

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7663171号
(P7663171)

(45)発行日 令和7年4月16日(2025.4.16)

(24)登録日 令和7年4月8日(2025.4.8)

(51)国際特許分類 F I
G 0 6 F 40/44 (2020.01) G 0 6 F 40/44

請求項の数 6 (全54頁)

(21)出願番号	特願2020-137323(P2020-137323)	(73)特許権者	301022471 国立研究開発法人情報通信研究機構 東京都小金井市貫井北町4-2-1
(22)出願日	令和2年8月17日(2020.8.17)	(74)代理人	100143498 弁理士 中西 健
(65)公開番号	特開2022-33437(P2022-33437A)	(74)代理人	100136319 弁理士 北原 宏修
(43)公開日	令和4年3月2日(2022.3.2)	(74)代理人	100148275 弁理士 山内 聡
審査請求日	令和5年7月4日(2023.7.4)	(74)代理人	100142745 弁理士 伊藤 世子
(出願人による申告) 令和元年度、総務省、グローバルコミュニケーション計画の推進 - 多言語音声翻訳技術の研究開発及び社会実証 - 産業技術力強化法第17条の適用を受ける特許出願		(72)発明者	マリ バンジャマン 東京都小金井市貫井北町4-2-1 国立研究開発法人情報通信研究機構内
		(72)発明者	藤田 篤

最終頁に続く

(54)【発明の名称】 疑似対訳データ生成用機械翻訳モデルの学習方法、疑似対訳データ取得方法、および、機械翻訳モデルの学習方法

(57)【特許請求の範囲】

【請求項1】

パラメータを設定することで学習処理を行うことができ、入力データ埋込部と機械学習処理部とを含む疑似対訳データ生成用機械翻訳モデルの学習方法であって、

生成する疑似対訳データの対象とする分野である適応先分野以外の分野の第1言語のデータである他分野第1言語データと、当該他分野第1言語データの第2言語の翻訳データである他分野第2言語データとからなる対訳データを複数含む他分野対訳データ集合 $D_{setp}(L1 - L2)$ と、

前記適応先分野以外の分野の第1言語のデータを複数含む他分野単言語データ集合 $D_{setm}(L1)$ と、

前記適応先分野以外の分野の第2言語のデータを複数含む他分野単言語データ集合 $D_{setm}(L2)$ と、

前記適応先分野の第1言語のデータを複数含む適応先分野単言語データ集合 $D_{setm}(R1)$ と、

前記適応先分野の第2言語のデータを複数含む適応先分野単言語データ集合 $D_{setm}(R2)$ と、

を用いて、パラメータを設定することで学習処理を行うことができ、前記入力データ埋込部とXLM処理部とを含むXLMモデルに対して、

(A) 前記他分野単言語データ集合 $D_{setm}(L1)$ 、前記他分野単言語データ集合 $D_{setm}(L2)$ 、前記適応先分野単言語データ集合 $D_{setm}(R1)$ 、および、前記

適応先分野単言語データ集合 D s e t m (R 2) に含まれるデータである単言語用入力データの
一部をマスクしたマスク化データを入力とし、前記単言語用入力データを正解データとし、
当該正解データと前記 X L M モデルの出力との損失が小さくなるように学習する
処理であるマスク化処理による学習処理と、

(B) 前記他分野対訳データ集合 D s e t p (L 1 - L 2) に含まれる対訳データの
前記他分野第 1 言語データおよび前記他分野第 2 言語データのうちの一方のデータを入力とし、
他方のデータを正解データとし、当該正解データと前記 X L M モデルの出力との損失が
小さくなるように学習する処理である教師ありデータによる学習処理と
を行うことで、前記 X L M モデルの最適パラメータを取得し、前記最適パラメータが設定
された前記 X L M モデルの前記入力データ埋込部に設定されている最適パラメータを、前
記疑似対訳データ生成用機械翻訳モデルの前記入力データ埋込部のパラメータの初期パラ
メータとして設定する初期化ステップと、

前記初期パラメータが設定されている状態の前記入力データ埋込部と、前記機械学習処理
部とを含む前記疑似対訳データ生成用機械翻訳モデルに対して、

(1) 正解データを入力データと同一にして前記疑似対訳データ生成用機械翻訳モデルの
学習処理を行う自己符号化処理、

(2) 入力データに対する前記疑似対訳データ生成用機械翻訳モデルの出力データを、再
度、前記疑似対訳データ生成用機械翻訳モデルに入力し、その前記疑似対訳データ生成用
機械翻訳モデルの出力が前記入力データと同一となるように前記疑似対訳データ生成用機
械翻訳モデルの学習処理を行うゼロショット折り返し機械翻訳処理、

(3) 他分野対訳データ集合 D s e t p (L 1 - L 2) に含まれる第 1 言語のデータおよ
び第 2 言語データのいずれか一方を前記疑似対訳データ生成用機械翻訳モデルの入力とし
、他方を正解データとして、前記疑似対訳データ生成用機械翻訳モデルの学習処理を行う
教師あり機械翻訳処理、

の少なくとも 1 つを用いて、学習処理を行うことで、前記疑似対訳データ生成用機械翻訳
モデルの最適パラメータを取得する最適化ステップと、
を備え、

前記疑似対訳データ生成用機械翻訳モデルは、

制御信号により指定された種別のデータを出力することができ、前記制御信号により指
定された、(1) 前記他分野第 1 言語データ、(2) 前記他分野第 2 言語データ、(3)
前記適応先分野の第 1 言語のデータ、および、(4) 前記適応先分野の第 2 言語のデータ
のいずれかを出力するように設定される、

疑似対訳データ生成用機械翻訳モデルの学習方法。

【請求項 2】

請求項 1 に記載の疑似対訳データ生成用機械翻訳モデルの学習方法により取得された疑似
 対訳データ生成用機械翻訳モデルを用いて、適応先分野の疑似対訳データを取得する疑似
 対訳データ取得方法であって、

前記他分野対訳データ集合 D s e t p (L 1 - L 2) から取得した第 1 言語のデータに
 対して、前記制御信号により前記疑似対訳データ生成用機械翻訳モデルの出力が前記適応
 先分野の第 2 言語のデータとなるように設定した前記疑似対訳データ生成用機械翻訳モ
 デルを用いた機械翻訳処理を行うことで、前記他分野第 1 言語データの機械翻訳結果デー
 タである前記適応先分野の第 2 言語疑似翻訳データを取得する第 1 機械翻訳ステップと、

前記他分野対訳データ集合 D s e t p (L 1 - L 2) から取得した第 2 言語のデータに
 対して、前記制御信号により前記疑似対訳データ生成用機械翻訳モデルの出力が前記適応
 先分野の第 1 言語のデータとなるように設定した前記疑似対訳データ生成用機械翻訳モ
 デルを用いた機械翻訳処理を行うことで、前記他分野第 2 言語データの機械翻訳結果デー
 タである前記適応先分野の第 1 言語疑似翻訳データを取得する第 2 機械翻訳ステップと、

前記第 1 機械翻訳ステップで取得された前記適応先分野の第 2 言語疑似翻訳データと、
 前記第 2 機械翻訳ステップで取得された前記適応先分野の第 1 言語疑似翻訳データと、を
 対応づけることで、前記適応先分野の疑似対訳データを取得する疑似対訳データ取得ステ

10

20

30

40

50

ップと、
を備える疑似対訳データ取得方法。

【請求項 3】

請求項 1 に記載の疑似対訳データ生成用機械翻訳モデルの学習方法により取得された疑似対訳データ生成用機械翻訳モデルを用いて、適応先分野の疑似対訳データを取得する疑似対訳データ取得方法であって、

前記適応先分野単言語データ集合 $Dsetm(R1)$ から取得した第 1 言語のデータまたは前記適応先分野単言語データ集合 $Dsetm(R2)$ から取得した第 2 言語のデータに対して、前記制御信号により前記疑似対訳データ生成用機械翻訳モデルの出力が前記適応先分野の第 2 言語のデータまたは第 1 言語のデータとなるように設定した前記疑似対訳データ生成用機械翻訳モデルを用いた機械翻訳処理を行うことで、前記適応先分野の第 1 言語のデータの機械翻訳結果データである前記適応先分野の第 2 言語疑似翻訳データまたは前記適応先分野の第 2 言語のデータの機械翻訳結果データである前記適応先分野の第 1 言語疑似翻訳データを取得する単言語データ機械翻訳ステップと、

10

前記単言語データ機械翻訳ステップで、前記疑似対訳データ生成用機械翻訳モデルの入力とした前記適応先分野の第 1 言語のデータと、前記単言語データ機械翻訳ステップで取得された前記適応先分野の第 2 言語疑似翻訳データと、を対応づける、または、前記疑似対訳データ生成用機械翻訳モデルの入力とした前記適応先分野の第 2 言語のデータと、前記単言語データ機械翻訳ステップで取得された前記適応先分野の第 1 言語疑似翻訳データと、を対応づけることで、前記適応先分野の疑似対訳データを取得する疑似対訳データ取得ステップと、

20

を備える疑似対訳データ取得方法。

【請求項 4】

前記疑似対訳データ取得ステップが取得した前記適応先分野の疑似対訳データの各文対に対して、機械翻訳処理の結果の精度を示す信頼度を取得し、取得した前記信頼度が所定の値以上である文対を含む前記疑似対訳データのみを選択して出力するフィルター処理ステップをさらに備える、

請求項 2 または 3 に記載の疑似対訳データ取得方法。

【請求項 5】

パラメータを設定することで学習処理を行うことができる機械翻訳モデルであって、適応先分野の第 1 言語のデータに対して機械翻訳を行い第 2 言語のデータを取得するための前記機械翻訳モデルの学習方法であって、

30

請求項 2 から 4 のいずれかに記載の疑似対訳データ取得方法により取得された適応先分野の疑似翻訳データと、

前記適応先分野以外の分野の第 1 言語のデータである他分野第 1 言語データと、当該他分野第 1 言語データの第 2 言語の翻訳データである他分野第 2 言語データとからなる対訳データを複数含む他分野対訳データ集合 $Dsetp(L1 - L2)$ と、

を用いて、前記機械翻訳モデルに対して、

(A) 前記他分野対訳データ集合 $Dsetp(L1 - L2)$ に含まれる対訳データの前記他分野第 1 言語データを入力とし、前記他分野対訳データ集合 $Dsetp(L1 - L2)$ に含まれる対訳データの前記他分野第 2 言語データを正解データとし、当該正解データと前記機械翻訳モデルの出力との損失が小さくなるようにする学習処理と、

40

(B) 前記疑似翻訳データに含まれる前記適応先分野の第 1 言語のデータを入力とし、前記疑似翻訳データに含まれる前記適応先分野の第 2 言語のデータを正解データとし、当該正解データと前記機械翻訳モデルの出力との損失が小さくなるようにする学習処理と、

を行うことで、前記機械翻訳モデルの最適パラメータを取得し、当該最適パラメータを機械翻訳モデルに設定することで、学習済み機械翻訳モデルを取得する処理であるの学習処理を行う機械翻訳モデル学習ステップ、

を備える機械翻訳モデルの学習方法。

【請求項 6】

50

請求項 1 に記載の疑似対訳データ生成用機械翻訳モデルの学習方法により取得された学習済みの機械翻訳モデル、または、請求項 5 に記載の機械翻訳モデルの学習方法により取得された学習済みの機械翻訳モデルを用いて機械翻訳処理を行う機械翻訳装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ニューラル機械翻訳の技術に関する。

【背景技術】

【0002】

ニューラル機械翻訳 (NMT: Neural Machine Translation) とは、多層ニューラルネットワークを用いた機械翻訳技術である。ニューラル機械翻訳システムは、起点言語の各語および文全体をベクトルまたはテンソルに変換するニューラルネットワーク (エンコーダ) と、得られたベクトルまたはテンソルから目標言語の語の系列を生成するニューラルネットワーク (デコーダ) とで構成され、対訳データに基づいてニューラルネットワークのモデルを学習する (パラメータを最適化する) ものである。ニューラル機械翻訳システムで達成できる翻訳性能は、所与の起点言語、目標言語、対象分野についての対訳データの規模に強く依存する。実用的なニューラル機械翻訳システムの構築には大規模な対訳データが必要である (言語対や分野によるが少なくとも数十万個の対訳データが必要である)。対訳データが小規模にしか存在しない場合、ニューラル機械翻訳システムにおいて、高品質な機械翻訳処理を実現することは困難である。

10

【0003】

実現したい機械翻訳の対象とする言語対 (対象言語対) と同一言語対の他分野 (機械翻訳の対象とする分野以外の分野) の対訳データ、あるいは他の言語対であるが同一分野の対訳データが存在する場合、それらの対訳データを用いた機械翻訳分野適用技術、多言語機械翻訳技術によって対象分野の機械翻訳を実現できる場合がある。

20

【0004】

また、対訳データとは異なり、安価かつ大量に入手できる単言語データを活用する技術も生み出されてきている。

【0005】

また、近年、他分野の対訳データや他の言語対の対訳データを用いることなく、単言語データのみから機械翻訳システムを構築する教師なし機械翻訳技術も提案されている (例えば、非特許文献 1 を参照)。

30

【先行技術文献】

【非特許文献】

【0006】

【文献】 Mikel Artetxe, Gorka Labaka, Eneko Agirre (2018). Unsupervised Statistical Machine Translation. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 3632-3642.

【発明の概要】

【発明が解決しようとする課題】

40

【0007】

しかしながら、上記の技術のいずれにおいても、対象言語対・対象分野における数千文程度の対訳データの存在を前提としており、この前提が成り立たない場合、上記の技術により実現される機械翻訳システムにおいて、高い機械翻訳の性能を達成することはできない。また、上記の技術が必要とする対象言語対・対象分野における数千文程度の対訳データを作成するためには、無視できない規模の人的・金銭的成本が必要となる。

【0008】

このように、ニューラル機械翻訳システムを新規分野向けに適応させる際には、当該ニューラル機械翻訳システムの機械翻訳の対象とする分野 (対象分野) における対訳データが一定量必要であり、そのような対訳データがない場合、機械翻訳の性能を改善すること

50

は困難である。

【 0 0 0 9 】

そこで、本発明は、上記課題に鑑み、適応先分野（機械翻訳の対象とする分野）の対訳データが一切ない場合であっても、当該適応先分野における機械翻訳を精度良く実行するための疑似対訳データを生成する疑似対訳データ生成装置、および、当該疑似対訳データ生成装置で生成された疑似対訳データを用いて適応先分野における機械翻訳を精度良く実行する機械翻訳システム、および、当該機械翻訳システムで用いられる疑似対訳データ生成用機械翻訳モデルの学習方法、疑似対訳データ取得方法、および、機械翻訳モデルの学習方法を実現することを目的とする。

【課題を解決するための手段】

【 0 0 1 0 】

上記課題を解決するための第1の発明は、疑似対訳データを生成するための疑似対訳データ生成用機械翻訳モデルの学習方法であって、初期化ステップと、最適化ステップと、を備える。

【 0 0 1 1 】

初期化ステップは、

生成する疑似対訳データの対象とする分野である適応先分野以外の分野の第1言語のデータである他分野第1言語データと、当該他分野第1言語データの第2言語の翻訳データである他分野第2言語データとからなる対訳データを複数含む他分野対訳データ集合 $D s e t p (L 1 - L 2)$ と、

適応先分野以外の分野の第1言語のデータを複数含む他分野単言語データ集合 $D s e t m (L 1)$ と、

適応先分野以外の分野の第2言語のデータを複数含む他分野単言語データ集合 $D s e t m (L 2)$ と、

適応先分野の第1言語のデータを複数含む適応先分野単言語データ集合 $D s e t m (R 1)$ と、

適応先分野の第2言語のデータを複数含む適応先分野単言語データ集合 $D s e t m (R 2)$ と、

を用いて、疑似対訳データ生成用機械翻訳モデルの学習処理を行い、当該学習処理を実行した後の疑似対訳データ生成用機械翻訳モデルに設定されているパラメータを初期パラメータに設定する。

【 0 0 1 2 】

最適化ステップは、

初期パラメータが設定されている状態の疑似対訳データ生成用機械翻訳モデルに対して、
 (1) 正解データを入力データと同一にして疑似対訳データ生成用機械翻訳モデルの学習処理を行う自己符号化処理、

(2) 入力データに対する疑似対訳データ生成用機械翻訳モデルの出力データ（この「出力データ」は、例えば、入力データとは異なる言語のデータであり、かつ、ゼロショット（学習したことのない分野）のデータである。）を、再度、疑似対訳データ生成用機械翻訳モデルに入力し、その疑似対訳データ生成用機械翻訳モデルの出力（この「出力」は、入力データと同じ言語である。）が入力データと同一となるように疑似対訳データ生成用機械翻訳モデルの学習処理を行うゼロショット折り返し機械翻訳処理、

(3) 他分野対訳データ集合 $D s e t p (L 1 - L 2)$ に含まれる第1言語のデータおよび第2言語データのいずれか一方を疑似対訳データ生成用機械翻訳モデルの入力とし、他方を正解データとして、疑似対訳データ生成用機械翻訳モデルの学習処理を行う教師あり機械翻訳処理、

の少なくとも1つを用いて、学習処理を行うことで、疑似対訳データ生成用機械翻訳モデルの最適パラメータを取得する。

【 0 0 1 3 】

この疑似対訳データ生成用機械翻訳モデルの学習方法では、

10

20

30

40

50

(1) 精度の高い他分野の対訳データ集合 (D s e t p (L 1 - L 2) (大規模 (対訳データ数が多数) であることが好ましい)) と、
 (2) 他分野の第1言語データ集合 (D s e t m (L 1)) と、
 (3) 他分野の第2言語データ集合 (D s e t m (L 2)) と、
 (4) 適応先分野の第1言語データ集合 (D s e t m (R 1)) と、
 (5) 適応先分野の第2言語データ集合 (D s e t m (R 2)) と、
 を用いて、疑似対訳データ生成用機械翻訳モデル (例えば、ニューラルネットワークモデル) を事前学習処理 (初期化ステップによる処理) により初期化し、さらに、パラメータ最適化処理 (最適化ステップによる処理) を行うことで、適応先分野 (機械翻訳の対象とする分野) の対訳データが一切ない場合であっても、適応先分野の第1言語および第2言語の疑似対訳データを生成するための疑似対訳データ生成用機械翻訳モデルを学習させることができる。

10

【0014】

そして、学習後の (学習済みの) 疑似対訳データ生成用機械翻訳モデルにより、適応先分野 (対象分野) における対訳データが存在しない場合であっても、適応先分野の疑似対訳データを生成することができる。さらに、生成された適応先分野の疑似対訳データを用いて、機械翻訳モデルを学習させ、学習済みの機械翻訳モデルにより機械翻訳処理を行うことで、適応先分野 (対象分野) における対訳データが存在しない場合であっても、適応先分野における機械翻訳を精度良く行うことができる。

【0015】

20

なお、事前学習処理により初期化する対象は、入力データから分散表現データを取得する部分のみとしてもよい。例えば、機械翻訳モデルが、入力データから分散表現データを取得する埋込層 (例えば、入力データ埋込部により実現) と、機械翻訳用のニューラルネットワークモデル (MT用ニューラルネットワークモデル) とからなる場合、事前学習処理として、入力データから分散表現データを取得する埋込層 (例えば、入力データ埋込部により実現) と、言語横断言語モデル用のニューラルネットワークモデル (XLM用ニューラルネットワークモデル、XLM: Cross-lingual language model) とからなる言語横断言語モデルを設定し、埋込層および当該言語横断言語モデルのパラメータ最適化処理を行う。そして、機械翻訳モデルを、当該言語横断言語モデルの最適化処理後の埋込層 (例えば、入力データ埋込部により実現) と、MT用ニューラルネットワークモデルとからなるモデルとし、当該言語横断言語モデルの最適化処理後の状態を初期状態として、埋込層および機械翻訳モデルのパラメータ最適化処理を行うようにしてもよい。

30

【0016】

第2の発明は、第1の発明である疑似対訳データ生成用機械翻訳モデルの学習方法により取得された疑似対訳データ生成用機械翻訳モデルを用いて、適応先分野の疑似対訳データを取得する疑似対訳データ取得方法であって、第1機械翻訳ステップと、第2機械翻訳ステップと、疑似対訳データ取得ステップと、を備える。

【0017】

第1機械翻訳ステップは、他分野対訳データ集合 D s e t p (L 1 - L 2) から取得した第1言語のデータに対して、出力を適応先分野の第2言語に設定して疑似対訳データ生成用機械翻訳モデルを用いた機械翻訳処理を行うことで、他分野の第1言語のデータの機械翻訳結果データである適応先分野の第2言語疑似翻訳データを取得する。

