

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2012-113459

(P2012-113459A)

(43) 公開日 平成24年6月14日(2012.6.14)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 17/28 (2006.01)	G06F 17/28 Z	5B075
G06F 17/30 (2006.01)	G06F 17/30 170Z	5B091
	G06F 17/30 350C	

審査請求 有 請求項の数 6 O L (全 13 頁)

(21) 出願番号	特願2010-260845 (P2010-260845)	(71) 出願人	000003078 株式会社東芝 東京都港区芝浦一丁目1番1号
(22) 出願日	平成22年11月24日(2010.11.24)	(71) 出願人	301063496 東芝ソリューション株式会社 東京都港区芝浦一丁目1番1号
		(74) 代理人	100100516 弁理士 三谷 恵
		(72) 発明者	中村 寛爾 東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内
		(72) 発明者	澁谷 貴志 東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内

最終頁に続く

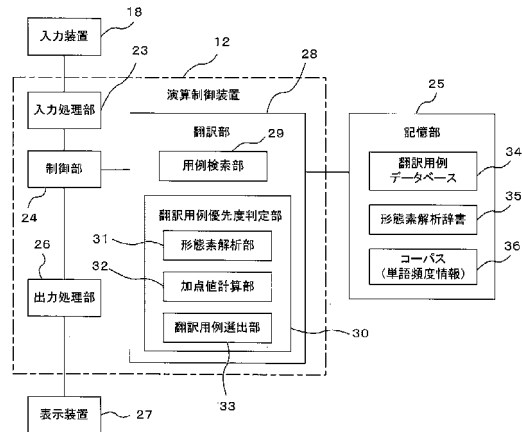
(54) 【発明の名称】 用例翻訳システム、用例翻訳方法及び用例翻訳プログラム

(57) 【要約】

【課題】類似度の計算方法に翻訳対象の分野情報を指標の一つとして加えることで、利用者の求める翻訳結果により近い用例訳文を提供することである。

【解決手段】用例検索部は翻訳対象原文と翻訳用例データベースの翻訳用例原文との類似度を計算し類似度が予め定めた閾値以上の翻訳用例を翻訳用例データベースから検索する。形態素解析部は用例検索部により複数の翻訳用例が検索されたとき複数の翻訳用例のそれぞれの訳文を形態素解析辞書の形態素解析情報を参照して形態素解析し単語を抽出する。加点値計算部は形態素解析部で抽出された前記単語につきコーパスの単語頻度情報を参照し単語の出現頻度に応じて翻訳用例の類似度の加算値を計算する。翻訳用例選出部は用例検索部で計算された類似度に加算値計算部で計算された加算値を加算して最も大きい類似度の翻訳用例を選出する。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

翻訳対象の第 1 言語の原文と翻訳目的の第 2 言語の訳文とを対にした翻訳用例を格納した翻訳用例データベースと、

文を形態素解析する際に参照する形態素解析辞書と、

所定の分野の文書で用いられた単語の出現頻度を格納したコーパスとを記憶した記憶装置と、

前記翻訳対象原文と前記翻訳用例データベースの翻訳用例原文との類似度を計算し、その類似度が予め定めた閾値以上の翻訳用例を前記翻訳用例データベースから検索する用例検索部と、

前記用例検索部により複数の翻訳用例が検索されたとき、当該複数の翻訳用例の各翻訳用例訳文を前記形態素解析辞書を参照して形態素解析し単語を抽出する形態素解析部と、

前記形態素解析部で抽出された単語につき前記コーパスに格納された当該単語の出現頻度に応じて前記類似度に加算する加算値計算部と、

前記加算後の類似度に基づいて翻訳用例を選出する翻訳用例選出部と、
を備えた用例翻訳システム。

【請求項 2】

前記コーパスは単語の出現頻度の更新日時情報をさらに格納し、前記加算値計算部は前記加算について前記更新日時情報の新しいものほど大きな重み付けをする請求項 1 記載の用例翻訳システム。

【請求項 3】

翻訳対象の第 1 言語の原文と翻訳目的の第 2 言語の訳文とを対にした翻訳用例を格納した翻訳用例データベースと、文を形態素解析する際に参照する形態素解析辞書と、所定の分野の文書で用いられた単語の出現頻度を単語頻度情報として格納したコーパスとを予め記憶装置に記憶しておき、

前記翻訳対象原文と前記翻訳用例データベースの翻訳用例原文との類似度を計算し、その類似度が予め定めた閾値以上の翻訳用例を前記翻訳用例データベースから検索し、

複数の翻訳用例が検索されたとき、当該複数の翻訳用例の各翻訳用例訳文を前記形態素解析辞書を参照して形態素解析し単語を抽出し、

抽出された前記単語につき前記コーパスに格納された当該単語の出現頻度に応じて前記類似度に加算し、

前記加算後の類似度に基づいて翻訳用例を選出して翻訳用例とする用例翻訳方法。

【請求項 4】

前記コーパスに前記単語頻度情報に加え単語の出現頻度の更新日時情報を予め格納しておき、前記単語の出現頻度に応じて計算した加算値に前記更新日時情報の新しいものほど大きな重み付け係数を乗算した加算値を計算する請求項 3 記載の用例翻訳方法。

