

(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2012-506596  
(P2012-506596A)

(43) 公表日 平成24年3月15日(2012.3.15)

(51) Int.Cl.		F I		テーマコード (参考)
<b>G06F 17/30</b>	<b>(2006.01)</b>	G06F 17/30	350C	5B091
<b>G06F 17/28</b>	<b>(2006.01)</b>	G06F 17/28	C	
		G06F 17/30	170A	

審査請求 未請求 予備審査請求 未請求 (全 16 頁)

(21) 出願番号	特願2011-533276 (P2011-533276)	(71) 出願人	500046438 マイクロソフト コーポレーション アメリカ合衆国 ワシントン州 9805 2-6399 レッドモンド ワン マイ クロソフト ウエイ
(86) (22) 出願日	平成21年10月20日 (2009.10.20)	(74) 代理人	110001243 特許業務法人 谷・阿部特許事務所
(85) 翻訳文提出日	平成23年4月21日 (2011.4.21)	(72) 発明者	ラガベンドラ ウドゥパ ユー アメリカ合衆国 98052-6399 ワシントン州 レッドモンド ワン マイ クロソフト ウエイ マイクロソフト コ ーポレーション エルシーエーインター ナショナル パテント内
(86) 国際出願番号	PCT/US2009/061352		
(87) 国際公開番号	W02010/048204		
(87) 国際公開日	平成22年4月29日 (2010.4.29)		
(31) 優先権主張番号	12/255,372		
(32) 優先日	平成20年10月21日 (2008.10.21)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 コンパラブルコーパスを使用する固有表現の翻字

(57) 【要約】

第1の言語の文書、および第2の言語の追加文書を検査することができる。追加文書が、第1の言語の文書に十分類似しているかどうか判定することができる。追加文書が第1の言語の文書に十分類似していると判定される場合には、文書内の固有表現を選択することができる。この方法では、文書内の固有表現と追加文書内の単語とを比較し、文書内の固有表現と単語が十分類似しているかどうか判定することにより、類似している固有表現を調査することができる。文書内の固有表現に類似する単語が検出されると、文書内の固有表現および類似する固有表現を、固有表現の翻字として記憶することができる。

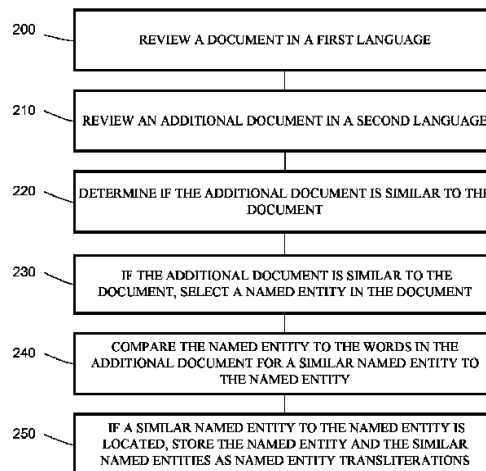


FIGURE 2

## 【特許請求の範囲】

## 【請求項 1】

多言語の固有表現の翻字を探索する方法であって、  
 第 1 の言語の文書を検査するステップと、  
 第 2 の言語の追加文書を検査するステップと、  
 前記追加文書が、前記文書に十分類似しているかどうか判定するステップと、  
 前記追加文書が、前記文書に十分類似している場合に、  
 前記文書内の固有表現を選択するステップと、  
 前記文書内の固有表現を前記追加文書内の単語と比較することを含む、十分類似して  
 いる固有表現を探索するステップと、  
 前記文書内の固有表現に十分類似している単語が検索される場合に、前記文書内の固有  
 表現および前記類似している単語を、固有表現の翻字として記憶するステップと  
 を含むことを特徴とする方法。

10

## 【請求項 2】

前記追加文書が前記文書に十分類似しているかを判定するステップが、言語横断の類似  
 度スコアを計算するステップをさらに含むことを特徴とする請求項 1 に記載の方法。

## 【請求項 3】

前記言語横断の類似度スコアを、Kullback-Leibler ダイバージェンス  
 を使用して計算することを特徴とする請求項 2 に記載の方法。

## 【請求項 4】

前記言語横断の類似度スコアを、複数の文書と追加文書のペアに対して計算することを  
 特徴とする請求項 2 に記載の方法。

20

## 【請求項 5】

最も高い類似度スコアを有する前記文書のペアを選択するステップをさらに含むことを  
 特徴とする請求項 4 に記載の方法。

## 【請求項 6】

類似している固有表現を探索するステップが、前記単語および前記文書内の固有表現に  
 対して言語横断の類似度スコアを計算するステップを含むことを特徴とする請求項 1 に記  
 載の方法。

## 【請求項 7】

前記言語横断の類似度スコアが、前記文書内の固有表現と前記単語の間の翻字の同等性  
 の度合いを測定することを特徴とする請求項 1 に記載の方法。

30

## 【請求項 8】

前記言語横断の類似度スコアが、複数の固有表現のペアに対して計算され、固有表現の  
 ペアが、前記文書内の固有表現および前記追加文書内の前記単語を含むことを特徴とする  
 請求項 7 に記載の方法。

## 【請求項 9】

前記追加文書内の前記単語が、前記追加文書内の単語のグループから順次選択され、前  
 記グループが、前記追加文書内の前置詞、動詞および形容詞を含まないことを特徴とする  
 請求項 8 に記載の方法。