40

【0018】

第2機械翻訳ステップは、他分野対訳データ集合 D s e t p (L 1 - L 2) から取得した第2言語のデータに対して、出力を適応先分野の第1言語に設定して疑似対訳データ生成用機械翻訳モデルを用いた機械翻訳処理を行うことで、他分野の第2言語のデータの機械翻訳結果データである適応先分野の第1言語疑似翻訳データを取得する。

【0019】

疑似対訳データ取得ステップは、第1機械翻訳ステップで取得された適応先分野の第2

50

言語疑似翻訳データと、第2機械翻訳ステップで取得された適応先分野の第1言語疑似翻訳データと、を対応づけることで、適応先分野の疑似対訳データを取得する。

【0020】

これにより、この疑似対訳データ取得方法では、適応先分野（対象分野）における対訳データが存在しない場合であっても、適応先分野における疑似対訳データを取得することができる。

【0021】

なお、第1機械翻訳ステップおよび第2機械翻訳ステップは、並列に実行されるものであってもよい。

【0022】

第3の発明は、第1の発明である疑似対訳データ生成用機械翻訳モデルの学習方法により取得された疑似対訳データ生成用機械翻訳モデルを用いて、適応先分野の疑似対訳データを取得する疑似対訳データ取得方法であって、単言語データ機械翻訳ステップと、疑似対訳データ取得ステップと、を備える。

【0023】

単言語データ機械翻訳ステップは、適応先分野単言語データ集合 $Dsetm(R1)$ から取得した第1言語のデータまたは適応先分野単言語データ集合 $Dsetm(R2)$ から取得した第2言語のデータに対して、出力を適応先分野の第2言語または第1言語に設定して疑似対訳データ生成用機械翻訳モデルを用いた機械翻訳処理を行うことで、適応先分野の第1言語のデータの機械翻訳結果データである適応先分野の第2言語疑似翻訳データまたは適応先分野の第2言語のデータの機械翻訳結果データである適応先分野の第1言語疑似翻訳データを取得する。

【0024】

疑似対訳データ取得ステップは、単言語データ機械翻訳ステップにおいて、疑似対訳データ生成用機械翻訳モデルの入力とした適応先分野の第1言語のデータまたは第2言語のデータと、単言語データ機械翻訳ステップで取得された適応先分野の第2言語疑似翻訳データまたは第1言語疑似翻訳データと、を対応づける、または、疑似対訳データ生成用機械翻訳モデルの入力とした適応先分野の第2言語のデータと、単言語データ機械翻訳ステップで取得された適応先分野の第1言語疑似翻訳データと、を対応づけることで、適応先分野の疑似対訳データを取得する。

【0025】

これにより、この疑似対訳データ取得方法では、適応先分野（対象分野）における対訳データが存在しない場合であっても、適応先分野における疑似対訳データを取得することができる。

【0026】

第4の発明は、第2または第3の発明であって、疑似対訳データ取得ステップが取得した適応先分野の疑似対訳データの各文対に対して、機械翻訳処理の結果の精度を示す信頼度を取得し、取得した信頼度が所定の値以上である文対を含む疑似対訳データのみを選択して出力するフィルター処理ステップをさらに備える。

【0027】

これにより、この疑似対訳データ取得方法において、所定の精度以上の疑似対訳データが取得されることが保証される。

【0028】

なお、「信頼度」は、その値が高い程、信頼度合いが高いことを示す指標である。また、信頼度は、機械翻訳処理により取得されるデータ（処理結果データ）が、単語列やトークンのデータである場合、個々の単語、トークンごとに付与されることが好ましい。

【0029】

第5の発明は、適応先分野の第1言語のデータに対して機械翻訳を行い第2言語のデータを取得するための機械翻訳モデルの学習方法であって、機械翻訳モデル学習ステップを備える。

10

20

30

40

50

【 0 0 3 0 】

機械翻訳モデル学習ステップは、第 2 から第 4 のいずれかの発明である疑似対訳データ取得方法により取得された適応先分野の疑似翻訳データと、適応先分野以外の分野の第 1 言語のデータである他分野第 1 言語データと、当該他分野第 1 言語データの第 2 言語の翻訳データである他分野第 2 言語データとからなる対訳データを複数含む他分野対訳データ集合 $Dsetp(L1 - L2)$ と、を用いて、機械翻訳モデルの学習を行う。

【 0 0 3 1 】

これにより、この機械翻訳モデルの学習方法では、適応先分野（対象分野）における対訳データが存在しない場合であっても、第 2 から第 5 のいずれかの発明である疑似対訳データ取得方法により取得された適応先分野の疑似翻訳データを用いて機械翻訳モデルを学習させることができる。そして、学習済みの機械翻訳モデルにより機械翻訳処理を行うことで、適応先分野（対象分野）における対訳データが存在しない場合であっても、適応先分野における機械翻訳を精度良く行うことができる。

10

【 0 0 3 2 】

第 6 の発明は、第 1 の発明である疑似対訳データ生成用機械翻訳モデルの学習方法により取得された学習済みの機械翻訳モデル、または、第 5 の発明である機械翻訳モデルの学習方法により取得された学習済みの機械翻訳モデルを用いて機械翻訳処理を行う機械翻訳装置である。

【 0 0 3 3 】

これにより、この機械翻訳装置では、第 1 の発明である疑似対訳データ生成用機械翻訳モデルの学習方法により取得された学習済みの機械翻訳モデル、または、第 5 の発明である機械翻訳モデルの学習方法により取得された学習済みの機械翻訳モデルを用いて機械翻訳処理を行うことができる。

20

【 0 0 3 4 】

なお、第 1 の発明である疑似対訳データ生成用機械翻訳モデルの学習方法をコンピュータに実行させるためのプログラムも本発明に含まれる。

【 0 0 3 5 】

また、第 2 から第 4 のいずれかの発明である疑似対訳データ取得方法をコンピュータに実行させるためのプログラムも本発明に含まれる。

【 0 0 3 6 】

また、第 5 の発明である機械翻訳モデルの学習方法をコンピュータに実行させるためのプログラムも本発明に含まれる。

30

【 0 0 3 7 】

さらに、第 5 の発明である機械翻訳モデルの学習方法により取得された学習済みの機械翻訳モデルを用いて機械翻訳処理を行う機械翻訳方法、当該機械翻訳方法をコンピュータに実行させるためのプログラムも本発明に含まれる。

【発明の効果】

【 0 0 3 8 】

本発明によれば、適応先分野（機械翻訳の対象とする分野）の対訳データが一切ない場合であっても、当該適応先分野における機械翻訳を精度良く実行するための疑似対訳データを生成する疑似対訳データ生成装置、および、当該疑似対訳データ生成装置で生成された疑似対訳データを用いて適応先分野における機械翻訳を精度良く実行する機械翻訳システム、および、当該機械翻訳システムで用いられる疑似対訳データ生成用機械翻訳モデルの学習方法、疑似対訳データ取得方法、および、機械翻訳モデルの学習方法を実現することができる。

40

【図面の簡単な説明】

【 0 0 3 9 】

【図 1】第 1 実施形態に係る機械翻訳システム 1000 の概略構成図。

【図 2】第 1 実施形態に係る疑似対訳データ生成装置 100 の概略構成図。

【図 3】第 1 実施形態に係る疑似対訳データ生成装置 100 の入力データ取得部 1、入力

50

データ埋込部 2、および X L M 処理部 3 の概略構成図。

【図 4】第 1 実施形態に係る疑似対訳データ生成装置 1 0 0 の入力データ取得部 1、入力データ埋込部 2、および機械翻訳処理部 5 の概略構成図。

【図 5】機械翻訳システム 1 0 0 0 で実行される処理のフローチャート。

【図 6】疑似対訳データ生成装置 1 0 0 で実行される処理のフローチャート。

【図 7】疑似対訳データ生成装置 1 0 0 で実行される処理を説明するための図。

【図 8】疑似対訳データ生成装置 1 0 0 で実行される処理を説明するための図。

【図 9】疑似対訳データ生成装置 1 0 0 で実行される処理を説明するための図。

【図 1 0】疑似対訳データ生成装置 1 0 0 で実行される処理を説明するための図。

【図 1 1】疑似対訳データ生成装置 1 0 0 で実行される処理を説明するための図。

10

【図 1 2】疑似対訳データ生成装置 1 0 0 で実行される処理を説明するための図。

【図 1 3】疑似対訳データ生成装置 1 0 0 で実行される処理を説明するための図。

【図 1 4】疑似対訳データ生成装置 1 0 0 で実行される処理を説明するための図。

【図 1 5】機械翻訳装置 M T 1 により取得されたデータの一例を示す図。

【図 1 6】C P U バス構成を示す図。

【発明を実施するための形態】

【0 0 4 0】

[第 1 実施形態]

第 1 実施形態について、図面を参照しながら、以下説明する。

【0 0 4 1】

20

< 1 . 1 : 機械翻訳システムの構成 >

図 1 は、第 1 実施形態に係る機械翻訳システム 1 0 0 0 の概略構成図である。

【0 0 4 2】

図 2 は、第 1 実施形態に係る疑似対訳データ生成装置 1 0 0 の概略構成図である。

【0 0 4 3】

図 3 は、第 1 実施形態に係る疑似対訳データ生成装置 1 0 0 の入力データ取得部 1、入力データ埋込部 2、および X L M 処理部 3 の概略構成図である。なお、図 3 では、第 1 セレクタ S E L 1 を省略して図示している。

【0 0 4 4】

図 4 は、第 1 実施形態に係る疑似対訳データ生成装置 1 0 0 の入力データ取得部 1、入力データ埋込部 2、および機械翻訳処理部 5 の概略構成図である。なお、図 4 では、第 1 セレクタ S E L 1 を省略して図示している。

30

【0 0 4 5】

機械翻訳システム 1 0 0 0 は、図 1 に示すように、対訳データ記憶部 D B p (L 1 - L 2) と、第 1 単言語データ記憶部 D B m (L 1) と、第 2 単言語データ記憶部 D B m (L 2) と、第 3 単言語データ記憶部 D B m (R 1) と、第 4 単言語データ記憶部 D B m (R 2) と、疑似対訳データ生成装置 1 0 0 と、疑似対訳データ格納部 D B 1 と、機械翻訳装置 M T 1 とを備える。

【0 0 4 6】

なお、以下では、疑似対訳データ生成装置 1 0 0 で生成する対訳データの対象とする分野を「適応先分野」といい、当該適応先分野以外の分野を「他分野」という。

40

【0 0 4 7】

対訳データ記憶部 D B p (L 1 - L 2) は、適応先分野以外の分野（他分野）の第 1 言語のデータである他分野第 1 言語データと、当該他分野第 1 言語データの第 2 言語の翻訳データである他分野第 2 言語データとからなる対訳データを複数含む他分野対訳データ集合 D s e t p (L 1 - L 2) を記憶する。そして、対訳データ記憶部 D B p (L 1 - L 2) は、疑似対訳データ生成装置 1 0 0 からの読み出し指令に基づいて、対訳データ記憶部 D B p (L 1 - L 2) に記憶している所定の対訳データ D 0 (L 1 - L 2) を読み出し、当該対訳データ D 0 (L 1 - L 2) を疑似対訳データ生成装置 1 0 0 へ出力する。対訳データ記憶部 D B p (L 1 - L 2) は、例えば、データベースにより実現される。なお、対

50

訳データ記憶部 DBp (L 1 - L 2) に記憶保持されている対訳データの数が多いたことが好ましい (対訳データ記憶部 DBp (L 1 - L 2) が大規模な対訳データを記憶保持していることが好ましい) 。

【 0 0 4 8 】

第 1 単言語データ記憶部 DBm (L 1) は、適応先分野以外の分野 (他分野) の第 1 言語の単言語データである他分野第 1 言語データを複数含む他分野第 1 言語データ集合 D s e t m (L 1) を記憶する。第 1 単言語データ記憶部 DBm (L 1) は、疑似対訳データ生成装置 1 0 0 からの読み出し指令に基づいて、第 1 単言語データ記憶部 DBm (L 1) に記憶している所定の第 1 言語の単言語データ D 0 (L 1) を読み出し、当該単言語データ D 0 (L 1) を疑似対訳データ生成装置 1 0 0 に出力する。第 1 単言語データ記憶部 DBm (L 1) は、例えば、データベースにより実現される。

10

【 0 0 4 9 】

第 2 単言語データ記憶部 DBm (L 2) は、適応先分野以外の分野 (他分野) の第 2 言語の単言語データである他分野第 2 言語データを複数含む他分野第 2 言語データ集合 D s e t m (L 2) を記憶する。第 2 単言語データ記憶部 DBm (L 2) は、疑似対訳データ生成装置 1 0 0 からの読み出し指令に基づいて、第 2 単言語データ記憶部 DBm (L 2) に記憶している所定の第 2 言語の単言語データ D 0 (L 2) を読み出し、当該単言語データ D 0 (L 2) を疑似対訳データ生成装置 1 0 0 に出力する。第 2 単言語データ記憶部 DBm (L 2) は、例えば、データベースにより実現される。

【 0 0 5 0 】

第 3 単言語データ記憶部 DBm (R 1) は、適応先分野の第 1 言語の単言語データである適応先分野第 1 言語データを複数含む適応先分野第 1 言語データ集合 D s e t m (R 1) を記憶する。第 3 単言語データ記憶部 DBm (R 1) は、疑似対訳データ生成装置 1 0 0 からの読み出し指令に基づいて、第 3 単言語データ記憶部 DBm (R 1) に記憶している所定の第 1 言語の単言語データ D 0 (R 1) を読み出し、当該単言語データ D 0 (R 1) を疑似対訳データ生成装置 1 0 0 に出力する。第 3 単言語データ記憶部 DBm (R 1) は、例えば、データベースにより実現される。

20

【 0 0 5 1 】

第 4 単言語データ記憶部 DBm (R 2) は、適応先分野の第 2 言語の単言語データである適応先分野第 2 言語データを複数含む適応先分野第 2 言語データ集合 D s e t m (R 2) を記憶する。第 4 単言語データ記憶部 DBm (R 2) は、疑似対訳データ生成装置 1 0 0 からの読み出し指令に基づいて、第 4 単言語データ記憶部 DBm (R 2) に記憶している所定の第 2 言語の単言語データ D 0 (R 2) を読み出し、当該単言語データ D 0 (R 2) を疑似対訳データ生成装置 1 0 0 に出力する。第 4 単言語データ記憶部 DBm (R 2) は、例えば、データベースにより実現される。

30

【 0 0 5 2 】

疑似対訳データ生成装置 1 0 0 は、図 2 に示すように、データ入力インターフェース I F 1 と、入力データ取得部 1 と、第 1 セレクタ S E L 1 と、入力データ埋込部 2 と、X L M 処理部 3 と、第 1 出力データ評価部 4 と、機械翻訳処理部 5 と、第 2 セレクタ S E L 2 と、第 2 出力データ評価部 6 と、第 1 バッファ B u f 1 と、第 2 入力データ取得処理部 7 と、フィルター処理部 8 とを備える。

40

【 0 0 5 3 】

データ入力インターフェース I F 1 は、疑似対訳データ生成装置 1 0 0 の各機能部を制御する制御部 (不図示) からの制御信号 C T L 1 を入力し、当該制御信号 C T L 1 に従い、対訳データ記憶部 DBp (L 1 - L 2)、第 1 単言語データ記憶部 DBm (L 1)、第 2 単言語データ記憶部 DBm (L 2)、第 3 単言語データ記憶部 DBm (R 1)、および、第 4 単言語データ記憶部 DBm (R 2) のいずれかから、所定のデータを読み出し、読み出したデータをデータ D 1 として入力データ取得部 1 に出力する。

【 0 0 5 4 】

入力データ取得部 1 は、図 3 に示すように、第 1 入力データ取得処理部 1 1 と、マスク

50

化処理部 1 2 と、正解データ取得部 1 3 と、入力データ出力部 1 4 と、を備える。入力データ取得部 1 は、入力されるデータの全部または一部を所定の期間、記憶保持することができるバッファ（不図示）を備えており、処理に応じて、当該バッファに記憶保持されているデータを使用することができる。

【 0 0 5 5 】

第 1 入力データ取得処理部 1 1 は、データ入力インターフェース I F 1 から出力されるデータ D 1 と、制御部から出力される制御信号 C T L 2 とを入力する。第 1 入力データ取得処理部 1 1 は、制御信号 C T L 2 に従い、データ D 1 から、マスク化処理部 1 2 および入力データ埋込部 2 に入力するためのデータを取得（生成）する。具体的には、第 1 入力データ取得処理部 1 1 は、入力されたデータ D 1（例えば、サブワード列（または単語列）のデータ、または、各サブワード（または各単語）に対応するインデックスの列のデータ）から、（ 1 ）トークン（トークンは、文字列を表すものであってもよい）を表すデータであるトークンデータ $x_{i0\ token}$ と、（ 2 ）当該トークンの位置を特定するための位置データ $x_{i\ pos}$ と、（ 3 ）当該トークンの言語を特定するための言語データ $x_{i\ lang}$ と、を取得する。そして、第 1 入力データ取得処理部 1 1 は、上記により取得した、トークンデータ $x_{i0\ token}$ をマスク化処理部 1 2 に出力し、位置データ $x_{i\ pos}$ と、言語データ $x_{i\ lang}$ とを第 1 セレクタ S E L 1 に出力する。

10

【 0 0 5 6 】

マスク化処理部 1 2 は、第 1 入力データ取得処理部 1 1 から出力されるトークンデータ $x_{i0\ token}$ と、制御部から出力される制御信号 C T L 2 とを入力する。（ 1 ）制御信号 C T L 2 がマスク化処理の実行を指示している場合、マスク化処理部 1 2 は、トークンデータ $x_{i0\ token}$ に対してマスク化処理を実行し、マスク化処理後のデータをトークンデータ $x_{i\ token}$ として、第 1 セレクタ S E L 1 に出力する。（ 2 ）制御信号 C T L 2 がマスク化処理の実行を指示していない場合、マスク化処理部 1 2 は、入力したトークンデータ $x_{i0\ token}$ をトークンデータ $x_{i\ token}$ として、第 1 セレクタ S E L 1 に出力する。

20

【 0 0 5 7 】

なお、入力データ取得部 1 から第 1 セレクタ S E L 1 に出力される（ 1 ）トークンデータ $x_{i\ token}$ と、（ 2 ）位置データ $x_{i\ pos}$ と、（ 3 ）言語データ $x_{i\ lang}$ と、を含むデータをデータ D 2 a と表記する。

30

【 0 0 5 8 】

正解データ取得部 1 3 は、データ入力インターフェース I F 1 から出力されるデータ D 1 と、制御部から出力される制御信号 C T L 2 とを入力する。正解データ取得部 1 3 は、制御信号 C T L 2 に従い、データ D 1 から、言語横断言語モデル（ X L M ）（入力データ埋込部 2（埋込層に対応）および X L M 処理部 3 の X L M 用ニューラルネットワークモデル 3 1 からなるモデル）、または、機械翻訳モデル（入力データ埋込部 2 および機械翻訳処理部 5 により構成されるニューラルネットワークモデル）の学習処理に用いる正解データ D _ c o r r e c t を生成し、当該正解データ D _ c o r r e c t を第 1 出力データ評価部 4 および第 2 出力データ評価部 6 に出力する。

【 0 0 5 9 】

入力データ出力部 1 4 は、データ入力インターフェース I F 1 から出力されるデータ D 1 と、制御部から出力される制御信号 C T L 2 とを入力する。入力データ出力部 1 4 は、制御信号 C T L 2 に従い、入力されたデータ D 1 を、データ D 1 _ o r g としてフィルター処理部 8 に出力する。

40

【 0 0 6 0 】

第 1 セレクタ S E L 1 は、入力データ取得部 1 の第 1 入力データ取得処理部 1 1 から出力されるデータ D 2 a と、第 2 入力データ取得処理部 7 から出力されるデータ D 2 b と、制御部から出力される選択信号 s e l 1 とを入力する。第 1 セレクタ S E L 1 は、選択信号 s e l 1 に従い、データ D 2 a、または、データ D 2 b を選択して、データ D 3 として入力データ埋込部 2 に出力する。

