

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5362353号
(P5362353)

(45) 発行日 平成25年12月11日 (2013.12.11)

(24) 登録日 平成25年9月13日 (2013.9.13)

(51) Int.Cl.

F I

G 0 6 F 17/21 (2006.01)

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/21 5 5 0 J

G 0 6 F 17/21 5 7 0 N

G 0 6 F 17/30 1 7 0 A

G 0 6 F 17/30 3 3 0 C

請求項の数 6 (全 12 頁)

(21) 出願番号 特願2008-520339 (P2008-520339)
 (86) (22) 出願日 平成18年6月30日 (2006.6.30)
 (65) 公表番号 特表2009-500754 (P2009-500754A)
 (43) 公表日 平成21年1月8日 (2009.1.8)
 (86) 国際出願番号 PCT/US2006/026012
 (87) 国際公開番号 W02007/008492
 (87) 国際公開日 平成19年1月18日 (2007.1.18)
 審査請求日 平成21年6月19日 (2009.6.19)
 (31) 優先権主張番号 11/177, 136
 (32) 優先日 平成17年7月8日 (2005.7.8)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 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 ウェイ
 (74) 代理人 100107766
 弁理士 伊東 忠重
 (74) 代理人 100070150
 弁理士 伊東 忠彦
 (74) 代理人 100091214
 弁理士 大貫 進介

最終頁に続く

(54) 【発明の名称】 文書中のコロケーション誤りを処理すること

(57) 【特許請求の範囲】

【請求項 1】

文書中のコロケーション誤りを訂正するためのコンピュータにより実施される方法であって、

前記文書中のテキストの文にアクセスするステップと、

前記テキストの前記文において1または複数の検査語を検出するステップであって、前記検査語は、コロケーション誤りを生じさせる可能性がある語であって、特定のコロケーション誤りタイプを示すこと、

検出された検査語に基づいて前記コロケーション誤りタイプを識別するステップと、

前記検出された検査語により示された特定のコロケーション誤りタイプに基づいてクエリを生成するステップであって、メモリ内の複数の異なる組のクエリ生成ルールから一組のルールを選択し、前記検査語が含まれた特定の組のクエリを生成するステップと、

前記特定の組のクエリを、検索モジュールにサブミットして検索結果を得るステップであって、前記検索結果は、ウェブ検索エンジンを用いて抽出された文書のテキストのサマリーを含むこと、

前記検索結果のテキストと前記特定の組のクエリのテキストとを比較して、1または複数のクエリのテキストが一致した場合には、コロケーション誤りが存在しないと判定し、全てのクエリのテキストが一致しなかった場合には、コロケーション誤りがあると判定するステップと、

コロケーション誤りがあると判定した場合に、前記コロケーション誤りタイプに基づい

10

20

て、前記検査語に対応する語をプレースホルダーに置換したクエリテンプレートを生成するステップと、

前記クエリテンプレートを、前記検索モジュールにサブミットして文字列を得るステップと、

前記クエリテンプレートと一致する文字列であって、前記プレースホルダーが任意の語と一致する文字列を識別するステップと、

前記識別された文字列を表示して、前記検出されたコロケーション誤りに対する訂正の選択肢を提供するステップと

を備えたことを特徴とする方法。

【請求項 2】

前記検査語を検出するステップは、前記文を構文解析して、前記文に含まれる品詞を識別するステップをさらに含み、

前記クエリを生成するステップは、識別された前記品詞に基づくことを特徴とする請求項 1 に記載の方法。

【請求項 3】

前記識別された文字列を、前記クエリテンプレートに対応する重みに基づいてランク付け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 4】

前記コロケーション誤りタイプは、動詞 - 名詞、前置詞 - 名詞、形容詞 - 名詞、および動詞 - 副詞のうちの少なくとも 1 つを含むことを特徴とする請求項 1 に記載の方法。

【請求項 5】

前記特定の組のクエリは、

前記テキストの文を含む文クエリと、

前記テキストの文のチャンクを含むチャンククエリと、

前記テキストの文の主要語ペアを含む語クエリと

を含むことを特徴とする請求項 1 に記載の方法。

【請求項 6】

前記検索モジュールは、ウェブベースの検索エンジンであることを特徴とする請求項 1 に記載の方法。

【発明の詳細な説明】

【背景技術】

【0001】

以下の説明は、単に、一般的な背景情報として提供され、主張される主題の範囲を確定する助けとして使用されることを意図するものではない。

【0002】

成長の一途をたどるグローバル経済、およびインターネットの急速な発展とともに、世界中の人々が、その人々の母語ではない言語で書くことにますます慣れ親しんでいる。残念ながら、大きく異なる文化および文体を有する、いくつかの社会に関して、いくつかの母語ではない言語で書く能力は、常に存在する障壁である。母語ではない言語（例えば、英語）で書く際、言語使用の誤りが、ノンネイティブスピーカ（例えば、中国語、日本語、韓国語その他の英語ではない他の言語を話す人々）によって頻繁に犯される。この種の誤りには、文法上の誤りと、動詞 - 目的語、形容詞 - 名詞、副詞 - 動詞等のコロケーションの不適切な使用の両方が含まれる。