40

## 【請求項 10】

最大の前記言語横断のスコアを有する前記固有表現のペアが、互いの翻字として選択さ  
 れることを特徴とする請求項 9 に記載の方法。

## 【請求項 11】

多言語の固有表現の翻字を探索するためのコンピュータ実行可能な命令を含むコンピュ  
 ータ記憶媒体であって、前記コンピュータ実行可能な命令が、  
 第 1 の言語の文書を検査する命令と、  
 第 2 の言語の追加文書を検査する命令と、  
 言語横断の類似度スコアを計算することにより、前記追加文書が、前記文書に十分類似  
 しているかどうか判定する命令と、

50

前記追加文書が、前記文書に十分類似している場合に、  
 前記文書内の固有表現を選択する命令と、  
 前記文書内の固有表現を前記追加文書内の単語と比較することを含む、十分に類似している固有表現を探索する命令と、  
 前記文書内の固有表現に十分類似している単語が検出される場合に、前記文書内の固有表現および前記類似している単語を、固有表現の翻字として記憶する命令と  
 を含むことを特徴とするコンピュータ記憶媒体。

【請求項 12】

前記言語横断の類似度スコアを、Kullback-Leiblerダイバージェンスを使用して計算することを特徴とする請求項 11 に記載のコンピュータ記憶媒体。

10

【請求項 13】

前記言語横断の類似度スコアを複数の文書と追加文書のペアに対して計算して、最も高い類似度スコアを有する前記文書のペアを選択することを特徴とする請求項 12 に記載のコンピュータ記憶媒体。

【請求項 14】

類似している固有表現を探索する命令が、前記単語および前記文書内の固有表現に対して言語横断の類似度スコアを計算する命令を含み、前記言語横断の類似度スコアが、前記文書内の固有表現と前記単語の間の翻字の同等性の度合いを測定することを特徴とする請求項 11 に記載のコンピュータ記憶媒体。

【請求項 15】

前記言語横断の類似度スコアが、複数の固有表現のペアに対して計算され、固有表現のペアが、前記文書内の固有表現および前記追加文書内の前記単語を含むことを特徴とする請求項 14 に記載のコンピュータ記憶媒体。

20

【請求項 16】

前記追加文書内の前記単語が、前記追加文書内の単語のグループから順次選択され、前記グループが、前記追加文書内の前置詞、動詞および形容詞を含まないことを特徴とする請求項 15 に記載のコンピュータ記憶媒体。

【請求項 17】

最大の前記言語横断のスコアを有する前記固有表現のペアが、互いの翻字として選択されることを特徴とする請求項 16 に記載のコンピュータ記憶媒体。

30

【請求項 18】

多言語の固有表現の翻字を探索するためのコンピュータ実行可能な命令を実行するプロセッサ、前記プロセッサと通信するメモリ、および入力/出力回路を備えるコンピュータ・システムであって、前記コンピュータ実行可能な命令が、

第 1 の言語の文書を検査する命令と、

第 2 の言語の追加文書を検査する命令と、

言語横断の類似度スコアを計算することにより、前記追加文書が、前記文書に十分類似しているかどうか判定する命令であって、前記言語横断の類似度スコアを、Kullback-Leiblerダイバージェンスを使用して計算する命令と、

40

前記追加文書が、前記文書に十分類似している場合に、

前記文書内の固有表現を選択する命令と、

前記文書内の固有表現を前記追加文書内の単語と比較することを含む、十分に類似している固有表現を探索する命令と、

前記文書内の固有表現に十分類似している単語が検出される場合に、前記文書内の固有表現および前記類似している単語を、固有表現の翻字として記憶する命令と

を含むことを特徴とするコンピュータ・システム。

【請求項 19】

前記言語横断の類似度スコアが、複数の文書と追加文書のペアに対して計算され、最も高い類似度スコアを有する前記文書のペアを選択することを特徴とする請求項 18 に記載

50

のコンピュータ・システム。

【請求項 20】

類似している固有表現を探索する命令が、前記単語および前記文書内の固有表現に対して言語横断の類似度スコアを計算する命令を含み、

前記言語横断の類似度スコアが、前記文書内の固有表現と前記単語の間の翻字の同等性の度合いを測定し、

前記追加文書内の前記単語が、前記追加文書内の単語のグループから順次選択され、前記グループが、前記追加文書内の前置詞、動詞および形容詞を含まず、

前記言語横断の類似度スコアが、複数の固有表現のペアに対して計算され、固有表現のペアが、前記文書内の固有表現および前記追加文書内の前記単語を含むことを特徴とする請求項 18 に記載のコンピュータ・システム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンパラブルコーパスを使用する固有表現の翻字に関する。

【背景技術】

【0002】

この背景は、本特許出願の基本的な事実関係を提供するものであり、解決すべき具体的な問題を説明するものではない。

【0003】

20

固有表現は、言語横断情報検索 (CLIR) における質問語の重要な部分を形成し、言語横断情報検索システムの性能に大きな影響を及ぼす。機械翻訳 (MT) では、未知語の多くが固有表現である。しかし、対訳辞書では十分な範囲の固有表現に乏しく、機械翻訳システムは、しばしば間違った翻字を生成する。

【発明の概要】

【0004】

この概要は、発明を実施するための形態において以下でさらに説明する選ばれた概念を、簡略化した形態で紹介するために提供される。この概要は、特許請求される主題の重要な特徴または本質的な特徴を特定するものではなく、特許請求される主題の範囲を限定するために使用するものでもない。

30

【0005】

多言語の固有表現の翻字を探索するための方法およびシステムが開示されている。第 1 の言語の文書を検査し、そして第 2 の言語の追加文書を検査する。追加文書が、第 1 の言語の文書に十分類似しているかどうか判定することができる。追加文書が第 1 の言語の文書に十分類似していると判定される場合には、この文書内の固有表現を選択することができる。この方法では、固有表現を追加文書内の単語と比較し、この固有表現と単語が十分類似しているかどうか判定することにより、類似している固有表現を検索することができる。固有表現に類似している単語が検出されると、この固有表現および類似している固有表現を、固有表現の翻字として記憶することができる。