50

【0061】

入力データ埋込部2は、図3、図4に示すように、トークン埋込部21と、位置埋込部22と、言語埋込部23とを備える。

【0062】

トークン埋込部21は、データD3に含まれるトークンデータ $x_{i\ token}$ を入力し、入力したトークンデータ $x_{i\ token}$ の分散表現データを取得し、取得した分散表現データを分散表現データ $x_{i\ 'token}$ として、XLM処理部3および機械翻訳処理部5に出力する。なお、トークン埋込部21は、例えば、トークンデータ $x_{i\ token}$ に対して、分散表現データを取得するための行列による行列演算を行うことで、分散表現データ $x_{i\ 'token}$ を取得する。例えば、トークン埋込部21は、下記の行列演算による処理を行うことで、分散表現データ $x_{i\ 'token}$ を取得する。なお、行列 W_{token} の各要素（重み付け係数に相当）は、パラメータ emb の一部である。パラメータ emb は、XLM処理部3または機械翻訳処理部5から入力データ埋込部2に入力されるパラメータ更新データ $update(emb)$ により更新される。

$$x_{i\ 'token} = x_{i\ token} \cdot W_{token}$$

$x_{i\ token}$ ：各トークン（入力データ）文字列を表すベクトル（例えば、 $1 \times n_1$ の行列（ n_1 次元ベクトル）（ n_1 ：自然数））

W_{token} ：分散表現データを取得するための行列（例えば、 $n_1 \times m_1$ の行列（ n_1, m_1 ：自然数））

$x_{i\ 'token}$ ：入力データ $x_{i\ token}$ の分散表現データ（例えば、 $1 \times m_1$ の行列（ m_1 次元ベクトル）（ m_1 ：自然数））

位置埋込部22は、データD3に含まれる位置データ $x_{i\ pos}$ を入力し、入力した位置データ $x_{i\ pos}$ の分散表現データを取得し、取得した分散表現データを分散表現データ $x_{i\ 'pos}$ として、XLM処理部3および機械翻訳処理部5に出力する。なお、位置埋込部22は、例えば、位置データ $x_{i\ pos}$ に対して、分散表現データを取得するための行列による行列演算を行うことで、分散表現データ $x_{i\ 'pos}$ を取得する。例えば、位置埋込部22は、下記の行列演算による処理を行うことで、分散表現データ $x_{i\ 'pos}$ を取得する。なお、行列 W_{pos} の各要素（重み付け係数に相当）は、パラメータ emb の一部である。パラメータ emb は、XLM処理部3または機械翻訳処理部5から入力データ埋込部2に入力されるパラメータ更新データ $update(emb)$ により更新される。

$$x_{i\ 'pos} = x_{i\ pos} \cdot W_{pos}$$

$x_{i\ pos}$ ：各トークン（入力データ）の位置を表すベクトル（例えば、 $1 \times n_2$ の行列（ n_2 次元ベクトル）（ n_2 ：自然数））

W_{pos} ：分散表現データを取得するための行列（例えば、 $n_2 \times m_2$ の行列（ n_2, m_2 ：自然数））

$x_{i\ 'pos}$ ：入力データ $x_{i\ pos}$ の分散表現データ（例えば、 $1 \times m_2$ の行列（ m_2 次元ベクトル）（ m_2 ：自然数））

言語埋込部23は、データD3に含まれる言語データ $x_{i\ lang}$ を入力し、入力した言語データ $x_{i\ lang}$ の分散表現データを取得し、取得した分散表現データを分散表現データ $x_{i\ 'lang}$ として、XLM処理部3および機械翻訳処理部5に出力する。なお、言語埋込部23は、例えば、言語データ $x_{i\ lang}$ に対して、分散表現データを取得するための行列による行列演算を行うことで、分散表現データ $x_{i\ 'lang}$ を取得する。例えば、言語埋込部23は、下記の行列演算による処理を行うことで、分散表現データ $x_{i\ 'lang}$ を取得する。なお、行列 W_{lang} の各要素（重み付け係数に相当）は、パラメータ emb の一部である。パラメータ emb は、XLM処理部3または機械翻訳処理部5から入力データ埋込部2に入力されるパラメータ更新データ $update(emb)$ により更新される。

$$x_{i\ 'lang} = x_{i\ lang} \cdot W_{lang}$$

$x_{i\ lang}$ ：各トークン（入力データ）の言語を表すベクトル（例えば、 $1 \times n_3$ の行列（ n_3 次元ベクトル）（ n_3 ：自然数））

10

20

30

40

50

W_{lang} : 分散表現データを取得するための行列 (例えば、 $n_3 \times m_3$ の行列 (n_3, m_3 : 自然数))

$x_{i'lang}$: 入力データ $x_{i'lang}$ の分散表現データ (例えば、 $1 \times m_3$ の行列 (m_3 次元ベクトル)) (m_3 : 自然数)

入力データ埋込部 2 は、上記により取得された分散表現データをデータ D_4 として、XLM 処理部 3 および機械翻訳処理部 5 に出力する。

【0063】

XLM 処理部 3 は、図 3 に示すように、XLM 用ニューラルネットワークモデル 31 を備える。

【0064】

XLM 用ニューラルネットワークモデル 31 は、例えば、下記文献 A に開示されているニューラルネットワークモデルであり、例えば、下記文献 B に開示されているトランスフォーマーモデルのアーキテクチャを採用したニューラルネットワークモデルにより実現される。

(文献 A) : Alexis Conneau and Guillaume Lample (2019). Cross-Lingual Language Model Pretraining. In Proceedings of the 32nd Neural Information Processing Systems Conference (NeurIPS), pp. 7057-7067.

(文献 B) : Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin (2017). Attention is All You Need. In Proceedings of the 30th Neural Information Processing Systems Conference (NeurIPS), pp. 5998-6008.

XLM 用ニューラルネットワークモデル 31 は、入力データ埋込部 2 から出力される分散表現データ D_4 ($= \{x_{i'token}, x_{i'pos}, x_{i'lang}\}$) を入力する。

【0065】

XLM 用ニューラルネットワークモデル 31 は、第 1 出力データ評価部 4 から出力されるパラメータ更新データ $update(x_{LM})$ に基づいて、XLM 用ニューラルネットワークモデル 31 のパラメータを更新する。また、XLM 用ニューラルネットワークモデル 31 は、XLM 用ニューラルネットワークモデルのパラメータを更新した後、埋込層 (入力データ埋込部 2 のトークン埋込部 21、位置埋込部 22、および、言語埋込部 23 に対応) のパラメータを更新するためのパラメータ更新データ $update(emb)$ を生成し、当該パラメータ更新データ $update(emb)$ を入力データ埋込部 2 に出力する。

【0066】

また、XLM 用ニューラルネットワークモデル 31 は、入力データ埋込部 2 から出力されるデータ D_4 に対して、XLM 処理部 3 を実行し、データ D_{5x} を取得する。そして、XLM 処理部 3 は、上記により取得したデータ D_{5x} を第 1 出力データ評価部 4 に出力する。

【0067】

第 1 出力データ評価部 4 は、図 2 に示すように、XLM 処理部 3 から出力されるデータ D_{5x} と、入力データ取得部 1 から出力される正解データ $D_{correct}$ とを入力する。第 1 出力データ評価部 4 は、損失を評価するために、XLM 処理部 3 の出力であるデータ D_{5x} に対する正解データ $D_{correct}$ から損失評価用のデータ $D_{correct'}$ を取得し、データ D_{5x} とデータ $D_{correct'}$ とから損失を取得する (詳細については後述)。そして、第 1 出力データ評価部 4 は、所定の学習データに対する損失 (学習損失) に基づいて XLM 処理部 3 の XLM 用ニューラルネットワークモデル 31 のパラメータ x_{LM} を更新するためのデータ $update(x_{LM})$ を生成し、当該データ $update(x_{LM})$ を XLM 処理部 3 に出力する。

【0068】

機械翻訳処理部 5 は、図 4 に示すように、MT 用ニューラルネットワークモデル 51 を備える。

【0069】

10

20

30

40

50

MT用ニューラルネットワークモデル51は、エンコーダ/デコーダ方式のニューラルネットワークモデルであり、例えば、上記文献Bに開示されているトランスフォーマーモデルのアーキテクチャによるエンコーダ、デコーダの構成を採用したニューラルネットワークモデルである。

【0070】

MT用ニューラルネットワークモデル51は、入力データ埋込部2から出力される分散表現データD4(= { x_i 'token, x_i 'pos, x_i 'lang})と、制御部から出力される制御信号CTL3とを入力する。MT用ニューラルネットワークモデル51は、制御信号CTL3で指示された言語のデータを出力する。

【0071】

また、MT用ニューラルネットワークモデル51は、第2出力データ評価部6から出力されるパラメータ更新データupdate(M_T)に基づいて、MT用ニューラルネットワークモデル51のパラメータを更新する。

【0072】

また、MT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4に対して、機械翻訳処理を実行し、データD5を取得する。なお、MT用ニューラルネットワークモデル51は、制御信号CTL3によって指定された種別のデータがデータD5として出力されるように、機械翻訳処理を実行する。なお、制御信号CTL3によって指定される種別は、(1)L1(他分野、第1言語)、(2)L2(他分野、第2言語)、(3)R1(適応先分野、第1言語)、および、(4)R2(適応先分野、第2言語)のいずれかである。

【0073】

機械翻訳処理部5は、MT用ニューラルネットワークモデル51により取得したデータD5を第2セクタSEL2に出力する。

【0074】

なお、機械翻訳処理部5は、疑似対訳データを取得する場合、入力データ埋込部2に入力されるデータ(例えば、他分野第1言語のデータD3(L1))に対して機械翻訳処理を行い取得したデータ(例えば、R2のデータ)と、入力データ埋込部2に入力されるデータ(他分野第1言語のデータD3(L1)の他分野第2言語の対訳文に相当するデータD3(L2))に対して機械翻訳処理を行い取得したデータ(例えば、R1のデータ)とをペアにして、データD5(例えば、データD5(R1-R2)(適応先分野第1言語の機械翻訳文(入力されたL2のデータの機械翻訳結果)と、適応先分野の第2言語の機械翻訳文(入力されたL1のデータの機械翻訳結果)とをペアとして構成する対訳データD5(R1-R2)))として、第2セクタSEL2に出力する。

【0075】

第2セクタSEL2は、1入力2出力の切替器であり、機械翻訳処理部5から出力されるデータD5と、制御部から出力される選択信号sel2とを入力する。第2セクタSEL2は、選択信号sel2に従い、データD5を、(1)データD6aとして第2出力データ評価部6およびフィルター処理部8に出力する、または、(2)データD6bとして第1バッファBuf1に出力する。

【0076】

第2出力データ評価部6は、図2に示すように、第2セクタSEL2から出力されるデータD6aと、入力データ取得部1から出力される正解データD__correctとを入力する。第2出力データ評価部6は、損失を評価するために、機械翻訳処理部5の出力であるデータD6aの正解データD__correctから損失評価用のデータD__correct'を取得し、データD6aとデータD__correct'とから損失を取得する(詳細については後述)。そして、第2出力データ評価部6は、所定の学習データに対する損失(学習損失)に基づいて機械翻訳処理部5のMT用ニューラルネットワークモデル51のパラメータ M_T を更新するためのデータupdate(M_T)を生成する。

【0077】

10

20

30

40

50

そして、第2出力データ評価部6は、データupdate (MT)を機械翻訳処理部5に出力する。

【0078】

第1バッファBuf1は、第2セレクタSEL2から出力されるデータD6bを記憶保持するためのバッファである。第1バッファBuf1は、記憶保持しているデータをデータD6b'として第2入力データ取得処理部7に出力する。

【0079】

第2入力データ取得処理部7は、第1入力データ取得処理部11と同様の機能部であり、第1バッファBuf1から出力されるデータD6b'を入力し、データD6b'から入力データ埋込部2に入力するためのデータを取得(生成)する。そして、第2入力データ取得処理部7は、取得したデータをデータD2bとして、第1セレクタSEL1に出力する。

10

【0080】

フィルター処理部8は、第2セレクタSEL2から出力されるデータD6aと、入力データ取得部1から出力されるデータD1_orgと、制御部から出力される制御信号CTL4とを入力する。

【0081】

また、フィルター処理部8は、(1)データD6aが対訳データ(疑似対訳データ)である場合、各対訳データに信頼度を付与し、当該信頼度によるフィルタリング処理を行い、(2)データD6aが機械翻訳データである場合、当該機械翻訳データとデータD1_orgとをペアリングして取得した対訳データ(疑似対訳データ)について、信頼度を付与し、当該信頼度によるフィルタリング処理を行う。

20

【0082】

そして、フィルター処理部8は、フィルター処理後の対訳データをデータDpsd1として、疑似対訳データ格納部DB1に出力する。

【0083】

疑似対訳データ格納部DB1は、図1に示すように、疑似対訳データ生成装置100から出力されるデータDpsd1を記憶保持するための記憶部である。また、疑似対訳データ格納部DB1は、機械翻訳装置MT1から読み出し指令に基づいて、記憶保持しているデータを読み出し、読み出したデータをデータDpsd2として機械翻訳装置MT1に出力する。

30

【0084】

機械翻訳装置MT1は、図1に示すように、対訳データ記憶部DBp(L1-L2)および/または疑似対訳データ格納部DB1から、データを読み出し、読み出した当該データを用いて、機械翻訳モデルの学習処理を行う。そして、機械翻訳装置MT1は、当該学習処理により取得した機械翻訳モデルの学習済みモデルを用いて、機械翻訳処理を行う。つまり、機械翻訳装置MT1は、起点言語データDin_eを入力し、当該起点言語データDin_eに対して機械翻訳処理を行い、目標言語データDout_jを取得する。

【0085】

<1.2: 機械翻訳システムの動作>

以上のように構成された機械翻訳システム1000の動作について、説明する。

40

【0086】

図5は、機械翻訳システム1000で実行される処理のフローチャートである。

【0087】

図6は、疑似対訳データ生成装置100で実行される処理(事前学習処理)のフローチャートである。

【0088】

図7~図11は、疑似対訳データ生成装置100で実行される処理(事前学習処理)を説明するための図である。

【0089】

図12~図14は、疑似対訳データ生成装置100で実行される処理(疑似対訳データ

50

生成処理)を説明するための図である。

【0090】

以下では、図面を参照しながら、機械翻訳システム1000の動作について、説明する。

【0091】

(1.2.1:事前学習処理)

まず、機械翻訳システム1000の疑似対訳データ生成装置100において実行される事前学習処理(疑似対訳データ生成用NMTモデル(入力データ埋込部2(埋込層に相当)と機械翻訳処理部5のMT用ニューラルネットワークモデル51とにより実現されるモデル)の事前学習処理)(図5のフローチャートのステップS1)について、XLMの学習処理と、疑似対訳データ生成用NMTモデルの学習処理とに分けて説明する。

10

【0092】

(1.2.1.1:XLMの学習処理)

まず、XLMの学習処理について説明する。

【0093】

XLMの学習処理は、図6に示すように、

(A)マスク化言語モデル(MLM)の処理(ステップS111)と、

(B)翻訳言語モデル(TLM)の処理(ステップS112)と、

(C)損失の計算処理(ステップS113)と、

(D)パラメータ($x_{LM, emb}$)の更新処理(ステップS114)と、

から構成される。

20

【0094】

なお、図6のフローチャートに示すように、ステップS111~ステップS114の処理は、終了条件が満たされるまで繰り返し実行される(図6のループ処理(ループ1)S110~S115)。

【0095】

(1.2.1.1A:マスク化言語モデル(MLM)の処理)

マスク化言語モデル(MLM)の処理において、疑似対訳データ生成装置100は、図7に示すように、入力データ埋込部2への入力をマスク化データとし、正解データを原データ(マスク化していないデータ)として、学習処理を行う。

【0096】

30

M1:MLMの処理(L1.mask L1)

具体的には、データ入力インターフェースIF1は、第1単言語データ記憶部DBm(L1)から他分野第1言語の単言語データD0(L1)を読み出し、読み出した単言語データをデータD1(L1)として、入力データ取得部1に出力する。

【0097】

入力データ取得部1の第1入力データ取得処理部11は、データD1(他分野第1言語のデータD1(L1))から、(1)トークンデータ $x_{i0 token}(=x_{i0 token}(L1))$ と、(2)当該トークンの位置を特定するための位置データ $x_{i pos}(=x_{i pos}(L1))$ と、(3)当該トークンの言語を特定するための言語データ $x_{i lang}(=x_{i lang}(L1))$ と、を取得する。

40

【0098】

そして、入力データ取得部1は、上記のように取得したトークンデータ $x_{i0 token}$ をマスク化処理部12に出力する。

【0099】

制御部(不図示)は、入力データ取得部1のマスク化処理部12に対して、マスク化処理を実行することを指示する制御信号CTL2を生成し、当該制御信号CTL2をマスク化処理部12に出力する。

【0100】

マスク化処理部12は、制御信号CTL2に従い、第1入力データ取得処理部11から入力したトークンデータ $x_{i0 token}$ の一部のトークンをマスクトークン(例えば、ラ

50

ベル「[mask]」を付与したデータ(トークン)に設定する)に置換するマスク化処理を行う。そして、マスク化処理部12は、当該マスク化処理後のトークンデータをトークンデータ $x_{i\text{token}}$ として、第1セレクタSEL1に出力する。

【0101】

また、第1入力データ取得処理部11は、位置データ $x_{i\text{pos}} (= x_{i\text{pos}}(L1))$ を第1セレクタSEL1に出力し、言語データ $x_{i\text{lang}} (= x_{i\text{lang}}(L1))$ を第1セレクタSEL1に出力する。

【0102】

つまり、入力データ取得部1は、上記により取得したデータをデータ $D2a (= D2a(L1.\text{mask}) = \{x_{i\text{token}}(L1), x_{i\text{pos}}(L1), x_{i\text{lang}}(L1)\})$ として第1セレクタSEL1に出力する。

10

【0103】

なお、データを示す変数の末尾に括弧書きを追加し、そのデータの種別を表すものとする。例えば、 $x_{i\text{token}}(L1)$ は、他分野第1言語(L1)のデータから導出されたデータであることを示している(以下、同様)。また、上記括弧の中の「.mask」という表記は、マスク化処理されたトークンデータを含むデータであることを示しており、例えば、「 $L1.\text{mask}$ 」は、他分野第1言語(L1)のデータから導出されたデータにおけるトークンデータに対してマスク化処理を行ったデータを含むデータであることを示している(以下、同様)。

【0104】

20

正解データ取得部13は、制御信号CTL2に従い、データ $D1 (= D1(L1))$ から、入力データ埋込部2およびXLM処理部3での学習処理(XLMの学習処理)に用いる正解データ $D_correct (= D_correct(L1))$ を生成する。そして、正解データ取得部13は、生成したデータを正解データ $D_correct (= D_correct(L1))$ として第1出力データ評価部4に出力する。

【0105】

制御部は、第1セレクタSEL1の端子「0」を選択する選択信号sel1を生成し、当該選択信号sel1を第1セレクタSEL1に出力する。

【0106】

第1セレクタSEL1は、選択信号sel1に従い、端子「0」を選択し、入力データ取得部1から出力されるデータ $D2a(L1.\text{mask})$ を、データ $D3 (= D3(L1.\text{mask}))$ として入力データ埋込部2に出力する。

30

【0107】

入力データ埋込部2のトークン埋込部21は、データ $D3(L1.\text{mask})$ に含まれるトークンデータ $x_{i\text{token}}$ を入力し、入力したトークンデータ $x_{i\text{token}}$ の分散表現データを、例えば、下記数式に相当する処理を実行して取得する。

$$x_{i\text{token}}' = x_{i\text{token}} \cdot W_{\text{token}}$$

$x_{i\text{token}}$: 各トークン(入力データ)文字列を表すベクトル(例えば、 $1 \times n1$ の行列($n1$ 次元ベクトル)($n1$:自然数))

W_{token} : 分散表現データを取得するための行列(例えば、 $n1 \times m1$ の行列($n1, m1$:自然数))

40

$x_{i\text{token}}'$: 入力データ $x_{i\text{token}}$ の分散表現データ(例えば、 $1 \times m1$ の行列($m1$ 次元ベクトル)($m1$:自然数))

入力データ埋込部2の位置埋込部22は、データ $D3(L1.\text{mask})$ に含まれる位置データ $x_{i\text{pos}}$ を入力し、入力した位置データ $x_{i\text{pos}}$ の分散表現データを、例えば、下記数式に相当する処理を実行して取得する。

$$x_{i\text{pos}}' = x_{i\text{pos}} \cdot W_{\text{pos}}$$

$x_{i\text{pos}}$: 各トークン(入力データ)の位置を表すベクトル(例えば、 $1 \times n2$ の行列($n2$ 次元ベクトル)($n2$:自然数))

W_{pos} : 分散表現データを取得するための行列(例えば、 $n2 \times m2$ の行列($n2,$

50

m2 : 自然数))

x i ' p o s : 入力データ x i p o s の分散表現データ (例えば、 1 x m 2 の行列 (m 2 次元ベクトル) (m 2 : 自然数))

入力データ埋込部 2 の言語埋込部 2 3 は、データ D 3 (L 1 . m a s k) に含まれる言語データ x i l a n g を入力し、入力した言語データ x i l a n g の分散表現データを、例えば、下記数式に相当する処理を実行して取得する。

$$x i ' l a n g = x i l a n g \cdot W l a n g$$

x i l a n g : 各トークン (入力データ) の言語を表すベクトル (例えば、 1 x n 3 の行列 (n 3 次元ベクトル) (n 3 : 自然数))