【請求項 5】

前記用例翻訳プログラム、翻訳対象の第 1 言語の原文と翻訳目的の第 2 言語の訳文とを対にした翻訳用例を格納した翻訳用例データベース、翻訳用例訳文を形態素解析する際に参照する形態素解析辞書、所定の分野の文書で用いられた単語の出現頻度を単語頻度情報として格納したコーパスを予め記憶した記憶装置と、前記翻訳対象原文を入力するとともに操作に必要な情報を入力する入力装置と、前記翻訳対象原文や前記翻訳用例を表示する表示装置と、前記用例翻訳プログラムを演算実行する演算制御装置とを備えた用例翻訳システムとして機能させるためのコンピュータにおいて、

前記コンピュータを、

前記翻訳対象原文と前記翻訳用例データベースの翻訳用例原文との類似度を計算し、その類似度が予め定めた閾値以上の翻訳用例を前記翻訳用例データベースから検索する用例検索手段と、

前記用例検索部により複数の翻訳用例が検索されたとき、当該複数の翻訳用例の各翻訳用例訳文を前記形態素解析辞書を参照して形態素解析し単語を抽出する形態素解析手段と

10

20

30

40

50

前記形態素解析部で抽出された前記単語につき前記コーパスに格納された当該単語の出現頻度に応じて前記類似度に加点する加点値計算手段と、

前記加点後の類似度に基づいて翻訳用例を選出する翻訳用例選出手段として機能させるための用例翻訳プログラム。

【請求項 6】

前記コーパスに前記単語頻度情報に加え単語の出現頻度の更新日時情報を予め格納しておき、前記加点値計算手段は単語の出現頻度に応じて計算した加算値に、前記更新日時情報の新しいものほど大きな重み付け係数を乗算した加算値を計算する請求項 5 記載の用例翻訳システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、翻訳対象の第 1 言語の原文と翻訳目的の第 2 言語の訳文とを対にした翻訳用例を用いて原文を訳文に翻訳する用例翻訳システム、用例翻訳方法及び用例翻訳プログラムに関する。

【背景技術】

【0002】

外国語を用いた情報交換のツールとして機械翻訳の重要性が高まっている。機械翻訳の技術の一つとして翻訳用例を用いた翻訳方法が知られている。これは予め原文と訳文とを対にした複数の翻訳用例を翻訳用例データベースに登録しておき、翻訳対象原文が入力された際に、翻訳対象原文と類似した翻訳用例原文を翻訳用例データベースから検索し、得られた翻訳用例原文を、対となる翻訳用例訳文とともに利用者に提示するものである。利用者は必要に応じて提示された翻訳用例訳文を部分的に修正することで、希望する訳文を比較的少ない作業量で得ることができる。

【0003】

ここで翻訳対象原文の類似文を検索する際、翻訳対象原文と各々の翻訳用例原文との類似度が計算される。類似度計算においては、翻訳対象原文と翻訳用例原文との間の一致する単語の割合を計算する方法が一般的な計算方法として知られている。類似度が同じ翻訳用例が複数見つかった場合には、検索で先に見つかったものを優先したり、登録時期が新しいものを優先したりすることで順位付けがなされることが多い。例として " I pass by the house every day. " という原文が与えられた場合、いま、翻訳用例データベースから類似度の高いものが検索され、下記のような用例が得られたとする。

【0004】

用例 1

原文 : I pass by the shop every day. (類似度 : 86%)

訳文 : 私は毎日その店のそばを通る

用例 2

原文 : I pass behind the shop every day. (類似度 : 71%)

訳文 : 私は毎日その店の後ろを通る

この結果、翻訳者は最も類似度の高い用例 1 の訳文を部分的に修正することで、希望の訳文を得ることができる。

【0005】

この方法により類似度が計算された場合、類似度が高いからといって必ずしも翻訳者の望む訳文に近い用例が優先的に検出されるとは限らない。例として、下記の用例 1、2 が翻訳用例データベースに登録されており、双方の原文と似た次の翻訳対象が入力されたとする。

【0006】

用例 1 (登録日 : 2009/08/11)

原文 : The stocks of this brand bring about profits.

10

20

30

40

50

訳文：この銘柄の株式は利益をもたらします

用例 2（登録日：2009/11/30）

原文：The fans of this brand bring about profits.

訳文：このブランドのファンたちは利益をもたらします

翻訳対象

原文：The shares of this brand bring about damage.

