

(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2005-216126

(P2005-216126A)

(43) 公開日 平成17年8月11日(2005.8.11)

(51) Int. Cl.⁷
G06F 17/28

F I
G O 6 F 17/28

テーマコード(参考)
5 B O 9 1

審査請求 有 請求項の数 14 O L (全 23 頁)

(21) 出願番号 特願2004-23913 (P2004-23913)
(22) 出願日 平成16年1月30日(2004.1.30)

(71) 出願人 301022471
独立行政法人情報通信研究機構
東京都小金井市貫井北町4-2-1
(74) 代理人 100130111
弁理士 新保 齋
(74) 代理人 100090893
弁理士 渡邊 敏
(72) 発明者 内元 清貴
東京都小金井市貫井北町4-2-1 独立
行政法人通信総合研究所内
(72) 発明者 井佐原 均
東京都小金井市貫井北町4-2-1 独立
行政法人通信総合研究所内
Fターム(参考) 5B091 AA05 CA21

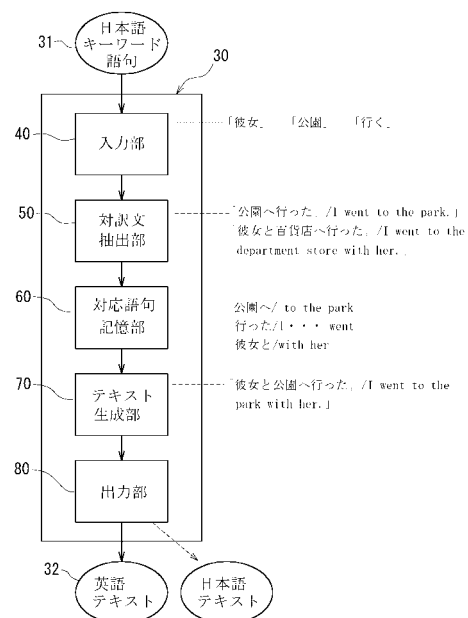
(54) 【発明の名称】 他言語のテキスト生成方法及びテキスト生成装置

(57) 【要約】

【課題】 使用者が適当なキーワード語句を与えることによりごく自然な他言語のテキスト生成を実現するテキスト生成方法及び装置を提供すること。

【解決手段】 原言語の単語をキーワード31として入力することにより、原言語・他言語間の対訳コーパスデータベースから抽出する対訳文抽出50し、該対訳文の部分対応情報から、各原言語のキーワード語句を含む原言語対応語句に対応する他言語の各他言語対応語句で構成する対応語句群テーブルを記憶60する。テキスト生成手段70では、該対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補32を生成する。

【選択図】 図7



【特許請求の範囲】**【請求項 1】**

原言語の語句をキーワード語句として入力することにより、原言語とは異なる他言語のテキストを生成する他言語テキスト生成方法であって、

入力手段から、単数又は複数の該原言語のキーワード語句を入力する入力ステップ、

対訳文中の語句間対訳関係に係る部分対応情報を含む原言語・他言語間の対訳コーパスデータベースを用い、対訳文抽出手段が、該キーワード語句を含む対訳文を、該対訳コーパスデータベースから抽出する対訳文抽出ステップ、

該対訳文の部分対応情報から、各原言語のキーワード語句を含む原言語対応語句に対応する他言語の各他言語対応語句で構成する対応語句群テーブルを記憶手段に記憶する対応語句記憶ステップ、

テキスト候補生成手段が、該対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補を生成するテキスト候補生成ステップ、

出力手段から、少なくとも1つのテキスト候補を出力する出力ステップ

の各ステップを含むことを特徴とする他言語テキスト生成方法。

【請求項 2】

前記他言語テキスト生成方法の対訳文抽出ステップにおいて、

前記入力ステップで入力したキーワード語句に対して、複数の対訳文が抽出され、部分対応情報から原言語対応語句が複数の種類存在するときに、

該対訳文抽出ステップの次に、

複数の原言語対応語句を使用者に選択可能に提示する原言語語句候補提示ステップを備え、

対応語句記憶ステップにおいて、

使用者が選択した場合に、その原言語対応語句に対応する他言語対応語句を対応語句群記憶テーブルに記憶する

ことを特徴とする請求項 1 に記載の他言語テキスト生成方法。

【請求項 3】

前記入力ステップにおいて、1個のキーワード語句を入力する毎に、

前記対訳文抽出ステップ及び、前記対応語句記憶ステップの各処理を行うと共に、

抽出された対訳文中において該キーワード語句と共起する共起語句を抽出し共起語句テーブルに記憶する共起語句抽出ステップ、

該共起語句テーブル中の共起語句を使用者に選択可能に提示する共起語句提示ステップの各ステップを備え、

該入力ステップにおいて、使用者が共起語句を選択した場合には、該共起語句を新たなキーワード語句として該入力ステップにおいて入力し、

全てのキーワード語句の入力が終了した後に、前記テキスト候補生成ステップに進む

ことを特徴とする請求項 1 又は 2 に記載の他言語テキスト生成方法。

【請求項 4】

前記他言語テキスト生成方法の対訳文抽出ステップにおいて、各処理に先だち、

入力ステップで入力されたキーワード語句について、構成する一部の形態素の加減、又は類語への置換を行う

ことを特徴とする請求項 1 ないし 3 に記載の他言語テキスト生成方法。

【請求項 5】

前記他言語テキスト生成方法において、

他言語テキストが複数の言語であって、

対訳文抽出ステップ、対応語句記憶ステップ、テキスト候補生成ステップにおいて、前記原言語と、各他言語との間についてそれぞれ処理を行い、

出力ステップにおいて、複数の言語のテキスト候補を出力する

ことを特徴とする請求項 1 ないし 4 に記載の他言語テキスト生成方法。

【請求項 6】

10

20

30

40

50

前記他言語テキスト生成方法のテキスト候補生成ステップにおいて、
テキスト候補生成手段が、該対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補を生成すると共に、
原言語テキスト候補生成手段が、該対応語句群テーブルに含まれる各原言語対応語句間の係り受け関係を仮定して原言語のテキスト候補を生成し、
出力ステップにおいて、
出力手段から、少なくとも1組の原言語及び他言語の対訳テキスト候補を共に出力することを特徴とする請求項1ないし5に記載の他言語テキスト生成方法。

【請求項7】

前記他言語テキスト生成方法において、
テキスト候補生成ステップの次に、
評価手段が、該テキスト候補を評価付けする評価ステップを有し、
出力ステップにおいては、該評価に基づいて少なくとも1つのテキスト候補を出力することを特徴とする請求項1ないし6に記載の他言語テキスト生成方法。

10

【請求項8】

原言語の単語をキーワードとして入力することにより、原言語とは異なる他言語のテキストを生成する他言語テキスト生成装置であって、
単数又は複数の該原言語のキーワード語句を入力する入力手段と、
対訳文中の語句間対訳関係に係る部分対応情報を含む原言語・他言語間の対訳コーパスデータベースと、
該キーワード語句を含む対訳文を、該対訳コーパスデータベースから抽出する対訳文抽出手段と、
該対訳文の部分対応情報から、各原言語のキーワード語句を含む原言語対応語句に対応する他言語の各他言語対応語句で構成する対応語句群テーブルを記憶可能な対応語句記憶手段と、
該対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補を生成するテキスト候補生成手段と、
少なくとも1つのテキスト候補を出力する出力手段と
を少なくとも備えたことを特徴とする他言語テキスト生成装置。

20

30

【請求項9】

前記他言語テキスト生成装置が、
入力したキーワード語句に対して前記対訳文抽出手段により複数の対訳文が抽出され、その部分対応情報から原言語対応語句が複数の種類存在するか否か判定し、複数の種類存在する場合には、使用者に該各原言語対応語句を提示する原言語語句候補提示手段を備えると共に、
前記入力手段から、使用者が提示された原言語対応語句の1個を選択可能であり、使用者が選択した場合には、前記対応語句記憶手段がその原言語対応語句に対応する他言語対応語句を対応語句群記憶テーブルに記憶する
請求項8に記載の他言語テキスト生成装置。

40

【請求項10】

前記他言語テキスト生成装置が、
入力手段から1個のキーワード語句を入力する毎に、前記対訳文抽出手段及び、前記対応語句記憶手段が作用する構成において、
抽出された対訳文中において該キーワード語句と共起する共起語句を抽出し共起語句テーブルに記憶する共起語句抽出手段と、
該共起語句テーブル中の共起語句を使用者に選択可能に提示する共起語句提示手段とを備え、
該入力手段から使用者が共起語句を選択した場合には、該共起語句を新たなキーワード語句として入力し、

50

全てのキーワード語句の入力が終了した後に、前記テキスト候補生成手段が作用することを特徴とする請求項 8 又は 9 に記載の他言語テキスト生成装置。

【請求項 1 1】

前記他言語テキスト生成装置において、

前記入力手段から入力されたキーワード語句について、構成する一部の形態素の加減、又は類語への置換を行うキーワード整形手段を備え、対訳文抽出手段において処理を行うことを特徴とする請求項 8 ないし 1 0 に記載の他言語テキスト生成方法。

