

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7194759号
(P7194759)

(45)発行日 令和4年12月22日(2022.12.22)

(24)登録日 令和4年12月14日(2022.12.14)

(51)国際特許分類 F I
 G 0 6 F 40/44 (2020.01) G 0 6 F 40/44
 G 0 6 N 20/00 (2019.01) G 0 6 N 20/00 1 3 0

請求項の数 6 (全17頁)

(21)出願番号	特願2020-572078(P2020-572078)	(73)特許権者	392026693 株式会社NTTドコモ 東京都千代田区永田町二丁目11番1号
(86)(22)出願日	令和1年10月4日(2019.10.4)	(74)代理人	100088155 弁理士 長谷川 芳樹
(86)国際出願番号	PCT/JP2019/039337	(74)代理人	100113435 弁理士 黒木 義樹
(87)国際公開番号	WO2020/166125	(74)代理人	100121980 弁理士 沖山 隆
(87)国際公開日	令和2年8月20日(2020.8.20)	(74)代理人	100128107 弁理士 深石 賢治
審査請求日	令和3年6月8日(2021.6.8)	(74)代理人	100183438 弁理士 内藤 泰史
(31)優先権主張番号	特願2019-22411(P2019-22411)	(72)発明者	村上 聡一郎 東京都千代田区永田町二丁目11番1号
(32)優先日	平成31年2月12日(2019.2.12)		最終頁に続く
(33)優先権主張国・地域又は機関	日本国(JP)		

(54)【発明の名称】 翻訳用データ生成システム

(57)【特許請求の範囲】

【請求項1】

原言語テキストにノイズを付与してノイズ付与原言語テキストを得るノイズ付与部と、
 前記ノイズ付与原言語テキストと、該ノイズ付与原言語テキストのノイズ付与前の原言語テキストに対応する目的言語テキストとを対応付けた疑似対訳コーパスを構築するコーパス構築部と、

ノイズを含んだ原言語テキスト群である訓練データを用いて、原言語テキストの各単語の次にノイズが入る場合に各単語に対してノイズのタイプを示すノイズラベルを予測するように学習されたノイズモデルを学習するノイズモデル学習部と、を備え、

前記ノイズ付与部は、

前記ノイズモデルを用いて、原言語テキストの各単語の特徴に応じて、ノイズのタイプを示すノイズラベルを付与し、該ノイズラベルを該ノイズラベルに対応する単語へ置き換えることにより、原言語テキストにノイズを付与し、

前記ノイズラベルの付与について、原言語テキストの各単語の特徴を入力として前記ノイズモデルから出力される各ノイズラベルのスコアに基づく各ノイズラベルの確率分布に従ってノイズラベルをサンプリングし、原言語テキストに付与するノイズラベルを決定する、翻訳用データ生成システム。

【請求項2】

前記疑似対訳コーパスを用いて翻訳モデルを学習する翻訳モデル学習部を更に備える、請求項1記載の翻訳用データ生成システム。

【請求項 3】

前記ノイズ付与部は、1つの前記ノイズラベルに対して置き換える単語を複数パターン導出し、1つの原言語テキストから複数パターンの前記ノイズ付与原言語テキストを得る、請求項 1 又は 2 記載の翻訳用データ生成システム。

【請求項 4】

前記ノイズ付与部は、前記ノイズモデルを用いて、各単語に対応する前記ノイズラベルを複数パターン導出し、1つの原言語テキストから複数パターンの前記ノイズ付与原言語テキストを得る、請求項 1～3 のいずれか一項記載の翻訳用データ生成システム。

【請求項 5】

前記ノイズ付与部は、前記ノイズモデルを用いて、原言語テキストの各単語の特徴である、形態素、品詞、及び単語の読みの少なくとも一つに応じて、前記ノイズラベルを付与する、請求項 4 記載の翻訳用データ生成システム。

10

【請求項 6】

前記ノイズモデルは、条件付き確率場又はニューラルネットワークを用いた手法により構築されている、請求項 1～5 のいずれか一項記載の翻訳用データ生成システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の一態様は、翻訳用データ生成システムに関する。

【背景技術】

20

【0002】

機械翻訳システムにおいて、利用者の自然発話入力に、言い淀み、言い直し、又はフィルター等（以下、これらを総称して「ノイズ」と記載する場合がある）が含まれることによって、翻訳精度が低下する場合がある。

【0003】

このような課題に対して、例えば特許文献 1 及び特許文献 2 等に示されるように、発話における言い直し箇所等を特定し利用者の発話内容を修正する技術が知られている。

【先行技術文献】

【特許文献】

【0004】

30

【文献】特開 2010 - 079647 号公報
特開 2007 - 057844 号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、ノイズ箇所を特定して修正することは容易ではなく、上述した技術によっても翻訳精度を十分に担保することは困難である。

【0006】

本発明の一態様は上記実情に鑑みてなされたものであり、ノイズが含まれる自然発話に対しても高精度に翻訳を行うことを目的とする。

40

【課題を解決するための手段】

【0007】

本発明の一態様に係る翻訳用データ生成システムは、原言語テキストにノイズを付与してノイズ付与原言語テキストを得るノイズ付与部と、ノイズ付与原言語テキストと、該ノイズ付与原言語テキストのノイズ付与前の原言語テキストに対応する目的言語テキストとを対応付けた疑似対訳コーパスを構築するコーパス構築部と、を備える。

【0008】

本発明の一態様に係る翻訳用データ生成システムでは、原言語テキストにノイズが付与され、ノイズ付与原言語テキストとノイズ付与前の原言語テキストに対応する目的言語テキストとを対応付けた疑似対訳コーパスが構築される。このように、ノイズ付与原言語テ

50

キストがノイズ付与前の原言語テキストに対応する目的言語テキストに対応付けられた対訳コーパスが構築されることにより、このような対訳コーパスを利用して、例えば自然発話入力にフィラー等のノイズが含まれている場合においても、ノイズ付与前の原言語テキストに対応する目的言語テキストを適切に導出することが可能となる。すなわち、本発明の一態様に係る翻訳用データ生成システムによれば、ノイズが含まれる自然発話に対して頑健なコーパス（疑似対訳コーパス）を構築することができ、ノイズが含まれる自然発話に対しても高精度に翻訳を行うことができる。

【0009】

上記翻訳用データ生成システムは、疑似対訳コーパスを用いて翻訳モデルを学習する翻訳モデル学習部を更に備えていてもよい。構築したコーパスに基づいて翻訳モデルが学習されることにより、ノイズが含まれる自然発話に対してより高精度に翻訳を行うことができる。

10

【0010】

上記翻訳用データ生成システムは、ノイズを含んだ原言語テキスト群である訓練データを用いて、原言語テキストに対するノイズの付与に係るノイズモデルを学習するノイズモデル学習部を更に備え、ノイズ付与部は、ノイズモデルを用いて、原言語テキストにノイズを付与してもよい。予めノイズが含まれている原言語テキスト群に基づきノイズモデルが学習され、該ノイズモデルに基づいてノイズの付与が行われることによって、実際に含まれる可能性が高いノイズが付与され易くなり、翻訳精度をより向上させることができる。