ここで、この従来 of の計算方法により翻訳対象原文と各々の翻訳用例原文との類似度を計算した場合、用例 1、用例 2 とともに 8 単語中 6 単語が原文と一致するため同じ類似度になり、登録日の新しい用例 2 の訳文が類似文の訳文として利用者に提示される。

【0007】

しかし、もし、この翻訳対象原文が株式関連の文書中に現れた文であった場合、推測される訳文は「この銘柄の株式は損害をもたらします」となり、提示された用例 2 より用例 1 の訳文の方が近いということになる。このように、従来 of の計算方法により選ばれた類似文の訳文は必ずしも最適というわけではなく、類似度が同じ、またはやや低い別の用例の訳文の方が有用である場合も少なくない。

【先行技術文献】

【特許文献】

【0008】

【特許文献 1】特開 2006 - 24114 号公報

【発明の概要】

【発明が解決しようとする課題】

【0009】

従来 of の類似度の計算方法に翻訳対象の分野情報を指標の一つとして加えることで、利用者の求める翻訳結果により近い用例訳文を提供することである。

【課題を解決するための手段】

【0010】

実施形態 of の翻訳用例システムは、翻訳対象の第 1 言語の原文と翻訳目的の第 2 言語の訳文とを対にした翻訳用例を格納した翻訳用例データベースと、文を形態素解析する際に参照する形態素解析辞書と、所定の分野 of の文書で用いられた単語の出現頻度を格納したコーパスとを記憶した記憶装置を備える。用例検索部は翻訳対象原文と翻訳用例データベース of の翻訳用例原文との類似度を計算し、その類似度が予め定めた閾値以上の翻訳用例を翻訳用例データベースから検索する。形態素解析部は用例検索部により複数の翻訳用例が検索されたとき、当該複数の翻訳用例 of のそれぞれの各翻訳用例訳文を形態素解析辞書を参照して形態素解析し単語を抽出する。加点値計算部は形態素解析部で抽出された前記単語につきコーパスに格納された当該単語 of の出現頻度に応じて類似度に加算する。翻訳用例選出部は加算後の類似度に基づいて翻訳用例を選出する。

【図面の簡単な説明】

【0011】

【図 1】実施形態に係る用例翻訳システムの機能ブロック構成図。

【図 2】実施形態に係る用例翻訳システムのハードウェア構成を示すブロック構成図。

【図 3】実施形態に係るコーパスの一例の説明図。

【図 4】実施形態でコーパスを作成する場合 of の処理内容を示すフローチャート。

【図 5】実施形態でコーパスを作成する場合 of の特定分野 of の文書及びコーパス of の説明図。

【図 6】実施形態に係る用例翻訳システムに翻訳対象原文が入力されてからコーパスを利用して翻訳用例を選出するまでの処理内容を示すフローチャート。

【図 7】実施形態に係るコーパス of の他の一例 of の説明図。

【発明を実施するための形態】

【0012】

以下、実施形態を図面に基づいて説明する。図 1 は、実施形態に係る用例翻訳システムの機能ブロック構成図、図 2 は実施形態に係る用例翻訳システムのハードウェア構成を示

10

20

30

40

50

すブロック構成図である。

【0013】

図2において、用例翻訳システム11は、例えば一般的なコンピュータに用例翻訳プログラムなどのソフトウェアプログラムがインストールされ、そのソフトウェアプログラムが演算制御装置12のプロセッサ13において実行されることにより実現される。

【0014】

演算制御装置12は機械翻訳に関する各種演算を行うものであり、演算制御装置12はプロセッサ13とメモリ14とを有し、メモリ14にはプログラム15が記憶され、プロセッサ13により処理が実行される際には作業エリア16が用いられる。演算制御装置12の演算結果等は表示装置17に表示出力される。

10

【0015】

入力装置18は演算制御装置12に情報を入力するものであり、例えば、マウス19、キーボード20、読み取り装置21a、読み込み装置21bから構成される。読み取り装置21aは、例えばOCR(光学式文字読み取り装置)等であり、読み込み装置21bは、例えば磁気テープ、磁気ディスク、光ディスク等、コンピュータ可読媒体からの読み込み装置である。

【0016】

例えば、マウス19やキーボード20は表示装置17を介して演算制御装置12に各種指令を入力し、キーボード20、読み取り装置21a、読み込み装置21bは、翻訳対象の文書を入力する。すなわち、読み取り装置21a、読み込み装置21bは、翻訳対象の文書のファイルを記憶媒体に入出力するものである。さらに、演算制御装置12の演算結果や用例翻訳に必要な知識・規則を蓄積した辞書等を記憶するハードディスクドライブ(HDD)22が設けられている。

20

【0017】

図1において、演算制御装置12内の各機能ブロックは、用例翻訳プログラムを構成する各プログラム15の機能に対応する。すなわち、プロセッサ13が用例翻訳プログラムを構成する各プログラム15を実行することで、演算制御装置12は、各機能ブロックとして機能することとなる。また、記憶装置25の各ブロックは、演算制御装置12内のメモリ14及びハードディスクドライブ22の記憶領域に対応する。