【請求項 1 2】

前記他言語テキスト生成装置において、

対訳コーパスデータベースに、原言語と、複数の他言語との間の対訳文中の語句間対訳関係に係る部分対応情報を含み、

対訳文抽出手段と、対応語句記憶手段と、テキスト候補生成手段において、該原言語と、各他言語との間についてそれぞれ処理を行うと共に、

出力手段から、複数の言語のテキスト候補を出力する

ことを特徴とする請求項 8 ないし 1 1 に記載の他言語テキスト生成装置。

【請求項 1 3】

前記他言語テキスト生成装置において、

前記テキスト候補生成手段が、前記対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補を生成すると共に、

該対応語句群テーブルに含まれる各原言語対応語句間の係り受け関係を仮定して原言語のテキスト候補を生成する原言語テキスト候補生成手段を備え、

出力手段から、少なくとも 1 組の原言語及び他言語の対訳テキスト候補を共に出力することを特徴とする請求項 8 ないし 1 2 に記載の他言語テキスト生成装置。

【請求項 1 4】

前記他言語テキスト生成装置において、

前記テキスト候補を評価付けする評価手段を備えた

ことを特徴とする請求項 8 ないし 1 3 に記載の他言語テキスト生成装置。

【発明の詳細な説明】

【技術分野】

30

【0 0 0 1】

本発明は自然言語処理の方法及び装置に関する。特に、原言語の単数又は複数のキーワード語句から他言語のテキストを生成する手法に関わる。

【背景技術】

【0 0 0 2】

計算機を用いてテキストを解析、生成するための方法は従来から数多く提案されている。それらを大別すると、人間が作成した規則に基づく方法と統計的学習に基づく方法に分けることができる。前者の方法では、多様な知識を利用することで処理精度を向上させようとしてきた。一方、後者の方法では、単純な知識を大量に利用することで処理精度を向上させようとしてきた。

40

テキストを精度良く解析、生成するためには、文内、文間に現われる表層的情報から得られる様々な知識をはじめとして、辞書的な知識、言語学的な知見など、できるだけ多様な知識を利用するのが良いと考えられる。

しかし、前者の方法では、多様な知識を扱うためには規則を精緻化しなくてはならず、必然的に規則が競合しやすくなり、規則同士の優先順位を決めるのが困難になる。

一方、後者の方法では、多様な知識を利用しようとする学習データに過学習する傾向があるため、過学習を避けるためにさらに多くの学習データが必要となることが多い。後者の方法で多様な知識を利用することができればより良い精度が期待できる。しかし、後者の方法では、これまで知識を充実させるという方向の研究はほとんどなされてこなかった。

50

【0003】

本件発明者らは、後者の統計的学習に基づく方法を採用し、テキスト解析・生成のための新しいモデルを提案しており、例えば特許文献1において開示している。このモデルは、主に最大エントロピー原理に基づくもので、過学習の問題を避けつつ、多様な知識を効率良く扱うことができる。実験により、既存の統計的方法に比べて高い精度が得られることを示すとともに、学習データから得られる知識や、辞書的な知識、言語学的な知見などの多様な知識を効率的に利用する方法、および、テキスト解析・生成に有効な知識とはどのようなものであるかが明らかになっている。

【0004】

【特許文献1】特許公開2002-334076号公報

10

【0005】

一方、具体的なテキスト生成の処理方法としては、例えば本件出願人による特許文献2に開示されるテキスト生成のシステムがある。該システムでは、キーワードを入力してそれを含むテキストをデータベースから抽出し、該テキストを形態素解析・構文解析した後、もとのキーワードをテキストに組み合わせることでテキストの生成を行うように処理している。

また、特許文献3に開示されたシステムでは、キーワードとなる単語を入力して、文字単位候補を生成し、文字単位候補の係り受け関係を仮定してテキスト候補を生成するテキスト生成方法を開示している。本方法によると、キーワードが十分でない場合にも自然なテキストを生成できる長所がある。

20

【0006】

【特許文献2】特許公開2003-196280号公報

【特許文献3】特許公開2003-271592号公報

【0007】

これらはいずれも、例えば日本語のキーワードから日本語のテキストを生成するものであって、異なる言語のテキストを生成する手法ではない。すなわち従来技術では単言語のコーパスを用いて、単言語のキーワードからテキスト生成する方法が提供されているだけであり、上記特許文献3の方法を他言語に適用する方法は実現できていなかった。

【0008】

また、入力する言語と出力する言語が異なる言語処理としては機械翻訳が知られている。機械翻訳の一般的な手法は、翻訳元言語のテキストを入力し、それを解析、その解析結果から翻訳先言語を生成する。

30

しかし、入力時に必ずしもテキストを入力せず、適当なキーワードを与えることで、より自然なテキストを出力できるのであれば、使用者にとって他者とのコミュニケーションをより図りやすくなることも考えられる。

【0009】

例えば近年、ネットワークを通じて世界中の人々が容易に情報を交換できるようになったが、依然として言語バリアが存在しており、異文化間のコミュニケーションは容易ではない。これまでに、機械翻訳の技術は向上してきたが、商用の機械翻訳システムを用いてもなお異文化間のコミュニケーションは難しいということが指摘されている。

40

そこで、異文化間コミュニケーションにおける言語バリアを克服するために、システムに対する人間の協調をうまく引き出し、異文化間コミュニケーションを可能とするような他言語のテキスト生成方法が求められている。

【発明の開示】

【発明が解決しようとする課題】

【0010】

本発明は、上記従来技術の有する問題点に鑑みて創出されたものであり、その目的は、使用者が適当なキーワード語句を与えることによりごく自然な他言語のテキスト生成を実現するテキスト生成方法及び装置を提供することである。

【課題を解決するための手段】

50

【0011】

本発明は、上記の課題を解決するために、次のようなテキスト生成方法を創出する。

すなわち、請求項1に記載の発明は、原言語の語句をキーワード語句として入力することにより、原言語とは異なる他言語のテキストを生成する他言語テキスト生成方法である。

そして、入力ステップにおいて、入力手段から、単数又は複数の該原言語のキーワード語句を入力する。対訳文中の語句間対訳関係に係る部分対応情報を含む原言語・他言語間の対訳コーパスデータベースを用い、対訳文抽出手段が、入力されたキーワード語句を含む対訳文を、該対訳コーパスデータベースから抽出する対訳文抽出ステップ、該対訳文の部分対応情報から、各原言語のキーワード語句を含む原言語対応語句に対応する他言語の各他言語対応語句で構成する対応語句群テーブルを記憶手段に記憶する対応語句記憶ステップ、テキスト候補生成手段が、該対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補を生成するテキスト候補生成ステップを備える。

10

最後に出力ステップで、出力手段から、少なくとも1つのテキスト候補を出力する。

【0012】

請求項2に記載の発明は、上記他言語テキスト生成方法の対訳文抽出ステップにおいて、入力ステップで入力したキーワード語句に対して、複数の対訳文が抽出され、部分対応情報から原言語対応語句が複数の種類存在するときに、対訳文抽出ステップの次に、複数の原言語対応語句を使用者に選択可能に提示する原言語語句候補提示ステップを備える。

20

そして、対応語句記憶ステップにおいて、使用者が選択した場合に、その原言語対応語句に対応する他言語対応語句を対応語句群記憶テーブルに記憶することを特徴とする。

【0013】

請求項3に記載の発明は、入力ステップにおいて、1個のキーワード語句を入力する毎に、対訳文抽出ステップ及び、前記対応語句記憶ステップの各処理を行うと共に、抽出された対訳文中において該キーワード語句と共起する共起語句を抽出し共起語句テーブルに記憶する共起語句抽出ステップ、該共起語句テーブル中の共起語句を使用者に選択可能に提示する共起語句提示ステップの各ステップを備える。入力ステップにおいて、使用者が共起語句を選択した場合には、その共起語句を新たなキーワード語句として入力し、全てのキーワード語句の入力が終了した後に、テキスト候補生成ステップに進むように構成する。

30

【0014】

請求項4に記載の発明は、前記他言語テキスト生成方法の対訳文抽出ステップにおいて、各処理に先だち、入力ステップで入力されたキーワード語句について、構成する一部の形態素の加減、又は類語への置換を行うことを特徴とする。

【0015】

請求項5に記載の発明は、上記の他言語テキストが複数の言語であって、対訳文抽出ステップ、対応語句記憶ステップ、テキスト候補生成ステップにおいて、原言語と、各他言語との間についてそれぞれ処理を行うものである。例えば、他言語テキストとして、第1言語、第2言語、第3言語がある場合には、原言語と第1言語、原言語と第2言語、原言語と第3言語の間でそれぞれ上記処理を行う。これにより、出力ステップにおいては、全ての他言語のテキスト候補を出力するように構成する。

40

【0016】

請求項6に記載の発明は、テキスト候補生成ステップにおいて、テキスト候補生成手段が、該対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補を生成すると共に、原言語テキスト候補生成手段が、該対応語句群テーブルに含まれる各原言語対応語句間の係り受け関係を仮定して原言語のテキスト候補を生成する。

そして、出力ステップにおいて出力手段から、少なくとも1組の原言語及び他言語の対訳テキスト候補を共に出力する。

50

【0017】

請求項7に係る他言語テキスト生成方法は、テキスト候補生成ステップの次に、評価手段が、該テキスト候補を評価付けする評価ステップを有し、出力ステップにおいては、該評価に基づいて少なくとも1つのテキスト候補を出力する構成をとることができる。