【0011】

上記翻訳用データ生成システムにおいて、ノイズ付与部は、原言語テキストの各単語に、ノイズのタイプを示すノイズラベルを付与し、該ノイズラベルを該ノイズラベルに対応する単語へ置き換えることにより、原言語テキストにノイズを付与してもよい。原言語テキストの各単語に応じたノイズラベルが付与された後に該ノイズラベルに応じた単語（ノイズ）が導出されることにより、ノイズ付与の容易性及び妥当性を担保することができる。

20

【0012】

上記翻訳用データ生成システムにおいて、ノイズ付与部は、1つのノイズラベルに対して置き換える単語を複数パターン導出し、1つの原言語テキストから複数パターンのノイズ付与原言語テキストを得てもよい。これにより、1つの原言語テキストから効率的に疑似対訳コーパスを充実させ、翻訳精度をより向上させることができる。

30

【0013】

上記翻訳用データ生成システムにおいて、ノイズ付与部は、各単語に対応するノイズラベルを複数パターン導出し、1つの原言語テキストから複数パターンのノイズ付与原言語テキストを得てもよい。これにより、1つの原言語テキストから効率的に疑似対訳コーパスを充実させ、翻訳精度をより向上させることができる。

【0014】

上記翻訳用データ生成システムにおいて、ノイズ付与部は、原言語テキストの各単語の特徴に応じて、ノイズラベルを付与してもよい。これにより、各単語に関連して含まれやすいノイズに係るノイズラベルを、各単語に適切に付与することができる。

【0015】

上記翻訳用データ生成システムにおいて、ノイズ付与部は、原言語テキストの各単語の特徴である、形態素、品詞、及び単語の読みの少なくとも一つに応じて、ノイズラベルを付与してもよい。これにより、各単語に関連して含まれやすいノイズに係るノイズラベルを、各単語に適切に付与することができる。

40

【0016】

上記翻訳用データ生成システムにおいて、ノイズ付与部は、原言語テキストの各単語の特徴を入力としてノイズモデルから出力される各ノイズラベルのスコアに基づく各ノイズレベルの確率分布に従ってノイズラベルをサンプリングし、原言語テキストに付与するノイズラベルを決定してもよい。これにより、例えばノイズモデルから出力されたスコアが高いノイズラベルを付与することが可能となり、各単語に関連して含まれやすいノイズに

50

係るノイズラベルを、各単語に適切に付与することができる。

【0017】

上記翻訳用データ生成システムにおいて、ノイズモデルは、条件付き確率場又はニューラルネットワークを用いた手法により構築されていてもよい。これにより、機械学習によってノイズモデルを適切に構成することができる。

【発明の効果】

【0018】

本発明の一態様によれば、ノイズが含まれる自然発話に対しても高精度に翻訳を行うことができる。

【図面の簡単な説明】

10

【0019】

【図1】本実施形態に係る翻訳用データ生成システムの処理イメージを模式的に示す図である。

【図2】本実施形態に係る翻訳用データ生成システムの機能構成を示す図である。

【図3】ノイズモデルの概要を説明する図である。

【図4】ノイズラベルを説明する図である。

【図5】疑似対訳コーパスの構築イメージを示す図である。

【図6】翻訳用データ生成システムが実行する処理を示すフローチャートである。

【図7】本実施形態及び比較例の翻訳例を示す表である。

【図8】翻訳用データ生成装置のハードウェア構成を示す図である。

20

【発明を実施するための形態】

【0020】

以下、添付図面を参照しながら本発明の実施形態を詳細に説明する。図面の説明において、同一又は同等の要素には同一符号を用い、重複する説明を省略する。

【0021】

最初に、図1を参照して、本実施形態に係る翻訳用データ生成システム1の処理イメージを説明する。図1は、本実施形態に係る翻訳用データ生成システム1の処理イメージを模式的に示す図である。翻訳用データ生成システム1は、既存の対訳コーパス（一般的に用いられる対訳コーパス）の原言語側のテキスト（原言語テキスト）に対してノイズを付与すると共に、ノイズが付与された原言語テキスト（ノイズ付与原言語テキスト）と、ノイズ付与前の原言語テキストに対応する目的言語側のテキスト（目的言語テキスト）とを対応付けた疑似対訳コーパスを構築し、該疑似対訳コーパスを用いて機械翻訳モデル（例えばNMT（Neural Machine Translation）モデル）を学習（構築）するシステムである。ここでのノイズとは、利用者の自然発話入力に含まれ得る言い淀み、言い直し、又はフィルター等である。

30

【0022】

図1に示される例では、既存の対訳コーパスにおける「主要な高速道路よりも観光ルートの方を走りたいです」との原言語テキストに対して、所定のルールに従って複数パターンのノイズを付与し（詳細は後述）、「えー主要な高速道路よりもまー観光ルートの方を走りたいです」「あ主要な高速道路よりもえー観光ルートの方を走りたいです」「えーっと主要な高速道路よりもまー観光ルートの方を走りたいです」という3パターンのノイズ付与原言語テキストを得ている。そして、ノイズ付与前の原言語テキストに対応する目的言語テキストである「I would rather take a scenic route than a main highway.」と、上述した3パターンのノイズ付与原言語テキストとを対応付けた疑似対訳コーパスが構築されて、該疑似対訳コーパスを用いて機械翻訳モデルが学習（構築）されている。このように、ノイズ付与原言語テキストがノイズ付与前の原言語テキストに対応する目的言語テキストに対応付けられた疑似対訳コーパスが構築されることにより、このような疑似対訳コーパスを利用して、例えば自然発話入力にフィルター等のノイズが含まれている場合においても、ノイズ付与前の原言語テキストに対応する目的言語テキストを適切に導出することが可能となる。以下、翻訳用データ生成システム1の機能の詳細について説明する。

40

50

【 0 0 2 3 】

図 2 は、本実施形態に係る翻訳用データ生成システム 1 の機能構成を示す図である。図 2 に示されるように、翻訳用データ生成システム 1 は、翻訳用データ生成装置 1 0 と、対訳コーパス DB 2 0 と、訓練情報 DB 3 0 と、ノイズモデル学習装置 4 0 (ノイズモデル学習部) と、翻訳モデル学習装置 5 0 (翻訳モデル学習部) と、を備えている。なお、翻訳用データ生成システム 1 は、必ずしも上記の各構成を備えるものでなくてもよく、例えば翻訳用データ生成装置 1 0 のみで構成されていてもよいし、翻訳用データ生成装置 1 0 及びノイズモデル学習装置 4 0 のみで構成されていてもよいし、翻訳用データ生成装置 1 0、ノイズモデル学習装置 4 0、及び翻訳モデル学習装置 5 0 のみで構成されていてもよいし。

10

【 0 0 2 4 】

対訳コーパス DB 2 0 は、対訳コーパスを記憶しているデータベースである。対訳コーパスとは、原言語テキスト及び目的言語テキストの組み合わせを構造化したものである。対訳コーパス DB 2 0 が記憶する対訳コーパスは、通常利用されるものでよく、例えば K F T T (Kyoto Free Translation Task) 又は B T E C 等の日本語・英語の対訳コーパスである。本実施形態では、翻訳用データ生成装置 1 0 によって、対訳コーパス DB 2 0 が記憶する対訳コーパスの原言語テキストにノイズが付与され、疑似対訳コーパスが生成される(詳細は後述)。