【0018】

入力装置18は、翻訳対象原文の文書の電子データを入力するものであり、利用者の入力操作に基づく文書の入力が可能である。また、入力装置18は、入力処理部23を介して制御部24に対して各種コマンドを与える。入力装置18によって入力された翻訳対象原文の文書は、演算処理部12の入力処理部23により入力処理されて取り込まれ、制御部24を介して記憶装置25の図示省略の文書記憶エリアに記憶される。制御部24は、入力処理部23、出力処理部26、翻訳部28を制御するとともに、記憶装置25とのデータの授受の制御も行う。そして、演算制御装置12の演算結果は表示装置27に表示出力される。

30

【0019】

翻訳部28は、用例検索部29及び翻訳用例優先度判定部30を有し、翻訳用例優先度判定部30は、形態素解析部31、加価値計算部32、翻訳用例選出部33を有している。これらの詳細については、後述する。

40

【0020】

また、記憶部25には、複数の翻訳用例が予め登録されている翻訳用例データベース34が格納されている。翻訳用例は第1言語の翻訳用例原文と第2言語の翻訳用例訳文とが対となって格納されている。また、記憶部25には、翻訳対象原文や翻訳用例の原文及び訳文を形態素解析をする際に参照される形態素解析辞書35が格納されている。形態素解析辞書35には、形態素解析の対象となる第1言語や第2言語の文法の知識(文法のルールの集まり)や辞書(品詞等の情報付きの単語リスト)が形態素解析情報として格納されている。

50

【0021】

さらに、記憶部25には、翻訳に関連する分野の単語頻度情報が登録されているコーパス36が格納されている。コーパス36は、大量のテキストデータを翻訳システムで利用可能な形式にして登録したものであり、本実施形態では、翻訳対象原文と同じ分野の文書で用いられた名詞単語及びその出現頻度が単語頻度情報として格納され、また出現頻度の更新日時情報が格納されている。コーパス36の詳細は後述する。

【0022】

翻訳部28の用例検索部29は、翻訳対象原文と翻訳用例データベース34の翻訳用例原文との類似度を計算し、類似度が予め定めた閾値以上の翻訳用例を翻訳用例データベース34から検索するものである。

10

【0023】

翻訳用例優先度判定部30は、用例検索部29により複数の翻訳用例が検索されたとき、検索された複数の翻訳用例のうち、どの翻訳用例を優先して選出するかを判定するものである。

【0024】

翻訳用例優先度判定部30の形態素解析部31は、用例検索部29により複数の翻訳用例が検索されたときは、複数の翻訳用例のそれぞれの翻訳用例訳文を形態素解析辞書35の形態素解析情報を参照して形態素解析し、名詞単語を抽出する。

【0025】

翻訳用例優先度判定部30の加点値計算部32は、形態素解析部31で抽出された名詞単語につきコーパス36の単語頻度情報を参照し、名詞単語の出現頻度に応じて翻訳用例の類似度の加算値を計算する。

20

【0026】

翻訳用例優先度判定部30の翻訳用例選出部33は、用例検索部29で計算された翻訳用例の類似度に、加点値計算部32で計算された加算値を加算して、類似度の合計値が最も大きい翻訳用例を選出し、出力処理部26を介して表示装置27に表示出力する。

【0027】

図3はコーパス36の説明図である。コーパス36は特定分野の文書で用いられた名詞単語の出現頻度を単語頻度情報として格納するとともに、出現頻度の更新日時を更新日時情報として格納している。図3では株式関連分野の場合のコーパスを示している。

30

【0028】

例えば、株式という名詞単語は出現頻度が30で更新日時は2010年10月29日であり、証券という名詞単語は出現頻度が27で更新日時は2010年10月09日であり、以下、同様に株式関連分野の文書に用いられた名詞単語の出現頻度と更新日時とを情報として格納している。

【0029】

図4は実施形態でコーパスを作成する場合の処理内容を示すフローチャートである。これは、図示は省略するが、コーパス作成プログラムをコンピューターにインストールし、そのソフトウェアプログラムを演算制御装置12のプロセッサ13において実行することにより実現される。

40

【0030】

いま、図5(a)に示す株式分野の文書が入力装置18から入力処理部23を介して記憶部25の図示省略の文書記憶エリアに記憶されたとする。図4に示すように、まず、コーパスの作成機能は、読み込まれた文書を文単位に切り出す(S1)。図5(a)の株式分野の文書の場合、「株式とは、株式会社における社員権、持分のことである。」という文と、「通常、持分が社員の出資額などに応じて不均一な形態を取るのに対して、均一的な細分化された割合的な構成単位を取る点に特徴がある。」という文との二つの文からなっているので、この二つの文を切り出す。

【0031】

次に、一つ目の文「株式とは、株式会社における社員権、持分のことである。」につき

50

、形態素解析により単語分割をする（S2）。そして、分割した単語から名詞を識別する（S3）。この場合の名詞は、図5（a）の下線を引いた単語であり、「株式」、「株式会社」、「社員権」、「持分」の4個の名詞である。