【0018】

本発明は次のような他言語のテキスト生成装置を提供することもできる。

すなわち請求項8に記載の発明は、原言語の単語をキーワードとして入力することにより、原言語とは異なる他言語のテキストを生成する他言語テキスト生成装置であって、単数又は複数の該原言語のキーワード語句を入力する入力手段と、対訳文中の語句間対訳関係に係る部分対応情報を含む原言語・他言語間の対訳コーパスデータベースと、該キーワード語句を含む対訳文を、該対訳コーパスデータベースから抽出する対訳文抽出手段と、該対訳文の部分対応情報から、各原言語のキーワード語句を含む原言語対応語句に対応する他言語の各他言語対応語句で構成する対応語句群テーブルを記憶可能な対応語句記憶手段と、該対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補を生成するテキスト候補生成手段と、少なくとも1つのテキスト候補を出力する出力手段とを少なくとも備えたことを特徴とする。

10

【0019】

請求項9に記載の他言語テキスト生成装置は、入力したキーワード語句に対して前記対訳文抽出手段により複数の対訳文が抽出され、その部分対応情報から原言語対応語句が複数の種類存在するか否かが判定し、複数の種類存在する場合には、使用者に該各原言語対応語句を提示する原言語語句候補提示手段を備え、前記入力手段から、使用者が提示された原言語対応語句の1個を選択可能であり、使用者が選択した場合には、前記対応語句記憶手段がその原言語対応語句に対応する他言語対応語句を対応語句群記憶テーブルに記憶する構成を提供する。

20

【0020】

請求項10に記載の発明によると、入力手段から1個のキーワード語句を入力する毎に、前記対訳文抽出手段及び、前記対応語句記憶手段が作用する構成において、抽出された対訳文中において該キーワード語句と共に共起する共起語句を抽出し共起語句テーブルに記憶する共起語句抽出手段と、該共起語句テーブル中の共起語句を使用者に選択可能に提示する共起語句提示手段とを備える。そして、入力手段から使用者が共起語句を選択した場合には、該共起語句を新たなキーワード語句として入力し、全てのキーワード語句の入力が終了した後に、前記テキスト候補生成手段が作用することを特徴とする。

30

【0021】

請求項11に記載の発明は、前記入力手段から入力されたキーワード語句について、構成する一部の形態素の加減、又は類語への置換を行うキーワード整形手段を備え、対訳文抽出手段において処理を行うものである。

【0022】

請求項12に記載の発明によると、対訳コーパスデータベースに、原言語と、複数の他言語との間の対訳文中の語句間対訳関係に係る部分対応情報を含み、対訳文抽出手段と、対応語句記憶手段と、テキスト候補生成手段において、該原言語と、各他言語との間についてそれぞれ処理を行うと共に、出力手段から、複数の言語のテキスト候補を出力する。

40

【0023】

請求項13に記載の他言語テキスト生成装置は、前記テキスト候補生成手段が、前記対応語句群テーブルに含まれる各他言語対応語句間の係り受け関係を仮定して他言語のテキスト候補を生成すると共に、該対応語句群テーブルに含まれる各原言語対応語句間の係り受け関係を仮定して原言語のテキスト候補を生成する原言語テキスト候補生成手段を備え、出力手段から、少なくとも1組の原言語及び他言語の対訳テキスト候補を共に出力することを特徴とするものである。

【0024】

請求項14に記載の発明は、テキスト候補を評価付けする評価手段を備えたことを特徴

50

とする他言語テキスト生成装置を提供する。

【発明の効果】

【0025】

以上の発明により次の効果を奏する。

すなわち、請求項1又は8に記載のテキスト生成方法又は装置によると、キーワード語句を与えることによって、対訳コーパスから他言語のテキストを生成することができるので、自然な他言語を出力することができる。また、キーワードを入力することにより、原言語がテキストである場合に比して処理が容易であると共に、原言語テキストの解析誤りによる他言語テキストの誤りがなく、より正確なニュアンスのテキスト生成に寄与する。

【0026】

請求項2、3、9、10に記載のテキスト生成方法又は装置では、使用者に対してキーワード語句を提示することにより、使用者においては原言語で提示されるために理解が容易で指示が簡便に行える一方、本方法を備えた装置では正確なキーワード語句を用いて処理が行えるため、高精度な他言語テキストの生成が可能になる。

【0027】

請求項4、11に記載のテキスト生成方法又は装置によれば、対訳文抽出の際に、キーワード語句を変形させることにより、効率的な対訳文抽出処理が行える。この際、複数の形態素から成る場合には例えば語尾の助詞を削除したり、変形させたりして、対訳コーパス中に完全に一致するキーワード語句がなくとも抽出が行えるようにする。また、同義語、狭義語、広義語などの類語に置き換えることもできる。

【0028】

請求項5、12に記載のテキスト生成方法又は装置によれば、原言語1言語のキーワード語句を入力するだけで、同時に複数の言語のテキストを生成することができるので、効率の向上が図れるだけでなく、同時に多くの言語の使用者とのコミュニケーションにも寄与する。

【0029】

請求項6、13に記載のテキスト生成方法又は装置によると、原言語のテキスト候補を他言語のテキスト候補と共に出力することができるので、使用者が生成された他言語のテキストの意味を正確に把握することが可能になる。

【0030】

請求項7、14に記載のテキスト生成方法又は装置は、評価する処理を行うことにより、テキスト候補が複数ある場合にも、自動的に1個又は特定の候補数だけテキストを出力できる。例えば、後述の学習モデルによる確率値に応じて確率の高いものから所定数だけ順序付けして出力することもできる。

【発明を実施するための最良の形態】

【0031】

以下、本発明の最良と考えられる実施形態を、図面に示す実施例を基に説明する。なお、実施形態は下記に限定されるものではない。

まず、本発明の要部につき説述する。従来から母国語などを入力して異なる言語のテキスト(文章又はその集合)を出力する機械翻訳技術は知られており、近年高精度な機械翻訳が可能になりつつある。しかしながら、原言語のテキストを解析する過程と、他言語のテキストを生成する過程それぞれで、それぞれの言語が有する自然な言い回しや語順などが崩れてしまう場合があり、翻訳としては誤りではなくとも、コミュニケーションを図るために最適なテキストを得ることは難しい問題があった。

【0032】

また、機械翻訳の性能が十分に高くないと、原言語の入力時に機械翻訳に適する言い回しに直して入力しなければならなかったり、必要な言葉を過不足無く入力文に盛り込まなければならなかったりして、誰にでも簡便に使用することは難しい。一方で、インターネットの普及により世界中の誰とでも気軽にコミュニケーションをとれるようになった昨今において、正しいニュアンスの他言語を生成し、コミュニケーションを図れるような支援

10

20

30

40

50

方法の提供は急務である。

【 0 0 3 3 】

そこで本発明では母国語などのキーワード語句をいくつか入力することで、該キーワード語句の対訳語句を用いる他言語のテキストを生成する方法を創出した。使用者は母国語で伝えたい内容のうち重要な単語等を入力することにより、装置があらかじめ備えている対訳テキストのデータベースからそれらを用いる他言語テキストが生成される。その上、伝えている内容は原言語で確認できるため、使用者は正確なニュアンスの他言語テキストが生成されているかを確認することができる。

【 0 0 3 4 】

この方法で用いる対訳コーパスと呼ばれるデータベースは、原言語と他言語のそれぞれの文が対訳関係を持って格納されており、最初は人手によって正確な翻訳文を作成することが望ましい。そして、それぞれの文には構文情報も付与されており、句のレベルでの言語間の対応も付与されている。

【 0 0 3 5 】

本件発明者らが開発している対訳コーパスの1つとして、日本語と英語の対訳コーパスが完成しており、該コーパスは新聞記事を基にプロの翻訳家により作成したもので、日英文数は現在約4万である。

本コーパスは、英訳は日本文1文に対して1つの訳文(1文)とし、自然な英文に訳出している。日本文で主語が省略されている場合は、前文章の流れで必要に応じて主語を補い、主語に代名詞を持ってくるか、固有名詞かは前文からの自然な流れで決定する。このように作出するため、本コーパスは日本文・英文共に自然な言葉で表現されている。

【 0 0 3 6 】

コーパスのデータ形式を簡単に説述する。例えば日本文で、「また、一九九五年中の衆院解散・総選挙の可能性に否定的な見解を表明、二十日招集予定の通常国会前の内閣改造を明確に否定した。」に対して、図1のような依存構造木を定義し、依存構造木の左側に文節毎に付したIDを用いて

* 0 12D

また また * 接続詞 * * *

、 * 特殊 読点 * *

* 1 2D

一九九五 いちきゅうきゅうご * 名詞 数詞 * *

年 ねん * 接尾辞 名詞性名詞助数辞 * *

中 ちゅう * 接尾辞 名詞性名詞接尾辞 * *

の の * 助詞 接続助詞 * *

というように順に文節の番号、係り受け先、形態素、読み、品詞などを定義する。

【 0 0 3 7 】

さらに、この対訳文「He also responded negatively to the possibility of dissolution of the House of Representatives and general elections before the end of 1995, and clearly denied a cabinet reshuffle would take place prior to the ordinary Diet session scheduled to be convened on the 20th.」について、「

<P id="6,7">He<¥P> <P id="1">also<¥P> <P id="6,7">responded<¥P> <P id="5">negatively<¥P> <P id="4">to the possibility<¥P> <P id="3">of dissolution of the House of Representatives and general elections<¥P> . . . 」