【0003】

多くの人々は、適切な文法を使用して母語ではない言語で書く能力を有するが、それらの人々は依然として、2つの語の間のコロケーションの誤りに苦悩する可能性がある。さらに他の人々は、文法と、2つの語の間のコロケーションなどの他の誤りとの両方に苦悩する。スペルチェックプログラムおよび文法チェックプログラムが文法上の誤りを訂正するのに役立つとはいえ、2つの語の間のコロケーションの誤りの検出および/または訂正は、特にこれらの誤りがその他の点で文法的に正しい可能性があるので、困難である可能

10

20

30

40

50

性がある。したがって、文法チェッカは通常、語の間のコロケーションと関係する誤りを検出するのに全くと言っていいほど助けとならない。以下の説明において、母語ではない言語の例として英語が使用されるが、以上の問題は、他の言語境界にも存在する。

【0004】

例えば、その他の点で文法的に正しい場合でも、文をネイティブらしい英語でなくするコロケーション誤りを含む、以下の文を考える。すなわち、

1. Open the light .

2. Everybody hates the crowded traffic on weekends .

3. This is a check of US\$500 .

4. I congratulate you for your success .

以上の文のネイティブらしい英語のバージョンは、以下のようでなければならない。すなわち、

1. Turn on the light .

2. Everybody hates the heavy traffic on weekends .

3. This is a check for US\$500 .

4. I congratulate you on your success .

である。

【0005】

英語を母語としない話者が直面する障壁の例として、中国人ユーザの窮状を検討されたい。文化、背景、および考え方の習慣により、中国人はしばしば、文法的ではあるが自然ではない英文を作る。例えば、中国人は、中国語の主語を英語の主語に直接に翻訳し、目的語および動詞に関しても、同じことを行う傾向がある。英語で書く際、中国人は、しばしば、動詞と前置詞の間、形容詞と名詞の間、動詞と名詞の間等のコロケーションを決定するのに困難を覚える。さらに、ビジネス分野等の特定の分野では、特別なライティングスキル及びライティングスタイルが必要とされる。

【発明の開示】

【発明が解決しようとする課題】

【0006】

一般的な辞書は、リーディング（一種の復号プロセス）の目的でノンネイティブスピーカーによって主に使用されるが、これらの辞書は、ライティング（一種の符号化プロセス）のための十分な支援を提供しない。これらの辞書は、単一の語の説明を提供するだけであり、通常、関係のある句およびコロケーションを説明する十分な情報を提供しない。さらに、この種の情報を、その情報のいくらかが辞書において提供されている場合でさえ、辞書から簡単に得る方法が存在しない。他方、現在広く使用されている文法チェックツールは、犯しやすい文法上の誤りを検出する、いくらかの限られた能力を有するが、コロケーション誤りを検出することができない。

【課題を解決するための手段】

【0007】

この概要は、発明を実施するための最良の形態において後述されるいくつかの概念を簡略化された形で導入するために提供される。この概要は、主張される主題の重要な特徴または不可欠な特徴を特定することを意図しておらず、また、主張される主題の範囲を確定する助けとして使用されることも意図していない。

【0008】

文がアクセスされ、少なくとも1つのクエリが、前記文に基づいて生成される。前記少なくとも1つのクエリは、例えばウェブ検索エンジンを使用して、ある文書コレクション内のテキストと比較されることが可能である。前記文中のコロケーション誤りが、前記少なくとも1つのクエリと前記文書コレクション内のテキストとの比較に基づき、検出され、かつ/または訂正されることが可能である。

【発明を実施するための最良の形態】

【0009】

図1は、本発明を実施することができる適切なコンピューティングシステム環境100の例を示している。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一例に過ぎず、本発明の用法または機能の範囲について何ら限定を示唆することを意図していない。また、コンピューティング環境100が、例示的な動作環境100に示されるコンポーネントのいずれの1つ又は組合せに関連する依存関係または要件も有すると解釈してはならない。