【 0 0 2 5 】

訓練情報 DB 3 0 は、ノイズモデル(詳細は後述)を学習するための訓練情報(訓練データ)を記憶しているデータベースである。訓練情報とは、予めノイズがアノテーションされた原言語テキスト群(自然発話の書き起こしコーパス。学習用発話データ)である。このような訓練情報は、例えば通常のコーパスに含まれる原言語テキストにノイズがアノテーションされることによって構築されている。

20

【 0 0 2 6 】

ノイズモデル学習装置 4 0 は、訓練情報 DB 3 0 に記憶されている訓練情報(ノイズを含んだ原言語テキスト群である訓練データ)を用いて、原言語テキストに対するノイズの付与に係るノイズモデルを学習する。ノイズモデルの学習データ(訓練データ)としては、例えば、日本語はなし言葉コーパス(C S J)又は S w i t c h B o a r d C o r p u s 等の自然発話コーパスの書き起こしコーパスが用いられてもよい。ノイズモデルは、原言語テキストが入力された場合に、該原言語テキストに係るノイズラベルの情報を出力するものである。ノイズラベルとは、ノイズのタイプ(種別)を示す情報である。図 4 は、ノイズラベルを説明する図である。図 4 に示されるように、本実施形態では、ノイズラベルとして、< F >、< D >、0 の 3 種類がある。< F > は、フィルターを示すノイズラベルである。< D > は言い淀み又は言い直しを示すノイズラベルである。0 はノイズ無しを示すノイズラベルである。ノイズラベルの情報とは、ノイズラベルの種類(上述した< F >、< D >、0)と各ノイズラベルが対応付けられる単語(詳細には形態素)とが紐づいた情報であり、例えば後述するノイズラベル系列である。

30

【 0 0 2 7 】

図 3 は、ノイズモデルの概要を説明する図である。図 3 に示されるように、ノイズモデルは、例えば、品詞タグ付けや固有表現抽出タスク等で広く用いられている双方向再帰的ニューラルネットワーク(B i R N N : Bi-directional Recurrent Neural Networks)を用いて構築されている。なお、ノイズモデルは、R N N 等のその他のニューラルネットワークを用いた手法や、C R F (Conditional random field)等の条件付き確率場を用いた手法により構築されていてもよい。ノイズモデルは、入力された原言語テキストの各入力要素(単語詳細には形態素)の次にノイズが入る場合、その要素に対して適当なノイズラベルを予測するように学習されている。ノイズモデルを用いたノイズ付与においては、原言語テキストの形態素系列 $w = (w_0, w_1, \dots, w_n)$ からノイズラベル系列 $l = (l_0, l_1, \dots, l_n)$ を予測する系列ラベリング問題として考える。

40

【 0 0 2 8 】

50

いま、「 $\langle F \text{えー} \rangle$ それでは会議を $\langle D \text{を} \rangle$ 始め $\langle F \text{あー} \rangle$ ます」という学習用発話データを例にノイズモデルを学習する方法を説明する。ここで、学習用発話データに含まれる「 $\langle F \text{えー} \rangle$ 」は「えー」がフィルター $\langle F \rangle$ に相当することを表している。この場合、まず、学習用発話データから形態素系列 $w = (\langle B O S \rangle, \text{それでは}, \text{会議}, \text{を}, \text{始め}, \text{ます}, \langle E O S \rangle)$ が抽出される。図3に示されるように形態素系列は、 $t = 0 \sim 6$ までのタイムステップに各形態素($\langle B O S \rangle$ 及び $\langle E O S \rangle$ を含む)が対応している。次に、同じ学習用発話データであってノイズがアノテーションされた情報に基づいて、ノイズラベル系列 $l = (\langle F \rangle, 0, 0, \langle D \rangle, \langle F \rangle, 0, 0)$ が生成される。最後に、形態素系列 w からノイズラベル系列 l を予測する系列ラベリング問題としてBiRNNを学習する。BiRNNでは、入力系列に対する出力系列の予測誤差が用いられ、パラメータ学習が行われる。

10

【0029】

図2に戻り、翻訳モデル学習装置50は、翻訳用データ生成装置10において構築された疑似対訳コーパスを用いて翻訳モデルを学習する。翻訳モデルとしては、Transformer又はRNN-based Sequence-to-Sequenceモデル等を用いてもよい。

【0030】

翻訳用データ生成装置10は、その機能として、解析部11と、ノイズ付与部12と、コーパス構築部13と、記憶部14とを備えている。

【0031】

解析部11は、対訳コーパスDB20から原言語テキストを取得し、取得した原言語テキストに対して形態素解析を行う。すなわち、例えば、解析部11は、「主要な高速道路よりも観光ルートの方を走りたいです。」という原言語テキストを取得すると、該原言語テキストについて形態素系列 $w = (\text{主要}, \text{な}, \text{高速}, \text{道路}, \text{より}, \text{も}, \text{観光}, \text{ルート}, \text{の}, \text{方}, \text{を}, \text{走}, \text{り}, \text{たい}, \text{です})$ を抽出する。

20

【0032】

ノイズ付与部12は、原言語テキスト(詳細には解析部11が抽出した形態素系列)にノイズを付与してノイズ付与原言語テキストを得る。ノイズ付与部12は、ノイズモデル学習装置40によって学習されたノイズモデルを用いて、原言語テキストにノイズを付与する。ノイズ付与部12は、原言語テキストの各単語の特徴(具体的には形態素)に応じて、各形態素にノイズラベルを付与し、該ノイズラベルを該ノイズラベルに対応する単語(ノイズとしての単語)へ置き換えることにより、原言語テキストにノイズを付与する。ノイズ付与部12は、ノイズモデルを用いることにより、入力された原言語テキストの形態素系列に対応するノイズラベル系列を予測し、対応する形態素系列の次にノイズラベルを挿入する。そして、ノイズ付与部12は、挿入したノイズラベルを、ノイズを表す単語に置換し、最終的な出力であるノイズが付与された原言語テキストであるノイズ付き原言語テキストを得る。なお、ノイズ付与部12は、原言語テキストの形態素に応じてノイズラベルを付与するとして説明したがこれに限定されず、原言語テキストの各単語の品詞や読み(発音)に応じてノイズラベルを付与してもよい。また、ノイズ付与部12は、単語の形態素、品詞、及び読み等の2つ以上の情報に応じて、ノイズラベルを付与してもよい。