と、上記日本文の文節IDをタグ(<P id=" " >と<¥P>で囲まれた部分)で表示しながら、各ワードの部分対応情報としている。

【 実施例 1 】

【 0 0 3 8 】

図2には本発明による第1の実施形態に係る他言語テキスト生成方法のフローチャートを示す。図のように、原言語(日本語)のキーワード語句(1)を入力し、そのキーワード語句(1)を含む対訳文を、対訳コーパスデータベース(10)から抽出(2)する。

10

20

30

40

50

そして、対訳文中からキーワード語句に関係する対応語句を部分対応情報(11)から抽出し、対応語句群テーブル(12)として記憶する。なお、該部分対応情報(11)は実際には対訳コーパスデータベース(10)中に含まれている情報であるから、両データは一体である。

ここまでの処理によって入力したキーワード語句に対応する他言語の語句が得られる。この後、これらの語句間の係り受け関係を仮定して他言語のテキスト候補を生成(4)する。

得られたテキスト候補はそのまま出力する構成でもよいが、本実施例ではこの後これらを一覧(5)し、候補の中から最も適当な他言語(英語)テキスト(6)を出力する。

【0039】

次に、本発明によるテキスト生成方法を実現するテキスト生成装置の構成を図7に示す。本装置(30)は、例えば「彼女」「公園」「行く」などの日本語キーワード語句を入力すると、入力部(40)で装置(30)内への取り込み処理を行い、対訳文抽出部(50)において「公園へ行った/I went to the park」「彼女と百貨店へ行った/I went to the department store with her」などの対訳文が抽出される。

【0040】

さらに対応語句記憶部(60)で、部分対応情報から上記対訳文の中でキーワード語句に関する「公園へ/to the park」「行った/I went ...」「彼女と/with her」などが抽出され、記憶する。

テキスト生成部(70)では、これらの対応語句から「I went to the park with her」という英語のテキストを生成し、出力部(80)から英語テキスト(32)を出力する。

次に各部(40)ないし(80)の詳細を説述する。

【0041】

入力部(40)は図8に示すようにCPU(41)とそれに接続されたマウス(42)やキーボード(43)、CDドライブ、ハードディスクドライブ、MOドライブ、フロッピー(登録商標)ディスクドライブなどの記憶装置(44)等から構成される。また、CPU(41)の動作に伴い、必要に応じて公知のメモリを用いることもできる。

使用者はマウス(42)やキーボード(43)により直接キーワード語句を入力することができる。

【0042】

また、本発明はインターネットやイントラネットのネットワーク(45)を介して他のコンピュータサーバー等からキーワード語句を受信することも可能である。

公知のタッチパネルモニタ(46)を設けてより簡便な入力方法を提供してもよい。

入力部(40)により日本語キーワード語句(31)は図9に示される対訳文抽出(50)・対応語句記憶(60)部に送られる。

【0043】

本実施例では、対訳文抽出(50)・対応語句記憶(60)部は1個の処理部(51)として図示する。ここでもCPU及びメモリが協働して各処理を行う。

まず対訳文抽出部(50)は外部記憶装置に格納された対訳コーパスデータベース(52)から日本語キーワード語句(31)を文中に含む対訳文を抽出する。

このとき、日本語キーワード語句(31)として使用者が形容詞や助詞を含めた場合や、複数のキーワード語句を1個のキーワード語句として入力した場合には、周知の処理方法によって基本形に変形したり、分割して複数のキーワード語句にしてもよい。この際、形態素解析等の言語処理方法が用いられることは公知である。

【0044】

もっとも本発明において対応語句記憶部(60)が最適な対応語句を抽出する上で、助詞や形容詞が重要な働きを果たす場合が多く、なるべくそれらを含めた形で対訳コーパスデータベース(52)から対訳文を抽出するのが望ましい。助詞は後述する係り受け関係を特定するのに有効であるし、形容詞が含まれることで対応語句の多義性の解消などに寄

10

20

30

40

50

与することも考えられる。

【0045】

また、上記対訳文抽出の際に、入力したキーワードに対応する対訳文が対訳コーパスデータベース(52)に見つけれない場合には、再び入力部(40)に処理を戻して、使用者に再入力を求めるようにしてもよい。或いは、シソーラスを用いて自動的に他のキーワード語句に置き換えるように構成してもよい。

【0046】

具体的には処理部(51)に図示しないキーワード整形部を設け、入力部(40)で入力されたキーワード語句を、整形処理する。該処理では、キーワード語句を公知の形態素解析処理により形態素に分割し、キーワード語句が複数の形態素から成る場合には、上記コーパスにおける接続助詞や格助詞を適宜削除したり、或いは対訳コーパス中に存在する形に合わせて加えたりする。形容詞に含む語尾、例えば「否定的な」の「な」を削除・変形させてもよい。

また、記憶手段にシソーラスを格納した上で、該キーワード語句の全形態素又は一部形態素を置換してもよい。

【0047】

次の対応語句記憶部(60)では、対訳文抽出部(50)で抽出された対訳文から、日本語のキーワード語句を含む日本語対応語句に対応する英語対応語句を、部分対応情報に基づいて抽出し、対応語句群テーブル(53)として記憶手段に記憶する。

すなわち、図7の例では「to the park」「I went ...」「with her」が記憶される。

【0048】

次に、以上により形成された対応語句群テーブル(53)を、図10に示すテキスト生成部(70)に入力し、英語テキストを生成する。

いくつかの語句を入力し、その語句を含むテキストを生成する方法としては次のような手法がある。すなわち、本件出願人が前記の特許文献3で開示するテキスト生成方法を、翻訳先言語である英語に適用して用いる。

【0049】

本テキスト生成部(70)の具体的な構成例として図10に示す各部を備える。テキスト生成部(70)は、例えばCPUとメモリ、ハードディスクなどの外部記憶媒体を備えるパーソナルコンピュータなどにより構成することができ、主な処理をCPUにおいて行い、処理の結果を随時メモリ、外部記憶媒体に記録する。

【0050】

本実施例で、入力された英語対応語句が、単語列ではなく単語列の主辞となる内容語である場合には、テキスト候補生成部(73)における処理の前に、単語列の候補を生成する。これは英語対応語句が内容語だけの場合、テキスト候補生成部(73)において係り受け関係を決定しただけではテキストが形成されない場合があるからである。

【0051】

該処理において、入力された英語対応語句(53)は2つの処理に用いられる。その1つは単語列生成規則獲得部(71)であり、もう1つは単語列候補生成部(72)である。以下では英語対応語句(53)のうち、単語列の主辞となる内容語であるものを特に英語対応単語と呼び、英語対応語句(53)が英語対応単語である場合には単語列候補生成部(72)で処理する一方、該当しない場合にはテキスト候補生成部(73)に英語対応語句(53)を送る。

内容語は、その語の品詞が、動詞、形容詞、名詞、指示詞、副詞、接続詞、連体詞、感動詞、未定義語である形態素の見出し語であるとし、それ以外の形態素の見出し語を機能語とする。

【0052】

単語列生成規則獲得部(71)では、英語対応単語が与えられたとき、それぞれを含む文を対訳コーパス(75)から検索し、形態素解析、構文解析(係り受け解析)をする。そして、そこから英語対応単語を含む単語列を抽出して、英語対応単語から英語対応語句

10

20

30

40

50

(53) を生成する単語列生成規則(76)を獲得し、記録する。このとき、対訳コーパス(75)を用いて、英語と日本語の対応付けをした単語列生成規則とするので、上記英語対応単語に対応する日本語の単語も同時に単語列として生成することができる。

例えば、「1995」「before the end of 1995 / 一九九五年中の」、「possibility」「to the possibility / 可能性に」などの単語列生成規則(76)を獲得し、記録する。

なお、ここでは英語対応単語に着目して英語と日本語の対応語句の組を生成したが、日本語キーワードから英語と日本語の対応語句の組を生成することも可能である。

【0053】

ここで、生成規則の自動獲得には次の手法を用いる。英語対応語句の集合をVとし、英語対応語句k(V)から単語列を生成する規則の集合をR_kとすると、規則rk(R_k)は次の形式で表現されるものと定義する。

$k \quad hk \quad m^*$

hk は英語対応語句を含む主辞形態素、m* は同じ単語列内でhkに連続する任意個の形態素とする。英語対応単語が与えられると、この形式を満たす規則を翻訳先言語のコーパス(75)から自動獲得する。

【0054】

一方、単語列候補生成部(72)では、単語列生成規則(76)を参照しながら、入力された英語対応語句(53)から出力する英語テキスト(32)を構成する単語列の候補を生成する。日本語テキストも同時に出力する場合には、このときに合わせて日本語対応語句についても単語列の候補を生成する。

例えば、「1995」では自然なテキストを構成する単語列とはなりにくい、「before the end of 1995」あるいは「in 1995」のように「1995」という単語と極めて密接な関連性を有する語句を付加し、後段の処理によるテキスト生成に備える。

【0055】

本実施例のように、単語列生成規則獲得部(71)により対訳コーパス(75)から入力する英語対応語句(53)(及び日本語対応語句)の単語列規則を生成することで、最小限の計算量で効果的に単語列生成規則を得ることができ、処理速度の向上に寄与する。