【0032】

次に、変数*i*に「1」をセットし（S4）、*i*個目の名詞を取り出し（S5）、*i*個目の名詞の出現頻度に1を加算する（S6）。そして、更新日時を更新する（S7）。最初は*i* = 1であるから、1個目の名詞である「株式」が取り出される。1個目の名詞「株式」は、図5（b）に示すように、「株式」の出現頻度29に1を加算し、更新日時を本日の2010年10月29日に更新する。

【0033】

次に、すべての名詞を取り出したか否かを判定し（S8）、すべての名詞を取り出していないときは、変数*i*に1を加算し（S9）、ステップS5に戻る。一方、すべての名詞を取り出しているときは、次の文はあるかどうかを判定し（S10）、次の文があるときはステップS2に戻り、次の文がないときは処理を終了する。

【0034】

このように、一つ目の文につき、ステップS5～ステップS9の処理により、2個目～4個目の名詞「株式会社」、「社員権」、「持分」についても、出現頻度に1を加算し、更新日時を本日の2010年10月29日に更新する。4個目の「持分」について処理が終了すると、ステップS10により、二つ目の文についてステップS2～S10までの処理が開始される。

【0035】

次に、二つ目の文「通常の持分が社員の出資額などに応じて不均一な形態を取るのに対して、均一的な細分化された割合的な構成単位を取る点に特徴がある。」につき、形態素解析により単語分割をし（S2）、分割した単語から名詞を識別する（S3）。この場合の名詞は、図5（a）の下線を引いた単語であり、「通常」、「持分」、「社員」、「出資額」、「不均一」、「形態」、「均一」、「割合」、「構成単位」、「点」、「特徴」の11個の名詞である。

【0036】

一つ目の文の場合と同様に、変数*i*に「1」をセットし（S4）、*i*個目の名詞を取り出し（S5）、*i*個目の名詞の出現頻度に1を加算する（S6）。そして、更新日時を更新する（S7）。

【0037】

最初は*i* = 1であるから、1個目の名詞である「通常」が取り出され、「通常」の出現頻度5に1を加算し、更新日時を本日の2010年10月29日に更新する。以下同様に、2個目～11個目の名詞「持分」、「社員」、「出資額」、「不均一」、「形態」、「均一」、「割合」、「構成単位」、「点」、「特徴」についても、出現頻度に1を加算し、更新日時を本日の2010年10月29日に更新する。2個目の「持分」については、一つ目の文にも出現しているので、1が2回加算されることになる。そして、11個目の名詞「特徴」の処理が終了すると、図5（a）の場合には、次の文はないので処理を終了する。

【0038】

このようにして、コーパス36には特定分野の名詞単語の出現頻度や更新日時が更新されて格納される。

【0039】

次に、図6は、実施形態に係る用例翻訳システムに翻訳対象原文が入力されてからコーパスを利用して翻訳用例を選出するまでの処理内容を示すフローチャートである。

【0040】

用例翻訳システムの利用者により入力装置18から翻訳対象原文が入力されると、入力処理部23により入力処理されて取り込まれ、制御部24を介して記憶装置25の図示省略の文書記憶エリアに記憶される。そして、制御部24は翻訳部28を起動する。

10

20

30

40

50

【 0 0 4 1 】

翻訳部 2 8 は起動がかけられると、まず翻訳対象原文に対し形態素解析を行う (S 1 1)。用例検索部 2 9 は、その結果をもとに翻訳用例データベース 3 4 から翻訳用例を検索する (S 1 2)。このとき翻訳対象原文と、翻訳用例データベース 3 4 に登録されている翻訳用例原文との類似度を計算することになるが、この類似度は双方の文中に同じ単語がどれだけ含まれるかという割合で決定される。そして、用例検索部 2 9 は、類似度が予め定めた閾値以上を満たす翻訳用例は検索できたか否かを判定し (S 1 3)、翻訳用例が検索できない場合は処理を終了する。この場合は、翻訳用例を用いない通常の翻訳処理を行うことになる。

【 0 0 4 2 】

一方、用例検索部 2 9 は、翻訳用例が検索できたときは、複数の翻訳用例か否かを判定する (S 1 4)。複数の翻訳用例でない場合、つまり一つの翻訳用例である場合には、その翻訳用例を選出する (S 1 5)。

【 0 0 4 3 】

ステップ S 1 4 の判定で、複数の翻訳用例が検索されたときは、翻訳用例優先度判定部 3 0 は変数 j に「 1 」をセットする (S 1 6)。これにより、形態素解析部 3 1 は j 個目の翻訳用例訳文を形態素解析し名詞単語を抽出する (S 1 7)。