W l a n g : 分散表現データを取得するための行列 (例えば、 n 3 x m 3 の行列 (n 3 , m 3 : 自然数))

10

x i ' l a n g : 入力データ x i l a n g の分散表現データ (例えば、 1 x m 3 の行列 (m 3 次元ベクトル) (m 3 : 自然数))

入力データ埋込部 2 は、上記により取得された分散表現データをデータ D 4 (L 1 . m a s k) として、X L M 処理部 3 に出力する。なお、上記のトークン文字列や言語識別子の埋め込み方法、位置の埋め込み方法は、上記に限定されることなく、他の方法を用いてもよい。

【 0 1 0 8 】

X L M 処理部 3 の X L M 用ニューラルネットワークモデル 3 1 は、入力データ埋込部 2 から出力されるデータ D 4 (L 1 . m a s k) に対して、X L M 処理を実行し、データ D 5 x (= D 5 x (L 1 . m a s k)) を取得する。そして、X L M 処理部 3 は、上記により取得したデータ D 5 x を第 1 出力データ評価部 4 に出力する。

20

【 0 1 0 9 】

第 1 出力データ評価部 4 は、X L M 処理部 3 から出力されるデータ D 5 x (L 1 . m a s k) と、入力データ取得部 1 から出力される正解データ D _ c o r r e c t (= D _ c o r r e c t (L 1)) とを入力する。第 1 出力データ評価部 4 は、X L M 処理部 3 の出力であるデータ D 5 (L 1 . m a s k) と、その正解データ D _ c o r r e c t (L 1) とから損失を取得する。

【 0 1 1 0 】

M 2 : M L M の処理 (L 2 . m a s k L 2)

30

次に、疑似対訳データ生成装置 1 0 0 のデータ入力インターフェース I F 1 は、第 2 単言語データ記憶部 D B m (L 2) から他分野の第 2 言語の単言語データ D 0 (L 2) を読み出し、読み出した単言語データをデータ D 1 (L 2) として、入力データ取得部 1 に出力する。そして、上記と同様の処理を行う。つまり、入力データ埋込部 2 への入力データを、データ D 3 (L 2 . m a s k) 、すなわち、

$$D 3 (L 2 . m a s k) = \{ x i t o k e n (L 2) , x i p o s (L 2) , x i l a n g (L 2) \}$$

x i t o k e n (L 2) : 第 2 単言語データ記憶部 D B m (L 2) から読み出した単言語データ D 1 (L 2) のトークンデータ x i 0 t o k e n (L 2) の一部のトークンをマスクトークンに置換するマスク化処理を行うことで取得したマスク化トークンデータ

40

x i p o s (L 2) : 第 2 単言語データ記憶部 D B m (L 2) から読み出した単言語データ D 1 (L 2) のトークンデータの位置を示すデータ (位置データ)

x i l a n g (L 2) : 第 2 単言語データ記憶部 D B m (L 2) から読み出した単言語データ D 1 (L 2) のトークンの言語を示すデータ (言語データ)

とし、正解データを D _ c o r r e c t (L 2) (= 原データ (マスク化していないデータ)) として、上記と同様の処理を行う。

【 0 1 1 1 】

M 3 : M L M の処理 (R 1 . m a s k R 1)

次に、疑似対訳データ生成装置 1 0 0 のデータ入力インターフェース I F 1 は、第 3 単言語データ記憶部 D B m (R 1) から適応先分野の第 1 言語の単言語データ D 0 (R 1)

50

を読み出し、読み出した単言語データをデータD1(R1)として、入力データ取得部1に出力する。そして、上記と同様の処理を行う。つまり、入力データ埋込部2への入力データを、データD3(R1.mask)、すなわち、

$D3(R1.mask) = \{x_{i token}(R1), x_{i pos}(R1), x_{i lang}(R1)\}$

$x_{i token}(R1)$: 第3単言語データ記憶部DBm(R1)から読み出した単言語データD1(R1)のトークンデータ $x_{i0 token}(R1)$ の一部のトークンをマスクトークンに置換するマスク化処理を行うことで取得したマスク化トークンデータ

$x_{i pos}(R1)$: 第3単言語データ記憶部DBm(R1)から読み出した単言語データD1(R1)のトークンデータの位置を示すデータ(位置データ) 10

$x_{i lang}(R1)$: 第3単言語データ記憶部DBm(R1)から読み出した単言語データD1(R1)のトークンデータの言語を示すデータ(言語データ)

とし、正解データをD__correct(R1)(=原データ(マスク化していないデータ))として、上記と同様の処理を行う。

【0112】

M4:MLMの処理(R2.mask R2)

次に、疑似対訳データ生成装置100のデータ入力インターフェースIF1は、第4単言語データ記憶部DBm(R2)から適応先分野の第2言語の単言語データD0(R2)を読み出し、読み出した単言語データをデータD1として、入力データ取得部1に出力する。そして、上記と同様の処理を行う。つまり、入力データ埋込部2への入力データを、データD3(R2.mask)、すなわち、 20

$D3(R2.mask) = \{x_{i token}(R2), x_{i pos}(R2), x_{i lang}(R2)\}$

$x_{i token}(R2)$: 第4単言語データ記憶部DBm(R2)から読み出した単言語データD1(R2)のトークンデータ $x_{i0 token}(R2)$ の一部のトークンをマスクトークンに置換するマスク化処理を行うことで取得したマスク化トークンデータ

$x_{i pos}(R2)$: 第4単言語データ記憶部DBm(R2)から読み出した単言語データD1(R2)のトークンデータの位置を示すデータ(位置データ)

$x_{i lang}(R2)$: 第4単言語データ記憶部DBm(R2)から読み出した単言語データD1(R2)のトークンデータの言語を示すデータ(言語データ) 30

とし、正解データをD__correct(R2)(=原データ(マスク化していないデータ))として、上記と同様の処理を行う。

【0113】

(1.2.1.1B:翻訳言語モデル(TLM)の処理)

次に、翻訳言語モデル(TLM)の処理において、疑似対訳データ生成装置100は、図7に示すように、入力データ埋込部2への入力(対訳データ)をマスク化データとし、正解データを原データ(マスク化していないデータ)として、学習処理を行う。

【0114】

T1:TLMの処理((L1-L2).mask L1-L2)

具体的には、データ入力インターフェースIF1は、対訳データ記憶部DBp(L1-L2)から他分野の対訳データD0(L1-L2)を読み出し、読み出した対訳データをデータD1(L1-L2)として、入力データ取得部1に出力する。 40

【0115】

入力データ取得部1の第1入力データ取得処理部11は、データD1(対訳データD1(L1-L2)(他分野第1言語の文と、当該文の第2言語の対訳文をペアとして構成する対訳データD1(L1-L2))から、(1)トークンデータ $x_{i0 token}(=x_{i0 token}(L1-L2))$ と、(2)当該トークンの位置を特定するための位置データ $x_{i pos}(=x_{i pos}(L1-L2))$ と、(3)当該トークンの言語を特定するための言語データ $x_{i lang}(=x_{i lang}(L1-L2))$ と、を取得する。

【0116】

そして、入力データ取得部 1 は、上記のように取得したトークンデータ $x_{i0\ token}$ ($L1 - L2$) をマスク化処理部 12 に出力する。

【0117】

制御部 (不図示) は、入力データ取得部 1 のマスク化処理部 12 に対して、マスク化処理を実行することを指示する制御信号 $CTL2$ を生成し、当該制御信号 $CTL2$ をマスク化処理部 12 に出力する。

【0118】

マスク化処理部 12 は、制御信号 $CTL2$ に従い、第 1 入力データ取得処理部 11 から入力したトークンデータ $x_{i0\ token}$ ($L1 - L2$) の一部のトークンをマスクトークン (例えば、文字列を人工的なトークン「[MASK]」としたデータ (トークン) に設定する) に置換するマスク化処理を行う。そして、マスク化処理部 12 は、当該マスク化処理後のトークンデータをトークンデータ $x_{i\ token}$ ($(L1 - L2).mask$) として、第 1 セレクタ $SEL1$ に出力する。なお、対訳データ $L1 - L2$ (他分野第 1 言語の文と、当該文の第 2 言語の対訳文をペアとして構成する対訳データ) に対してマスク化処理を行い取得されるデータを「 $(L1 - L2).mask$ 」と表記する (以下、同様)。

【0119】

また、第 1 入力データ取得処理部 11 は、位置データ $x_{i\ pos}$ ($=x_{i\ pos}$) ($=x_{i\ pos}$ ($L1 - L2$)) を第 1 セレクタ $SEL1$ に出力し、言語データ $x_{i\ lang}$ ($=x_{i\ lang}$) ($=x_{i\ lang}$ ($L1 - L2$)) を第 1 セレクタ $SEL1$ に出力する。

【0120】

つまり、入力データ取得部 1 は、上記により取得したデータをデータ $D2a$ ($=D2a$ ($(L1 - L2).mask$) $=\{x_{i\ token}$ ($(L1 - L2).mask$), $x_{i\ pos}$ ($L1 - L2$), $x_{i\ lang}$ ($L1 - L2$)}) として第 1 セレクタ $SEL1$ に出力する。

【0121】

正解データ取得部 13 は、制御信号 $CTL2$ に従い、データ $D1$ から、入力データ埋込部 2 および MLM 処理部 3 での学習処理 (MLM の学習処理) に用いる正解データ $D_correct$ ($=D_correct$ ($L1 - L2$)) を生成する。そして、正解データ取得部 13 は、生成したデータを正解データ $D_correct$ ($L1 - L2$) として第 1 出力データ評価部 4 に出力する。

【0122】

制御部は、第 1 セレクタ $SEL1$ の端子「0」を選択する選択信号 $sel1$ を生成し、当該選択信号 $sel1$ を第 1 セレクタ $SEL1$ に出力する。

【0123】

第 1 セレクタ $SEL1$ は、選択信号 $sel1$ に従い、端子「0」を選択し、入力データ取得部 1 から出力されるデータ $D2a$ ($(L1 - L2).mask$) を、データ $D3$ ($(L1 - L2).mask$) として、入力データ埋込部 2 に出力する。

【0124】

入力データ埋込部 2 のトークン埋込部 21 は、データ $D3$ ($(L1 - L2).mask$) に含まれるトークンデータ $x_{i\ token}$ ($=x_{i\ token}$ ($(L1 - L2).mask$)) を入力し、入力したトークンデータ $x_{i\ token}$ の分散表現データ $x_{i\ 'token}$ を、上記の MLM の処理で説明したのと同じ処理を実行して取得する。

【0125】

入力データ埋込部 2 の位置埋込部 22 は、データ $D3$ ($(L1 - L2).mask$) に含まれる位置データ $x_{i\ pos}$ ($=x_{i\ pos}$ ($L1 - L2$)) を入力し、入力した位置データ $x_{i\ pos}$ の分散表現データ $x_{i\ 'pos}$ を、上記の MLM の処理で説明したのと同じ処理を実行して取得する。

【0126】

入力データ埋込部 2 の言語埋込部 23 は、データ $D3$ ($(L1 - L2).mask$) に含まれる言語データ $x_{i\ lang}$ ($=x_{i\ lang}$ ($L1 - L2$)) を入力し、入力した

10

20

30

40

50

言語データ $x_{i l a n g}$ の分散表現データ $x_{i ' l a n g}$ を、上記の M L M の処理で説明したのと同じ処理を実行して取得する。

【0127】

入力データ埋込部 2 は、上記により取得された分散表現データをデータ $D 4 ((L 1 - L 2) . m a s k)$ として、X L M 処理部 3 に出力する。

【0128】

X L M 処理部 3 の X L M 用ニューラルネットワークモデル 3 1 は、入力データ埋込部 2 から出力されるデータ $D 4 ((L 1 - L 2) . m a s k)$ に対して、X L M 処理を実行し、データ $D 5 x (= D 5 x ((L 1 - L 2) . m a s k))$ を取得する。そして、X L M 処理部 3 は、上記により取得したデータ $D 5 x ((L 1 - L 2) . m a s k)$ を第 1 出力データ評価部 4 に出力する。

10

【0129】

第 1 出力データ評価部 4 は、X L M 処理部 3 から出力されるデータ $D 5 x ((L 1 - L 2) . m a s k)$ と、入力データ取得部 1 から出力される正解データ $D _ c o r r e c t (L 1 - L 2)$ とを入力する。第 1 出力データ評価部 4 は、X L M 処理部 3 の出力であるデータ $D 5 x ((L 1 - L 2) . m a s k)$ と、その正解データ $D _ c o r r e c t (L 1 - L 2)$ とから損失を取得する。

【0130】

T 2 : T L M の処理 $((L 2 - L 1) . m a s k \quad L 2 - L 1)$

次に、疑似対訳データ生成装置 1 0 0 のデータ入力インターフェース I F 1 は、対訳データ記憶部 $D B p (L 1 \quad L 2)$ から他分野の対訳データ $D 0 (L 1 - L 2)$ を読み出し、読み出した対訳データの第一言語のデータと第二言語のデータを入れ替えた対訳データをデータ $D 1 (L 2 - L 1)$ として、入力データ取得部 1 に出力する。

20

【0131】

入力データ取得部 1 の第 1 入力データ取得処理部 1 1 は、データ $D 1$ (対訳データ $D 1 (L 2 - L 1)$ (他分野第 2 言語の文と、当該文の第 1 言語の対訳文をペアとして構成する対訳データ $D 1 (L 2 - L 1)$)) から、(1) トークンデータ $x_{i t o k e n} (= x_{i t o k e n} (L 2 - L 1))$ と、(2) 当該トークンの位置を特定するための位置データ $x_{i p o s} (= x_{i p o s} (L 2 - L 1))$ と、(3) 当該トークンの言語を特定するための言語データ $x_{i l a n g} (= x_{i l a n g} (L 2 - L 1))$ と、を取得する。

30

【0132】

そして、入力データ取得部 1 は、上記のように取得したトークンデータ $x_{i t o k e n} (L 2 - L 1)$ をトークンデータ $x_{i 0 t o k e n} (L 2 - L 1)$ としてマスク化処理部 1 2 に出力する。

【0133】

制御部 (不図示) は、入力データ取得部 1 のマスク化処理部 1 2 に対して、マスク化処理を実行することを指示する制御信号 $C T L 2$ を生成し、当該制御信号 $C T L 2$ をマスク化処理部 1 2 に出力する。

【0134】

マスク化処理部 1 2 は、制御信号 $C T L 2$ に従い、第 1 入力データ取得処理部 1 1 から入力したトークンデータ $x_{i 0 t o k e n} (L 2 - L 1)$ の一部のトークンをマスクトークン (例えば、文字列を人工的なトークン「[M A S K]」としたデータ (トークン) に設定する) に置換するマスク化処理を行う。そして、マスク化処理部 1 2 は、当該マスク化処理後のトークンデータをトークンデータ $x_{i t o k e n} ((L 2 - L 1) . m a s k)$ として、第 1 セレクタ $S E L 1$ に出力する。なお、対訳データ $L 2 - L 1$ (他分野第 2 言語の文と、当該文の第 1 言語の対訳文をペアとして構成する対訳データ) に対してマスク化処理を行い取得されるデータを「 $(L 2 - L 1) . m a s k$ 」と表記する (以下、同様)。

40

【0135】

また、第 1 入力データ取得処理部 1 1 は、位置データ $x_{i p o s} (= x_{i p o s}) (= x_{i p o s} (L 2 - L 1))$ を第 1 セレクタ $S E L 1$ に出力し、言語データ $x_{i l a n g} (=$

50

x_{i1ang} ($= x_{i1ang}(L2 - L1)$) を第1セクタSEL1に出力する。

【0136】

つまり、入力データ取得部1は、上記により取得したデータをデータD2a ($= D2a((L2 - L1).mask) = \{x_{itoken}((L2 - L1).mask), x_{ipos}(L2 - L1), x_{i1ang}(L2 - L1)\}$) として第1セクタSEL1に出力する。

【0137】

正解データ取得部13は、制御信号CTL2に従い、データD1から、入力データ埋込部2およびXLM処理部3での学習処理(XLMの学習処理)に用いる正解データD__correct ($= D_correct(L2 - L1)$) を生成する。そして、正解データ取得部13は、生成したデータを正解データD__correct(L2 - L1)として第1出力データ評価部4に出力する。

10

【0138】

制御部は、第1セクタSEL1の端子「0」を選択する選択信号sel1を生成し、当該選択信号sel1を第1セクタSEL1に出力する。

【0139】

第1セクタSEL1は、選択信号sel1に従い、端子「0」を選択し、入力データ取得部1から出力されるデータD2a ($(L2 - L1).mask$) を、データD3 ($(L2 - L1).mask$) として、入力データ埋込部2に出力する。

【0140】

20

入力データ埋込部2のトークン埋込部21は、データD3 ($(L2 - L1).mask$) に含まれるトークンデータ x_{itoken} ($= x_{itoken}(L2 - L1)$) を入力し、入力したトークンデータ x_{itoken} の分散表現データ $x_{i'token}$ を、上記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0141】

入力データ埋込部2の位置埋込部22は、データD3 ($(L2 - L1).mask$) に含まれる位置データ x_{ipos} ($= x_{ipos}(L2 - L1)$) を入力し、入力した位置データ x_{ipos} の分散表現データ $x_{i'pos}$ を、上記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0142】

30

入力データ埋込部2の言語埋込部23は、データD3 ($(L2 - L1).mask$) に含まれる言語データ x_{i1ang} ($= x_{i1ang}(L2 - L1)$) を入力し、入力した言語データ x_{i1ang} の分散表現データ $x_{i'1ang}$ を、上記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0143】

入力データ埋込部2は、上記により取得された分散表現データをデータD4 ($(L2 - L1).mask$) として、XLM処理部3に出力する。

【0144】

XLM処理部3のXLM用ニューラルネットワークモデル31は、入力データ埋込部2から出力されるデータD4 ($(L2 - L1).mask$) に対して、XLM処理を実行し、データD5x ($= D5x((L2 - L1).mask)$) を取得する。そして、XLM処理部3は、上記により取得したデータD5x ($(L2 - L1).mask$) を第1出力データ評価部4に出力する。

40

【0145】

第1出力データ評価部4は、XLM処理部3から出力されるデータD5x ($(L2 - L1).mask$) と、入力データ取得部1から出力される正解データD__correct(L2 - L1) とを入力する。第1出力データ評価部4は、XLM処理部3の出力であるデータD5x ($(L2 - L1).mask$) と、その正解データD__correct(L2 - L1) とから損失を取得する。

【0146】

50

(1 . 2 . 1 . 1 C : 損失の計算処理)

上記の処理、すなわち、(A) M L M の処理 (4 種類の単言語データを使用) (ステップ S 1 1 1) と、(B) T L M の処理 (L 1 - L 2 の対訳データを使用) (ステップ S 1 1 2) とを実行した後、第 1 出力データ評価部 4 は、損失の計算処理を行う。なお、損失の計算を M 文 (M 個の文、M : 自然数) ごとに行う場合、M 文の中の第 i 番目の文に含まれるサブワード数を N_i ($1 \leq i \leq M$) とし、第 i 番目の文についての X L M 処理部 3 からの出力データ $D5x(X_{in})$ (入力データを X_{in} で表す) の j 番目 (j : 自然数、 $1 \leq j \leq N_i$) のサブワードに相当するデータを $D5x((X_{in}, X_{out}), i, j)$ (「 (X_{in}, X_{out}) 」 は、入力データが X_{in} であり、出力データが X_{out} であることを表す) とすると、第 1 出力データ評価部 4 は、下記数式のように、X L M 処理部 3 から出力されるデータと、正解データとから損失 $Loss$ を取得する。

10

【数 1】

$$Loss = Loss_MLM + Loss_TLM$$

【数 2】

$$\begin{aligned} Loss_MLM = & \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(L1, i, j), D5x((L1.mask \rightarrow L1), i, j)) \\ & + \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(L2, i, j), D5x((L2.mask \rightarrow L2), i, j)) \\ & + \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(R1, i, j), D5x((R1.mask \rightarrow R1), i, j)) \\ & + \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(R2, i, j), D5x((R2.mask \rightarrow R2), i, j)) \end{aligned}$$

20

30

【数 3】

$$\begin{aligned} Loss_TLM = & \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(L1 - L2, i, j), D5x(((L1 - L2).mask \rightarrow (L1 - L2)), i, j)) \\ & + \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(L2 - L1, i, j), D5x(((L2 - L1).mask \rightarrow (L2 - L1)), i, j)) \end{aligned}$$