30

40

【0033】

ノイズ付与部12は、具体的には、まず、原言語テキストの形態素系列をノイズモデルに入力し、各タイムステップ(各形態素系列)におけるノイズモデルの出力ベクトル h_t を取得する。本実施形態では、各タイムステップにおけるノイズラベルについて、単純にノイズラベルの事後確率が最大となるものを推定結果とするのではなく、出力ベクトル h_t に指数をとった値 $\exp(h_t / \quad)$ で定義される多項分布に基づくサンプリングにより決定する。すなわち、各タイムステップにおけるノイズラベル l_t は以下の(1)式に基づき推定される。

$$l_t \sim \exp(h_t / \quad) \cdot \dots \cdot (1)$$

上記(1)式において、 l_t はノイズラベルの推定結果、 h_t はノイズモデルの出力ベク

50

トル、 θ は温度パラメータである。出力ベクトル h_t は、3種類のラベルタイプ ($\langle F \rangle$, $\langle D \rangle$, 0) についての3次元ベクトルで示される。温度パラメータ θ は、ノイズラベルのバリエーションの強弱を操作するためのパラメータである。温度パラメータ θ の値を大きく ($\theta \rightarrow \infty$) するとノイズラベルの確率分布は一様分布に近づき、小さく ($\theta \rightarrow 0$) すると最も高い確率のノイズラベルが選択されるようになる。

【0034】

例えば温度パラメータ θ が比較的小さい場合のノイズラベルの決定について説明する。いま、ノイズモデルの出力ベクトル $h_t = (-0.1$ (0 の重みスコア), 0.3 ($\langle F \rangle$ の重みスコア), -0.3 ($\langle D \rangle$ の重みスコア)) であり、温度パラメータ $\theta = 0.15$ であるとする。この場合、 $h_t / \theta = (-0.6666\dots, 2, -2)$ となる。各ノイズラベルの重みスコアを0以上とすべく指数をとると、 $\exp(h_t / \theta) = (0.51, 7.39, 0.13)$ となる。重みスコアを確率値として扱うべく値域が $[0, 1]$ 且つ全ての値を足して1になるように正規化すると、確率分布は $(0.06$ (0 がノイズラベルとして選ばれる確率), 0.92 ($\langle F \rangle$ がノイズラベルとして選ばれる確率), 0.02 ($\langle D \rangle$ がノイズラベルとして選ばれる確率)) となる。このような確率分布 (多項分布) に基づきノイズラベルを1回だけサンプリング (試行) することは、カテゴリカル分布からのサンプリングに相当する。この場合、ノイズラベル $\langle F \rangle$ の確立が92%と極めて高く、サンプリング結果として選択される可能性が極めて高い。

10

【0035】

例えば温度パラメータ θ が比較的大きい場合のノイズラベルの決定について説明する。いま、ノイズモデルの出力ベクトル $h_t = (-0.1$ (0 の重みスコア), 0.3 ($\langle F \rangle$ の重みスコア), -0.3 ($\langle D \rangle$ の重みスコア)) であり、温度パラメータ $\theta = 1.0$ であるとする。この場合、 $h_t / \theta = (-0.1, 0.3, -0.3)$ となる。各ノイズラベルの重みスコアを0以上とすべく指数をとると、 $\exp(h_t / \theta) = (0.90, 1.35, 0.74)$ となる。重みスコアを確率値として扱うべく値域が $[0, 1]$ 且つ全ての値を足して1になるように正規化すると、確率分布は $(0.30$ (0 がノイズラベルとして選ばれる確率), 0.45 ($\langle F \rangle$ がノイズラベルとして選ばれる確率), 0.25 ($\langle D \rangle$ がノイズラベルとして選ばれる確率)) となる。このように、温度パラメータ θ を大きくすると、上述した温度パラメータ $\theta = 0.15$ の場合と比較して、ノイズラベル 0 及びノイズラベル $\langle D \rangle$ が選択されやすくなっていることがわかる。温度パラメータ θ が ∞ に近づくほど、各ノイズラベルの確立は $33.333\dots\%$ に近づき、確率分布が一様分布に近づく。

20

30

【0036】

このように、ノイズ付与部12は、原言語テキストの各単語の特徴 (形態素系列) を入力としてノイズモデルから出力される各ノイズラベルのスコアに基づく確率分布に従ってノイズラベルをサンプリングし、原言語テキストに付与するノイズラベルを決定している。なお、上述した説明においては、ノイズモデルの出力値を基に定義される確率分布が多項分布を表すとして説明したが、これに限定されず、確率分布はポアソン分布又は正規分布等を表すものであってもよい。

【0037】

ノイズ付与部12は、つづいて、ノイズモデルを用いて予測したノイズラベル系列を、ノイズを表す単語に置き換える。ノイズ付与部12は、例えば、各ノイズラベルに対応する語彙集合 V_{type} からユニグラム確率に基づきサンプリングを行う。例えば、フィルターのノイズラベル $\langle F \rangle$ を、フィルターを表す単語へ置換する場合、以下の (2) 式に基づきフィルターを表す単語が決定される。

$$w_t' \sim V_{\langle F \rangle} \dots (2)$$

上記 (2) 式において、 $V_{\langle F \rangle}$ はノイズラベル $\langle F \rangle$ の語彙集合、 w_t' はタイムステップ t に挿入されるフィルター (ノイズ) を表す単語である。以上によって、原言語テキストの形態素系列 $w = (w_0, w_1, \dots, w_n)$ からノイズを表す単語を含む系列 $w' = (w_0, w_1, w_1', w_2, w_2', \dots, w_n)$ を得る。

40

50

【 0 0 3 8 】

ノイズ付与部 1 2 は、1つの原言語テキストから複数パターンのノイズ付与原言語テキストを得る。ノイズ付与部 1 2 は、例えば、1つのノイズラベルに対して置き換える単語（ノイズを表す単語）を複数パターン導出し、1つの原言語テキストから複数パターンのノイズ付与原言語テキストを得てもよい。また、ノイズ付与部 1 2 は、例えば各形態素に対応するノイズラベルを複数パターン導出し、1つの原言語テキストから複数パターンのノイズ付与原言語テキストを得てもよい。

【 0 0 3 9 】

コーパス構築部 1 3 は、ノイズ付与原言語テキストと、該ノイズ付与原言語テキストのノイズ付与前の原言語テキストに対応する目的言語テキストとを対応付けた疑似対訳コーパスを構築する。図 5 は、疑似対訳コーパスの構築イメージを示す図である。図 5 に示される例では、ノイズ付与前の原言語テキストである「主要な高速道路よりも観光ルートの方を走りたいです」についてのノイズ付与原言語テキスト（「えー主要な高速道路よりも観光ルートの方を走りたいです」等の 7 つのノイズ付与原言語テキスト）と、ノイズ付与前の原言語テキストに対応する目的言語テキストである「I would rather take a scenic route than a main highway.」とが対応付けられた（対訳ペアとした）疑似対訳コーパスが構築されている。

10

【 0 0 4 0 】

記憶部 1 4 は、コーパス構築部 1 3 によって構築された疑似対訳コーパスを記憶する DB である。翻訳モデル学習装置 5 0 は、記憶部 1 4 に記憶されている疑似対訳コーパスを用いて翻訳モデルを学習する。

20

【 0 0 4 1 】

次に、図 6 を参照して、翻訳用データ生成システム 1 が実行する処理を説明する。図 6 は、翻訳用データ生成システム 1 が実行する処理を示すフローチャートである。なお、図 6 に示される処理が実行される前提として、ノイズモデル学習装置 4 0 によってノイズモデルが構築（学習）されているものとする。