【図面の簡単な説明】

40

【0006】

【図 1】携帯用コンピューティング装置の図である。

【図 2】多言語の固有表現の翻字を探索する方法を示す図である。

【図 3】コンパラブルコーパスを示す図である。

【発明を実施するための形態】

【0007】

以下の説明では、数多くの様々な実施形態の詳細な説明を行っているが、この説明の法的な範囲は、この特許の最後で述べる特許請求の範囲の説明によって規定されるものと理解すべきである。実現可能な実施形態をすべて説明することは、不可能ではないとしても実際的ではないので、詳細な説明は例示的なものに過ぎないと解釈すべきであり、必ずし

50

もあらゆる実現可能な実施形態を説明するものではない。現行技術、またはこの発明の出願日以降に開発される技術のいずれかを使用して、数多くの代替実施形態を実施することもできるが、こうした技術もやはり、特許請求の範囲に記載の範囲の中にある。

【0008】

「本明細書では、用語「\_\_」は...を意味するものとしてここに定義する」といった文、または同様の文を使用して、この特許において用語が明示的に定義されない場合には、明示的にせよ暗示的にせよ、そのありのままの意味または通常の意味を超えて、その用語の意味を限定する意図はなく、こうした用語は、(特許請求の範囲の言葉以外の)この発明の任意のセクションでなされる任意の記述に基づいて範囲を限定的されるものと解釈すべきではないことも理解されたい。この特許の最後にある特許請求の範囲に記載された任意の用語が、この特許において単一の意味と合致するように参照される限りにおいて、それは、読み手を混乱させないよう明確にするためになされるものであり、こうした特許請求の範囲の用語が、暗にまたは他の方法で、その単一の意味に限定されるものではない。最後に、請求項の要素が、任意の構造を列挙することなく、用語「手段」と機能を列挙することによって定義されない限り、任意の請求項の要素の範囲は、米国特許法第112条第6項の適用に基づいて解釈されるものではない。

10

【0009】

図1に、本明細書で説明するユーザ・インターフェースを表示し提供するように動作することのできる、適切なコンピューティング・システム環境100の一例を示す。コンピューティング・システム環境100は、適切なコンピューティング環境のほんの一例に過ぎず、特許請求の範囲に記載の方法および装置の、使用法または機能の範囲に関して、何ら制限を示唆するものではないことに留意されたい。コンピューティング環境100は、例示的な動作環境100に示した構成要素のうち任意の1つの構成要素もしくはそれらの組合せに関して、いかなる依存性または要求をも、有すると解釈すべきではない。

20

【0010】

図1を参照すると、特許請求された方法および装置の各ブロックを実施するための例示的なシステムは、コンピュータ110の形態の汎用コンピューティング装置を備える。コンピュータ110の構成部品には、それだけには限らないが、処理ユニット120、システム・メモリ130、および、システム・メモリを含む様々なシステム構成部品を処理ユニット120に結合するシステム・バス121が含まれ得る。

30

【0011】

コンピュータ110は、モデム172または他のネットワーク・インターフェース170を用いて、ローカル・エリア・ネットワーク(LAN)171および/またはワイド・エリア・ネットワーク(WAN)173を介して、リモート・コンピュータ180など、1つまたは複数のリモート・コンピュータとの論理的な接続を使用して、ネットワーク環境で動作することができる。

【0012】

コンピュータ110は、通常、コンピュータ110がアクセスすることのできる任意の使用可能な媒体でよい様々なコンピュータ読取り可能な媒体を備え、揮発性および不揮発性の媒体、取外し可能および取外し不可能な媒体の両方を含む。システム・メモリ130には、読取り専用メモリ(ROM)131およびランダム・アクセス・メモリ(RAM)132など、揮発性および/または不揮発性のメモリの形態のコンピュータ記憶媒体が含まれる。ROMには、基本入出力システム(BIOS)133が含まれ得る。RAM132は、通常、オペレーティング・システム134、アプリケーション・プログラム135、他のプログラム・モジュール136、およびプログラム・データ137を含む、データおよび/またはプログラム・モジュールを含む。コンピュータ110にはまた、ハード・ディスク・ドライブ141、磁気ディスク152との間で読み書きする磁気ディスク・ドライブ151、光ディスク156との間で読み書きする光ディスク・ドライブ155など、他の取外し可能/取外し不可能な、揮発性/不揮発性コンピュータ記憶媒体が含まれ得る。ハード・ディスク・ドライブ141、151、および155は、インターフェース1

40

50

40、150を介して、システム・バス121とインターフェースすることができる。

【0013】

ユーザは、キーボード162や、マウス、トラックボールまたはタッチ・パッドと普通呼ばれているポインティング装置161などの入力装置を用いて、コンピュータ20にコマンドおよび情報を入力することができる。他の入力装置(図示せず)は、マイクロホン、ジョイスティック、ゲーム・パッド、衛星用パラボラ・アンテナ、スキャナまたは同様のものを含んでもよい。上記その他の入力装置は、システム・バスに結合されたユーザ入力インターフェース160を介して、処理ユニット120にしばしば接続されているが、パラレル・ポート、ゲーム・ポート、またはユニバーサル・シリアル・バス(USB)など、他のインターフェースおよびバス構造によって接続されてもよい。ビデオ・インターフェース190などのインターフェースを介して、モニタ191または他のタイプの表示装置もシステム・バス121に接続してもよい。モニタに加えて、コンピュータにはまた、出力周辺インターフェース190を介して接続してもよいスピーカ197やプリンタ196など、他の周辺出力装置が含まれ得る。