40

【数 4】

$$\text{loss}(p, q) = - \sum_{k=1}^V p(k) \log(q(k))$$

50

V : サブワード語彙のサイズ (各トークン (入力データ) 文字列を表すベクトルの次元数)

p : 確率分布 ($p(k)$ は、k 番目 (第 k 次元) の要素の確率を示す)

q : 確率分布 ($q(k)$ は、k 番目 (第 k 次元) の要素の確率を示す)

なお、上記数式において、 $D_correct'(x, i, j)$ は、第 1 出力データ評価部 4 により正解データ $D_correct(x)$ の第 i 番目の文の第 j 番目のサブワードから取得されるデータ (ベクトル) であり、サブワード語彙のサイズ (= 各トークン (入力データ) 文字列を表すベクトルの次元数 (これを n_1 とする)) と同じ次元数のベクトル (n_1 次元のベクトル) である。そして、例えば、 $D_correct'(x, i, j)$ は、 n_1 次元のうち 1 次元のみ値が「1」であり (n_1 次元ベクトルの要素のうち、当該正解サブワードに対応する要素のみが「1」)、それ以外は値が「0」である one-hot ベクトルである。

10

【0147】

また、上記数式において、例えば、 $D_{5 \times (L_1 \text{ mask } L_1), i, j}$ は、 n_1 次元の実数ベクトル (各次元のデータ (ベクトルの要素) が、対応するサブワードである確率を示す実数ベクトル) であり、softmax 関数により確率化 (実数ベクトルの各要素の総和が「1」となるように正規化) されている。そして、 $D_{5 \times (L_1 \text{ mask } L_1), i, j}$ は、入力データを $L_1 \text{ mask}$ としたときの XLM 処理部 3 からの出力データであり、第 i 番目の文の第 j 番目のサブワードに相当するデータである。

【0148】

20

また、上記数式において、 $loss(p, q)$ は、交差エントロピー誤差を求める数式であり、これにより、確率分布間 (上記数式では、確率分布 p、q) の相違を定量化できる。

【0149】

このように、第 1 出力データ評価部 4 は、上記数式に相当する処理により、損失 $loss$ を取得する。

【0150】

(1.2.1.D : パラメータの更新処理)

第 1 出力データ評価部 4 は、所定の学習データに対して上記で算出した損失 (学習損失) に基づいて XLM 処理部 3 の XLM 用ニューラルネットワークモデル 31 のパラメータ x_{LM} を更新するためのデータ $update(x_{LM})$ を生成し、当該データ $update(x_{LM})$ を XLM 処理部 3 に出力する。

30

【0151】

XLM 処理部 3 の XLM 用ニューラルネットワークモデル 31 は、データ $update(x_{LM})$ に基づいて、パラメータ x_{LM} を更新する。

【0152】

また、XLM 用ニューラルネットワークモデル 31 は、XLM 用ニューラルネットワークモデル 31 のパラメータを更新した後、埋込層 (入力データ埋込部 2 のトークン埋込部 21、位置埋込部 22、および、言語埋込部 23 に対応) のパラメータを更新するためのパラメータ更新データ $update(emb)$ を生成し、当該パラメータ更新データ $update(emb)$ を入力データ埋込部 2 に出力する。

40

【0153】

入力データ埋込部 2 は、パラメータ更新データ $update(emb)$ に基づいて、埋込層 (入力データ埋込部 2 のトークン埋込部 21、位置埋込部 22、および、言語埋込部 23 に対応) のパラメータを更新する。例えば、入力データ埋込部 2 は、パラメータ更新データ $update(emb)$ に基づいて、変換行列 (例えば、行列 W_{token} 、 W_{pos} 、 W_{lang}) の要素 (値) を更新する。

【0154】

疑似対訳データ生成装置 100 は、所定の終了条件を満たすまで、上記処理 (図 6 のループ 1 の処理) を繰り返し実行する。その際、例えば、学習に用いるデータとは別の調整

50

用データに対して第1出力データ評価部4で算出される損失(テスト損失)を評価値として参照する。

【0155】

そして、疑似対訳データ生成装置100は、上記処理(図6のループ1の処理)の終了条件が満たされたときをもって、XLMの学習処理を終了させる。これにより、入力データ埋込部2の初期化が完了する。

【0156】

つまり、XLMの学習処理が終了した時点において、入力データ埋込部2で設定されている分散表現データを取得するためのパラメータ(例えば、変換行列 W_{token} 、 W_{pos} 、 W_{lang})を、疑似対訳データ生成用NMTモデル(入力データ埋込部2(埋込層に相当)と機械翻訳処理部5のMT用ニューラルネットワークモデル51とにより実現されるモデル)の学習処理の初期値(初期パラメータ)に設定する。

10

【0157】

つまり、疑似対訳データ生成装置100において、XLMの学習処理が終了した時点におけるパラメータが設定されている状態(入力データ埋込部2の状態)を、疑似対訳データ生成用NMTモデルの学習処理の初期状態として、疑似対訳データ生成用NMTモデルの学習処理を開始させる。

【0158】

なお、図6のフローチャートのループ処理(ループ1)の終了条件は、例えば、以下のよう設定される。

20

(1) 事前に定めた反復回数だけループ処理(ループ1)が実行された。

(2) 言語横断言語モデル(XLM)の学習処理において、第1出力データ評価部4における評価値が一定以上(事前に定めた値以上)の変化を示さなかった。

(3) 言語横断言語モデル(XLM)の学習処理において、第1出力データ評価部4における評価値が事前に定めた値を下回った。

(4) 言語横断言語モデル(XLM)の学習処理において、第1出力データ評価部4における評価値が事前に定めた回数更新されなかった。

【0159】

上記の終了条件を満たす場合、疑似対訳データ生成装置100は、XLMの学習処理が完了した判断し、XLMの学習処理を終了させる。

30

【0160】

そして、疑似対訳データ生成装置100は、上記のようにして、パラメータが上記初期値に設定された状態で、以下に説明する、疑似対訳データ生成用NMTモデル(入力データ埋込部2(埋込層に相当)と機械翻訳処理部5のMT用ニューラルネットワークモデル51とにより実現されるモデル)の学習処理を行う。

【0161】

(1.2.1.2: 疑似対訳データ生成用NMTモデルの学習処理)

次に、疑似対訳データ生成装置100では、疑似対訳データ生成用NMTモデルの学習処理(ステップS120~S126)が実行される。

【0162】

40

疑似対訳データ生成用NMTモデルの学習処理は、図6に示すように、

(A) 自己符号化処理(ステップS121)と、

(B) ゼロショット折り返し機械翻訳処理(ステップS122)と、

(C) 教師データあり機械翻訳処理(ステップS123)と、

(D) 損失の計算処理(ステップS124)と、

(E) パラメータ(M_T , emb)の更新処理(ステップS125)と、

から構成される。

【0163】

なお、図6のフローチャートに示すように、ステップS121~S125の処理は、終了条件が満たされるまで繰り返し実行される(図6のループ処理(ループ2)S120~

50

S 1 2 6)。

【 0 1 6 4 】

(1 . 2 . 1 . 2 A : 自己符号化処理)

自己符号化処理において、疑似対訳データ生成装置 1 0 0 は、図 8 に示すように、入力データ埋込部 2 への入力データと同一の出力データが出力されるように学習処理 (疑似対訳データ生成用 N M T モデルの学習処理) を行う。つまり、自己符号化処理において、正解データは、入力データと同一のデータに設定される。すなわち、自己符号化処理において、疑似対訳データ生成装置 1 0 0 は、

(1) 入力データ埋込部 2 への入力を D 3 (L 1) とし、機械翻訳処理部 5 の出力データ D 5 (L 1) の正解データを D _ c o r r e c t (L 1) とする、

10

(2) 入力データ埋込部 2 への入力を D 3 (L 2) とし、機械翻訳処理部 5 の出力データ D 5 (L 2) の正解データを D _ c o r r e c t (L 2) とする、

(3) 入力データ埋込部 2 への入力を D 3 (R 1) とし、機械翻訳処理部 5 の出力データ D 5 (R 1) の正解データを D _ c o r r e c t (R 1) とする、および、

(4) 入力データ埋込部 2 への入力を D 3 (R 2) とし、機械翻訳処理部 5 の出力データ D 5 (R 2) の正解データを D _ c o r r e c t (R 2) とし、学習処理 (自己符号化処理によるパラメータ最適化処理) を実行する。

【 0 1 6 5 】

A 1 : 自己符号化処理 (L 1 L 1)

入力データ埋込部 2 への入力を D 3 (L 1) とし、機械翻訳処理部 5 の出力データ D 5 (L 1) の正解データを D _ c o r r e c t (L 1) とする場合について、説明する。

20

【 0 1 6 6 】

データ入力インターフェース I F 1 は、第 1 単言語データ記憶部 D B m (L 1) から他分野第 1 言語の単言語データ D 0 (L 1) を読み出し、読み出した単言語データをデータ D 1 (= D 1 (L 1)) として、入力データ取得部 1 に出力する。

【 0 1 6 7 】

入力データ取得部 1 の第 1 入力データ取得処理部 1 1 は、データ D 1 (他分野第 1 言語のデータ D 1 (L 1)) (他分野第 1 言語の文を構成するデータ D 1 (L 1)) から、(1) トークンデータ x i 0 t o k e n と、(2) 当該トークンの位置を特定するための位置データ x i p o s と、(3) 当該トークンの言語を特定するための言語データ x i l a n g と、を取得する。

30

【 0 1 6 8 】

そして、入力データ取得部 1 は、上記のように取得したトークンデータ x i 0 t o k e n をマスク化処理部 1 2 に出力する。

【 0 1 6 9 】

制御部 (不図示) は、入力データ取得部 1 のマスク化処理部 1 2 に対して、マスク化処理を実行しないことを指示する制御信号 C T L 2 を生成し、当該制御信号 C T L 2 をマスク化処理部 1 2 に出力する。

【 0 1 7 0 】

マスク化処理部 1 2 は、制御信号 C T L 2 に従い、第 1 入力データ取得処理部 1 1 から入力したトークンデータ x i 0 t o k e n をトークンデータ x i t o k e n として、第 1 セレクタ S E L 1 に出力する。

40

【 0 1 7 1 】

また、第 1 入力データ取得処理部 1 1 は、位置データ x i p o s (= x i p o s (L 1)) を第 1 セレクタ S E L 1 に出力し、言語データ x i l a n g (= x i l a n g) (= x i l a n g (L 1)) を第 1 セレクタ S E L 1 に出力する。

【 0 1 7 2 】

つまり、入力データ取得部 1 は、上記により取得したデータをデータ D 2 a (= D 2 a (L 1) = { x i t o k e n (L 1) , x i p o s (L 1) , x i l a n g (L 1) }) として第 1 セレクタ S E L 1 に出力する。

50

【0173】

正解データ取得部13は、制御信号CTL2に従い、データD1(=D1(L1))から、疑似対訳データ生成用NMTモデル(入力データ埋込部2(埋込層に相当)と機械翻訳処理部5のMT用ニューラルネットワークモデル51とにより実現されるモデル)の学習処理に用いる正解データD__correct(=D__correct(L1))を生成する。具体的には、正解データ取得部13は、第1セクタSEL1への入力データD2a(L1)と同一のデータを正解データD__correct(L1)に設定する。そして、正解データ取得部13は、上記により生成した正解データD__correct(L1)を第2出力データ評価部6に出力する。

【0174】

制御部は、第1セクタSEL1の端子「0」を選択する選択信号sel1を生成し、当該選択信号sel1を第1セクタSEL1に出力する。

【0175】

第1セクタSEL1は、選択信号sel1に従い、端子「0」を選択し、入力データ取得部1から出力されるデータD2a(L1)を、データD3(L1)として、入力データ埋込部2に出力する。

【0176】

入力データ埋込部2のトークン埋込部21は、データD3(L1)に含まれるトークンデータx_i'tokenを入力し、入力したトークンデータx_i'tokenの分散表現データx_i'tokenを、上記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0177】

入力データ埋込部2の位置埋込部22は、データD3(L1)に含まれる位置データx_i'posを入力し、入力した位置データx_i'posの分散表現データx_i'posを、上記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0178】

入力データ埋込部2の言語埋込部23は、データD3(L1)に含まれる言語データx_i'langを入力し、入力した言語データx_i'langの分散表現データx_i'langを、上記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0179】

入力データ埋込部2は、上記により取得された分散表現データをデータD4(L1)として、機械翻訳処理部5に出力する。

【0180】

機械翻訳処理部5のMT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4(L1)(分散表現データD4(L1)(={x_i'token, x_i'pos, x_i'lang}))と、制御部から出力される制御信号CTL3とを入力する。

【0181】

制御部は、MT用ニューラルネットワークモデル51からL1のデータを出力することを指示する制御信号CTL3を生成し、当該制御信号CTL3を機械翻訳処理部5に出力する。

【0182】

MT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4(L1)に対して、機械翻訳処理を実行し、データD5(L1)を取得する。なお、MT用ニューラルネットワークモデル51は、制御信号CTL3に従い、出力データの種別をL1(他分野、第1言語)とする。

【0183】

MT用ニューラルネットワークモデル51は、上記により取得したデータD5(L1)を第2セクタSEL2に出力する。

【0184】

制御部は、第2セクタSEL2の端子「0」を選択する選択信号sel2を生成し、

10

20

30

40

50

当該選択信号 $s e l 2$ を第 2 セレクタ $S E L 2$ に出力する。

【 0 1 8 5 】

第 2 セレクタ $S E L 2$ は、選択信号 $s e l 2$ に従い、端子「0」を選択し、機械翻訳処理部 5 から出力されるデータ $D 5$ (データ $D 5 (L 1)$) を、第 2 出力データ評価部 6 にデータ $D 6 a$ (データ $D 6 a (L 1)$) として出力する。

【 0 1 8 6 】

第 2 出力データ評価部 6 は、第 2 セレクタ $S E L 2$ を介して機械翻訳処理部 5 から出力されるデータ $D 5 (L 1)$ (=データ $D 6 a (L 1)$) と、入力データ取得部 1 から出力される正解データ $D_c o r r e c t (L 1)$ とを入力する。第 2 出力データ評価部 6 は、機械翻訳処理部 5 の出力であるデータ $D 5 (L 1)$ と、その正解データ $D_c o r r e c t (L 1)$ とから損失を取得する。

10

【 0 1 8 7 】

A 2 : 自己符号化処理 (L 2 L 2)

入力データ埋込部 2 への入力を $D 3 (L 2)$ とし、機械翻訳処理部 5 の出力データ $D 5 (L 2)$ の正解データを $D_c o r r e c t (L 2)$ とする場合についても、疑似対訳データ生成装置 100 は、上記と同様の処理を実行する。つまり、疑似対訳データ生成装置 100 は、上記処理において、 $D 3 (L 1)$ を $D 3 (L 2)$ に置換し、正解データ $D_c o r r e c t (L 1)$ を正解データ $D_c o r r e c t (L 2)$ に置換して、上記処理と同様の処理を行う。

【 0 1 8 8 】

20

A 3 : 自己符号化処理 (R 1 R 1)

入力データ埋込部 2 への入力を $D 3 (R 1)$ とし、機械翻訳処理部 5 の出力データ $D 5 (R 1)$ の正解データを $D_c o r r e c t (R 1)$ とする場合についても、疑似対訳データ生成装置 100 は、上記と同様の処理を実行する。つまり、疑似対訳データ生成装置 100 は、上記処理において、 $D 3 (L 1)$ を $D 3 (R 1)$ に置換し、正解データ $D_c o r r e c t (L 1)$ を正解データ $D_c o r r e c t (R 1)$ に置換して、上記処理と同様の処理を行う。

【 0 1 8 9 】

A 4 : 自己符号化処理 (R 2 R 2)

入力データ埋込部 2 への入力を $D 3 (R 2)$ とし、機械翻訳処理部 5 の出力データ $D 5 (R 2)$ の正解データを $D_c o r r e c t (R 2)$ とする場合についても、疑似対訳データ生成装置 100 は、上記と同様の処理を実行する。つまり、疑似対訳データ生成装置 100 は、上記処理において、 $D 3 (L 1)$ を $D 3 (R 2)$ に置換し、正解データ $D_c o r r e c t (L 1)$ を正解データ $D_c o r r e c t (R 2)$ に置換して、上記処理と同様の処理を行う。

30

【 0 1 9 0 】

(1 . 2 . 1 . 2 B : ゼロショット折り返し機械翻訳処理)

ゼロショット折り返し機械翻訳処理において、疑似対訳データ生成装置 100 は、図 9、図 10 に示すように、(1) 入力データ埋込部 2 への入力データに対して機械翻訳処理 (疑似対訳データ生成用 NMT モデル (入力データ埋込部 2 (埋込層に相当) と機械翻訳処理部 5 の MT 用ニューラルネットワークモデル 5 1 とにより実現されるモデル) による機械翻訳処理) を行い (1 回目の機械翻訳処理) 、入力データと異なる言語のデータ (ゼロショット機械翻訳のデータ) を出力させ、(2) その出力されたデータに対して機械翻訳処理 (疑似対訳データ生成用 NMT モデルにより機械翻訳処理) を行い (2 回目の機械翻訳処理) 、入力データと同一のデータが出力されるように学習処理を行う。つまり、ゼロショット折り返し機械翻訳処理において、正解データは、入力データと同一のデータに設定される。

40

【 0 1 9 1 】

具体的には、ゼロショット折り返し機械翻訳処理において、疑似対訳データ生成装置 100 は、以下の (1) ~ (6) の処理を行う。

50

(1) L 1 R 2 L 1 :

1 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(1)}(L 1)$ とし、機械翻訳処理部 5 からの出力を $D 5^{(1)}(R 2)$ とし、2 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(2)}(R 2) (= D 5^{(1)}(R 2))$ とし、機械翻訳処理部 5 からの出力を $D 5^{(2)}(L 1)$ とし、また、正解データを $D_correct(L 1)$ とする。

【 0 1 9 2 】

なお、1 回目の機械翻訳処理の入力データ埋込部 2 への入力データを $D 3^{(1)}(x)$ と、入力データ埋込部 2 からの出力データを $D 4^{(1)}(x)$ と、2 回目の機械翻訳処理の入力データ埋込部 2 への入力データを $D 3^{(2)}(x)$ と、入力データ埋込部 2 からの出力データを $D 4^{(2)}(x)$ と表記する(以下同様)。

10

(2) R 1 L 2 R 1 :

1 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(1)}(R 1)$ とし、機械翻訳処理部 5 からの出力を $D 5^{(1)}(L 2)$ とし、2 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(2)}(L 2) (= D 5^{(1)}(L 2))$ とし、機械翻訳処理部 5 からの出力を $D 5^{(2)}(R 1)$ とし、また、正解データを $D_correct(R 1)$ とする。

(3) R 1 R 2 R 1 :

1 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(1)}(R 1)$ とし、機械翻訳処理部 5 からの出力を $D 5^{(1)}(R 2)$ とし、2 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(2)}(R 2) (= D 5^{(1)}(R 2))$ とし、機械翻訳処理部 5 からの出力を $D 5^{(2)}(R 1)$ とし、また、正解データを $D_correct(R 1)$ とする。

20

(4) L 2 R 1 L 2 :

1 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(1)}(L 2)$ とし、機械翻訳処理部 5 からの出力を $D 5^{(1)}(R 1)$ とし、2 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(2)}(R 1) (= D 5^{(1)}(R 1))$ とし、機械翻訳処理部 5 からの出力を $D 5^{(2)}(L 2)$ とし、また、正解データを $D_correct(L 2)$ とする。

(5) R 2 L 1 R 2 :

1 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(1)}(R 2)$ とし、機械翻訳処理部 5 からの出力を $D 5^{(1)}(L 1)$ とし、2 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(2)}(L 1) (= D 5^{(1)}(L 1))$ とし、機械翻訳処理部 5 からの出力を $D 5^{(2)}(R 2)$ とし、また、正解データを $D_correct(R 2)$ とする。

30

(6) R 2 R 1 R 2 :

1 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(1)}(R 2)$ とし、機械翻訳処理部 5 からの出力を $D 5^{(1)}(R 1)$ とし、2 回目の機械翻訳処理において、入力データ埋込部 2 への入力を $D 3^{(2)}(R 1) (= D 5^{(1)}(R 1))$ とし、機械翻訳処理部 5 からの出力を $D 5^{(2)}(R 2)$ とし、また、正解データを $D_correct(R 1)$ とする。

40

【 0 1 9 3 】

B 1 : ゼロショット折り返し機械翻訳処理 (L 1 R 2 L 1)