【 0 0 4 4 】

次に、加点値計算部 3 2 は、j 個目の翻訳用例訳文の名詞単語の出現頻度に応じて加算値を計算する (S 1 8)。すなわち、加点値計算部 3 2 は、コーパス 3 6 の名詞単語の出現頻度を参照して名詞単語の出現頻度を取得し、その出現頻度が高いほど大きな加算値を算出する。加算値の算出の仕方については後述する。そして、加点値計算部 3 2 は j 個目の翻訳用例の類似度に加算値を加算し (S 1 9)、すべての翻訳用例を取り出したか否かを判定し (S 2 0)、すべての翻訳用例を取り出していないときは、変数 j に 1 を加算し (S 2 1)、ステップ S 1 7 に戻る。

【 0 0 4 5 】

一方、すべての翻訳用例を取り出しているときは、翻訳用例選出部 3 3 は、加点値計算部 3 2 により、翻訳用例の類似度に加算値を加算して得られた各々の翻訳用例の類似度合計値を比較し、最も大きい類似度合計値の翻訳用例を選出する (S 2 2)。

【 0 0 4 6 】

図 6 に示した処理内容につき具体例を用いて説明する。いま、翻訳対象原文として、下記の文が与えられたとする。

【 0 0 4 7 】

The shares of this brand bring about damage.

この翻訳対象原文の類似文を検索する場合を考える。なお、この翻訳対象原文は株式関連の文章中に現れた文であり、用例翻訳システムには予め株式の単語情報を登録したコーパス 3 6 を持っているものとする。

【 0 0 4 8 】

まず、翻訳部 2 8 は上記の翻訳対象原文を形態素解析により単語分割し、用例検索部 2 9 は、分割された単語をもとに翻訳用例データベース 3 4 中の翻訳用例を検索する。その結果、類似度の高い翻訳用例として、下記の二つの翻訳用例が得られたとする。

【 0 0 4 9 】

用例 1

原文 : The stocks of this brand bring about profits.

訳文 : この銘柄の株式は利益をもたらします

用例 2

原文 : The fans of this brand bring about profits.

訳文 : このブランドのファンたちは利益をもたらします

翻訳対象原文とこれら二つの翻訳用例原文との類似度は、どちらも $75 \{ (\text{一致する単語数} / \text{全単語数}) \times 100 \text{ で計算} \}$ で同じである。この場合、二つの翻訳用例が得られ

10

20

30

40

50

たので、図6のステップS16以降の処理に移ることになる。

【0050】

形態素解析部31は、用例1及び用例2のそれぞれの翻訳用例訳文に対し、形態素解析により単語分割を行い名詞単語を取り出す。

【0051】

用例1から、「銘柄」、「株式」、「利益」を取り出し、用例2から「ブランド」、「ファン」、「利益」を取り出す。

【0052】

これらすべての名詞単語について、加点値計算部32はコーパス36を参照し、出現頻度の高い名詞単語についてはポイントを加点する。ここでは、簡略化のため、出現頻度を10で割り小数点以下を切り捨てたものを加点するポイントとする。

10

【0053】

用例1の「銘柄」の出現頻度は、図3に示すように「16」であり、「株式」の出現頻度は「30」であり、「利益」は未登録の名詞単語であるので出現頻度は「0」である。従って、用例1の名詞単語の出現頻度の合計は、 $(16 + 30 + 0 = 46)$ であり、これを10で割り小数点以下を切り捨てると加算値は「4」と計算される。用例1の類似度は75であるので、これに加算値4を加算すると、用例1の類似度合計値は79となる。

【0054】

一方、用例2の「ブランド」、「ファン」、「利益」は、図3に示すように、すべて未登録の名詞単語であるので出現頻度は「0」である。従って、用例2の加算値は「0」と計算される。用例2の類似度は75であるので、これに加算値0を加算すると、用例2の類似度合計値は75となる。

20

【0055】

翻訳用例選出部33は、最も大きい類似度合計値の翻訳用例を選出する。この場合は、用例1の類似度合計値が79で用例2の類似度合計値が75であるので、用例1が選出される。翻訳対象原文の内容を考慮すると、用例1の方が用例2より有用である。以上のように、コーパス36を利用することによって翻訳対象原文と同じ分野の翻訳用例の類似度を高くすることで、より有用な翻訳用例を類似文として利用者に提示することができる。

【0056】

このように、類似度の近い翻訳用例が複数ある場合、翻訳者の指定する分野の単語の出現頻度を利用することにより、指定の分野に近い訳文を持つ翻訳用例ほど類似度が高くなるため、単語の出現頻度を利用しない場合に比べ、より翻訳者の希望に近い類似文が検出される。

30

【0057】

以上の説明では、類似度に加点するポイントについて、出現頻度をもとに計算したが、出現頻度だけでなく更新日時の情報も合わせて利用してもよい。一般的に、長い期間をかけてコーパス36を作成する場合、古い単語情報よりも新しい単語情報の方がより有用である場合が多い。そこで、更新日時が新しいものほど加点ポイントが大きくなるよう重みを付けることによって、単語の新鮮さを類似文検索における指標へ反映させる。更新日時と係数との例を以下に示す。

40

【0058】

現在から	係数
半年以内	: 1.0
1年以内	: 0.9
3年以内	: 0.8
3年以上経過	: 0.7

具体例として以下の翻訳対象原文及び類似度の高い翻訳用例として検出された用例1、用例2を考える。また、利用するコーパス36には、図7に示すような単語情報が登録されていたとする。

【0059】

50

翻訳対象原文

The government must fix a safety net immediately.