10

20

【0014】

図2は、多言語の固有表現の翻字のために文書を探査する方法を示すことができる。固有表現(NE)は、情報検索(IR)システムにおける質問語の重要な部分を形成し、その性能に大きな影響を及ぼすことがある。これら固有表現は、言語横断情報検索(CLIR)において、さらに重要になることがある。さらに、未知語の多くが実際には固有表現なので、固有表現は、機械翻訳(MT)システムの性能においても重要な役割を演じる。固有表現は、言語横断情報検索システムがうまく働くかどうかにとって重要であり、機械翻訳の性能に著しい影響を及ぼすが、対訳辞書は、固有表現を十分に網羅していないので、それが手作りであれ統計に基づくものであれ、限られたサポートしか提供しない。ニュース記事およびインターネットにより、毎日、新規の固有表現が言語の語彙に取り入れられる。機械翻字の代替手法では、スペルミスまたは間違っただけの翻字がしばしば生じるが、それらは、言語横断情報検索において雑音として働き、MTにおける翻訳品質を劣化させる。

【0015】

最近では、多言語で同時にニュース記事が利用可能になることにより、こうしたニュースのコーパスから、固有表現の翻字、特に固有表現の翻字の相当語句(NETE)を探査することへの見込みのある選択肢に関心が集まってきた。形式的には、ニュースのコンパラブルコーパスとは、適度に長い期間にわたる、一対の言語の時間通りに並んだニュース記事である。世界中の多くの報道機関が、毎日、こうしたニュース・コンテンツを多言語で作成する。ニュースのコンパラブルコーパスから探索された固有表現の翻字の相当語句は、対訳辞書および機械翻字システムを効果的に補完するのに、CLIRおよびMTを含む多くのタスクにおいて貴重なものになり得る。

30

【0016】

ニュース記事は、通常、固有表現310を多く含み、したがって、ニュースのコンパラブルコーパスには、固有表現の翻字の相当語句が豊富に存在する。多くの世界中の言語でのニュース・コーパスの膨大な量と永久の可用性により、こうしたコーパスから固有表現310の相当語句を探査する効果的な方式を考え出すことができる場合には、固有表現310の相当語句を探査するための、膨大で貴重で尽きることのないデータ・ソースが示される。この可能性が、説明する方法およびシステムの推進要因の1つである。

40

【0017】

大規模なコンパラブルコーパスから固有表現の翻字の相当語句を効果的に探索するには、いくつかの課題がある。第1に、固有表現の識別および検証には、多くの言語で利用可能でない言語ツールおよび言語学的資源が必要となることがある。第2に、コンパラブルコーパスにおける固有表現の大部分が散在しており、コーパス内の固有表現の署名の頻度にはほとんど依存する必要がない。第3に、探索する方法は、より大きいコーパスを探査するときに効果的に計算できるように、候補の生成を制限しなければならない。さらに、

50

候補を制限は、誤判定による精度の劣化が減少することになる。最後に、各言語にわたって探索を効果的にするために、言語特有の知識をほとんど使用しないことが重要である。

【0018】

本出願は、大規模なコンパラブルコーパスから、固有表現の翻字の相当語句を効果的に探索するための、MINTと呼ばれる新規な方法を紹介し、上記に掲載されたあらゆる課題に取り組む。MINTは、ただ1つの言語について固有表現認識装置(NER)の可用性を想定しており、したがって、固有表現認識装置が利用可能である言語とペアを組むときに、リソースが乏しい言語からさえ固有表現の翻字の相当語句を探索するのに適用可能である。それに応じて、本出願では以下のことを行う。

【0019】

類似しているコンテンツを有するニュース記事が先験的に知られている場合、これらのニュース記事を、効果的かつ徹底的に探索することができることを認識する。

【0020】

標準のコンパラブルコーパスの場合など、類似している記事が先験的に知られていない場合も、MINTが、上記洞察および言語横断情報検索技法を使用して、最新技術よりもはるかに良好に探索できることを実験的に証明する。

【0021】

様々な特性を有する多くのコーパスにわたって、その有効性を証明する。

【0022】

最後に、本出願は、3つの異なる言語ファミリー(すなわち、スラブ語、インドヨーロッパ語およびドラビダ語)からの互いに異なる言語セット(すなわち、ロシア語、ヒンディー語、カンナダ語およびタミル語)間でのいくつかのコンパラブルコーパス上でその性能を説明することにより、この方法が言語に依存しないことを証明する。

【0023】

MINT法は、ニュースが、人々、場所、組織、および他の固有表現を含む出来事に関するものなので、類似しているコンテンツを有する多言語のニュース記事が、非常にオーバーラップする固有表現のセットを含んでいるはずである、という重要な考えに基づく。同じニュースの出来事を報告する多言語のニュース記事は、それぞれの言語の固有表現に言及するはずであり、したがって、固有表現の翻字の相当語句を豊富に生成することが予想される。図3には、BBCによって公開された、ヒンディー語および英語での一対の類似している記事における固有表現の翻字の相当語句を示してある。1つのソースによって公開された、英語とタミル語での200対の類似しているニュース記事を分析すると、英語側の単一語の固有表現310のうち87%において、タミル語側で少なくとも1つの相当語句が存在したことが分かっている。MINT法は、この経験に裏打ちされた考えを用いて、こうしたコーパスから固有表現の翻字の相当語句を探索する。

【0024】

MINTは2つの段階を有することができる。第1の段階で、各文書と比較して、ソース側のあらゆる文書について、類似しているコンテンツを有するターゲット側の文書のセットを識別する。類似している文書が識別されると、それらの文書を入力として第2の段階に与え、そこで、それらの文書から固有表現の翻字の相当語句を探索する。