【0010】

本発明は、他の多数の汎用または専用のコンピューティングシステム環境またはコンピューティングシステム構成で動作する。本発明で使用するのに適する可能性がある周知のコンピューティングシステム、コンピューティング環境、および/またはコンピューティング構成の例には、これらに限定されないが、パーソナルコンピュータ、サーバコンピュータ、ハンドヘルドデバイスまたはラップトップデバイス、マルチプロセッサシステム、マイクロプロセッサベースのシステム、セットトップボックス、プログラマブル家庭用電子機器、ネットワークPC、ミニコンピュータ、メインフレームコンピュータ、電話システム、以上のシステム又はデバイスのいずれかを含む分散コンピューティング環境等が含まれる。

【0011】

本発明は、コンピュータによって実行される、プログラムモジュールなどのコンピュータ実行可能命令の一般的な文脈において説明することができる。一般に、プログラムモジュールには、特定のタスクを実行する又は特定の抽象データ型を実装する、ルーチン、プログラム、オブジェクト、コンポーネント、データ構造等が含まれる。また本発明は、通信ネットワークを介してリンクされたリモート処理装置によってタスクが実行される、分散コンピューティング環境で実施してもよい。分散コンピューティング環境では、プログラムモジュールは、メモリ記憶装置を含むローカルコンピュータ記憶媒体とリモートコンピュータ記憶媒体の両方の中に配置されることが可能である。プログラムおよびモジュールによって実行されるタスクについては、後段で、図の助けを借りて説明する。当業者は、任意の形のコンピュータ読み取り可能な媒体上に書き込むことが可能なプロセッサ実行可能命令として、説明および図を実施することができる。

【0012】

図1を参照すると、本発明を実施するための例示的システムは、コンピュータ110の形で汎用コンピューティングデバイスを含む。コンピュータ110のコンポーネントには、これらに限定されないが、処理装置120、システムメモリ130、及びシステムメモリを含む様々なシステムコンポーネントを処理装置120に結合するシステムバス121が含まれる。システムバス121は、様々なバスアーキテクチャのいずれかを使用する、メモリバスまたはメモリコントローラ、周辺バス、およびローカルバスを含め、いくつかのタイプのバス構造のいずれであってもよい。例として、限定としてではなく、そのようなアーキテクチャには、ISA (Industry Standard Architecture) バス、MCA (Micro Channel Architecture) バス、EISA (Enhanced ISA) バス、VESA (Video Electronics Standards Association) ローカルバス、およびメザニン (Mezzanine) バスとしても知られるPCI (Peripheral Component Interconnect) バスが含まれる。

【0013】

コンピュータ110は通常、様々なコンピュータ読み取り可能な媒体を備える。コンピュータ読み取り可能な媒体は、コンピュータ110がアクセスすることができる任意の利用可能な媒体であることが可能であり、揮発性媒体と不揮発性媒体、リムーバブルな媒体とノンリムーバブルな媒体がともに含まれる。例として、限定としてではなく、コンピュータ読み取り可能な媒体は、コンピュータ記憶媒体および通信媒体を含むことが可能であ

10

20

30

40

50

る。コンピュータ記憶媒体には、コンピュータ読み取り可能な命令、データ構造、プログラムモジュールその他のデータの情報の格納のために任意の方法または技術で実装された、揮発性および不揮発性の、リムーバブルおよびノンリムーバブルな媒体が含まれる。コンピュータ記憶媒体には、これらに限定されないが、RAM、ROM、EEPROM、フラッシュメモリその他のメモリ技術、CD-ROM、DVD（デジタルバーサタイルディスク）その他の光ディスクストレージ、磁気カセット、磁気テープ、磁気ディスクストレージその他の磁気記憶装置、または所望の情報を格納するのに使用することができ、コンピュータ110がアクセスすることができる他の任意の媒体が含まれる。通信媒体は通常、搬送波等の変調されたデータ信号の中に、または他のトランスポート機構で、コンピュータ読み取り可能な命令、データ構造、プログラムモジュールその他のデータを実体10
化し、通信媒体には、あらゆる情報配信媒体（information delivery media）が含まれる。「変調されたデータ信号」という用語は、信号内に情報を符号化するように特性の1つ又は複数が設定または変更されている信号を意味する。例として、限定としてではなく、通信媒体には、有線ネットワークまたは直接有線接続等の有線媒体、および音響媒体、RF媒体、赤外線媒体その他の無線媒体等の無線媒体が含まれる。また、以上の媒体のいずれかの媒体の組合せも、コンピュータ読み取り可能な媒体の範囲内に含まれなければならない。