【 0 0 4 2 】

図 6 に示されるように、翻訳用データ生成システム 1 では、まず、翻訳用データ生成装置 1 0 の解析部 1 1 が対訳コーパス DB 2 0 から原言語テキストを取得する（ステップ S 1）。つづいて、解析部 1 1 は、取得した原言語テキストに対して形態素解析を実行する（ステップ S 2）。

30

【 0 0 4 3 】

つづいて、翻訳用データ生成装置 1 0 のノイズ付与部 1 2 は、解析部 1 1 が抽出した形態素系列に対してノイズを付与し、ノイズ付与原言語テキストを得る（ステップ S 3）。詳細には、ノイズ付与部 1 2 は、ノイズモデルを用いることにより、入力された原言語テキストの形態素系列に対応するノイズラベル系列を予測し、対応する形態素系列の次にノイズラベルを挿入する。そして、ノイズ付与部 1 2 は、挿入したノイズラベルを、ノイズを表す単語に置換し、最終的な出力であるノイズが付与された原言語テキストであるノイズ付き原言語テキストを得る。

【 0 0 4 4 】

つづいて、翻訳用データ生成装置 1 0 のコーパス構築部 1 3 は、ノイズ付与原言語テキストと、該ノイズ付与原言語テキストのノイズ付与前の原言語テキストに対応する目的言語テキストとを対応付けた疑似対訳コーパスを構築する（ステップ S 4）。

40

【 0 0 4 5 】

最後に、翻訳モデル学習装置 5 0 は、コーパス構築部 1 3 によって構築された疑似対訳コーパスを用いて翻訳モデルを学習する（ステップ S 5）。以上が、翻訳用データ生成システム 1 が実行する処理の一例である。

【 0 0 4 6 】

次に、本実施形態の作用効果について説明する。

【 0 0 4 7 】

50

本実施形態に係る翻訳用データ生成システム1は、原言語テキストにノイズを付与してノイズ付与原言語テキストを得るノイズ付与部12と、ノイズ付与原言語テキストと、該ノイズ付与原言語テキストのノイズ付与前の原言語テキストに対応する目的言語テキストとを対応付けた疑似対訳コーパスを構築するコーパス構築部13と、を備える。

【0048】

本実施形態に係る翻訳用データ生成システム1では、原言語テキストにノイズが付与され、ノイズ付与原言語テキストとノイズ付与前の原言語テキストに対応する目的言語テキストとを対応付けた疑似対訳コーパスが構築される。このように、ノイズ付与原言語テキストがノイズ付与前の原言語テキストに対応する目的言語テキストに対応付けられた対訳コーパスが構築されることにより、このような対訳コーパスを利用して、例えば自然発話入力にフィラー等のノイズが含まれている場合においても、ノイズ付与前の原言語テキストに対応する目的言語テキストを適切に導出することが可能となる。すなわち、本実施形態に係る翻訳用データ生成システム1によれば、ノイズが含まれる自然発話に対して頑健なコーパス（疑似対訳コーパス）を構築することができ、ノイズが含まれる非流暢な自然発話に対しても高精度に翻訳を行うことができる。なお、このような翻訳用データ生成システム1により生成された情報が翻訳に用いられる場合には、利用者の発話内容を修正して翻訳モデルに入力する必要がなく、利用者の発話内容をそのまま翻訳モデルに入力することができる。また、例えば、特開2010-079647号公報及び特開2007-057844号公報に記載されたシステムでは、音声認識装置を用いて、逐次利用者の発話を受け取り言い直し判定を行っているが、本実施形態に係る翻訳用データ生成システム1では音声認識装置が不要であり、認識結果のテキスト情報のみが利用できればよい。このように、本実施形態に係る翻訳用データ生成システム1では、発話内容の修正処理や言い直し判定処理が実施されることを抑制できるため、CPU等の処理部における処理負荷を軽減するという技術的効果も併せて奏する。

【0049】

図7は、本実施形態及び比較例の翻訳例を示す表である。図7の上段に示されるように、ノイズが含まれる自然発話入力に対して、比較例では訳抜けが生じている。また、図7の下段に示されるように、ノイズが含まれる自然発話入力に対して、比較例ではノイズを含めた状態で翻訳しており、所望の翻訳を行うことができていない。比較例に示されるように、従来、ノイズが含まれる自然発話に対して高精度に翻訳を行うことは困難であった。この点、図7の上段及び下段に示されるように、本実施形態の翻訳用データ生成システム1によって構築された疑似対訳コーパスが考慮されて翻訳が行われた場合には、ノイズが含まれる自然発話に対しても翻訳誤りが起きにくく、高精度に翻訳を行うことができる。

【0050】

翻訳用データ生成システム1は、疑似対訳コーパスを用いて翻訳モデルを学習する翻訳モデル学習装置50を備えている。構築したコーパスに基づいて翻訳モデルが学習されることにより、ノイズが含まれる自然発話に対してより高精度に翻訳を行うことができる。

【0051】

翻訳用データ生成システム1は、ノイズを含んだ原言語テキスト群である訓練データを用いて、原言語テキストに対するノイズの付与に係るノイズモデルを学習するノイズモデル学習装置40を備え、ノイズ付与部12は、ノイズモデルを用いて、原言語テキストにノイズを付与する。予めノイズが含まれている原言語テキスト群に基づきノイズモデルが学習され、該ノイズモデルに基づいてノイズの付与が行われることにより、実際に含まれる可能性が高いノイズが付与され易くなり、翻訳精度をより向上させることができる。

【0052】

翻訳用データ生成システム1において、ノイズ付与部12は、原言語テキストの各単語に、ノイズのタイプを示すノイズラベルを付与し、該ノイズラベルを該ノイズラベルに対応する単語へ置き換えることにより、原言語テキストにノイズを付与する。原言語テキストの各単語に応じたノイズラベルが付与された後に該ノイズラベルに応じた単語（ノイズ）が導出されることにより、ノイズ付与の容易性及び妥当性を担保することができる。

10

20

30

40

50

【 0 0 5 3 】

翻訳用データ生成システム 1 において、ノイズ付与部 1 2 は、1 つのノイズラベルに対して置き換える単語を複数パターン導出し、1 つの原言語テキストから複数パターンのノイズ付与原言語テキストを得る。これにより、1 つの原言語テキストから効率的に疑似対訳コーパスを充実させ、翻訳精度をより向上させることができる。

【 0 0 5 4 】

翻訳用データ生成システム 1 において、ノイズ付与部 1 2 は、各単語に対応するノイズラベルを複数パターン導出し、1 つの原言語テキストから複数パターンのノイズ付与原言語テキストを得る。これにより、1 つの原言語テキストから効率的に疑似対訳コーパスを充実させ、翻訳精度をより向上させることができる。

10

【 0 0 5 5 】

翻訳用データ生成システム 1 において、ノイズ付与部 1 2 は、原言語テキストの各単語の特徴に応じて、ノイズラベルを付与する。これにより、各単語に関連して含まれやすいノイズに係るノイズラベルを、各単語に適切に付与することができる。