用例 1

原文：The city must fix a safety net immediately.

訳文：市は早急に安全網を整備しなければならない。

【0060】

用例 2

原文：The prefecture must fix a safety net immediately.

訳文：県は早急にセーフティーネットを整備しなければならない。

【0061】

このとき、類似度に加算するポイントは、出現頻度に更新日時から計算した係数をかけたものを利用する。単純に出現頻度のみを参照した場合、「セーフティーネット」より「安全網」の方が重要な単語となり、用例 2 より用例 1 が優先される。

【0062】

しかし、更新日時による係数をかけた場合、例えば、前述の係数を用いて「安全網」の出現頻度に 0.7、「セーフティーネット」の出現頻度に 1.0 をかけると、加点ポイントは用例 2 の方が大きくなる。実際、現在では「安全網」という言葉より「セーフティーネット」の方が使われることが多く、用例 2 を優先的に利用者へ提示することは妥当な判断といえる。

【0063】

以上のように、翻訳用例を用いて翻訳を行う場合に、類似した翻訳用例が複数検索された際、それら複数の翻訳用例に対して名詞単語の出現頻度や更新日時の情報をもとにポイントを加点することによって、翻訳する分野により近くより新しい訳し方を持つ翻訳用例が優先的に検出されるようになる。

【0064】

本発明のいくつかの実施形態を説明したが、これらの実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。これら新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

【符号の説明】

【0065】

1 1 ... 用例翻訳システム、1 2 ... 演算制御装置、1 3 ... プロセッサ、1 4 ... メモリ、1 5 ... プログラム、1 6 ... 作業エリア、1 7 ... 表示装置、1 8 ... 入力装置、1 9 ... マウス、2 0 ... キーボード、2 1 a ... 読み取り装置、2 1 b ... 読み込み装置、2 2 ... ハードディスクドライブ、2 3 ... 入力処理部、2 4 ... 制御部、2 5 ... 記憶部、2 6 ... 出力処理部、2 7 ... 表示装置、2 8 ... 翻訳部、2 9 ... 用例検索部、3 0 ... 翻訳用例優先度判定部、3 1 ... 形態素解析部、3 2 ... 加点値計算部、3 3 ... 翻訳用例選出部、3 4 ... 翻訳用例データベース、3 5 ... 形態素解析辞書、3 6 ... コーパス