【0014】

システムメモリ130は、ROM（読み取り専用メモリ）131やRAM（ランダムアクセスメモリ）132などの揮発性メモリおよび/または不揮発性メモリの形で、コンピュータ記憶媒体を備える。起動中等に、コンピュータ110内部の要素間で情報を転送するのを助ける基本ルーチンを含むBIOS（基本入出力システム）133が、通常ROM131の中に格納される。RAM132は通常、処理装置120が、即時にアクセスすることができ、かつ/または現在、処理しているデータおよび/またはプログラムモジュールを含む。限定ではなく例として、図1は、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その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を示している。20

【0015】

またコンピュータ110は、他のリムーバブル/ノンリムーバブル、揮発性/不揮発性のコンピュータ記憶媒体も備えることが可能である。単に例として、図1は、ノンリムーバブルな不揮発性の磁気媒体に対して読み取り又は書き込みを行うハードディスクドライブ141、リムーバブルな不揮発性の磁気ディスク152に対して読み取り又は書き込みを行う磁気ディスクドライブ151、およびCD-ROMその他の光媒体等のリムーバブルな不揮発性の光ディスク156に対して読み取り又は書き込みを行う光ディスクドライブ155を示している。例示的動作環境において使用することができる、他のリムーバブル/ノンリムーバブル、揮発性/不揮発性のコンピュータ記憶媒体には、これらに限定されないが、磁気テープカセット、フラッシュメモリカード、DVD、デジタルビデオテープ、ソリッドステートRAM、ソリッドステートROMなどが含まれる。ハードディスクドライブ141は通常、インタフェース140のようなノンリムーバブルなメモリインタフェースを介してシステムバス121に接続され、磁気ディスクドライブ151および30
光ディスクドライブ155は通常、インタフェース150のようなリムーバブルなメモリインタフェースでシステムバス121に接続される。40

【0016】

前述し図1に示したドライブ及び関連するコンピュータ記憶媒体は、コンピュータ読み取り可能な命令、データ構造、プログラムモジュールその他のデータのストレージをコンピュータ110に提供する。図1では例えば、ハードディスクドライブ141が、オペレーティングシステム144、アプリケーションプログラム145、他のプログラムモジュール146、およびプログラムデータ147を格納しているものとして示されている。これらのコンポーネントは、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と同一であ50

ることも異なることも可能であることに留意されたい。オペレーティングシステム 144、アプリケーションプログラム 145、他のプログラムモジュール 146、およびプログラムデータ 147 には、少なくともそれらが異なるコピーであることを示すために、異なる符号を与えている。

【0017】

ユーザは、キーボード 162、マイクロホン 163、および、マウス、トラックボール又はタッチパッドなどのポインティングデバイス 161 などの入力デバイスを介して、コマンドおよび情報をコンピュータ 110 に入力することができる。他の入力デバイス（図示せず）としては、ジョイスティック、ゲームパッド、サテライトディッシュ、スキャナなどを含めることが可能である。これら及びその他の入力デバイスはしばしば、システムバスに結合されたユーザ入力インタフェース 160 を介して処理装置 120 に接続されるが、パラレルポート、ゲームポート、または USB（ユニバーサルシリアルバス）などの、他のインタフェースおよびバス構造で接続してもよい。また、モニタ 191 その他のタイプのディスプレイデバイスも、ビデオインタフェース 190 のようなインタフェースを介してシステムバス 121 に接続される。モニタに加え、コンピュータは、出力周辺インタフェース 190 を介して接続することができる、スピーカ 197 やプリンタ 196 などの、他の周辺出力デバイスも備えることが可能である。

【0018】

コンピュータ 110 は、リモートコンピュータ 180 のような 1 つ又は複数のリモートコンピュータに対する論理接続を使用するネットワーク化された環境で動作することができる。リモートコンピュータ 180 は、パーソナルコンピュータ、ハンドヘルドデバイス、サーバ、ルータ、ネットワーク PC、ピアデバイスその他の一般的なネットワークノードであることが可能であり、通常、コンピュータ 110 に関連して前述した要素の多く又はすべてを備える。図 1 に示した論理接続には、LAN（ローカルエリアネットワーク）171 および WAN（ワイドエリアネットワーク）173 が含まれるが、他のネットワークを含むことも可能である。そのようなネットワーキング環境は、オフィス、企業全体のコンピュータ網、イントラネット及びインターネットで一般的である。