【0025】

再び図2を参照すると、ブロック200で、第1の言語の文書300(図3)を検査することができる。理想的には、この文書は、対象となるいくつかの固有表現の翻字の相当語句を含むものとする。この文書は、問題となる固有表現の翻字の相当語句を含むときに選択してもよく、または、特定の日向けに書かれたニュース記事を用いた逐次探索でもよい。もちろん、第1の言語の文書を選択する他の方法も可能であり、またそのように意図されている。

【0026】

ブロック210で、第2の言語の追加文書305を検査することができる。理想的には、追加文書305も、固有表現の翻字の相当語句を有するように選択することができる。

10

20

30

40

50

たとえば、第1の言語の文書がスポーツ記事である場合、第2の言語の追加文書として学術論文を検査することは、ほとんど意味がない。文書300と追加文書305の間で、類似している固有表現の翻字の相当語句が存在する確率が低いからである。

【0027】

ブロック220で、追加文書305が文書300に類似しているかどうか判定することができる。この判定は、様々な方式で行うことができる。実施形態によっては、言語横断の文書類似度モデルを使用して、言語横断の類似度スコアを計算する。言語横断の文書類似度モデルにより、ソース言語とターゲット言語の一对の文書間における類似度を測定することができる。文書の確率分布と追加文書の確率分布との間の負のKullback-Leibler (KL) ダイバージェンスを、類似度測定値として使用してもよい。

10

【0028】

確率論および情報理論では、Kullback-Leiblerダイバージェンス(また、情報ダイバージェンス、情報利得(information gain)、または相対エントロピー)は、2つの確率分布PとQの間の差の非可換尺度である。KLは、Pに基づく符号を使用するとき、およびQに基づく符号を使用するとき、Pからのサンプルを符号化するのに必要となるビットの数の予想される差を測定する。通常Pは、データの「真」の分布、観測値、または精密に計算された理論分布を表す。測定値Qは通常、Pの理論、モデル、記述、または近似を表す。

【0029】

本出願では、ソース言語およびターゲット言語の語彙を示す $V_s$ 、 $V_t$ を有する、それぞれソース言語およびターゲット言語の2つの文書 $D_S$ 300、 $D_T$ 305が与えられている場合、これら2つの文書300と305の間の類似度は、 $KL(D_S, D_T)$ で与えられ得る。

20

【0030】

$$\sum_{w_T \in V_T} P(w_T | D_S) \log \frac{p(w_T | D_T)}{p(w_T | D_S)}$$

【0031】

ここで、 $p(w | D)$ は、単語 $w$ が文書 $D$ 内に存在する確率である。所与のソース言語の文書300に類似しているターゲットの文書305を見つけることが関心であるので、分子は、ターゲット言語の文書と無関係であるときには無視してもよい。最後に、以下のように $p(w_t | D_s)$ を展開する。

30

【0032】

$$\sum_{w_s \in V_s} p(w_s | D_s) p(w_T | w_s)$$

【0033】

言語横断の類似度スコアは、以下のように指定することができる。

【0034】

40

$$\text{CrossLanguageDocumentSimilarity}(D_S, D_T, MD) =$$

【0035】

$$\sum_{w_T \in V_T} \sum_{w_S \in V_S} p(w_S | D_S) p(w_T | w_S) \log p(w_T | D_T)$$

【0036】

擬似コードでは、文書300、305の比較を、以下のように進めることができる。

【0037】

Input: Comparable news corpora ( $C_S, C_T$ ) in languages (S,T)  
 Crosslanguage Document Similarity Model MD for (S,T)  
 Threshold score  $\alpha$ .

Output: Set  $A_{S,T}$  of pairs of similar articles ( $D_S, D_T$ ) from ( $C_S, C_T$ ).

```

1  $A_{ST} \leftarrow \emptyset$ ; //Set of Similar articles ( $D_S, D_T$ )
2 for each article  $D_S$  in  $C_S$  do
3    $X_S \leftarrow \emptyset$ ; //Set of candidates for  $D_S$ .
4   for each article  $d_T$  in  $C_T$  do
5     score = CrossLanguageDocumentSimilarity( $D_S, d_T, MD$ );
6     if (score  $\geq \alpha$ ) then  $X_S \leftarrow X_S \cup (d_T, \text{score})$ ;
7   end
8    $D_T = \text{BestScoringCandidate}(X_S)$ ;
9 if ( $D_T \neq \emptyset$ ) then  $A_{ST} \leftarrow A_{ST} \cup (D_S, D_T)$ ;
10 end

```

10

【0038】

擬似コードから分かるように、複数の追加文書305を文書300と比較することができる。実施形態によっては、追加文書305は事前選別して、固有表現310に類似している単語315を有する可能性のある追加文書305のみを検査することを確実にする。例として、Michael Phelpsに着目した文書300は、おそらくスポーツに関連したものになる。この知識を使用して、検査される追加文書305のタイプを減らしてもよい。各追加文書/文書のペア320(元のソース文書300および個別の各追加文書305)に対して類似度スコアを計算することができ、最も高い類似度を有するペア320を、さらに分析すべき文書/追加文書のペア320として使用することができる。

20

【0039】

ブロック230で、追加文書305が文書300に十分類似している場合、文書内の固有表現310を選択することができる。たとえば、水泳選手Michael Phelpsは、アメリカ人の名前であり、多くのスポーツ記事で認識が容易である可能性がある。しかし、Michael Phelpsは、他の言語では作成が難しいことがある。したがって、Michael Phelpsは、この方法が探索しようとし得る固有表現310の一例であってよい。