【0056】

もっとも、必ずしも英語対応語句(53)に関連する語句をコーパスから抽出する構成を取る必要はなく、計算能力に応じて任意の語句を入力された英語対応語句(53)の前後に付加してもよい。あるいは、別に対訳辞書データベースを備えて、それに含まれる慣用表現の情報から単語列を生成することもできる。上記「possibility」「to the possibility」などは対訳辞書データベースに記載される表現であり、単語列の候補として生成することができる。

【0057】

また、日本語など主格を多く省略する言語を入力した場合には、「respond」「He responded」などのように主語を補って単語列候補を生成することができる。このとき、日本語などの多くの言語では主格が明らかな時や、形式主語であるときに省略されることに着目し、入力に主格が何であるかの情報だけでなく、主格がないという情報を用いることで、「respond」「He responded」を生成せず、「respond」「It is responded that」を生成するようにすることもできる。

【0058】

次に、テキスト候補生成部(73)においてテキスト候補を生成する。テキスト候補はグラフあるいは木の形で表現する。ここでは英語対応語句(53)のうち、「to the park」「I went ...」「with her」の3語句の関係を例として説述する。

すなわち、図11のように、各英語対応語句(53a)(53b)(53c)の間に係り受けの関係を仮定して、テキスト候補(54)のような英語対応語句を単位とした依存構造木の形でテキスト候補を生成する。このとき、3語の場合に全ての係り受け関係は3! × 2 = 12通りであるが、翻訳先言語の文法・特性に合わせて語順の固定などにより候

10

20

30

40

50

補の数を削減することができる。

【 0 0 5 9 】

生成されたテキスト候補 (5 4) は、評価部 (7 4) でコーパスから学習した英語対応語句生成モデル (7 7) や言語モデル (7 8) を用いて順序付けされる。

以下、英語対応語句生成モデル (7 7) と、言語モデル (7 8) として形態素モデル及び係り受けモデルについて説述する。

【 0 0 6 0 】

英語対応語句生成モデルでは、次の 5 種類の情報を素性として用いたモデル (K M 1 ないし 5) を考える。以下で、英語対応語句の集合 V は、ある回数以上コーパスに出現した主辞単語の集合とし、単語列は前記で表現されるものと仮定する。また、各英語対応語句 k_i は単語 w_j ($1 \leq j \leq m$) に対応していると仮定する。図 1 2 にモデルの説明図を示す。

10

【 0 0 6 1 】

[K M 1]

前方の二単語を考慮 (trigram)

k_i は前方の二単語 w_{j-1} と w_{j-2} のみに依存すると仮定する。

【 数 1 】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_{j-1}, w_{j-2})$$

20

【 0 0 6 2 】

[K M 2]

後方の二単語を考慮 (後方 trigram)

k_i は後方の二単語 w_{j+1} と w_{j+2} のみに依存すると仮定する。

【 数 2 】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_{j+1}, w_{j+2})$$

【 0 0 6 3 】

30

[K M 3]

係り単語列を考慮 (係り単語列)

k_i を含む単語列に係る単語列がある場合、 k_i はそのうち最も文末側の単語列の末尾から二単語 w_l と w_{l-1} のみに依存すると仮定する (図 1 2 参照) 。

【 数 3 】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_l, w_{l-1})$$

【 0 0 6 4 】

40

[K M 4]

受け単語列を考慮 (受け単語列)

k_i を含む単語列を受ける単語列がある場合、 k_i はその単語列内の主辞単語から二単語 w_s と w_{s+1} のみに依存すると仮定する (図 1 2 参照) 。

【 数 4 】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_s, w_{s+1})$$

【 0 0 6 5 】

50

[K M 5]

係り単語列を最大二単語列考慮(係り二単語列)

k_i を含む単語列に係る単語列がある場合、 k_i は、そのうち最も文末側の単語列の末尾から二単語 w_l 、 w_{l-1} と、最も文頭側の単語列の末尾から二単語 w_h 、 w_{h-1} のみに依存すると仮定する(図 1 2 参照)。

【数 5】

$$P(K|M, D, T) = \prod_{i=1}^n P(k_i | w_l, w_{l-1}, w_h, w_{h-1})$$

10

【 0 0 6 6 】

次に、形態素モデル(MM)について示す。形態素に付与すべき文法的属性が1個あると仮定する。テキストつまり文字列が与えられたとき、その文字列が形態素であり、かつ $j(1 \leq j \leq l)$ 番目の文法的属性を持つとしたときの尤もらしさを確率値として求めるモデルを用いる。

テキスト T が与えられたとき、順序付き形態素集合 M が得られる確率は、各形態素 $m_i(1 \leq i \leq n)$ が独立であると仮定し、

【数 6】

$$P(M|T) = \prod_{i=1}^n P(m_i | m_1^{i-1}, T)$$

20

と表す。ここで、 m_i は1から l までのいずれかの文法的属性を表わす。

【 0 0 6 7 】

一方、係り受けモデル(DM)は、テキスト T と順序付き形態素集合 M が与えられたとき、各単語列に対する係り受けの順序付き集合 D が得られる確率は、各々の係り受け $d_1 \cdots d_n$ が独立であると仮定し、

【数 7】

$$P(D|M, T) = \prod_{i=1}^n P(d_i | M, T)$$

30

と表わす。

【 0 0 6 8 】

例えば、「to the park」「I went ...」「with her」の3つの英語対応語句(5 3)から「I went with her to the park.」と「I went to the park with her」の2つの候補が生成されたとする。係り受けモデルにより、このうち尤もらしい係り受け構造を持つ候補が優先される。

【 0 0 6 9 】

以上に示すような各モデルを用い、本発明では評価部(7 4)においてテキスト候補(5 4)に評価付けを行う。 40

評価部(7 4)では上記手法により句と句の依存関係や、形態素の並びとしての尤もらしさなどが考慮されるため、例えば英語における3単現の s の有無などについても、適切なものが評価値が高くなるので、文法的な正確さにも寄与する。

そして、評価値が最大あるいは閾値を超えるテキスト候補、あるいは評価値の上位 N 個を表層文に変換して出力する。

【 0 0 7 0 】

出力部(8 0)における出力方法としては、モニタによる表示の他、音声合成を用いた発声、翻訳システムなど他の言語処理システムへのデータ出力などが可能である。また、ネットワーク接続された他のコンピュータなどにテキストデータを送出してもよい。 50

【 0 0 7 1 】

本発明は、以上のように英語テキスト(32)を生成するものであるが、最後に文法的な補正処理を加えてもよい。すなわち、上記のように文法的にもある程度正しい出力が可能であるが、本方法による生成では時制の誤りや前置詞・主語の欠落などが生じる可能性もある。その場合、公知のOCR(光学的文字読み取り認識)技術における誤り修正の手法を適用することが考えられる。

【 0 0 7 2 】

英語側のテンス(時制)、(相:完了形、進行形などで表わされる)、モダリティ(法相:may, can, mustなどで表わされる)に不整合がある場合は、本件出願人らによる特許文献4に開示した方法などにより修正することができる。

10

例えば、「彼女と公園に行った」なら時制が過去と推定して、英語でも過去形を用いる、「彼女と公園に行ってきたところだ」なら完了形を用いる、「彼女と公園に行くだろう」なら、英語で may を用いる、というように間違った英語が選択された場合に修正する。

【 0 0 7 3 】

【特許文献4】特許第3388393号

【 0 0 7 4 】

また、三単現のsや前置詞の間違いなどは、例えば、非特許文献1に開示されるような文法的誤りのパターンを機械学習させ、誤りの検出を行う手法などにより、修正することができる。

20

【 0 0 7 5 】

【非特許文献1】「Automatic Error Detection in the Japanese Learners' English Spoken Data」, Emi Izumi, Kiyotaka Uchimoto, Toyomi Saiga, Thepchai Supnithi, Hito shi Isahara, Proceedings of the ACL2003 Interactive Poster/Demo Sessions pp.145-148, 2003

【 実施例 2 】

【 0 0 7 6 】

本発明の第2の実施形態として、図3にフローチャートを示す処理がある。すなわち、日本語キーワード語句(1)を入力して対訳文を抽出(2)した際、複数の対訳文が抽出され、その部分対応情報から日本語対応語句が複数の種類存在する場合に、日本語キーワード語句の絞り込み処理(20)を行うようにする。

30

図9に従って説述すると、対訳文抽出部(50)で日本語キーワード(31)を含む対訳文を対訳コーパスデータベース(52)から抽出する。例えばキーワード語句として「彼女」を入力したとき、複数の対訳文中に「彼女が」「彼女と」「彼女に」が日本語対応語句として抽出されることがある。

本実施例に係る日本語語句候補提示部(61)は、これらの日本語対応語句を使用者にすべて提示し、使用者はいずれの日本語対応語句がキーワード語句として最適であるか選択するようにする。

【 0 0 7 7 】

選択にはマウス(62)、キーボード(63)などを用い、使用者への提示はモニタ(64)で表示する。また、タッチパネルモニタ(65)を用いて優れたユーザインタフェースを提供することもできる。

40

本実施形態では、同様に「公園」と入力した場合には「公園へ / to the park」「公園で / in the park」を、「行く」の場合には「行く / I will go」「行った / I went ...」などを候補とし提示する。このように使用者がキーワード語句を入力するたびに対訳コーパスデータベース(10)から選択できる対応語句を提示することで、使用者の介入を容易にしながら、より適切なテキスト生成を図るようにする。