【 0 0 5 6 】

翻訳用データ生成システム 1 において、ノイズ付与部 1 2 は、原言語テキストの各単語の特徴である、形態素、品詞、及び単語の読みの少なくとも一つに応じて、ノイズラベルを付与する。これにより、各単語に関連して含まれやすいノイズに係るノイズラベルを、各単語に適切に付与することができる。

【 0 0 5 7 】

翻訳用データ生成システム 1 において、ノイズ付与部 1 2 は、原言語テキストの各単語の特徴を入力としてノイズモデルから出力される各ノイズラベルのスコアに基づく各ノイズレベルの確率分布に従ってノイズラベルをサンプリングし、原言語テキストに付与するノイズラベルを決定する。これにより、例えばノイズモデルから出力されたスコアが高いノイズラベルを付与することが可能となり、各単語に関連して含まれやすいノイズに係るノイズラベルを、各単語に適切に付与することができる。

20

【 0 0 5 8 】

最後に、翻訳用データ生成装置 1 0 のハードウェア構成について、図 8 を参照して説明する。上述の翻訳用データ生成装置 1 0 は、物理的には、プロセッサ 1 0 0 1、メモリ 1 0 0 2、ストレージ 1 0 0 3、通信装置 1 0 0 4、入力装置 1 0 0 5、出力装置 1 0 0 6、バス 1 0 0 7 などを含むコンピュータ装置として構成されてもよい。

30

【 0 0 5 9 】

なお、以下の説明では、「装置」という文言は、回路、デバイス、ユニットなどに読み替えることができる。翻訳用データ生成装置 1 0 のハードウェア構成は、図に示した各装置を 1 つ又は複数含むように構成されてもよいし、一部の装置を含まずに構成されてもよい。

【 0 0 6 0 】

翻訳用データ生成装置 1 0 における各機能は、プロセッサ 1 0 0 1、メモリ 1 0 0 2 などのハードウェア上に所定のソフトウェア（プログラム）を読み込ませることで、プロセッサ 1 0 0 1 が演算を行い、通信装置 1 0 0 4 による通信や、メモリ 1 0 0 2 及びストレージ 1 0 0 3 におけるデータの読み出し及び/又は書き込みを制御することで実現される。

40

【 0 0 6 1 】

プロセッサ 1 0 0 1 は、例えば、オペレーティングシステムを動作させてコンピュータ全体を制御する。プロセッサ 1 0 0 1 は、周辺装置とのインターフェース、制御装置、演算装置、レジスタなどを含む中央処理装置（CPU：Central Processing Unit）で構成されてもよい。例えば、翻訳用データ生成装置 1 0 のノイズ付与部 1 2 等の制御機能はプロセッサ 1 0 0 1 で実現されてもよい。

【 0 0 6 2 】

また、プロセッサ 1 0 0 1 は、プログラム（プログラムコード）、ソフトウェアモジュールやデータを、ストレージ 1 0 0 3 及び/又は通信装置 1 0 0 4 からメモリ 1 0 0 2 に

50

読み出し、これらに従って各種の処理を実行する。プログラムとしては、上述の実施の形態で説明した動作の少なくとも一部をコンピュータに実行させるプログラムが用いられる。例えば、翻訳用データ生成装置10のノイズ付与部12等の制御機能は、メモリ1002に格納され、プロセッサ1001で動作する制御プログラムによって実現されてもよく、他の機能ブロックについても同様に実現されてもよい。上述の各種処理は、1つのプロセッサ1001で実行される旨を説明してきたが、2以上のプロセッサ1001により同時又は逐次に行われてもよい。プロセッサ1001は、1以上のチップで実装されてもよい。なお、プログラムは、電気通信回線を介してネットワークから送信されてもよい。

【0063】

メモリ1002は、コンピュータ読み取り可能な記録媒体であり、例えば、ROM(Read Only Memory)、EPROM(Erasable Programmable ROM)、EEPROM(Electrically Erasable Programmable ROM)、RAM(Random Access Memory)などの少なくとも1つで構成されてもよい。メモリ1002は、レジスタ、キャッシュ、メインメモリ(主記憶装置)などと呼ばれてもよい。メモリ1002は、本発明の一実施の形態に係る無線通信方法を実施するために実行可能なプログラム(プログラムコード)、ソフトウェアモジュールなどを保存することができる。

【0064】

ストレージ1003は、コンピュータ読み取り可能な記録媒体であり、例えば、CD-ROM(Compact Disc ROM)などの光ディスク、ハードディスクドライブ、フレキシブルディスク、光磁気ディスク(例えば、コンパクトディスク、デジタル多用途ディスク、Blu-ray(登録商標)ディスク)、スマートカード、フラッシュメモリ(例えば、カード、スティック、キードライブ)、フロッピー(登録商標)ディスク、磁気ストリップなどの少なくとも1つで構成されてもよい。ストレージ1003は、補助記憶装置と呼ばれてもよい。上述の記憶媒体は、例えば、メモリ1002及び/又はストレージ1003を含むデータベース、サーバその他の適切な媒体であってもよい。

【0065】

通信装置1004は、有線及び/又は無線ネットワークを介してコンピュータ間の通信を行うためのハードウェア(送受信デバイス)であり、例えばネットワークデバイス、ネットワークコントローラ、ネットワークカード、通信モジュールなどともいう。

【0066】

入力装置1005は、外部からの入力を受け付ける入力デバイス(例えば、キーボード、マウス、マイクロフォン、スイッチ、ボタン、センサなど)である。出力装置1006は、外部への出力を実施する出力デバイス(例えば、ディスプレイ、スピーカー、LEDランプなど)である。なお、入力装置1005及び出力装置1006は、一体となった構成(例えば、タッチパネル)であってもよい。

【0067】

また、プロセッサ1001やメモリ1002などの各装置は、情報を通信するためのバス1007で接続される。バス1007は、単一のバスで構成されてもよいし、装置間で異なるバスで構成されてもよい。

【0068】

また、翻訳用データ生成装置10は、マイクロプロセッサ、デジタル信号プロセッサ(DSP: Digital Signal Processor)、ASIC(Application Specific Integrated Circuit)、PLD(Programmable Logic Device)、FPGA(Field Programmable Gate Array)などのハードウェアを含んで構成されてもよく、当該ハードウェアにより、各機能ブロックの一部又は全てが実現されてもよい。例えば、プロセッサ1001は、これらのハードウェアの少なくとも1つで実装されてもよい。

【0069】

以上、本実施形態について詳細に説明したが、当業者にとっては、本実施形態が本明細書中に説明した実施形態に限定されるものではないということは明らかである。本実施形態は、特許請求の範囲の記載により定まる本発明の趣旨及び範囲を逸脱することなく修正

10

20

30

40

50

及び変更態様として実施することができる。したがって、本明細書の記載は、例示説明を目的とするものであり、本実施形態に対して何ら制限的な意味を有するものではない。例えば、本発明の一態様に係る翻訳用データ生成システムは、事前に定義したノイズ単語（フィルター、言い淀み、言い直し等）を原言語テキストのランダムな位置に付与するものであってもよい。ランダムな位置に付与する単語（ノイズ）は、例えばノイズ単語候補からランダムに選択されてもよい。このような構成においては、ノイズモデルの学習データ（ラベル付きデータ）がなくても、ノイズ単語が定義できさえすれば、ノイズをランダムに付与するノイズモデルを構築することができる。