【0019】

LAN ネットワーキング環境で使用される場合、コンピュータ 110 は、ネットワークインタフェース又はアダプタ 170 を介して LAN 171 に接続される。WAN ネットワーキング環境で使用される場合、コンピュータ 110 は通常、インターネットなどの WAN 173 を介して通信を確立するためのモデム 172 または他の手段を備える。内部にあることも外部にあることも可能なモデム 172 は、ユーザ入力インタフェース 160 その他の適切な機構を介してシステムバス 121 に接続することができる。ネットワーク化された環境では、コンピュータ 110 に関連して示したプログラムモジュール又はプログラムモジュールの部分は、リモートメモリ記憶装置の中に格納することができる。限定ではなく例として、図 1 は、リモートアプリケーションプログラム 185 が、リモートコンピュータ 180 上に存在するものとして示す。図示したネットワーク接続は、例示的であり、コンピュータ間で通信リンクを確立する他の手段も使用できることを理解されたい。

【0020】

図 2 は、テキスト内のコロケーション誤り（collocation error）を検出して訂正するためのシステム 200 の流れ図である。多くのタイプのコロケーション誤りが存在する。システム 200 の一態様では、4 つのタイプのコロケーション誤りが検出される。コロケーション誤りのタイプには、以下が含まれる。すなわち、

1. 動詞 - 名詞（VN、例えば、* learn / acquire knowledge）
2. 前置詞 - 名詞（PN、例えば、* on / in the morning）
3. 形容詞 - 名詞（AN、例えば、* social / socialist country）、および
4. 動詞 - 副詞（VA、例えば、situations change * largely / greatly）

前処理モジュール 202 が、テキストを処理して、テキストの品詞タグ付け (tagging) および構文解析をもたらす。多くの異なるタイプのパーサを、テキストを処理するのに使用することが可能である。以下は、例示的な文である。

I have recognized this person for years .

前処理モジュール 202 が、この文にタグ付けを行い、この文を以下のとおりに切り分ける。

[NP I / PRP] [VP have / VBP recognized / VBN] [NP this / DT person / NN] [PP for / IN] [NP years . < / s > / NNS]

この処理済みのテキストを使用して、クエリ生成モジュール 204 が、クエリを構築する。一例では、4つのクエリセットが、先に特定されたコロケーション誤りタイプの各タイプに関して生成される。例えば、コロケーション誤りタイプは、動詞 - 名詞、前置詞 - 名詞、形容詞 - 名詞、および動詞 - 副詞であることが可能である。生成されるクエリは、文のフルテキストと共に、補助が削除された文の短縮された部分も含むことが可能である。前掲の文に関する例示的な短縮されたクエリには、「have recognized this person」、「have recognized」、「this person」、および「recognized person」を含めることが可能である。

【0021】

クエリは、検索モジュール 206 にサブミットされる。一実施形態では、検索モジュールは、MSN (登録商標) Search (search.msn.com)、Google (登録商標) (www.google.com)、および/またはYahoo! (登録商標) (www.yahoo.com) などのウェブベースの検索エンジンであることが可能である。ウェブは、膨大な量のテキストを含むので、コロケーション誤りを検出する安価なリソースであることが可能である。誤り検出モジュール 208 が、クエリ生成モジュール 204 によって生成されたクエリを、検索モジュール 206 によって獲得された結果と比較する。誤り訂正モジュール 210 が、誤り検出モジュール 208 によって識別された誤りに関する候補訂正を提供する。

【0022】

図3は、図2に示されるシステム200において実施することが可能な方法220の流れ図である。ステップ222で、文にアクセスする。文には、ワードプロセッサ、例えば、ワシントン州リッチモンドのマイクロソフトコーポレーションから入手可能なMicrosoft Word (登録商標) に入力されているテキストが含まれてもよい。ステップ224で、その文を構文解析してチャンク (chunk) にし、文中の品詞が識別される。次に、ステップ226で、その構文解析に基づいてクエリを生成する。ステップ228で、クエリを、MSN (登録商標) Search、Google (登録商標) および/またはYahoo! (登録商標) 等の検索エンジンにサブミットする。ステップ230で、文中のコロケーション誤りが、それらのクエリと検索エンジンからの結果とを比較することにより検出される。誤りを検出した後、ステップ232で、コロケーション誤りの代替物のランク付けされた候補をユーザに提示する。

【0023】

図4は、図2のクエリ生成モジュール204のブロック図である。クエリ生成モジュール204は、構文解析済みの文240、例えば前処理モジュール202から受け取る構文解析済みの文を受け入れる。構文解析済みの文240に基づき、クエリ生成モジュール204は、文クエリ242、チャンククエリ244、および語クエリ246を生成する。先に特定された可能なコロケーション誤りのタイプを所与として、検査語 (checking word) (すなわち、コロケーション誤りを生じさせる可能性がある語) を、以下のとおり検出する。すなわち、タイプVNにおける動詞、タイプPNにおける前置詞、タイプANにおける形容詞、およびタイプVAにおける副詞である。タイプに応じて、クエリ生成モジュール204は、異なるクエリセットを以下のとおり生成する。すなわち、