【 0 0 7 8 】

さらに、周知の文字入力方法としてローマ字や読み仮名の最初の1文字を入力した時点から順にその文字から始まる単語列を表示する手法がある。これを本実施形態に適用する

50

と、例えばkと入力した時点で「彼は」「彼女は」「今日」・・・などが表示され、kanまで入力すると「彼女と」「彼女が」・・・と絞られるようになる。対訳コーパスデータベースからこれらの候補を漸次抽出するのが処理上困難である場合には、適当な辞書データベースを別に設けて該辞書で単語のレベルまで絞りかけた後に、対訳コーパスデータベースから日本語対応語句を抽出すると良い。

【実施例3】

【0079】

本発明の第3の実施形態として、図4にフローチャートを示す処理がある。ここでは、対訳文を抽出(2)した際に、該対訳文においてキーワード語句と共起する語句を抽出(21)する。抽出された共起語句は使用者に提示(22)し、使用者が選択した共起語句は新たなキーワード語句(1)として追加する。 10

図13に示すように、日本語キーワード語句で「彼女」「公園」と入力した時点で、対訳文抽出部(50)が「彼女と公園へ行った/I went to the park with her」を抽出し、共起語句抽出部(66)は「彼女と」「公園へ」と共起する語句として「行った」を抽出する。このような共起語句の抽出方法は公知である。

【0080】

そして、共起語句提示部(67)でモニタ(64)等から使用者に対して「行った」を提示し、使用者がそれをキーワード語句とするのが適当と判断した場合にはマウス(62)(63)から選択することによりこれを新たなキーワードとして再び対訳文抽出部(50)に入力するか、対応語句記憶部(60)において「公園へ/to the park」を対応語句群テーブル(53)に記憶する。 20

前者の場合にはさらに選択した共起語句と共起する語句を選択することができるが、対訳文の数が膨大になる可能性があるため、後者の方法でもよい。

【実施例4】

【0081】

本発明の実施形態4に係る構成は、図5に示すように、日本語キーワード語句を入力すると、同時に2つの言語についてテキスト生成を行うテキスト生成方法である。

すなわち、図示の例では日英対訳コーパスデータベース(10a)と日本語対語対訳コーパスデータベース(10b)を用いてそれぞれについて対訳文抽出(2a)(2b)、部分対応情報(11a)(11b)を用いた対訳語句記憶(3a)(3b)、得られた対応語句群記憶テーブル(12a)(12b)からテキスト候補生成(4a)(4b)、評価(5a)(5b)を行い、英語テキスト(6a)、タイ語テキスト(6b)を同時に出力する。 30

これらの各方法において、上記実施例1ないし3で述べたような処理方法を導入してもよい。

本構成では、複数の言語テキストを同時に出力できるため、ネットワーク上において複数の言語の使用者が共存する場合などに特に好適である。

【実施例5】

【0082】

第5の実施形態は、テキスト候補生成において、日本語テキスト候補と英語テキスト候補を同時に生成し、使用者に生成された他言語の内容把握を容易にするものである。 40

図6に示すように、対応語句を記憶(3)する際に、対応語句群テーブル(12)に日英の対訳語句を共に記憶しておき、英語テキスト候補生成(4)に合わせて日本語テキスト候補を生成(23)する。両言語における係り受け関係を対応させておくことにより、生成された両テキストは同内容の対訳テキストが得られていると考えられるため、これらを使用者に提示することで、使用者は日本語による生成内容の確認を行うことができる。

【0083】

また、日本語テキスト候補の中から適切な係り受け関係になっているものを使用者が選択するようにすることで、係り受け関係を特定することができるため、英語テキスト候補の中から、係り受け関係が正しくかつ自然なテキストを得ることができる。 50

【実施例 6】

【0084】

以上説述した実施例は、いずれも日本語キーワード語句を直接入力するものであったが、本発明を次のようなシステムに実装して利用することができる。すなわち、本システムでは図14に示すように、ユーザーは日本語テキストを入力する。例えば、「彼女は公園へ行った」と入力部(40') (前記入力部(40)と入力する対象がテキストである他は同様の構成である)で入力すると、次のようなキーワードの抽出処理を図15の構成図におけるキーワード抽出部(90)で行う。

【0085】

キーワード抽出部(90)の構成を図5に示す。ここでもCPU及びメモリが協働して各処理を行う。キーワード抽出部(90)では、入力された日本語入力テキストからそのテキストの内容を特徴的に表すキーワードを抽出する。

このような技術は、言語処理において文書を要約する技術や、文書検索などの要素技術として公知の多数の手法が知られており、それらを適宜用いることができるが、ここでは一例として非特許文献2に記載の方法を用いる。

【0086】

【非特許文献2】情報処理学会自然言語処理研究会 1999-NL-133, 1999「タームのrepresentativeness」を測る 久光徹、丹羽芳樹、辻井潤一

【0087】

本方法によると、特徴語を選ぶために文書中の単語の話題性もしくは分野代表性(representativeness、本明細書ではこれを特徴性と呼ぶ。)を測ることが可能であり、かつ数値的な評価によるため、本発明の実施に好適である。以下に、簡単に説述する。

まず、キーワード抽出部(90)では、公知の形態素解析技術を用いて、日本語テキストを形態素解析部(91)において形態素解析する。解析された形態素はメモリ又は図示しない外部記憶装置などに形態素テーブルとして記録する。

【0088】

そして、形態素テーブルから形態素を順次読み出し、その形態素(以下、これを着目タームと呼ぶ)毎に特徴性を測る。

まず文書抽出部(92)において、着目タームWについて、Wを含む文書すべてを任意の文書データベース(93)から抽出する。文書データベース(93)は複数の日本語(翻訳元言語)の文書が含まれたものであり、外部記憶装置などに記憶されている。日本語単言語のコーパスや日英の対訳コーパスの日本語部分を用いてもよい。

【0089】

次に、着目タームWが抽出された文書すべての集合における単語分布と、文書データベース(93)に含まれる全文書の単語分布とを、単語分布算出部(94)において算出し、各単語分布間の異なり度を測る。

具体的には異なり度合算部(95)において次のような計算処理を行う。

【0090】

すなわち、着目タームW、Wを含む文書すべての集合D(W)、全文書の集合D₀、D(W)における単語分布P_{D(W)}、D₀における単語分布P₀として、Wの特徴性Rep(W)を、2つの分布{P_{D(W)}, P₀}の距離Dist{P_{D(W)}, P₀}に基づいて定義する。

単語分布間の距離計測方法として、本実施例では対数尤度比を用いている。すなわち、全単語を{W₁, ..., W_n}、単語w_iがD(W)、D₀に出現する頻度をそれぞれk_i、K_iとすると、P_{D(W)}、P₀の距離Dist{P_{D(W)}, P₀}を、次のように定義する。

【0091】

10

20

30

40

【数 8】

$$Dist(P_{D(W)}, P_0) = \sum_{i=1}^n k_i \log \frac{k_i}{\#D(W)} - \sum_{i=1}^n k_i \log \frac{K_i}{\#D_0}$$

ここで、 $\#D(W)$ は着目ターム W について $D(W)$ の含む単語数、 $\#D_0$ は同様に全文書の含む単語数である。

【0092】

数 8 の定義によると、 $\#D(W)$ が離れた着目ターム同士の特徴性を有効に比較することが難しいため、数 9 のように正規化を行った特徴性 $Rep(W)$ を定義する。なお $B(\cdot)$ は $\#D(W)$ が適切な数となる範囲内 (例えば $1000 \leq \#D(W) \leq 20000$) で特徴性が精度よく求められるような指数関数を用いた近似関数である。

【0093】

(数 9)

$$Rep(W) = Dist\{P_{D(W)}, P_0\} / B(\#D(W))$$

【0094】

ここで、「する」などのように著しく $\#D(W)$ が大きい場合には、 $D(W)$ の抽出数を限定し、 $\#D(W) \leq 20000$ を満たすようにすることで、上記近似関数を有効に用いることができると共に計算量を削減できる。

キーワード抽出部 (90) では以上の方法により特徴性を算出すると共に、所定の閾値に従って、キーワード決定部 (96) により入力した日本語入力テキストのキーワードを抽出する。

【0095】

ここで、例えば「彼女」「公園」「行く」がキーワードとして抽出されるので、上記実施例と同様に、対訳文抽出部 (50') により対訳コーパスデータベース (10) から対訳文を抽出する。上記では説明のため省略したが、このとき例えば「彼女は動物園へ行った。/ She went to the zoo.」なども同時に抽出されている。

そして、対訳語句記憶部 (60') も同様であり、テキスト生成部 (70') に進む。以上、各処理部 (40') (50') (60') (70') は前記実施例の (40) (50) (60) (70) と同態様の処理部であって、特記しない構成は同一である。

【0096】

前記実施例のテキスト生成部 (70) は図 10 に示すような構成であったが、ここでは例えば評価部 (74) で閾値を超えるテキスト候補を、実施例 5 のように複数の対訳文の形で出力し、最後に類似度評価部 (100) において、対訳文のうち日本語テキストと、最初に入力した日本語入力テキストの類似度を評価する。