30

【0040】

複数の追加文書305が、文書300と比較された場合、最も高いと判定された類似度を有する文書300/追加文書305のペア320を、さらに分析すべき選択されたペア320として選ぶことができる。いずれのペア320も十分な類似度スコアに達しない場合、追加文書305を戻さなくてもよく、この方法は、終了してもよく、また新規の文書に関して再び開始してもよい。

【0041】

ブロック240で、固有表現310に類似している単語315を得るために、固有表現310を追加文書305内の単語315と比較することができる。想定できるように、単語は、句でも、文章の断片でも、実体名でもよい。この方法は、コレクション $A_{s,t}$ 内の記事( $D_s, D_t$ )の各ペア上で働くことができ、固有表現の翻字の相当語句のセット $P_{s,t}$ を生成する。 $P_{s,t}$ 内の各ペア( $e_s, e_t$ )は、言語Sの固有表現 $s_{310}$ および言語Tのトークン $e_t_{315}$ から構成され、これらは互いの翻字の相当語句である。さらに、翻字の類似度モデルMTによって測定される、 $s_{310}$ と $e_t_{315}$ の間の翻字の類似度は、少なくとも0とすることができる。

40

【0042】

擬似コードでは、この方法の一実施形態を以下のように進めることができる。

【0043】

Input:

Set  $A_{ST}$  of similar documents ( $D_S, D_T$ ) in languages (S,T)  
 Transliteration Similarity Model MT for (S,T)  
 Threshold score  $\beta$ .

Output: Set  $P_{S,T}$  of NETEs ( $s_S, s_T$ ) from  $A_{ST}$ ;

```

1  $P_{ST} \leftarrow \emptyset$ ;
2 for each pair of articles ( $D_S, D_T$ ) in  $A_{ST}$  do
3   for each named entity  $s_T$  in  $D_S$  do
4      $Y_S \leftarrow \emptyset$ ; // Set of candidates for  $s_S$ 
5     for each candidate  $e_T$  in  $D_T$  do
6       score = TransliterationSimilarity( $s_S, e_T, MT$ );
7       if (score  $\geq \beta$ ) then  $Y_S \leftarrow Y_S \cup (e_T, \text{score})$ ;
8     end
9      $s_T = \text{BestScoringCandidate}(Y_S)$ ;
10    if ( $s_T \neq \text{null}$ ) then  $P_{ST} \leftarrow P_{ST} \cup (s_S, e_T)$ ;
11  end
12 end

```

10

【 0 0 4 4 】

翻字の類似度モデルは、ソースの固有表現 3 1 0 とターゲット言語の単語 3 1 5 の間の翻字の同等性の度合いを測定する。以下のように、翻字の類似度モデル MT として、ロジスティック関数を利用してもよい。

20

【 0 0 4 5 】

$$\text{TransliterationSimilarity}(e_S, e_T, MT) = \frac{1}{1 + e^{-w \cdot \phi(e_S, e_T)}}$$

【 0 0 4 6 】

ここで、 $(e_S, e_T)$  は、ペア  $(e_S, e_T)$  についての特徴ベクトルであり、 $w$  は、重みベクトルである。翻字の類似度は、 $[0 \dots 1]$  の範囲の値をとるとし得る。このモデルで利用されている特徴は、ある種の文字配列の出現、 $e_S$  と  $e_T$  のサブストリングの結合、文字の配列の単調性、および 2 つのストリング内の文字数の差など、 $(e_S, e_T)$  において観察される関心を引く言語横断の関連性を取り込んでもよい。重みベクトル  $w$  は、既知の翻字の相当語句のトレーニング・コーパスにわたって、識別しながら学習される。もちろん、固有表現 3 1 0 と単語 3 1 5 の類似度を判定する他の方式が可能でもよく、また企図されている。

30

【 0 0 4 7 】

実施形態によっては、追加文書 3 0 5 内のすべての単語 3 1 5 を、文書 3 0 0 の固有表現 3 1 0 と比較する。他の実施形態では、追加文書 3 0 5 を詳しく調べて、多くの用語を分析から排除する。たとえば、「the」、「a」、「an」など英語の冠詞は、固有表現 3 1 0 の一部である可能性は非常に低いので、これらの単語は分析しなくてもよい。さらに、固有表現 3 1 0 内に動詞が存在する可能性は低いので、動詞は分析しなくてもよい。さらに他の例として、形容詞が固有表現 3 1 0 の一部である可能性は低いので、形容詞も分析しなくてもよい。追加文書 3 0 5 内の単語 3 1 5 をさらにふるい分けることが可能であり、また企図されている。その結果、固有表現 3 1 0 と比較される、追加文書 3 0 5 内の単語 3 1 5 の数は、極めて少ない可能性があり、まさに対象とすべきである。

40

【 0 0 4 8 】

ブロック 2 5 0 で、固有表現 3 1 0 に類似している単語 3 1 5 が検出される場合、固有表現 3 1 0 および類似している単語 3 1 5 を、固有表現の翻字として格納することができる。複数の単語 3 1 5 が固有表現と比較された場合、最も高いと判定された類似度を有する単語 3 1 5 / 固有表現 3 1 0 のペアを、固有表現の翻字として選択することができる。

50

いずれのペアも十分な類似度スコアに達しない場合、いずれの単語 3 1 5 も、固有表現 3 1 0 の翻字として戻さないとし得る。

【0049】

次いで、この翻字を様々な目的に使用してもよい。一実施形態では、翻訳ソフトウェアがこの翻字を使用して、翻訳を改善してもよい。他の実施形態では、この翻字を探索ソフトウェアで使用して、複数の言語の関連する結果を探索する助けとし得る。もちろん、他の使用法も可能であり、またそのように意図されている。