入力データ埋込部 2 への 1 回目の入力を $D 3^{(1)}(L 1)$ とし、1 回目の機械翻訳処理部 5 からの出力を $D 5^{(1)}(R 2)$ とし、入力データ埋込部 2 への 2 回目の入力を $D 3^{(2)}(R 2) (= D 5^{(1)}(R 2))$ とし、2 回目の機械翻訳処理部 5 からの出力を $D 5^{(2)}(L 1)$ とし、正解データを $D_correct(L 1)$ とする場合について、説明する。

【 0 1 9 4 】

データ入力インターフェース I F 1 は、第 1 単言語データ記憶部 $D B m(L 1)$ から他

50

分野の第1言語の単言語データD0(L1)を読み出し、読み出した単言語データをデータD1(=D1(L1))として、入力データ取得部1に出力する。

【0195】

入力データ取得部1の第1入力データ取得処理部11は、データD1(他分野第1言語のデータD1(L1))(他分野第1言語の文を構成するデータD1(L1))から、(1)トークンデータxi0tokenと、(2)当該トークンの位置を特定するための位置データxi_posと、(3)当該トークンの言語を特定するための言語データxi_langと、を取得する。

【0196】

そして、入力データ取得部1は、上記のように取得したトークンデータxi0tokenをマスク化処理部12に出力する。

10

【0197】

制御部(不図示)は、入力データ取得部1のマスク化処理部12に対して、マスク化処理を実行しないことを指示する制御信号CTL2を生成し、当該制御信号CTL2をマスク化処理部12に出力する。

【0198】

マスク化処理部12は、制御信号CTL2に従い、第1入力データ取得処理部11から入力したトークンデータxi0tokenをトークンデータxitokenとして、第1セクタSEL1に出力する。

【0199】

また、第1入力データ取得処理部11は、位置データxi_pos(=xi_pos(L1))を第1セクタSEL1に出力し、言語データxi_lang(=xi_lang(L1))を第1セクタSEL1に出力する。

20

【0200】

つまり、入力データ取得部1は、上記により取得したデータをデータD2a(=D2a(L1)={xitoken(L1),xi_pos(L1),xi_lang(L1)})として第1セクタSEL1に出力する。

【0201】

正解データ取得部13は、制御信号CTL2に従い、データD1(=D1(L1))から、疑似対訳データ生成用NMTモデル(入力データ埋込部2(埋込層に相当)と機械翻訳処理部5のMT用ニューラルネットワークモデル5.1とにより実現されるモデル)の学習処理に用いる正解データD__correct(=D__correct(L1))を生成する。具体的には、正解データ取得部13は、第1セクタSEL1への入力データD2a(L1)と同一のデータを正解データD__correct(L1)に設定する。そして、正解データ取得部13は、上記により生成した正解データD__correct(L1)を第2出力データ評価部6に出力する。

30

【0202】

制御部は、第1セクタSEL1の端子「0」を選択する選択信号sel1を生成し、当該選択信号sel1を第1セクタSEL1に出力する。

【0203】

第1セクタSEL1は、選択信号sel1に従い、端子「0」を選択し、入力データ取得部1から出力されるデータD2a(L1)を、データD3⁽¹⁾(L1)として、入力データ埋込部2に出力する。

40

【0204】

入力データ埋込部2のトークン埋込部21は、データD3⁽¹⁾(L1)に含まれるトークンデータxitokenを入力し、入力したトークンデータxitokenの分散表現データxi'tokenを、上記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0205】

入力データ埋込部2の位置埋込部22は、データD3⁽¹⁾(L1)に含まれる位置データxi_posを入力し、入力した位置データxi_posの分散表現データxi'posを、上

50

記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0206】

入力データ埋込部2の言語埋込部23は、データD3⁽¹⁾(L1)に含まれる言語データ*x i l a n g*を入力し、入力した言語データ*x i l a n g*の分散表現データ*x i l a n g*を、上記のMLMの処理で説明したのと同じ処理を実行して取得する。

【0207】

入力データ埋込部2は、上記により取得された分散表現データをデータD4⁽¹⁾(L1)として、機械翻訳処理部5に出力する。

【0208】

機械翻訳処理部5のMT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4⁽¹⁾(L1)(分散表現データD4⁽¹⁾(L1)(=*{ x i t o k e n , x i p o s , x i l a n g }*))と、制御部から出力される制御信号CTL3とを入力する。

10

【0209】

制御部は、MT用ニューラルネットワークモデル51からR2のデータを出力することを指示する制御信号CTL3を生成し、当該制御信号CTL3を機械翻訳処理部5に出力する。

【0210】

MT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4⁽¹⁾(L1)に対して、機械翻訳処理を実行し、データD5⁽¹⁾(R2)を取得する。なお、MT用ニューラルネットワークモデル51は、制御信号CTL3に従い、出力データの種別をR2(適応先分野、第2言語)とする。

20

【0211】

MT用ニューラルネットワークモデル51は、上記により取得したデータD5⁽¹⁾(R2)を第2セクタSEL2に出力する。

【0212】

制御部は、第2セクタSEL2の端子「1」を選択する選択信号sel2を生成し、当該選択信号sel2を第2セクタSEL2に出力する。

【0213】

第2セクタSEL2は、選択信号sel2に従い、端子「1」を選択し、機械翻訳処理部5から出力されるデータD5(データD5⁽¹⁾(R2))を、第1バッファBuf1にデータD6b(データD6b⁽¹⁾(R2))として記憶する(図9を参照)。

30

【0214】

次に、第2入力データ取得処理部7は、第1バッファBuf1から出力されるデータD6b'(データD6b⁽¹⁾(R2))を入力し、データD6b'(データD6b⁽¹⁾(R2))から入力データ埋込部2に入力するためのデータ(=*{ x i t o k e n (R 2) , x i p o s (R 2) , x i l a n g (R 2) }*)を取得(生成)する。そして、第2入力データ取得処理部7は、取得したデータをデータD2b(データD2b⁽²⁾(R2)(=*{ x i t o k e n (R 2) , x i p o s (R 2) , x i l a n g (R 2) }*))として、第1セクタSEL1に出力する。

40

【0215】

制御部は、第1セクタSEL1の端子「1」を選択する選択信号sel1を生成し、当該選択信号sel1を第1セクタSEL1に出力する。

【0216】

第1セクタSEL1は、選択信号sel1に従い、端子「1」を選択し、第2入力データ取得処理部7から出力されるデータD2b⁽²⁾(R2)を、データD3⁽²⁾(R2)として、入力データ埋込部2に出力する。

【0217】

入力データ埋込部2のトークン埋込部21は、データD3⁽²⁾(R2)に含まれるトークンデータ*x i t o k e n*を入力し、入力したトークンデータ*x i t o k e n*の分散表現デ

50

ータ $x_i \text{ 't o k e n}$ を、上記の M L M の処理で説明したのと同じ処理を実行して取得する。

【0218】

入力データ埋込部 2 の位置埋込部 2 2 は、データ $D 3^{(2)}(R 2)$ に含まれる位置データ $x_i \text{ 'p o s}$ を入力し、入力した位置データ $x_i \text{ 'p o s}$ の分散表現データ $x_i \text{ 'p o s}$ を、上記の M L M の処理で説明したのと同じ処理を実行して取得する。

【0219】

入力データ埋込部 2 の言語埋込部 2 3 は、データ $D 3^{(2)}(R 2)$ に含まれる言語データ $x_i \text{ 'l a n g}$ を入力し、入力した言語データ $x_i \text{ 'l a n g}$ の分散表現データ $x_i \text{ 'l a n g}$ を、上記の M L M の処理で説明したのと同じ処理を実行して取得する。

【0220】

入力データ埋込部 2 は、上記により取得された分散表現データをデータ $D 4^{(2)}(R 2)$ として、機械翻訳処理部 5 に出力する。

【0221】

機械翻訳処理部 5 の M T 用ニューラルネットワークモデル 5 1 は、入力データ埋込部 2 から出力されるデータ $D 4^{(2)}(R 2)$ (分散表現データ $D 4^{(2)}(R 2)$ ($= \{ x_i \text{ 't o k e n}, x_i \text{ 'p o s}, x_i \text{ 'l a n g} \}$)) と、制御部から出力される制御信号 C T L 3 とを入力する。

【0222】

制御部は、M T 用ニューラルネットワークモデル 5 1 から L 1 のデータを出力することを指示する制御信号 C T L 3 を生成し、当該制御信号 C T L 3 を機械翻訳処理部 5 に出力する。

【0223】

M T 用ニューラルネットワークモデル 5 1 は、入力データ埋込部 2 から出力されるデータ $D 4^{(2)}(R 2)$ に対して、機械翻訳処理を実行し、データ $D 5^{(2)}(L 1)$ を取得する。なお、M T 用ニューラルネットワークモデル 5 1 は、制御信号 C T L 3 に従い、出力データの種別を L 1 (他分野、第 1 言語) とする。

【0224】

M T 用ニューラルネットワークモデル 5 1 は、上記により取得したデータ $D 5^{(2)}(L 1)$ を第 2 セレクタ S E L 2 に出力する。

【0225】

制御部は、第 2 セレクタ S E L 2 の端子「0」を選択する選択信号 $s e l 2$ を生成し、当該選択信号 $s e l 2$ を第 2 セレクタ S E L 2 に出力する。

【0226】

第 2 セレクタ S E L 2 は、選択信号 $s e l 2$ に従い、端子「0」を選択し、機械翻訳処理部 5 から出力されるデータ $D 5^{(2)}(L 1)$ を、データ $D 6 a^{(2)}(L 1)$ として、第 2 出力データ評価部 6 に出力する。

【0227】

第 2 出力データ評価部 6 は、第 2 セレクタ S E L 2 から出力されるデータ $D 6 a^{(2)}(L 1)$ ($= D 5^{(2)}(L 1)$) と、入力データ取得部 1 から出力される正解データ $D _c o r r e c t(L 1)$ とを入力する。第 2 出力データ評価部 6 は、機械翻訳処理部 5 の出力であるデータ $D 5^{(2)}(L 1)$ と、その正解データ $D _c o r r e c t(L 1)$ とから損失を取得する。

【0228】

B 2 : ゼロショット折り返し機械翻訳処理 (R 1 L 2 R 1)

入力データ埋込部 2 への 1 回目の入力を $D 3^{(1)}(R 1)$ とし、1 回目の機械翻訳処理部 5 からの出力を $D 5^{(1)}(L 2)$ とし、入力データ埋込部 2 への 2 回目の入力を $D 3^{(2)}(L 2)$ ($= D 5^{(1)}(L 2)$) とし、2 回目の機械翻訳処理部 5 からの出力を $D 5^{(2)}(R 1)$ とし、正解データを $D _c o r r e c t(R 1)$ とする場合についても、疑似対訳データ生成装置 1 0 0 は、上記と同様の処理を実行する。

【0229】

10

20

30

40

50

つまり、疑似対訳データ生成装置100は、上記処理において、入力データ埋込部2への1回目の入力D3⁽¹⁾(L1)をD3⁽¹⁾(R1)に、機械翻訳処理部5からの出力D5⁽¹⁾(R2)をD5⁽¹⁾(L2)に置換し、入力データ埋込部2への2回目の入力D3⁽²⁾(R2)をD3⁽²⁾(L2)に、機械翻訳処理部5からの出力D5⁽²⁾(L1)をD5⁽²⁾(R1)に置換し、正解データD__correct(L1)を正解データD__correct(R1)に置換して、上記処理と同様の処理を行う。

【0230】

B3：ゼロショット折り返し機械翻訳処理(R1 R2 R1)

入力データ埋込部2への1回目の入力をD3⁽¹⁾(R1)とし、1回目の機械翻訳処理部5からの出力をD5⁽¹⁾(R2)とし、入力データ埋込部2への2回目の入力をD3⁽²⁾(R2)(=D5⁽¹⁾(R2))とし、2回目の機械翻訳処理部5からの出力をD5⁽²⁾(R1)とし、正解データをD__correct(R1)とする場合についても、疑似対訳データ生成装置100は、上記と同様の処理を実行する。

10

【0231】

つまり、疑似対訳データ生成装置100は、上記処理において、入力データ埋込部2への1回目の入力D3⁽¹⁾(L1)をD3⁽¹⁾(R1)に、機械翻訳処理部5からの出力D5⁽¹⁾(R2)を同じくD5⁽¹⁾(R2)とし、入力データ埋込部2への2回目の入力D3⁽²⁾(R2)を同じくD3⁽²⁾(R2)とし、機械翻訳処理部5からの出力D5⁽²⁾(L1)をD5⁽²⁾(R1)に置換し、正解データD__correct(L1)を正解データD__correct(R1)に置換して、上記処理と同様の処理を行う。

20

【0232】

B4：ゼロショット折り返し機械翻訳処理(L2 R1 L2)

入力データ埋込部2への1回目の入力をD3⁽¹⁾(L2)とし、1回目の機械翻訳処理部5からの出力をD5⁽¹⁾(R1)とし、入力データ埋込部2への2回目の入力をD3⁽²⁾(R1)(=D5⁽¹⁾(R1))とし、2回目の機械翻訳処理部5からの出力をD5⁽²⁾(L2)とし、正解データをD__correct(L2)とする場合についても、疑似対訳データ生成装置100は、上記と同様の処理を実行する。

【0233】

つまり、疑似対訳データ生成装置100は、上記処理において、入力データ埋込部2への1回目の入力D3⁽¹⁾(L1)をD3⁽¹⁾(L2)に、機械翻訳処理部5からの出力D5⁽¹⁾(R2)をD5⁽¹⁾(R1)に置換し、入力データ埋込部2への2回目の入力D3⁽²⁾(R2)をD3⁽²⁾(R1)に、機械翻訳処理部5からの出力D5⁽²⁾(L1)をD5⁽²⁾(L2)に置換し、正解データD__correct(L1)を正解データD__correct(L2)に置換して、上記処理と同様の処理を行う。

30

【0234】

B5：ゼロショット折り返し機械翻訳処理(R2 L1 R2)

入力データ埋込部2への1回目の入力をD3⁽¹⁾(R2)とし、1回目の機械翻訳処理部5からの出力をD5⁽¹⁾(L1)とし、入力データ埋込部2への2回目の入力をD3⁽²⁾(L1)(=D5⁽¹⁾(L1))とし、2回目の機械翻訳処理部5からの出力をD5⁽²⁾(R2)とし、正解データをD__correct(R2)とする場合についても、疑似対訳データ生成装置100は、上記と同様の処理を実行する。

40

【0235】

つまり、疑似対訳データ生成装置100は、上記処理において、入力データ埋込部2への1回目の入力D3⁽¹⁾(L1)をD3⁽¹⁾(R2)に、機械翻訳処理部5からの出力D5⁽¹⁾(R2)をD5⁽¹⁾(L1)に置換し、入力データ埋込部2への2回目の入力D3⁽²⁾(R2)をD3⁽²⁾(L1)に、機械翻訳処理部5からの出力D5⁽²⁾(L1)をD5⁽²⁾(R2)に置換し、正解データD__correct(L1)を正解データD__correct(R2)に置換して、上記処理と同様の処理を行う。

【0236】

B6：ゼロショット折り返し機械翻訳処理(R2 R1 R2)

50

入力データ埋込部 2 への 1 回目の入力を $D3^{(1)}(R2)$ とし、1 回目の機械翻訳処理部 5 からの出力を $D5^{(1)}(R1)$ とし、入力データ埋込部 2 への 2 回目の入力を $D3^{(2)}(R1)$ ($=D5^{(1)}(R1)$) とし、2 回目の機械翻訳処理部 5 からの出力を $D5^{(2)}(R2)$ とし、正解データを $D_correct(R2)$ とする場合についても、疑似対訳データ生成装置 100 は、上記と同様の処理を実行する。

【0237】

つまり、疑似対訳データ生成装置 100 は、上記処理において、入力データ埋込部 2 への 1 回目の入力 $D3^{(1)}(L1)$ を $D3^{(1)}(R2)$ に、機械翻訳処理部 5 からの出力 $D5^{(1)}(R2)$ を $D5^{(1)}(R1)$ に置換し、入力データ埋込部 2 への 2 回目の入力 $D3^{(2)}(R2)$ を $D3^{(2)}(R1)$ に、機械翻訳処理部 5 からの出力 $D5^{(2)}(L1)$ を $D5^{(2)}(R2)$ に置換し、正解データ $D_correct(L1)$ を正解データ $D_correct(R2)$ に置換して、上記処理と同様の処理を行う。

10

【0238】

以上のようにして、疑似対訳データ生成装置 100 では、

- (1) L1 R2 L1
- (2) R1 L2 R1
- (3) R1 R2 R1
- (4) L2 R1 L2
- (5) R2 L1 R2
- (6) R2 R1 R2

20

の 6 パターン (6 種類の場合) について、ゼロショット折り返し機械翻訳処理が実行される。

【0239】

(1.2.1.2C: 教師データあり機械翻訳処理)

教師データあり機械翻訳処理において、疑似対訳データ生成装置 100 は、図 11 に示すように、入力データ埋込部 2 への入力データを対訳データ記憶部 DBp (L1 - L2) から取得した対訳データ $D0(L1 - L2)$ の一方の言語のデータとし、当該入力データに対応する対訳データが出力されるように学習処理を行う。つまり、教師あり機械翻訳処理において、正解データは、対訳データ記憶部 DBp (L1 - L2) から読み出した対訳データ $D0(L1 - L2)$ に基づいて、設定される。

30

【0240】

なお、教師あり機械翻訳処理において、疑似対訳データ生成装置 100 は、

- (1) 入力データ埋込部 2 への入力を $D3(L1)$ とし、正解データを $D_correct(L2)$ とする、あるいは、
- (2) 入力データ埋込部 2 への入力を $D3(L2)$ とし、正解データを $D_correct(L1)$ として、

疑似対訳データ生成用 NMT モデル (入力データ埋込部 2 (埋込層に相当) と機械翻訳処理部 5 の MT 用ニューラルネットワークモデル 51 とにより実現されるモデル) の学習処理 (パラメータの最適化処理) を行う。

【0241】

C1: 教師あり機械翻訳処理 (L1 L2)

入力データ埋込部 2 への入力を $D3(L1)$ とし、正解データを $D_correct(L2)$ とする場合について、説明する。

【0242】

データ入力インターフェース IF1 は、対訳データ記憶部 DBp (L1 - L2) から他分野の対訳データ $D0(L1 - L2)$ を読み出し、読み出した対訳データをデータ $D1(L1 - L2)$ として、入力データ取得部 1 に出力する。

【0243】

入力データ取得部 1 の第 1 入力データ取得処理部 11 は、データ $D1$ (他分野対訳データ $D1(L1 - L2)$) から第 1 言語のデータをデータ $D1(L1)$ として抽出し、当該

40

50

第1言語のデータ(言語データ) $D1(L1)$ (他分野第1言語の文を構成するデータ $D1(L1)$) から、(1) トークンデータ $x_{i0\ token}$ と、(2) 当該トークンの位置を特定するための位置データ $x_{i\ pos}$ と、(3) 当該トークンの言語を特定するための言語データ $x_{i\ lang}$ と、を取得する。

【0244】

そして、入力データ取得部1は、上記のように取得したトークンデータ $x_{i0\ token}$ をマスク化処理部12に出力する。

【0245】

制御部(不図示)は、入力データ取得部1のマスク化処理部12に対して、マスク化処理を実行しないことを指示する制御信号 $CTL2$ を生成し、当該制御信号 $CTL2$ をマスク化処理部12に出力する。

10

【0246】

マスク化処理部12は、制御信号 $CTL2$ に従い、第1入力データ取得処理部11から入力したトークンデータ $x_{i0\ token}$ をトークンデータ $x_{i\ token}$ として、第1セクタ $SEL1$ に出力する。

【0247】

また、第1入力データ取得処理部11は、位置データ $x_{i\ pos}$ ($=x_{i\ pos}(L1)$) を第1セクタ $SEL1$ に出力し、言語データ $x_{i\ lang}$ ($=x_{i\ lang}(L1)$) を第1セクタ $SEL1$ に出力する。

【0248】

つまり、入力データ取得部1は、上記により取得したデータをデータ $D2a$ ($=D2a(L1) = \{x_{i\ token}(L1), x_{i\ pos}(L1), x_{i\ lang}(L1)\}$) として第1セクタ $SEL1$ に出力する。

20

【0249】

正解データ取得部13は、制御信号 $CTL2$ に従い、データ $D1$ (他分野対訳データ $D1(L1-L2)$) から第2言語のデータを $D1(L2)$ (他分野対訳データ $D1(L1-L2)$ から抽出したデータ $D1(L1)$ (他分野第1言語のデータ) の対訳データに相当するデータ) として抽出し、抽出したデータ $D1(L2)$ から、疑似対訳データ生成用 NMT モデルの学習処理に用いる正解データ $D_correct$ ($=D_correct(L2)$) を生成する。そして、正解データ取得部13は、当該正解データ $D_correct$ ($=D_correct(L2)$) を第2出力データ評価部6に出力する。