類似度の評価方法としては、例えばテキストに含まれる文字列の一致する割合がどの程度であるかを算出して求める方法、あるいは非特許文献 3 に開示されるような自動翻訳した結果と人間の翻訳結果を文字列の単位 (或いは単語単位) で比較してその一致度を基に計算する方法などを用いることができる。

【0097】

【非特許文献 3】「Bleu: a Method for Automatic Evaluation of Machine Translation」, Kishore Papineni, Salim Roukos, Todd Ward, Wei-Jing Zhu. IBM Research Report, RC22176(W0109-022) 2001

【0098】

類似度評価部 (100) では、「彼女は公園へ行った」という入力テキストと、テキスト生成部 (70') で生成された「彼女と公園へ行った」「彼女は公園へ行った」の類似度を比較し、より類似度の高い「彼女は公園へ行った。/ She went to the park」の対訳文を出力部 (80') から出力することができる。

【0099】

10

20

30

40

50

以上、本発明の実施形態を1から6まで説述した。上記では説明の便宜のために、各部(40)(50)(60)(70)(80)を別個に説述したが、これらは一体的に例えば1台のパーソナルコンピュータによって提供することができる。特に、CPU、メモリ、入出力装置、ネットワークに接続するためのネットワークアダプタ(図示していない)、外部記憶装置などは共用することが望ましく、装置の簡略化に寄与することができる。

【0100】

外部記憶装置に記録される対訳コーパスデータベース(10)、コーパス(75)はいずれも同一のデータベースの一部又は全部を用いることが可能である。

また、これらは外部記憶装置上に記録される場合にとどまらず、ネットワーク上の複数のサーバーに記録されたものを収集するように構成してもよい。

10

【図面の簡単な説明】

【0101】

【図1】本発明で用いるコーパスの依存構造木の説明図である。

【図2】本発明の第1の実施形態に係るテキスト生成方法のフローチャートである。

【図3】本発明の第2の実施形態に係るテキスト生成方法のフローチャートである。

【図4】本発明の第3の実施形態に係るテキスト生成方法のフローチャートである。

【図5】本発明の第4の実施形態に係るテキスト生成方法のフローチャートである。

【図6】本発明の第5の実施形態に係るテキスト生成方法のフローチャートである。

【図7】本発明のテキスト生成装置の構成図である。

【図8】本発明における入力部の構成図である。

20

【図9】本発明における対訳文抽出・対応語句記憶部の構成図である。

【図10】本発明におけるテキスト生成部の構成図である。

【図11】英語対応語句からのテキスト生成の例を示す説明図である。

【図12】英語対応語句と単語列との関係を示す説明図である。

【図13】本発明におけるテキスト生成部(実施例3)の構成図である。

【図14】本発明の第6の実施形態に係るテキスト生成方法のフローチャートである。

【図15】本発明におけるテキスト生成部(実施例6)の構成図である。

【符号の説明】

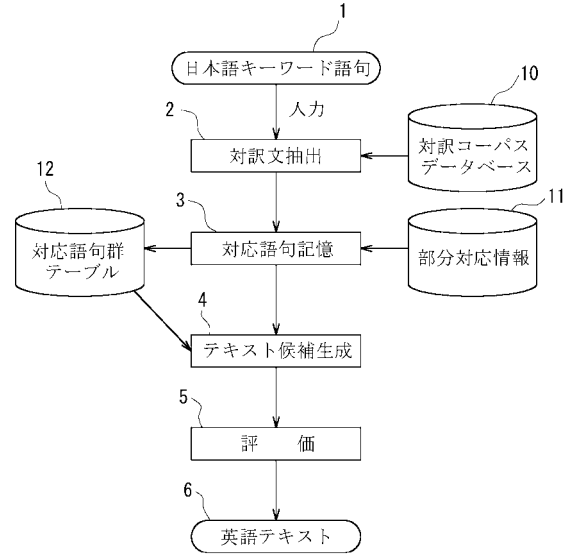
【0102】

| | | |
|----|------------|----|
| 30 | テキスト生成装置 | 30 |
| 31 | 日本語キーワード語句 | |
| 32 | 英語テキスト | |
| 40 | 入力部 | |
| 50 | 対訳文抽出部 | |
| 60 | 対応語句記憶部 | |
| 70 | テキスト生成部 | |
| 80 | 出力部 | |

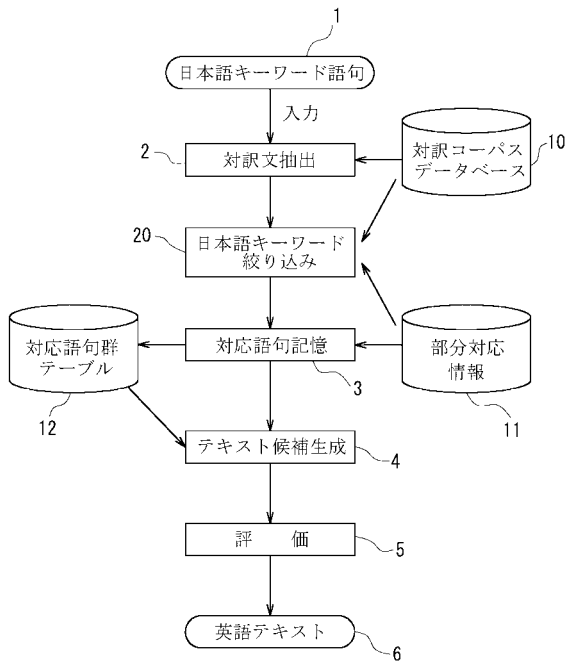
【 図 1 】

1 また、
 2 一九九五年中の「
 3 衆院解散・総選挙の「
 4 可能性に「
 5 否定的な「
 6 見解を「
 7 表明、
 8 二十日「
 9 召集予定の「
 10 通常国会前の「
 11 内閣改造を
 12 明確に
 13 否定した。