【0070】

本明細書で説明した各態様／実施形態は、LTE (Long Term Evolution)、LTE - A (LTE-Advanced)、SUPER 3G、IMT - Advanced、4G、5G、FRA (Future Radio Access)、W - CDMA (登録商標)、GSM (登録商標)、CDMA 2000、UMB (Ultra Mobile Broad-band)、IEEE 802.11 (Wi-Fi)、IEEE 802.16 (WiMAX)、IEEE 802.20、UWB (Ultra-Wide Band)、Bluetooth (登録商標)、その他の適切なシステムを利用するシステム及び／又はこれらに基づいて拡張された次世代システムに適用されてもよい。

10

【0071】

本明細書で説明した各態様／実施形態の処理手順、フローチャートなどは、矛盾の無い限り、順序を入れ替えてもよい。例えば、本明細書で説明した方法については、例示的な順序で様々なステップの要素を提示しており、提示した特定の順序に限定されない。

20

【0072】

入出力された情報等は特定の場所(例えば、メモリ)に保存されてもよいし、管理テーブルで管理してもよい。入出力される情報等は、上書き、更新、または追記され得る。出力された情報等は削除されてもよい。入力された情報等は他の装置へ送信されてもよい。

【0073】

判定は、1ビットで表される値(0か1か)によって行われてもよいし、真偽値(Boolean: trueまたはfalse)によって行われてもよいし、数値の比較(例えば、所定の値との比較)によって行われてもよい。

【0074】

本明細書で説明した各態様／実施形態は単独で用いてもよいし、組み合わせで用いてもよいし、実行に伴って切り替えて用いてもよい。また、所定の情報の通知(例えば、「Xであること」の通知)は、明示的に行うものに限られず、暗黙的(例えば、当該所定の情報の通知を行わない)ことによって行われてもよい。

30

【0075】

ソフトウェアは、ソフトウェア、ファームウェア、ミドルウェア、マイクロコード、ハードウェア記述言語と呼ばれるか、他の名称で呼ばれるかを問わず、命令、命令セット、コード、コードセグメント、プログラムコード、プログラム、サブプログラム、ソフトウェアモジュール、アプリケーション、ソフトウェアアプリケーション、ソフトウェアパッケージ、ルーチン、サブルーチン、オブジェクト、実行可能ファイル、実行スレッド、手順、機能などを意味するよう広く解釈されるべきである。

40

【0076】

また、ソフトウェア、命令などは、伝送媒体を介して送受信されてもよい。例えば、ソフトウェアが、同軸ケーブル、光ファイバケーブル、ツイストペア及びデジタル加入者回線(DSL)などの有線技術及び／又は赤外線、無線及びマイクロ波などの無線技術を使用してウェブサイト、サーバ、又は他のリモートソースから送信される場合、これらの有線技術及び／又は無線技術は、伝送媒体の定義内に含まれる。

【0077】

本明細書で説明した情報、信号などは、様々な異なる技術のいずれかを使用して表されてもよい。例えば、上記の説明全体に渡って言及され得るデータ、命令、コマンド、情報

50

、信号、ビット、シンボル、チップなどは、電圧、電流、電磁波、磁界若しくは磁性粒子、光場若しくは光子、又はこれらの任意の組み合わせによって表されてもよい。

【0078】

なお、本明細書で説明した用語及び/又は本明細書の理解に必要な用語については、同一の又は類似する意味を有する用語と置き換えてもよい。

【0079】

また、本明細書で説明した情報、パラメータなどは、絶対値で表されてもよいし、所定の値からの相対値で表されてもよいし、対応する別の情報で表されてもよい。

【0080】

ユーザ端末は、当業者によって、移動通信端末、加入者局、モバイルユニット、加入者ユニット、ワイヤレスユニット、リモートユニット、モバイルデバイス、ワイヤレスデバイス、ワイヤレス通信デバイス、リモートデバイス、モバイル加入者局、アクセス端末、モバイル端末、ワイヤレス端末、リモート端末、ハンドセット、ユーザエージェント、モバイルクライアント、クライアント、またはいくつかの他の適切な用語で呼ばれる場合もある。

10

【0081】

本明細書で使用する「判断(determining)」、「決定(determining)」という用語は、多種多様な動作を包含する場合がある。「判断」、「決定」は、例えば、計算(calculating)、算出(computing)、処理(processing)、導出(deriving)、調査(investigating)、探索(looking up) (例えば、テーブル、データベースまたは別のデータ構造での探索)、確認(ascertaining)した事を「判断」「決定」したとみなす事などを含み得る。また、「判断」、「決定」は、受信(receiving) (例えば、情報を受信すること)、送信(transmitting) (例えば、情報を送信すること)、入力(input)、出力(output)、アクセス(accessing) (例えば、メモリ中のデータにアクセスすること)した事を「判断」「決定」したとみなす事などを含み得る。また、「判断」、「決定」は、解決(resolving)、選択(selecting)、選定(choosing)、確立(establishing)、比較(comparing)などした事を「判断」「決定」したとみなす事を含み得る。つまり、「判断」「決定」は、何らかの動作を「判断」「決定」したとみなす事を含み得る。

20

【0082】

本明細書で使用する「に基づいて」という記載は、別段に明記されていない限り、「のみに基づいて」を意味しない。言い換えれば、「に基づいて」という記載は、「のみに基づいて」と「に少なくとも基づいて」の両方を意味する。

30

【0083】

本明細書で「第1の」、「第2の」などの呼称を使用した場合においては、その要素へのいかなる参照も、それらの要素の量または順序を全般的に限定するものではない。これらの呼称は、2つ以上の要素間を区別する便利な方法として本明細書で使用され得る。したがって、第1および第2の要素への参照は、2つの要素のみがそこで採用され得ること、または何らかの形で第1の要素が第2の要素に先行しなければならないことを意味しない。

【0084】

「含む(include)」、「含んでいる(including)」、およびそれらの変形が、本明細書あるいは特許請求の範囲で使用されている限り、これら用語は、用語「備える(comprising)」と同様に、包括的であることが意図される。さらに、本明細書あるいは特許請求の範囲において使用されている用語「または(or)」は、排他的論理和ではないことが意図される。

40

【0085】

本明細書において、文脈または技術的に明らかに1つのみしか存在しない装置である場合以外は、複数の装置をも含むものとする。

【0086】

本開示の全体において、文脈から明らかに単数を示したものでなければ、複数のもの

50

を含むものとする。

【符号の説明】

【0087】

1 ... 翻訳用データ生成システム、12 ... ノイズ付与部、13 ... コーパス構築部、40 ... ノイズモデル学習装置（ノイズモデル学習部）、50 ... 翻訳モデル学習装置（翻訳モデル学習部）。