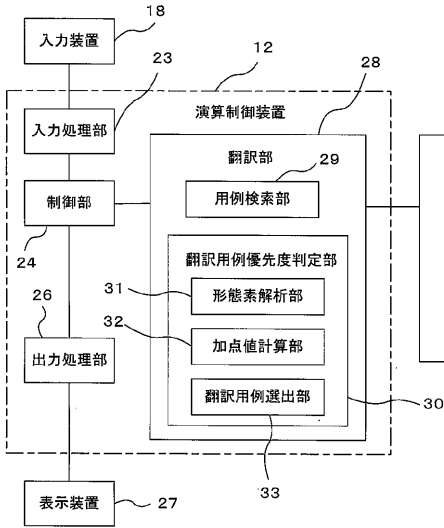
10

20

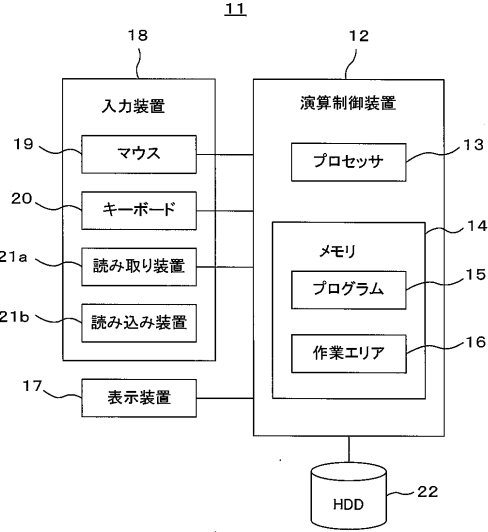
30

40

【図1】



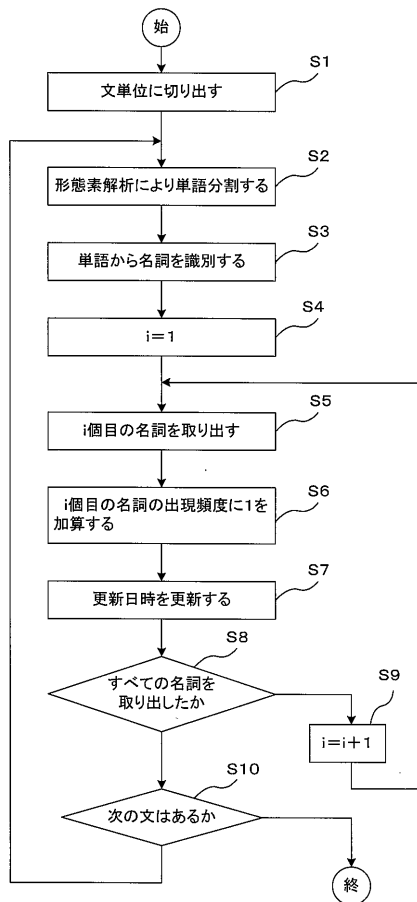
【図2】



【図3】

単語(名詞)	出現頻度	更新日時
株式	30	2010-10-29
証券	27	2010-10-09
銘柄	16	2010-10-01
配当	9	2010-09-09
株式会社	11	2010-10-29
社員権	3	2010-10-29
持分	9	2010-10-29
通常	6	2010-10-29
社員	4	2010-10-29
出資額	5	2010-10-29
役員	5	2009-07-23
平均	21	2009-08-31
不均一	4	2010-10-29
形態	12	2010-10-29
均一	7	2010-10-29
単語	2	2009-11-09
割合	26	2010-10-29
構成単位	9	2010-10-29
会員	3	2009-11-09
点	9	2010-10-29
特徴	13	2010-10-29
経営	14	2010-05-24
巨額	9	2010-09-03
少額	6	2010-01-15
⋮	⋮	⋮

【図4】



【 図 5 】

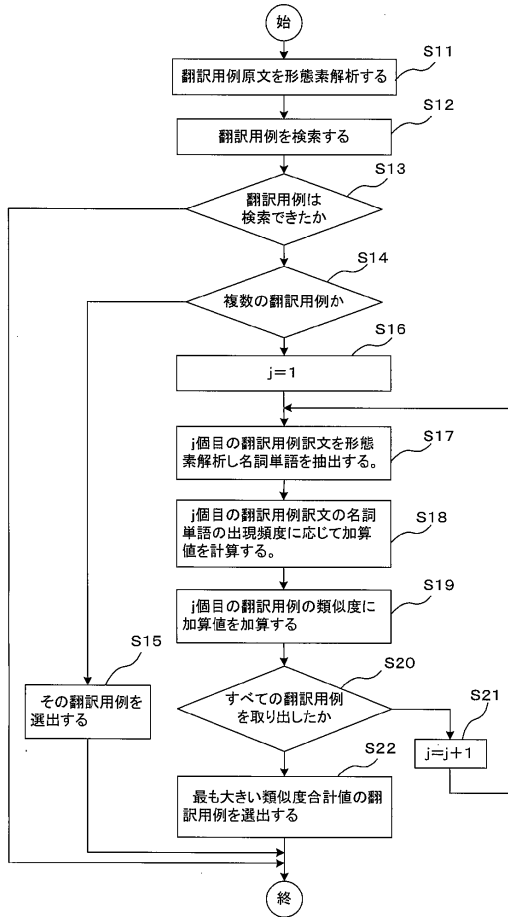
(a)

株式とは、株式会社における社員権、持分のことである。
 通常の持分が社員の出資額などに応じて不均一な形態を取るのに対して均一的な
 細分化された割合的な構成単位を取る点に特徴がある。

(b)

単語(名詞)	出現頻度	更新日時
株式	29+1	2010-10-29
証券	27	2010-10-09
銘柄	16	2010-10-01
配当	9	2010-09-09
株式会社	10+1	2010-10-29
社員権	2+1	2010-10-29
持分	7+1+1	2010-10-29
通常	5+1	2010-10-29
社員	3+1	2010-10-29
出資額	4+1	2010-10-29
役員	5	2009-07-23
平均	21	2009-08-31
不均一	3+1	2010-10-29
形態	11+1	2010-10-29
均一	6+1	2010-10-29
単語	2	2009-11-09
割合	25+1	2010-10-29
構成単位	8+1	2010-10-29
会員	3	2009-11-09
点	8+1	2010-10-29
特徴	12+1	2010-10-29
経営	14	2010-05-24
巨額	9	2010-09-03
少額	6	2010-01-15
⋮	⋮	⋮

【 図 6 】



【 図 7 】

単語(名詞)	出現頻度	更新日時
⋮	⋮	⋮
安全網	85	2005-03-28
⋮	⋮	⋮
セーフティネット	64	2010-02-19
⋮	⋮	⋮

フロントページの続き

(72)発明者 蔡 遠航

東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内

Fターム(参考) 5B075 ND03 ND20 PR04

5B091 BA04 BA13