1. 文クエリ 242: S - Query と呼ばれる、元の文と (各タイプに関して予め定められたされた補助 (auxiliary) を削除することにより) 短縮された文

2. チャンククエリ 244: C - Query と呼ばれる、文中の対応するチャンクペア、および

3. 語クエリ 246: W - Query と呼ばれる、文中の対応する主要語 (head word) ペア

である。

【0024】

タイプVN検出のための、文「I have recognized this person for years」に対する例示的クエリが以下に提示される。~は、2つの隣接する語が、互いに隣接しているか又は1語離れていることが可能であることを意味する。

S - Query: [“ I have recognized this person for years ”]

S - Query: [“ have recognized this person ”]

C - Query: [“ have recognized ” ~ “ this person ”]

W - Query: [“ recognized ” ~ “ person ”]

【0025】

各タイプのクエリを生成するための例示的規則は、以下のとおりである。すなわち、VN: 複数の S - Query、1つの C - Query V ~ N、および1つの W - Query V_h ~ N_h (N_h は、対応する名詞チャンクの主要語を表す。)、

PN: 前置詞を含む、PNの1つの C - Query、

AN: ANペアを含む、ANの1つの C - Query、および

VA: VAペアを含む複数の C - Query、およびVA主要語を含む複数の W - Query

である。

【0026】

図5は、ある文中の誤りを検出する方法250の流れ図である。クエリ生成モジュール204によって生成されたクエリが、ステップ251で、検索モジュール206にサブミットされる。検索モジュール206によって獲得された検索結果が、それらのクエリと比較される。一例では、それらの結果には、ウェブ検索エンジンを使用して取得された文書に関するテキストの要約が含まれる。ステップ252で、クエリ生成モジュール204からの複数のSクエリ242が、検索モジュールからの結果と比較される。次に、ステップ254で、Sクエリ242の1つ又は複数が、検索モジュール結果と一致するか否かについての判定が行われる。Sクエリの1つ又は複数が、検索モジュール結果と一致した場合、ステップ256で、コロケーション誤りは全く存在しないものと判定される。

【0027】

しかし一致が存在しない場合、方法250はステップ258に進み、複数のCクエリ244が、検索モジュール結果と比較される。ステップ260で、Cクエリの1つ又は複数、検索モジュール結果とよく一致するか否か、およびその比較に関するスコアが閾値より大きいかが否かが判定される。一例では、このスコアは、Cクエリのチャンクが、検索結果の中で出現する回数を、Cクエリの中の語がそれらの検索結果の中で同時に出現する回数で割ることによって計算される。そのスコアが、閾値より大きい場合、ステップ256で、コロケーション誤りは全く存在しないものと判定される。

【0028】

そのスコアが、閾値未満である場合、方法250はステップ262に進み、複数のW - Queryが、検索エンジンデータと比較される。ステップ264が、それらのWクエリと検索エンジンデータの間、よい一致が存在するか否か、および、その比較に関するスコアが、閾値より大きいかが否かを判定する。スコアが閾値より大きい場合、ステップ25

10

20

30

40

50

6で、コロケーション誤りは全く存在しないものと判定される。この比較に関するスコアは、C - Query比較スコアと同様であることが可能である。このため、W - Query比較スコアは、W - Queryが検索結果の中で出現する回数を、W - Queryの中の語ペアが同時に出現する合計回数で割ることによって計算することが可能である。そのスコアが閾値未満である場合、方法250はステップ266に進み、ユーザが可能なコロケーション誤りについての通知を受ける。

【0029】

図6は、可能な訂正済みのコロケーションをユーザに提示するための方法270の流れ図である。ステップ272で、クエリテンプレートが生成される。クエリテンプレートは、誤りとして識別された語に基づいて生成される（すなわち、前述した検査語は、図5の方法250により判定したコロケーション誤りを含む。）。クエリテンプレートは、コロケーション誤りを生じさせる検査語が「+」で置き換えられた後の入力文から導出される。前出の文において、「recognized」が、検査語として識別されており、このため、クエリテンプレートはこの語に基づいて展開される。例えば、VN検出のための、文「I have recognized this person for years」のクエリテンプレートは、以下のとおりである。+は、任意の語を表す。

S - QT: [" I have + this person for years "]

S - QT: [" I have + this person "]

S - QT: [" have + this person for years "]

S - QT: [" I have + this person "]

C - QT: [" + this person for years "]

C - QT: [" + this person "]