30

【0250】

制御部は、第1セクタ $SEL1$ の端子「0」を選択する選択信号 $sel1$ を生成し、当該選択信号 $sel1$ を第1セクタ $SEL1$ に出力する。

【0251】

第1セクタ $SEL1$ は、選択信号 $sel1$ に従い、端子「0」を選択し、入力データ取得部1から出力されるデータ $D2a(L1)$ を、データ $D3(L1)$ として、入力データ埋込部2に出力する。

【0252】

入力データ埋込部2のトークン埋込部21は、データ $D3(L1)$ に含まれるトークンデータ $x_{i\ token}$ を入力し、入力したトークンデータ $x_{i\ token}$ の分散表現データ $x_{i\ 'token}$ を、上記の MLM の処理で説明したのと同じ処理を実行して取得する。

40

【0253】

入力データ埋込部2の位置埋込部22は、データ $D3(L1)$ に含まれる位置データ $x_{i\ pos}$ を入力し、入力した位置データ $x_{i\ pos}$ の分散表現データ $x_{i\ 'pos}$ を、上記の MLM の処理で説明したのと同じ処理を実行して取得する。

【0254】

入力データ埋込部2の言語埋込部23は、データ $D3(L1)$ に含まれる言語データ $x_{i\ lang}$ を入力し、入力した言語データ $x_{i\ lang}$ の分散表現データ $x_{i\ 'lang}$ を、上記の MLM の処理で説明したのと同じ処理を実行して取得する。

50

【0255】

入力データ埋込部2は、上記により取得された分散表現データをデータD4(L1)として、機械翻訳処理部5に出力する。

【0256】

機械翻訳処理部5のMT用ニューラルネットワークモデル51は、入力データ埋込部2から出力される分散表現データD4(L1)($=\{x_i'token, x_i'pos, x_i'lang\}$)と、制御部から出力される制御信号CTL3とを入力する。

【0257】

制御部は、MT用ニューラルネットワークモデル51からL2のデータを出力することを指示する制御信号CTL3を生成し、当該制御信号CTL3を機械翻訳処理部5に出力する。

10

【0258】

MT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4(L1)に対して、機械翻訳処理を実行し、データD5(L2)を取得する。なお、MT用ニューラルネットワークモデル51は、制御信号CTL3に従い、出力データの種別をL2(他分野、第2言語)とする。

【0259】

MT用ニューラルネットワークモデル51は、上記により取得したデータD5(L2)を第2セクタSEL2に出力する。

【0260】

制御部は、第2セクタSEL2の端子「0」を選択する選択信号sel2を生成し、当該選択信号sel2を第2セクタSEL2に出力する。

20

【0261】

第2セクタSEL2は、選択信号sel2に従い、端子「0」を選択し、機械翻訳処理部5から出力されるデータD5(データD5(L2))を、第2出力データ評価部6にデータD6a(データD6a(L2))として出力する。

【0262】

第2出力データ評価部6は、第2セクタSEL2を介して機械翻訳処理部5から出力されるデータD5(L2)(=データD6a(L2))と、入力データ取得部1から出力される正解データD__correct(L2)とを入力する。第2出力データ評価部6は、機械翻訳処理部5の出力であるデータD5(L2)と、その正解データD__correct(L2)とから損失を取得する。

30

【0263】

C2: 教師あり機械翻訳処理(L2 L1)

入力データ埋込部2への入力をD3(L2)とし、正解データをD__correct(L1)とする場合についても、疑似対訳データ生成装置100は、上記と同様の処理を実行する。つまり、疑似対訳データ生成装置100は、上記処理において、D3(L1)をD3(L2)に置換し、正解データD__correct(L2)を正解データD__correct(L1)に置換して、上記処理と同様の処理を行う。

【0264】

(1.2.1.2D: 損失の計算処理)

上記の処理、すなわち、

(A) 自己符号化処理(ステップS121)、

(B) ゼロショット折り返し機械翻訳処理(ステップS122)、および

(C) 教師データあり機械翻訳処理(ステップS123)

を実行した後、第2出力データ評価部6は、損失の計算処理を行う。なお、損失の計算をM文(M個の文、M: 自然数)ごとに行う場合、M文の中の第i番目の文に含まれるサブワード数を $N_i(1 \leq i \leq M)$ とし、第i番目の文についての機械翻訳処理部5からの出力データD5(X_{in})(入力データを X_{in} で表す)のj番目(j : 自然数、 $1 \leq j \leq N_i$)のサブワードに相当するデータをD5($(X_{in}, X_{out}), i, j$)(「(X_{in}

40

50

X_{out})」は、入力データが X_{in} であり、出力データが X_{out} であることを表す)
 または $D5((X_{in} X_m X_{out}), i, j)$ (「 $(X_{in} X_m X_{out})$ 」は、
 入力データが X_{in} であり、1 回目の出力データが X_m であり、2 回目の入力
 が X_m であり、2 回目の出力が X_{out} であることを表す) とすると、第 2 出力データ評価部 6 は、下記
 数式のように、機械翻訳処理部 5 から出力されるデータと、正解データとから損失 $Loss$
 s を取得する。

【数 5】

$$Loss = Loss_self + Loss_zero + Loss_mt$$

10

【数 6】

$$\begin{aligned}
 Loss_self = & \sum_{i=1}^M \sum_{j=1}^{N_i} loss(D_correct'(L1, i, j), D5((L1 \rightarrow L1), i, j)) \\
 & + \sum_{i=1}^M \sum_{j=1}^{N_i} loss(D_correct'(L2, i, j), D5((L2 \rightarrow L2), i, j)) \\
 & + \sum_{i=1}^M \sum_{j=1}^{N_i} loss(D_correct'(R1, i, j), D5((R1 \rightarrow R1), i, j)) \\
 & + \sum_{i=1}^M \sum_{j=1}^{N_i} loss(D_correct'(R2, i, j), D5((R2 \rightarrow R2), i, j))
 \end{aligned}$$

20

【数 7】

30

40

50

$$\begin{aligned}
Loss_zero &= \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(L1, i, j), D5((L1 \rightarrow R2 \rightarrow L1), i, j)) \\
&+ \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(R1, i, j), D5((R1 \rightarrow L2 \rightarrow R1), i, j)) \\
&+ \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(R1, i, j), D5((R1 \rightarrow R2 \rightarrow R1), i, j)) & 10 \\
&+ \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(L2, i, j), D5((L2 \rightarrow R1 \rightarrow L2), i, j)) \\
&+ \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(R2, i, j), D5((R2 \rightarrow L1 \rightarrow R2), i, j)) \\
&+ \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(R2, i, j), D5((R2 \rightarrow R1 \rightarrow R2), i, j)) & 20
\end{aligned}$$

【数 8】

$$\begin{aligned}
Loss_mt &= \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(L2, i, j), D5((L1 \rightarrow L2), i, j)) \\
&+ \sum_{i=1}^M \sum_{j=1}^{N_i} \text{loss}(D_correct'(L1, i, j), D5((L2 \rightarrow L1), i, j)) & 30
\end{aligned}$$

【数 9】

$$\text{loss}(p, q) = - \sum_{k=1}^V p(k) \log(q(k))$$

40

V : サブワード語彙のサイズ (各トークン (入力データ) 文字列を表すベクトルの次元数)

p : 確率分布 (p(k) は、k 番目 (第 k 次元) の要素の確率を示す)

q : 確率分布 (q(k) は、k 番目 (第 k 次元) の要素の確率を示す)

なお、上記数式において、D_correct'(x, i, j) は、第 2 出力データ評価部 6 により正解データ D_correct(x) の第 i 番目の文の第 j 番目のサブワードから取得されるデータ (ベクトル) であり、サブワード語彙のサイズ (= 各トークン (入力データ) 文字列を表すベクトルの次元数 (これを n1 とする)) と同じ次元数のベクトル

50

ル ($n-1$ 次元のベクトル)である。そして、例えば、 $D_correct'(x, i, j)$ は、 $n-1$ 次元のうち1次元のみ値が「1」であり ($n-1$ 次元ベクトルの要素のうち、当該正解サブワードに対応する要素のみが「1」)、それ以外は値が「0」である one-hot ベクトルである。

【0265】

また、上記数式において、例えば、 $D5((L1 \ L2), i, j)$ は、 $n-1$ 次元の実数ベクトル (各次元のデータ (ベクトルの要素) が、対応するサブワードである確率を示す実数ベクトル) であり、softmax関数により確率化 (実数ベクトルの各要素の総和が「1」となるように正規化) されている。そして、 $D5((L1 \ L2), i, j)$ は、入力データを $L1$ としたときの機械翻訳処理部5からの出力データ種別を $L2$ としたときの出力データであり、第 i 番目の文の第 j 番目のサブワードに相当するデータである。

10

【0266】

なお、上記数式において、例えば、 $D5((L1 \ R2 \ L1), i, j)$ は、入力データの種別と翻訳処理結果データ (出力データ) の種別を明示した形式で示しており、 $(L1 \ R2 \ L1)$ は、1回目の処理において種別 $L1$ のデータが種別 $R2$ のデータとして出力され、2回目の処理において、1回目の出力を入力とし、種別 $L1$ のデータが出力されたときのデータであることを示している。

【0267】

また、上記数式において、 $loss(p, q)$ は、交差エントロピー誤差を求める数式であり、これにより、確率分布間 (上記数式では、確率分布 p, q) の相違を定量化できる。

20

【0268】

このように、第2出力データ評価部6は、上記数式に相当する処理により、損失 $Loss$ を取得する。

【0269】

(1.2.1.2E: パラメータの更新処理)

第2出力データ評価部6は、所定の学習データに対して上記で算出した損失 (学習損失) に基づいて MLM 処理部3の MT 用ニューラルネットワークモデル51のパラメータ MT を更新するためのデータ $update(MT)$ を生成し、当該データ $update(MT)$ を機械翻訳処理部5に出力する。

30

【0270】

機械翻訳処理部5の ML 用ニューラルネットワークモデル51は、データ $update(MT)$ に基づいて、パラメータ MT を更新する。

【0271】

また、 MT 用ニューラルネットワークモデル51は、 MT 用ニューラルネットワークモデル51のパラメータを更新した後、埋込層 (入力データ埋込部2のトークン埋込部21、位置埋込部22、および、言語埋込部23に対応) のパラメータを更新するためのパラメータ更新データ $update(emb)$ を生成し、当該パラメータ更新データ $update(emb)$ を入力データ埋込部2に出力する。

【0272】

入力データ埋込部2は、パラメータ更新データ $update(emb)$ に基づいて、埋込層 (入力データ埋込部2のトークン埋込部21、位置埋込部22、および、言語埋込部23に対応) のパラメータを更新する。例えば、入力データ埋込部2は、パラメータ更新データ $update(emb)$ に基づいて、変換行列 (例えば、行列 W_{token} 、 W_{pos} 、 W_{lang}) の要素 (値) を更新する。

40

【0273】

疑似対訳データ生成装置100は、所定の終了条件を満たすまで、上記処理 (図6のループ2の処理) を繰り返し実行する。その際、例えば、学習に用いるデータとは別の調整用データに対して第2出力データ評価部6で算出される損失 (テスト損失) を評価値として参照する。

50

【 0 2 7 4 】

そして、疑似対訳データ生成装置 1 0 0 は、上記処理（図 6 のループ 2 の処理）の終了条件が満たされたときをもって、疑似対訳データ生成用 N M T モデルの学習処理を終了させる。

【 0 2 7 5 】

なお、図 6 のフローチャートのループ処理（ループ 2）の終了条件は、例えば、以下のように設定される。

（ 1 ）事前に定めた反復回数だけループ処理（ループ 2）が実行された。

（ 2 ）疑似対訳データ生成用 N M T モデルの学習処理において、第 2 出力データ評価部 6 における評価値が一定以上（事前に定めた値以上）の変化を示さなかった。

10

（ 3 ）疑似対訳データ生成用 N M T モデルの学習処理において、第 2 出力データ評価部 6 における評価値が事前に定めた値を下回った。

（ 4 ）疑似対訳データ生成用 N M T モデルの学習処理において、第 2 出力データ評価部 6 における評価値が事前に定めた回数更新されなかった。

【 0 2 7 6 】

上記の終了条件を満たす場合、疑似対訳データ生成装置 1 0 0 は、疑似対訳データ生成用 N M T モデルの学習処理が完了した判断し、当該処理を終了させる。

【 0 2 7 7 】

以上のようにして、疑似対訳データ生成装置 1 0 0 は、事前学習処理を完了した時点におけるパラメータを疑似対訳データ生成用 N M T モデル（入力データ埋込部 2（埋込層に相当）と機械翻訳処理部 5 の M T 用ニューラルネットワークモデル 5 1 とにより実現されるモデル）の初期状態（初期パラメータを設定した状態（入力データ埋込部 2（埋込層に相当）のパラメータが X L M の学習により最適化された状態））として、上記処理により、疑似対訳データ生成用 N M T モデルの学習処理（パラメータ最適化処理）を行う。

20

【 0 2 7 8 】

（ 1 . 2 . 2 : 疑似対訳データ生成処理）

次に、疑似対訳データ生成処理（図 4 のステップ S 2）について説明する。

【 0 2 7 9 】

疑似対訳データ生成装置 1 0 0 は、上記の事前学習処理により取得された疑似対訳データ生成用 N M T モデル（入力データ埋込部 2（埋込層に相当）と機械翻訳処理部 5 の M T 用ニューラルネットワークモデル 5 1 とにより実現されるモデル）を用いて、適応先分野の疑似対訳データを自動的に生成する。

30

【 0 2 8 0 】

具体的には、疑似対訳データ生成装置 1 0 0 は、以下の（ 1 ）、（ 2 ）の方法を用いて、疑似対訳データ生成処理を行う。

（ 1 ）他分野の対訳データ（ L 1 - L 2 ）を適応先分野向けに改変する方法（第 1 の方法）

（ 2 ）適応先分野の単言語データ（ R 1 または R 2 ）を機械翻訳する方法（第 2 の方法）

以下、上記 2 つの方法による疑似対訳データ生成処理について、説明する。

【 0 2 8 1 】

（ 1 . 2 . 2 . 1 : 疑似対訳データ生成処理（第 1 の方法）（他分野対訳データを利用））

40

まず、第 1 の方法（他分野対訳データを利用する方法）について、説明する。

【 0 2 8 2 】

データ入力インターフェース I F 1 は、対訳データ記憶部 D B p（ L 1 - L 2 ）から他分野の対訳データ D 0（ L 1 - L 2 ）を読み出し、読み出した対訳データをデータ D 1（ L 1 - L 2 ）として、入力データ取得部 1 に出力する。

【 0 2 8 3 】

入力データ取得部 1 の第 1 入力データ取得処理部 1 1 は、データ D 1（他分野対訳データ D 1（ L 1 - L 2 ））から、 L 1 のデータをデータ D 1（ L 1 ）として抽出し、当該データ D 1（ L 1 ）から、（ 1 ）トークンデータ $x_{i0} \text{ token} (= x_{i0} \text{ token} (L$

50

1))と、(2)当該トークンの位置を特定するための位置データ $xipos (= xipos(L1))$ と、(3)当該トークンの言語を特定するための言語データ $xilang (= xilang(L1))$ と、を取得する。

【0284】

そして、入力データ取得部1は、上記により取得したトークンデータ $xitoken(L1)$ をマスク化処理部12に出力する。

【0285】

制御部(不図示)は、入力データ取得部1のマスク化処理部12に対して、マスク化処理を実行しないことを指示する制御信号CTL2を生成し、当該制御信号CTL2をマスク化処理部12に出力する。

【0286】

マスク化処理部12は、制御信号CTL2に従い、第1入力データ取得処理部11から入力したトークンデータ $xitoken(L1)$ を、トークンデータ $xitoken(L1)$ として第1セクタSEL1に出力する。

【0287】

また、第1入力データ取得処理部11は、位置データ $xipos (= xipos(L1))$ を第1セクタSEL1に出力し、言語データ $xilang (= xilang(L1))$ を第1セクタSEL1に出力する。

【0288】

つまり、入力データ取得部1は、上記により取得したデータをデータD2a(=D2a(L1)={ $xitoken(L1)$, $xipos(L1)$, $xilang(L1)$ })として第1セクタSEL1に出力する。

【0289】

制御部は、第1セクタSEL1の端子「0」を選択する選択信号sel1を生成し、当該選択信号sel1を第1セクタSEL1に出力する。

【0290】

第1セクタSEL1は、選択信号sel1に従い、端子「0」を選択し、入力データ取得部1から出力されるデータD2a(L1)を、データD3(L1)として、入力データ埋込部2に出力する。

【0291】

入力データ埋込部2は、データD3(L1)に含まれる(1)トークンデータ $xitoken(L1)$ 、(2)位置データ $xipos(L1)$ 、(3)言語データ $xilang(L1)$ から、分散表現データD4(L1)(={ $xitoken(L1)$, $xipos(L1)$, $xilang(L1)$ })を取得し、取得した分散表現データD4(L1)を機械翻訳処理部5に出力する。

【0292】

機械翻訳処理部5のMT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4(L1)(分散表現データD4(L1)(={ $xitoken(L1)$, $xipos(L1)$, $xilang(L1)$ }))と、制御部から出力される制御信号CTL3とを入力する。

【0293】

制御部は、MT用ニューラルネットワークモデル51からR2のデータを出力することを指示する制御信号CTL3を生成し、当該制御信号CTL3を機械翻訳処理部5に出力する。

【0294】

MT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4(L1)に対して、機械翻訳処理を実行し、データD5(R2)を取得し、データD5(R2)を保持する。なお、MT用ニューラルネットワークモデル51は、制御信号CTL3に従い、出力データの種別をR2(適応先分野、第2言語)とする。

【0295】

10

20

30

40

50

また、入力データ取得部 1 の第 1 入力データ取得処理部 1 1 は、データ D 1 (他分野対訳データ D 1 (L 1 - L 2)) から、L 2 のデータをデータ D 1 (L 2) として抽出し、当該データ D 1 (L 2) から、(1) トークンデータ $x_{i0\text{token}} (= x_{i0\text{token}}(L 2))$ と、(2) 当該トークンの位置を特定するための位置データ $x_{ipos} (= x_{ipos}(L 2))$ と、(3) 当該トークンの言語を特定するための言語データ $x_{ilang} (= x_{ilang}(L 2))$ と、を取得する。

【0296】

そして、入力データ取得部 1 は、上記により取得したトークンデータ $x_{i0\text{token}}(L 2)$ をマスク化処理部 1 2 に出力する。

【0297】

制御部 (不図示) は、入力データ取得部 1 のマスク化処理部 1 2 に対して、マスク化処理を実行しないことを指示する制御信号 CTL 2 を生成し、当該制御信号 CTL 2 をマスク化処理部 1 2 に出力する。

【0298】

マスク化処理部 1 2 は、制御信号 CTL 2 に従い、第 1 入力データ取得処理部 1 1 から入力したトークンデータ $x_{i0\text{token}}(L 2)$ を、トークンデータ $x_{itoken}(L 2)$ として第 1 セレクタ SEL 1 に出力する。

【0299】

また、第 1 入力データ取得処理部 1 1 は、位置データ $x_{ipos} (= x_{ipos}(L 2))$ を第 1 セレクタ SEL 1 に出力し、言語データ $x_{ilang} (= x_{ilang}(L 2))$ を第 1 セレクタ SEL 1 に出力する。

【0300】

つまり、入力データ取得部 1 は、上記により取得したデータをデータ D 2 a ($= D 2 a(L 2) = \{x_{itoken}(L 2), x_{ipos}(L 2), x_{ilang}(L 2)\}$) として第 1 セレクタ SEL 1 に出力する。

【0301】

制御部は、第 1 セレクタ SEL 1 の端子「0」を選択する選択信号 sel 1 を生成し、当該選択信号 sel 1 を第 1 セレクタ SEL 1 に出力する。

【0302】

第 1 セレクタ SEL 1 は、選択信号 sel 1 に従い、端子「0」を選択し、入力データ取得部 1 から出力されるデータ D 2 a (L 2) を、データ D 3 (L 2) として、入力データ埋込部 2 に出力する。

【0303】

入力データ埋込部 2 は、データ D 3 (L 2) に含まれる (1) トークンデータ $x_{itoken}(L 2)$ 、(2) 位置データ $x_{ipos}(L 2)$ 、(3) 言語データ $x_{ilang}(L 2)$ から、分散表現データ D 4 (L 2) ($= \{x_{i'token}(L 2), x_{i'pos}(L 2), x_{i'lang}(L 2)\}$) を取得し、取得した分散表現データ D 4 (L 2) を機械翻訳処理部 5 に出力する。