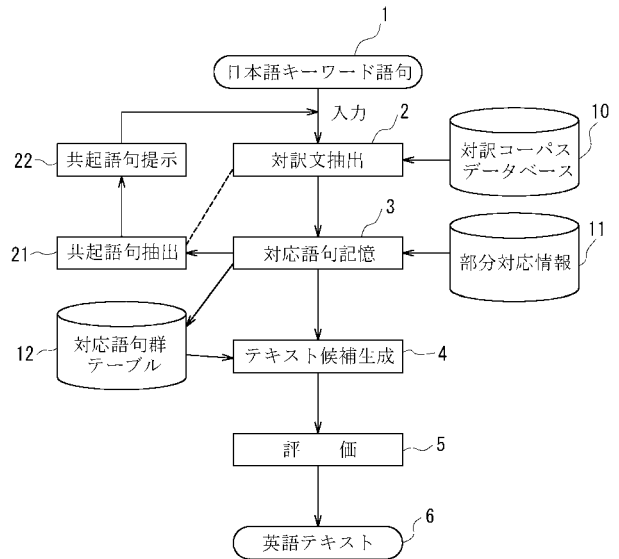
【 図 2 】



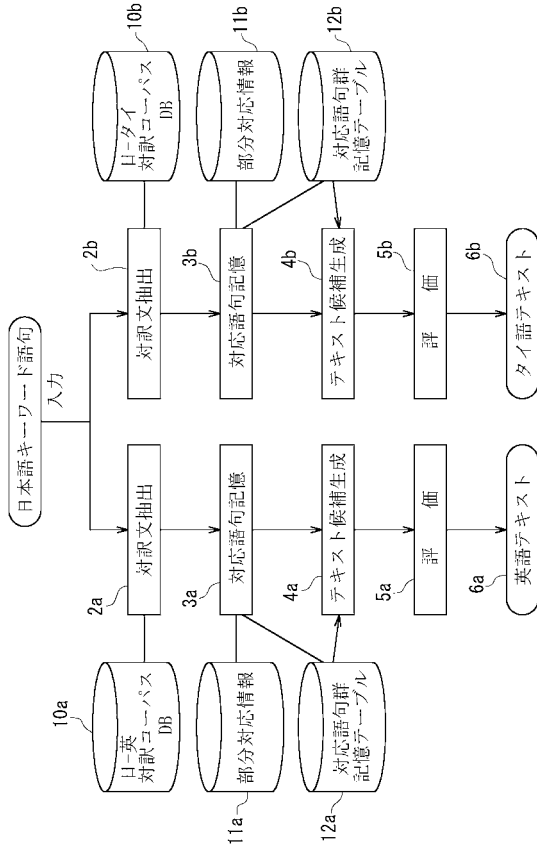
【 図 3 】



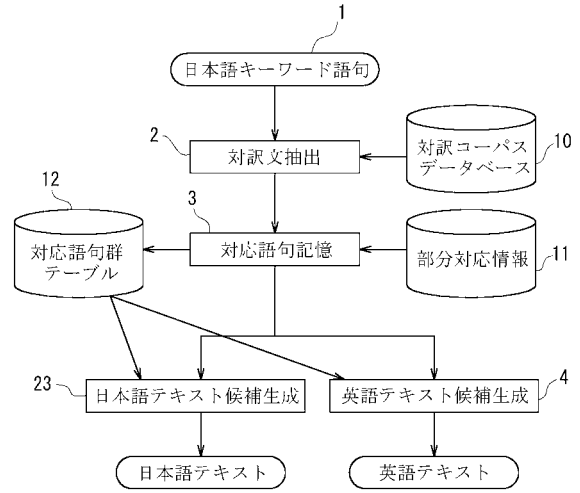
【 図 4 】



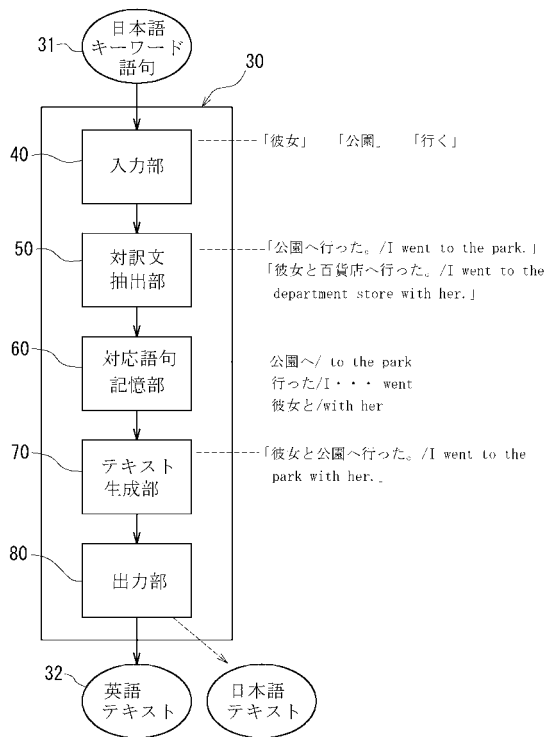
【 図 5 】



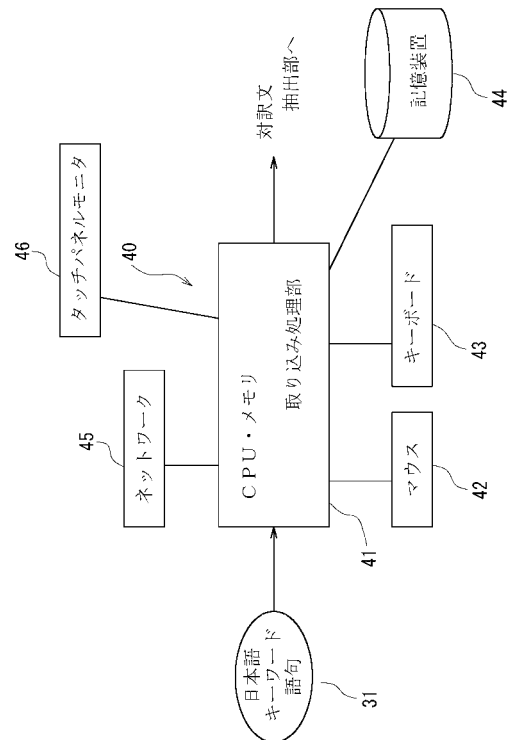
【 図 6 】



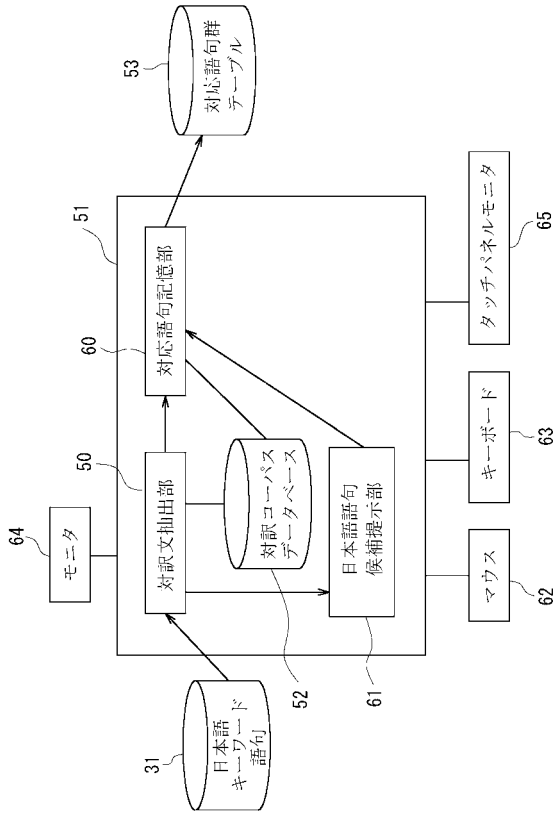
【 図 7 】



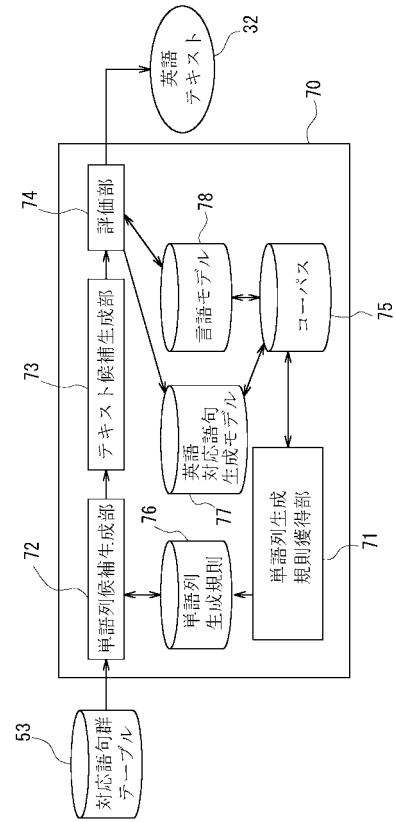
【 図 8 】



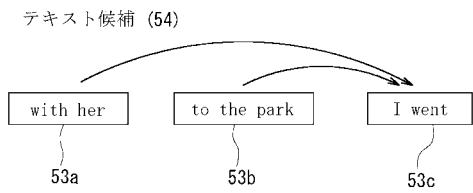
【 図 9 】



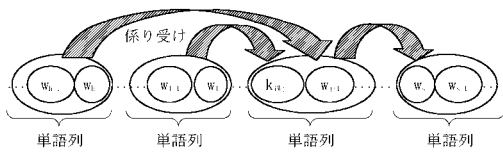
【 図 10 】



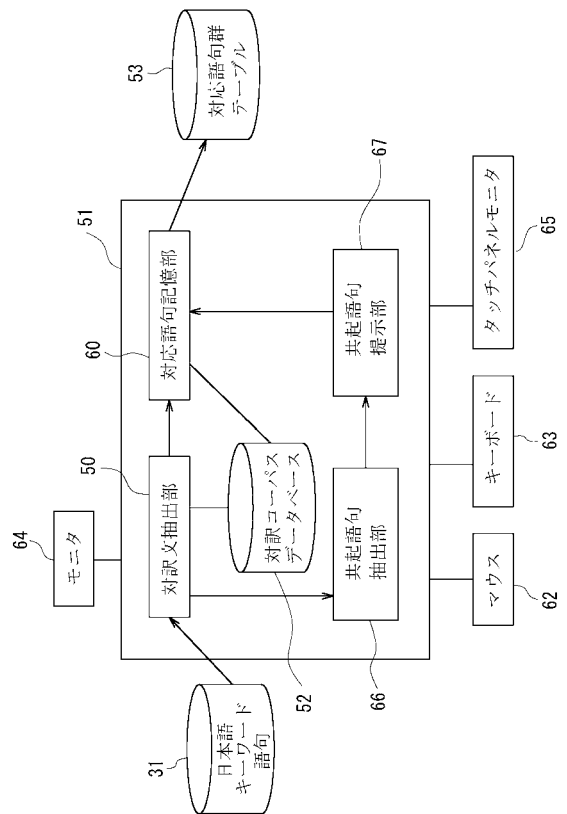
【 図 11 】



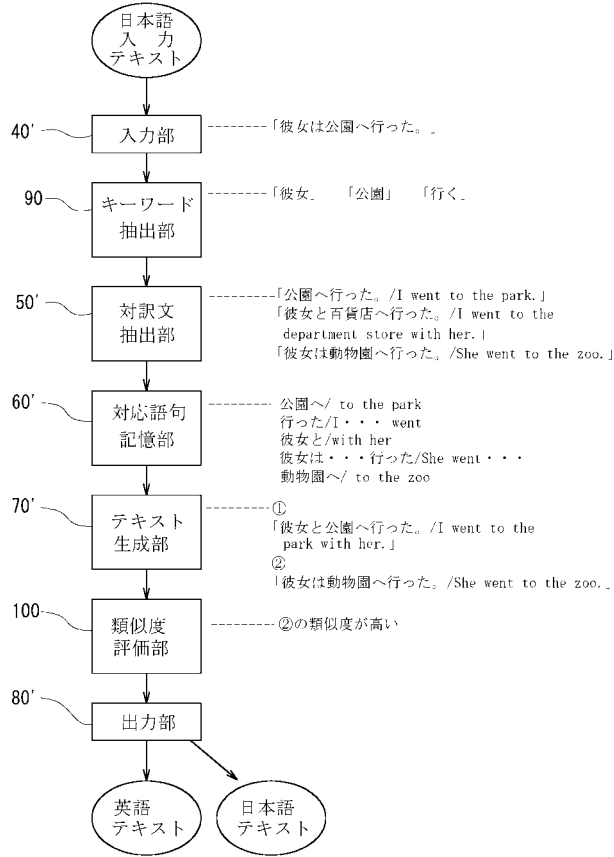
【 図 12 】



【 図 13 】



【 図 1 4 】



【 図 1 5 】