【0050】

前述の説明では、数多くの様々な実施形態の詳細な説明を行っているが、この特許の範囲は、この特許の最後で述べる特許請求の範囲の説明によって規定されるものと理解すべきである。あらゆる実現可能な実施形態を説明することは、不可能ではないとしても、実際的ではないはずなので、詳細な説明は例示的なものに過ぎないと解釈すべきであり、あらゆる実現可能な実施形態を説明するものではない。現行技術、またはこの発明の出願日以降に開発される技術のいずれかを使用して、数多くの代替実施形態を実施することもできるが、こうした技術もやはり、特許請求の範囲に記載の範囲内に収まるものとする。

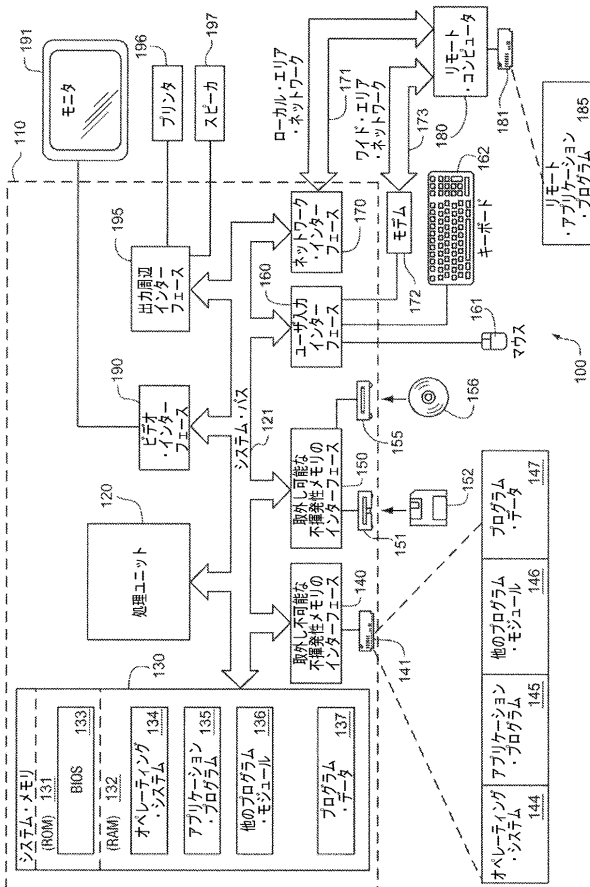
【0051】

このように、この特許請求の範囲に記載の精神および範囲から逸脱することなく、本明細書において説明し図示した技法および構造物において、多くの修正形態および変形形態を実施してもよい。したがって、本明細書において記載された方法および装置は、例示的なものに過ぎず、特許請求の範囲に記載の範囲を限定するものではないことを理解されたい。