10

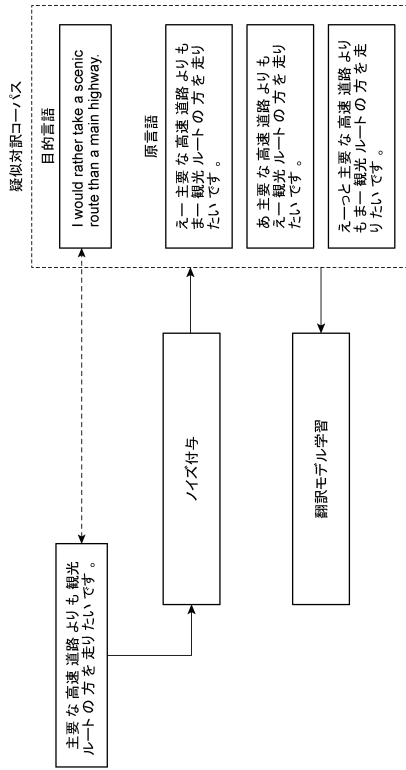
20

30

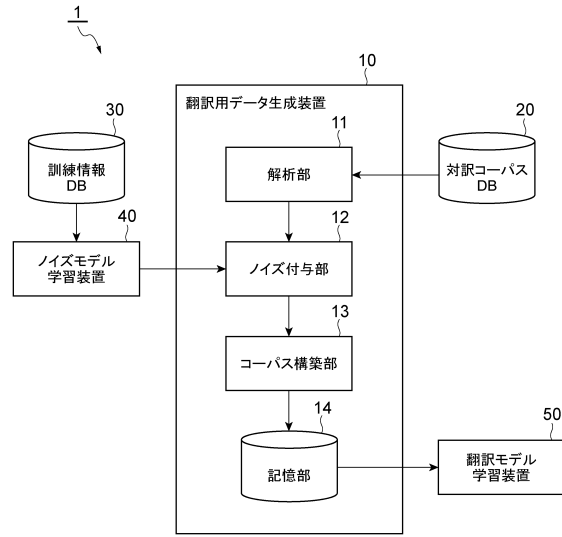
40

50

【図面】
【図 1】



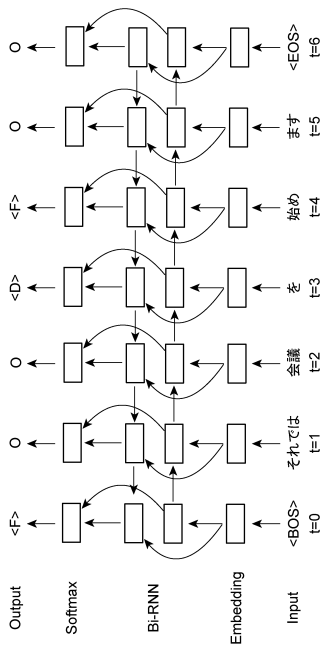
【図 2】



10

20

【図 3】



【図 4】

ノイズラベル	対象
<F>	フィラー
<D>	言い淀み、言い直し
O	ノイズ無し

30

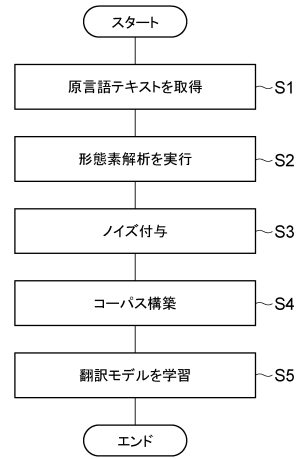
40

50

【図 5】

原言語テキスト	主要な高速道路よりも観光ルートの方を走りたいです。 えー主要な高速道路よりも観光ルートの方を走りたいです。 えと主要な高速道路よりも観光ルートの方を走りたいです。 えー主要な高速道路よりも観光ルートの方を走りたいです。 えーと主要な高速道路よりも観光ルートの方を走りたいです。流れいでたる。 ま主要な高速道路よりも観光ルートの方を走りたいです。 主要な高速道路よりも観光ルートの方を走りたいです。 えーと主要な高速道路よりも観光ルートの方を走りたいです。
目的言語テキスト	I would rather take a scenic route than a main highway.

【図 6】



10

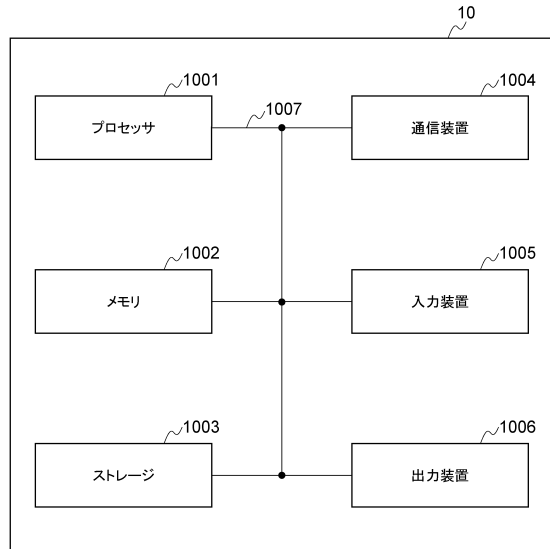
20

【図 7】

自然発語(入力)	比較例	本実施形態
えーとです。ね。えーちよつとえーキャンセルというところでキャンセル料がえー発生する形になるんですけど。えー本日。えー出発の三日前にはいぬの予約の際です。ね。あのキャンセル料のことば、こちらからさまで聞いてませんでしようか。	Let me see .Well , it seems like cancellation will cause a cancellation fee .	Well , a cancellation will cause a cancellation fee , but you made a reservation three days before your departure , did you confirm the cancellation ?
えーと。大阪にいや東京に行きたいんですが。どのように行けばえーいいでしょうか。	Well , I do not want to go to Osaka or Tokyo but how would you go if I go out ?	I would like to go to Tokyo , how can I go there ?

30

【図 8】



40

50

フロントページの続き

株式会社NTTドコモ内

審査官 長 由紀子

- (56)参考文献 特開2018-055671(JP,A)
今出 昌宏 外2名, ユーラルネット機械翻訳における自動コーパス生成適用, 一般社団法人人工知能学会 第31回全国大会論文集DVD [DVD-ROM], 一般社団法人人工知能学会, 2017年05月26日, pp.1-4
太田 健吾 外2名, フィラーの書き起こしのないコーパスからのフィラー付き言語モデルの構築, 情報処理学会研究報告, 日本, 社団法人情報処理学会, 2007年07月20日, 第2007巻 第75号, pp.1-6
増村 亮 外2名, Web上の言語資源を利用した大規模話し言葉データからの言語モデル作成, 日本音響学会 2011年 春季研究発表会講演論文集CD-ROM [CD-ROM], 社団法人日本音響学会, 2011年03月02日, pp.75-78
玉井 孝幸 外2名, 音声対話システムにおける発話予測を利用した音声認識, 情報処理学会研究報告, 日本, 社団法人情報処理学会, 2002年10月25日, 第2002巻 第98号, pp.1-6
- (58)調査した分野 (Int.Cl., DB名)
G06F 40/00-58
G06N 20/00