【0030】

クエリテンプレートを生成するための例示的な規則は、以下のとおりであることが可能である。すなわち、

VN: S - QT、C - QT（動詞が+で置き換えられている。）、

PN: S - QT、C - QT（前置詞が+で置き換えられている。）、

AN: S - QT、C - QT（形容詞が+で置き換えられている。）、および

VA: S - QT、C - QT（副詞が+で置き換えられている。）

である。

【0031】

ステップ274で、クエリテンプレートが、検索モジュールに、本明細書では検索エンジンにサブMITTされる。ステップ276で、検索エンジン結果からの文字列が取得される。それらの文字列は、取り巻く文脈の語を有するテキストの要約を含むことが可能である。クエリテンプレートと一致する文字列であって、+の位置は、任意の1つの語であることが可能な文字列が、文字列候補として識別される。そのコロケーション（+に取って代わる語、およびコロケーションタイプに応じた文字列の中の別の語によって形成される）を含まない候補が、ステップ278で削除される。残りの候補は、文字列候補と一致したクエリテンプレートの対応する重みに基づくスコアに応じて、ランク付けされる。例えば、クエリテンプレートの重みは、そのクエリテンプレートの中の語の数に基づくことが可能である。各候補に関するスコアは、それらの候補を含むすべての要約にわたって重みの合計をとることにより、計算される。候補を取得するクエリテンプレート（QT_s）に関するスコアは、以下によって表現することが可能である。すなわち、

【0032】

【数1】

スコア（候補） = $\sum_{QT_s} \text{重さ}(QT)$

【0033】

次に、ステップ280で、ランク付けされた候補リストがユーザに提示される。例えば、ポップアップ・メニューを使用して、ランク付けされたリストが提示されることが可能

10

20

30

40

50

である。ユーザは、そのリストから選択肢の1つを選び、コロケーション誤りを訂正することができる。

【0034】

主題を、構造上の特徴および／または方法上の動作に特有の言い回しで説明してきたが、添付の特許請求の範囲において規定される主題は、必ずしも前述した特定の特徴または動作に限定されないことを理解されたい。むしろ、前述した特定の特徴および動作は、請求項を実施する例示的な形態として開示されている。