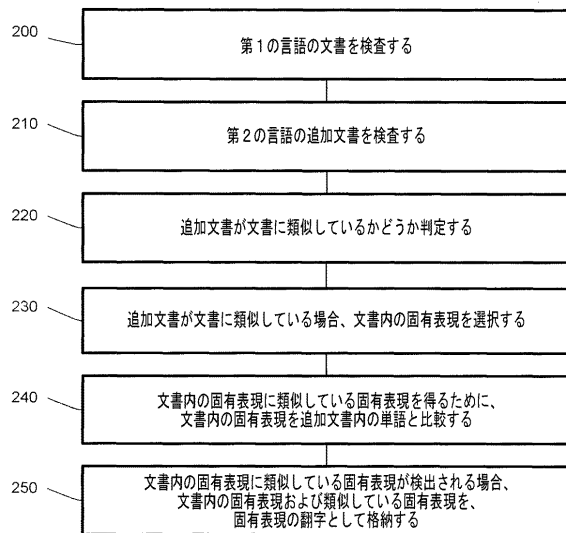
10

20

【図 1】



【図 2】



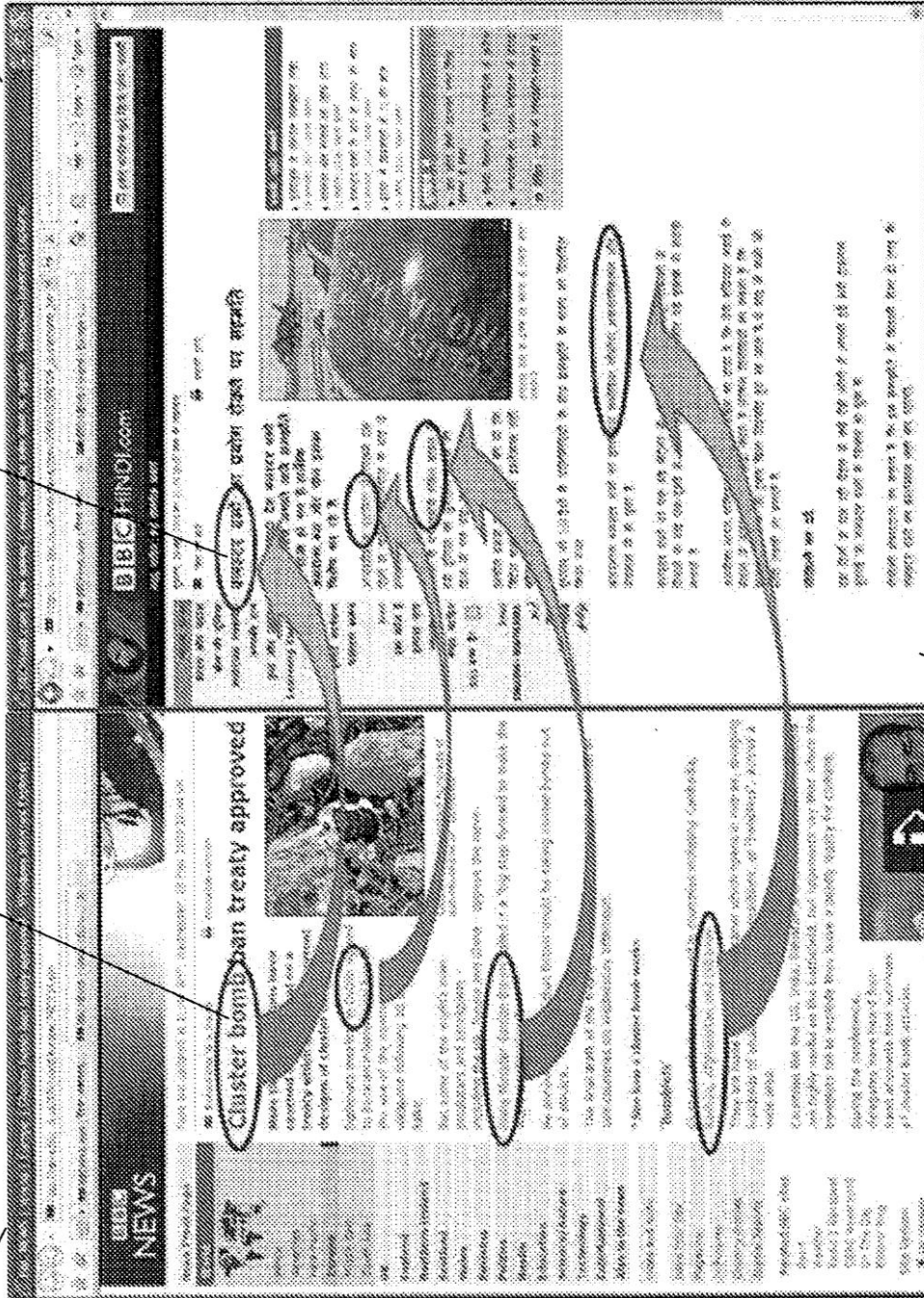
【 3 3 】

305

315



310

300



320

## 【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. <b>PCT/US2009/061352</b>
<b>A. CLASSIFICATION OF SUBJECT MATTER</b>		
<i>G06F 17/28(2006.01)i, G06F 17/21(2006.01)i</i>		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b>		
Minimum documentation searched (classification system followed by classification symbols) G06F 17/28		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Korean utility models and applications for utility models Japanese utility models and applications for utility models (Chinese Patents and application for patent)		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) eKOMPASS(KIPO internal) & Keywords: entity and document and (translate or interpret) and (pronunciation or phonetic or phoneme)		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	AHMED HASSAN et al., "IMPROVING NAMED ENTITY TRANSLATION BY EXPLOITING COMPARABLE AND PARALLEL CORPORA", Conference on Recent Advances in Natural Language Processing (RANLP 2007), ANML Workshop, 2007, pp.1-6.  See abstract; chapter 3; figs.2-3.  ( <a href="http://www.computing.dcu.ie/~hhasan/ranlp07.pdf">http://www.computing.dcu.ie/~hhasan/ranlp07.pdf</a> )	1-2,4-5,11,13 ,18-19
A	RICHARD SPROAT et al., "NAMED ENTITY TRANSLITERATION WITH COMPARABLE CORPORA", Proceeding of the 21th International Conference on Computational Linguistics and 44th Annual Meeting of the ACL, July 2006, pp.73-80.  See abstract; chapter 3 (3.1-3.2).  ( <a href="http://www.mt-archive.info/Coling-ACL-2006-Sproat.pdf">http://www.mt-archive.info/Coling-ACL-2006-Sproat.pdf</a> )	1-20
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* 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 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 03 JUNE 2010 (03.06.2010)		Date of mailing of the international search report <b>03 JUNE 2010 (03.06.2010)</b>
Name and mailing address of the ISA/KR  Korean Intellectual Property Office Government Complex-Daejeon, 139 Seonsa-ro, Seo-gu, Daejeon 302-701, Republic of Korea Facsimile No. 82-42-472-7140		Authorized officer KIM, SANG CHEOL Telephone No. 82-42-481-8521 

**INTERNATIONAL SEARCH REPORT**

International application No.

**PCT/US2009/061352**

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>DAN GOLDWASSER et al., "ACTIVE SAMPLE SELECTION FOR NAMED ENTITY TRANSLITERATION",            Proceeding of ACL-08: HLT, Short papers, June 2008, pp.53-56.</p> <p>See abstract; chapter 2; fig.2.</p> <p>(<a href="http://www.aclweb.org/anthology/P/P08/P08-2014.pdf">http://www.aclweb.org/anthology/P/P08/P08-2014.pdf</a>)</p>	1-20

**INTERNATIONAL SEARCH REPORT**  
Information on patent family members

International application No.  
**PCT/US2009/061352**

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
None			

## 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, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CL, 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, ST, SV, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW

(72)発明者 サラバナン クリシュナン

アメリカ合衆国 98052-6399 ワシントン州 レッドモンド ワン マイクロソフト  
ウェイ マイクロソフト コーポレーション エルシーイー - インターナショナル パテント内

(72)発明者 アルムガン クマラン

アメリカ合衆国 98052-6399 ワシントン州 レッドモンド ワン マイクロソフト  
ウェイ マイクロソフト コーポレーション エルシーイー - インターナショナル パテント内

Fターム(参考) 5B091 AA01 BA02 CC15 EA02