CONSIDERED TO BE RELEVANT <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 10%;">Category*</th> <th style="width: 70%;">Citation of document, with indication, where appropriate, of the relevant passages</th> <th style="width: 20%;">Relevant to claim No.</th> </tr> </thead> <tbody> <tr> <td style="text-align: center;">X</td> <td>US 2006/0080309 A1 (YACOUB et al.) 13 April 2006 (13.04.2006) Abstract, Para [0017]-[0141]</td> <td style="text-align: center;">1-18</td> </tr> <tr> <td style="text-align: center;">A</td> <td>US 2006/0184525 A1 (JONES et al.) 17 August 2006 (17.08.2006) Entire Document</td> <td style="text-align: center;">1-18</td> </tr> <tr> <td style="text-align: center;">A</td> <td>US 2004/0122811 A1 (PAGE) 24 June 2004 (24.06.2004) Entire Document</td> <td style="text-align: center;">1-18</td> </tr> <tr> <td style="text-align: center;">A</td> <td>US 2003/0229854 A1 (LEMAY) 11 December 2003 (11.12.2003) Entire Document</td> <td style="text-align: center;">1-18</td> </tr> </tbody> </table>			Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.	X	US 2006/0080309 A1 (YACOUB et al.) 13 April 2006 (13.04.2006) Abstract, Para [0017]-[0141]	1-18	A	US 2006/0184525 A1 (JONES et al.) 17 August 2006 (17.08.2006) Entire Document	1-18	A	US 2004/0122811 A1 (PAGE) 24 June 2004 (24.06.2004) Entire Document	1-18	A	US 2003/0229854 A1 (LEMAY) 11 December 2003 (11.12.2003) Entire Document	1-18
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.															
X	US 2006/0080309 A1 (YACOUB et al.) 13 April 2006 (13.04.2006) Abstract, Para [0017]-[0141]	1-18															
A	US 2006/0184525 A1 (JONES et al.) 17 August 2006 (17.08.2006) Entire Document	1-18															
A	US 2004/0122811 A1 (PAGE) 24 June 2004 (24.06.2004) Entire Document	1-18															
A	US 2003/0229854 A1 (LEMAY) 11 December 2003 (11.12.2003) Entire Document	1-18															
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/>																	
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family																	
Date of the actual completion of the international search 21 April 2008 (21.04.2008)		Date of mailing of the international search report <div style="text-align: center; font-size: 1.2em; font-weight: bold;">14 MAY 2008</div>															
Name and mailing address of the ISA/US Mail Stop PCT, Attn: ISA/US, Commissioner for Patents P.O. Box 1450, Alexandria, Virginia 22313-1450 Facsimile No. 571-273-3201		Authorized officer: <div style="text-align: center;">Lee W. Young</div> PCT Helpdesk: 571-272-4300 PCT OSP: 571-272-7774															

フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP(AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IS, IT, LT, LU, LV, MC, MT, NL, PL, PT, RO, SE, SI, SK, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RS, RU, SC, SD, SE, SG, SK, SL, SM, SV, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW

(72)発明者 スミス, レイ

アメリカ合衆国 カリフォルニア 94043, マウンテン ビュー, アンフィシアター パークウェイ 1600, グーグル インコーポレイテッド 気付

(72)発明者 ビンセント, ルーク

アメリカ合衆国 カリフォルニア 94043, マウンテン ビュー, アンフィシアター パークウェイ 1600, グーグル インコーポレイテッド 気付

(72)発明者 ブルームバーグ, ダン

アメリカ合衆国 カリフォルニア 94043, マウンテン ビュー, アンフィシアター パークウェイ 1600, グーグル インコーポレイテッド 気付

Fターム(参考) 5B029 AA01 AA02 BB02 CC26 CC27 CC29

5B075 NR03 NR12

5B109 NA03