【図面の簡単な説明】

【0035】

【図1】一般的なコンピューティング環境を示すブロック図である。

10

【図2】コロケーション誤りを検出して、訂正するためのシステムを示す流れ図である。

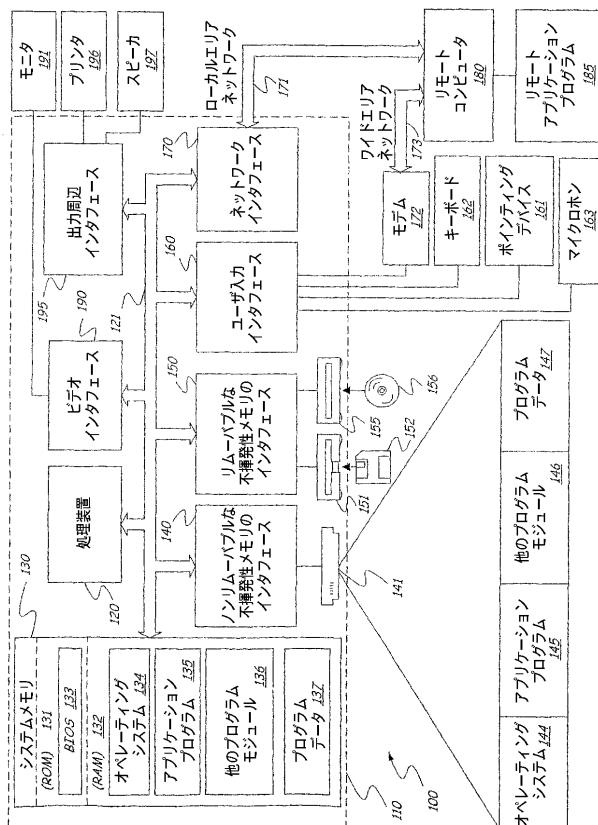
【図3】コロケーション誤りを検出して、訂正するための方法を示す流れ図である。

【図4】クエリ生成モジュールを示すブロック図である。

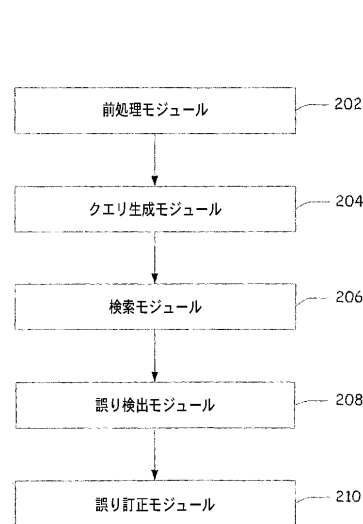
【図5】コロケーション誤りを検出するための方法を示す流れ図である。

【図6】候補コロケーション訂正を提示するための方法を示す流れ図である。

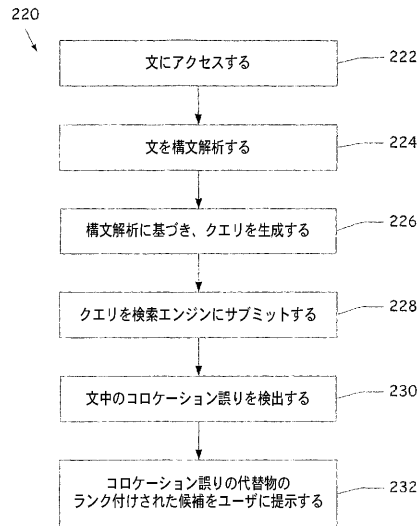
【図1】



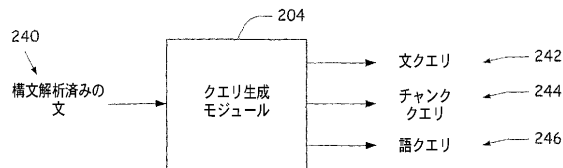
【図2】



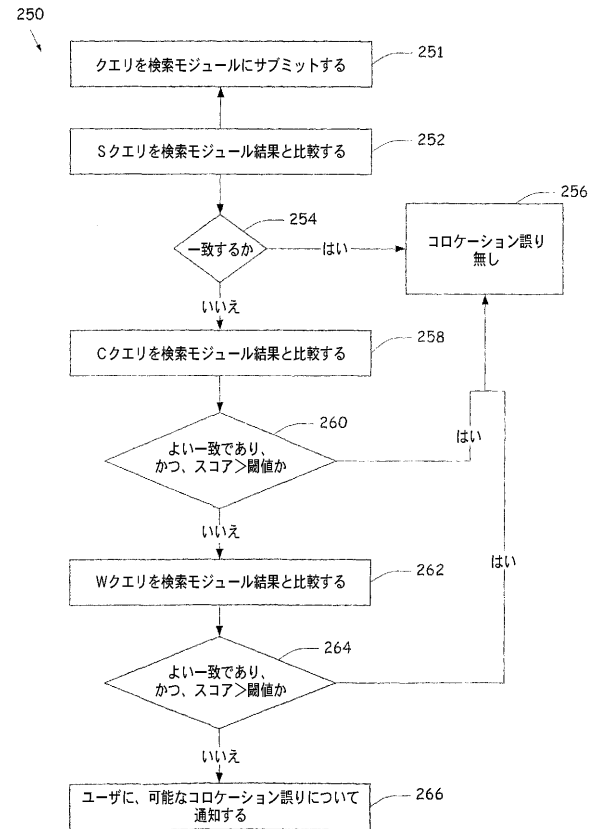
【図 3】



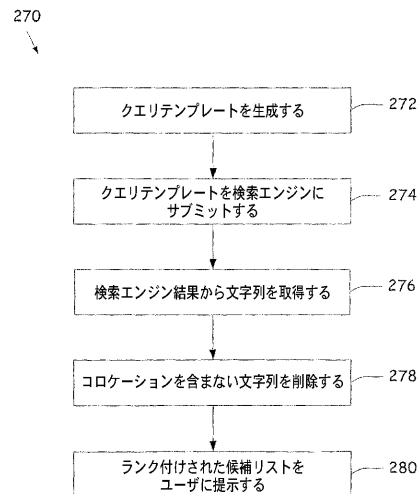
【図 4】



【図 5】



【図 6】



フロントページの続き

- (72)発明者 シャオ - ウーエン ホン
アメリカ合衆国 98052 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 インターナショナル パテント内
- (72)発明者 チェンフェン ガオ
アメリカ合衆国 98052 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 インターナショナル パテント内
- (72)発明者 ミン チョウ
アメリカ合衆国 98052 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 インターナショナル パテント内

審査官 長 由紀子

- (56)参考文献 特開2001-101186(JP,A)
大鹿 広憲 外3名, Googleを活用した英作文支援システムの構築, DEWS2005論
文集 [online], 日本, (社)電子情報通信学会データ工学研究専門委員会, 2005
年 5月 2日, DWS2005 4B-i8
大鹿 広憲 外3名, 検索エンジンを使った翻訳サポートシステムの構築, 情報処理学会研究報
告, 日本, 社団法人情報処理学会, 2004年 7月15日, 第2004号第72号, p.585-591

- (58)調査した分野(Int.Cl., DB名)
G06F 17/20 - 28
G06F 17/30