【0304】

機械翻訳処理部 5 の MT 用ニューラルネットワークモデル 5 1 は、入力データ埋込部 2 から出力されるデータ D 4 (L 2) (分散表現データ D 4 (L 2) ($= \{x_{i'token}(L 2), x_{i'pos}(L 2), x_{i'lang}(L 2)\}$)) と、制御部から出力される制御信号 CTL 3 とを入力する。

【0305】

制御部は、MT 用ニューラルネットワークモデル 5 1 から R 1 のデータを出力することを指示する制御信号 CTL 3 を生成し、当該制御信号 CTL 3 を機械翻訳処理部 5 に出力する。

【0306】

MT 用ニューラルネットワークモデル 5 1 は、入力データ埋込部 2 から出力されるデータ D 4 (L 2) に対して、機械翻訳処理を実行し、データ D 5 (R 1) を取得する。なお

10

20

30

40

50

、MT用ニューラルネットワークモデル51は、制御信号CTL3に従い、出力データの種別をR1（適応先分野、第1言語）とする。

【0307】

機械翻訳処理部5は、上記処理により取得したデータD5（R2）とデータD5（R1）とをペアリングして、データD5（R2 - R1）として、第2セクタSEL2に出力する。

【0308】

制御部は、第2セクタSEL2の端子「0」を選択する選択信号sel2を生成し、当該選択信号sel2を第2セクタSEL2に出力する。

【0309】

第2セクタSEL2は、選択信号sel2に従い、端子「0」を選択し、機械翻訳処理部5から出力されるデータD5（データD5（R2 - R1））を、データD6a（R2 - R1）として、フィルター処理部8に出力する（図12を参照）。

【0310】

制御部は、フィルター処理部8にて、R2 - R1の疑似対訳データに対してフィルター処理を実行することを指示する制御信号CTL4を生成し、当該制御信号をフィルター処理部8に出力する。

【0311】

フィルター処理部8は、制御信号CTL4に従い、入力されるデータD6a（R2 - R1）（R2およびR1の疑似対訳データ）に対してフィルター処理を行う。

【0312】

ここで、上記処理により取得された適応先分野の疑似対訳データ（D6a（R2 - R1））は、教師データがない適応先分野の対訳データを出力するように学習させたニューラルネットワークモデルを用いて処理して取得したデータであるため、対訳データとしての品質が低い可能性がある。

【0313】

そこで、フィルター処理部8は、上記処理により取得された適応先分野の疑似対訳データ（D6a（R2 - R1））に対して、フィルター処理を行う。例えば、フィルター処理部8は、上記処理により取得された適応先分野の疑似対訳データ（D6a（R2 - R1））の各文対に信頼度を付与し、付与した信頼度に基づいて、フィルタリングを行う。

【0314】

例えば、元のL1 - L2の対訳データ（対訳データ記憶部DBp（L1 - L2）から取得した対訳データD0（L1 - L2））を参照することが考えられる。より具体的には、フィルター処理部8は、対訳データL1 - L2（D0（L1 - L2））における文対a1 - a2を機械翻訳して取得した、適応先分野R1 - R2の疑似対訳データの文対b1 - b2の信頼度を、a1とb1との類似度、および、a2とb2の類似度に基づいて算出し、その信頼度が事前に定めた閾値Thよりも高いか否かを判断する。なお、フィルター処理部8は、対訳データ記憶部DBp（L1 - L2）から、対訳データL1 - L2（D0（L1 - L2））における文対a1 - a2を取得できるものとする（例えば、データ入力インターフェースIF1を介して取得する）。

【0315】

そして、フィルター処理部8は、上記により算出した信頼度が所定の閾値Thよりも高い疑似対訳データ（R1 - R2）のみをフィルタリングにより取得する。そして、当該フィルタリングにより取得した疑似対訳データを、疑似対訳データDpsd1（R1 - R2, #1）として、疑似対訳データ格納部DB1に出力する。

【0316】

（1.2.2.2：疑似対訳データ生成処理（第2の方法）（適応先分野単言語データを利用））

次に、第2の方法（適応先分野単言語データを利用する方法）について、説明する。

【0317】

10

20

30

40

50

図13に示すように、データ入力インターフェースIF1は、第3単言語データ記憶部DBm(R1)から適応先分野の第1言語の単言語データD0(R1)を読み出し、読み出した第1言語のデータをデータD1(R1)として、入力データ取得部1に出力する。

【0318】

入力データ取得部1の第1入力データ取得処理部11は、データD1(適応先分野第1言語データD1(R1))から、当該第1言語のデータ(言語データ)D1(R1)(適応先分野第1言語の文を構成するデータD1(R1))から、(1)トークンデータ $x_{i0\text{token}}$ (= $x_{i0\text{token}}(R1)$)と、(2)当該トークンの位置を特定するための位置データ x_{ipos} (= $x_{ipos}(R1)$)と、(3)当該トークンの言語を特定するための言語データ x_{ilang} (= $x_{ilang}(R1)$)と、を取得する。

10

【0319】

そして、入力データ取得部1は、上記により取得したトークンデータ $x_{i0\text{token}}$ (R1)をマスク化処理部12に出力する。

【0320】

制御部(不図示)は、入力データ取得部1のマスク化処理部12に対して、マスク化処理を実行しないことを指示する制御信号CTL2を生成し、当該制御信号CTL2をマスク化処理部12に出力する。

【0321】

マスク化処理部12は、制御信号CTL2に従い、第1入力データ取得処理部11から入力したトークンデータ $x_{i0\text{token}}$ (R1)を、トークンデータ x_{itoken} (R1)として第1セクタSEL1に出力する。

20

【0322】

また、第1入力データ取得処理部11は、位置データ x_{ipos} (= $x_{ipos}(R1)$)を第1セクタSEL1に出力し、言語データ x_{ilang} (= $x_{ilang}(R1)$)を第1セクタSEL1に出力する。

【0323】

つまり、入力データ取得部1は、上記により取得したデータをデータD2a(= $D2a(R1) = \{x_{itoken}(R1), x_{ipos}(R1), x_{ilang}(R1)\}$)として第1セクタSEL1に出力する。

【0324】

制御部は、第1セクタSEL1の端子「0」を選択する選択信号sel1を生成し、当該選択信号sel1を第1セクタSEL1に出力する。

30

【0325】

第1セクタSEL1は、選択信号sel1に従い、端子「0」を選択し、入力データ取得部1から出力されるデータD2a(R1)を、データD3(R1)として、入力データ埋込部2に出力する。

【0326】

入力データ埋込部2は、データD3(R1)に含まれる(1)トークンデータ x_{itoken} 、(2)位置データ x_{ipos} 、(3)言語データ x_{ilang} から、分散表現データD4(R1)(= $\{x_{i'token}(R1), x_{i'pos}(R1), x_{i'lang}(R1)\}$)を取得し、取得した分散表現データD4(R1)を機械翻訳処理部5に出力する。

40

【0327】

機械翻訳処理部5のMT用ニューラルネットワークモデル51は、入力データ埋込部2から出力される分散表現データD4(R1)(= $\{x_{i'token}(R1), x_{i'pos}(R1), x_{i'lang}(R1)\}$)と、制御部から出力される制御信号CTL3とを入力する。

【0328】

制御部は、MT用ニューラルネットワークモデル51からR2のデータを出力することを指示する制御信号CTL3を生成し、当該制御信号CTL3を機械翻訳処理部5に出力する。

50

【 0 3 2 9 】

MT用ニューラルネットワークモデル51は、入力データ埋込部2から出力されるデータD4(R1)に対して、機械翻訳処理を実行し、データD5(R2)を取得する。なお、MT用ニューラルネットワークモデル51は、制御信号CTL3に従い、出力データの種別をR2(適応先分野、第2言語)とする。

【 0 3 3 0 】

MT用ニューラルネットワークモデル51は、上記により取得したデータD5(R2)を第2セレクタSEL2に出力する。

【 0 3 3 1 】

制御部は、第2セレクタSEL2の端子「0」を選択する選択信号sel2を生成し、当該選択信号sel2を第2セレクタSEL2に出力する。

10

【 0 3 3 2 】

第2セレクタSEL2は、選択信号sel2に従い、端子「0」を選択し、機械翻訳処理部5から出力されるデータD5(データD5(R2))を、データD6a(R2)として、フィルター処理部8に出力する(図13を参照)。

【 0 3 3 3 】

制御部は、フィルター処理部8にて、R1-R2の疑似対訳データに対してフィルター処理を実行することを指示する制御信号CTL4を生成し、当該制御信号をフィルター処理部8に出力する。

【 0 3 3 4 】

また、入力データ取得部1は、第3単言語データ記憶部DBm(R1)から取得した適応先分野の言語データD0(R1)(機械翻訳処理の入力としたデータ)をデータD1__org(R1)として、フィルター処理部8に出力する。

20

【 0 3 3 5 】

そして、フィルター処理部8は、入力される、(1)データD1__org(R1)(機械翻訳処理の入力としたデータ)と、(2)データD6a(R2)(データD3(R1)の機械翻訳データ)とを対応づけることで適応先分野の疑似対訳データ(R1-R2)を取得する。そして、フィルター処理部8は、制御信号CTL4に従い、疑似対訳データ(R1-R2)にフィルター処理を行う。

【 0 3 3 6 】

ここで、上記処理により取得された適応先分野の疑似対訳データ(R1-R2)は、教師データがない適応先分野の単語データを出力するように学習させたニューラルネットワークモデルを用いて処理して取得したデータであるため、対訳データとしての品質が低い可能性がある。

30

【 0 3 3 7 】

そこで、フィルター処理部8は、上記処理により取得された適応先分野の疑似対訳データ(R1-R2)に対して、フィルター処理を行う。フィルター処理部8は、例えば、下記文献Dに開示されている機械翻訳の品質推定に関する技術を用いて、上記処理により取得された適応先分野の疑似対訳データ(R1-R2)に対して信頼度を付与する。

(文献D): Lucia Specia, Carolina Scarton, and Gustavo Henrique Paetzold (2018). Quality Estimation for Machine Translation. Morgan & Claypool.

40

そして、フィルター処理部8は、付与された信頼度に基づいて、フィルタリングを行う。例えば、フィルター処理部8は、上記により算出した信頼度が所定の閾値Thよりも高い疑似対訳データ(R1-R2)のみをフィルタリングにより取得する。そして、フィルター処理部8は、当該フィルタリングにより取得した疑似対訳データを、疑似対訳データDpsd1(R1-R2, #2)として、疑似対訳データ格納部DB1に出力する。

【 0 3 3 8 】

なお、上記では、データD1__orgを適応先分野の第1言語のデータ(R1のデータ(=D1__org(R1)))とし、データD3(R1)を機械翻訳処理により取得した翻訳データを適応先分野の第2言語のデータD5(R2)(データD6a(R2))とし

50

たが（図 13 の場合）、図 14 に示すように、データ D1__org を適応先分野の第 2 言語のデータ（R2 のデータ（= D1__org（R2）））とし、データ D3（R2）を機械翻訳処理により取得した翻訳データを適応先分野の第 2 言語のデータ D5（R1）（データ D6a（R2））として、上記と同様の処理により、疑似対訳データ（データ Dpsd1（R1 - R2, #2））を生成するようにしてもよい。

【0339】

まとめ

以上のように、疑似対訳データ生成装置 100 では、

(1) 大規模で（対訳データ数が多く）精度の高い他分野の対訳データ（対訳データ記憶部 DBp（L1 - L2）に記憶されている対訳データ）と、

10

(2) 他分野の第 1 言語データ（第 1 単言語データ記憶部 DBm（L1）に記憶されている単言語データ）と、

(3) 他分野の第 2 言語データ（第 2 単言語データ記憶部 DBm（L2）に記憶されている単言語データ）と、

(4) 適応先分野の第 1 言語データ（第 3 単言語データ記憶部 DBm（R1）に記憶されている単言語データ）と、

(5) 適応先分野の第 2 言語データ（第 4 単言語データ記憶部 DBm（R2）に記憶されている単言語データ）と、

を用いて、言語横断言語モデル（XLM）（入力データ埋込部 2（埋込層に相当）と XLM 用ニューラルネットワークモデル 31 とにより実現されるモデル）を最適化する（事前学習処理による初期化）。そして、疑似対訳データ生成装置 100 では、言語横断言語モデル（XLM）の最適化処理後（事前学習処理後）の入力データ埋込部 2（埋込層に相当）（XLM の最適化により取得されたパラメータ（例えば、変換行列 W_{token} 、 W_{pos} 、 W_{lang} が設定されている状態の入力データ埋込部 2））と、学習前の機械翻訳処理部 5 の MT 用ニューラルネットワークモデル 51 とからなる疑似対訳データ生成用 NMT モデルの状態を初期状態として、疑似対訳データ生成用 NMT モデルのパラメータ最適化処理を行う。これにより、疑似対訳データ生成装置 100 では、適応先分野（機械翻訳の対象とする分野）の対訳データが一切ない場合であっても、適応先分野の第 1 言語および第 2 言語の疑似対訳データ（R1 - R2）を生成できるモデル（学習済みモデル）を取得

20

できる。

30

【0340】

(1.2.3：翻訳モデル学習処理)

次に、翻訳モデルの学習処理（図 5 のステップ S3）について説明する。

【0341】

機械翻訳装置 MT1 は、NMT モデルを有しており、当該 NMT モデルにより機械翻訳を行う装置である。機械翻訳装置 MT1 は、疑似対訳データ生成装置 100 により生成され、疑似対訳データ格納部 DB1 に記憶された疑似対訳データを用いて、適応先分野の機械翻訳を行うための機械翻訳モデル（NMT モデル）の学習処理を行う。

【0342】

具体的には、例えば、機械翻訳装置 MT1 は、以下の (1) ~ (3) の方法のいずれかにより、適応先分野の機械翻訳を行うための機械翻訳モデル（NMT モデル）の学習処理を行う。

40

(1) 微調整法：

機械翻訳装置 MT1 は、対訳データ記憶部 DBp（L1 - L2）から、他分野の対訳データ D0'（L1 - L2）を取得し、他分野の対訳データ D0'（L1 - L2）を用いて、機械翻訳装置 MT1 の NMT モデルの学習処理を行う。その後、機械翻訳装置 MT1 は、疑似対訳データ生成装置 100 により生成された適応先分野の疑似対訳データ（R1 - R2 の対訳データ）を疑似対訳データ格納部 DB1 から疑似対訳データ Dpsd2（R1 - R2）として読み出し、読み出した適応先分野の疑似対訳データ Dpsd2（R1 - R2）により、機械翻訳装置 MT1 の NMT モデルの学習処理（パラメータの微調整）を行う。

50

(2) データ混合法：

機械翻訳装置MT1は、対訳データ記憶部DBp(L1-L2)から、他分野の対訳データD0'(L1-L2)を取得するとともに、疑似対訳データ格納部DB1から、適応分野の疑似対訳データDpsd2(R1-R2)を取得する。そして、機械翻訳装置MT1は、他分野(L1-L2)の対訳データD0'(L1-L2)と適応先分野(R1-R2)の疑似対訳データDpsd2(R1-R2)とを混合したデータを生成し、生成した当該データにより、機械翻訳装置MT1のNMTモデルを学習させる。なお、他分野(L1-L2)の対訳データD0'(L1-L2)と適応先分野(R1-R2)の疑似対訳データDpsd2(R1-R2)とを混合する際に、2種類の対訳データをタグで区別するようにしてもよい。また、2種類の対訳データを混合する前に、一方または両方をオーバーサンプリングまたはアンダーサンプリングすることで対訳データの混合比を変更してもよい。

10

(3) データ混合微調整法：

機械翻訳装置MT1は、対訳データ記憶部DBp(L1-L2)から、他分野の対訳データD0'(L1-L2)を取得し、他分野の対訳データD0'(L1-L2)を用いて、機械翻訳装置MT1のNMTモデルの学習処理を行う。その後、機械翻訳装置MT1は、機械翻訳装置MT1のNMTモデルの学習処理に用いた対訳データと同じ対訳データD0'(L1-L2)と、適応先分野の疑似対訳データ(R1-R2)(疑似対訳データ格納部DB1から読み出した適応先分野の疑似対訳データDpsd2(R1-R2))とを混合したデータを生成する。そして、機械翻訳装置MT1は、生成したデータ(2種類の対訳データを混合したデータ)を用いて、機械翻訳装置MT1のNMTモデルの学習処理(パラメータの微調整)を行う。なお、他分野(L1-L2)の対訳データD0'(L1-L2)と適応先分野(R1-R2)の疑似対訳データDpsd2(R1-R2)とを混合する際に、2種類の対訳データをタグで区別するようにしてもよい。また、2種類の対訳データを混合する前に、一方または両方をオーバーサンプリングまたはアンダーサンプリングすることで対訳データの混合比を変更してもよい。

20

【0343】

機械翻訳装置MT1は、上記(1)~(3)の方法のいずれかによるNMTモデルの学習処理を行い、NMTモデルの最適パラメータを取得する。そして、当該最適パラメータを設定したNMTモデル(学習済みモデル)を用いて、適応先分野の機械翻訳を行う。つまり、図1に示すように、機械翻訳装置MT1は、適応先分野の起点言語データDin__eを入力したとき、当該データDin__eに対して機械翻訳処理を行い、目標言語のデータDout__j(機械翻訳したデータ)を取得する。

30

【0344】

図15に、機械翻訳装置MT1により取得されたデータの一例を示す。図15の上段は、原文(起点言語データ)であり、図15の下段は、本発明の機械翻訳システム1000により取得された機械翻訳データである。また、図15の中段に、比較のため、従来技術の機械翻訳処理を行った結果データ(適応先分野の対訳データによる学習を行うことなく取得したNMTモデルによる機械翻訳の結果データ)を示す。なお、図15の下段の出力データは、他分野の対訳データとして、日英の対訳データ(適応先分野以外の対訳データ)を用い、適応先分野をSNS(Social Networking Service)分野としたときのデータ(機械翻訳結果データ)である。

40

【0345】

図15から分かるように、適応先分野なしの機械翻訳(従来技術の機械翻訳)により翻訳データでは、正しく翻訳できていないが、本発明の機械翻訳システム1000により取得された機械翻訳データは従来技術の機械翻訳に比べ良好なものとなっている。

【0346】

まとめ

このように、機械翻訳システム1000では、適応先分野(対象分野)における対訳データが存在しない場合であっても、疑似対訳データ生成装置100により、適応先分野の疑似対訳データを生成することができ、また、疑似対訳データ生成装置100により生成

50

された疑似対訳データを用いて、機械翻訳装置MT1のNMTモデルを学習させることができる。さらに、機械翻訳システム1000では、疑似対訳データ生成装置100により生成された疑似対訳データにより学習させたNMTモデルを用いて、機械翻訳処理を行うことで、適応先分野（対象分野）における対訳データが存在しない場合であっても、適応先分野における機械翻訳を精度良く行うことができる。

【0347】

機械翻訳システム1000では、XLMの学習、疑似対訳データ生成用機械翻訳モデルの学習、および、機械翻訳装置MT1による適応先分野のNMTモデルの学習のいずれにおいても適応先分野の正規の対訳データを一切使用していない。この点において、機械翻訳システム1000は、従来技術と大きく相違する。なお、機械翻訳システム1000において、疑似対訳データの中から最も信頼度が高い部分を抽出し、調整用の対訳データとして用いるようにしてもよい。

10

【0348】

また、機械翻訳システム1000では、適応先分野（対象分野）における対訳データが存在しないが、他分野の対訳データおよび適応先分野の単言語データが大規模に存在する場合に絶大な効果を発揮する。

【0349】

[他の実施形態]

上記実施形態で説明した機械翻訳システム1000、疑似対訳データ生成装置100、機械翻訳装置MT1において、各ブロックは、LSIなどの半導体装置により個別に1チップ化されても良いし、一部または全部を含むように1チップ化されても良い。

20

【0350】

なおここではLSIとしたが、集積度の違いにより、IC、システムLSI、スーパーLSI、ウルトラLSIと呼称されることもある。

【0351】

また集積回路化の手法はLSIに限るものではなく、専用回路または汎用プロセッサで実現してもよい。LSI製造後にプログラムすることが可能なFPGA(Field Programmable Gate Array)や、LSI内部の回路セルの接続や設定を再構成可能なりコンフィギュラブル・プロセッサを利用してよい。

【0352】

30

また上記各実施形態の各機能ブロックの処理の一部または全部は、プログラムにより実現されるものであってもよい。そして上記各実施形態の各機能ブロックの処理の一部または全部は、コンピュータにおいて、中央演算装置(CPU)により行われる。また、それぞれの処理を行うためのプログラムは、ハードディスク、ROMなどの記憶装置に格納されており、ROMにおいて、あるいはRAMに読み出されて実行される。

【0353】

また上記実施形態の各処理をハードウェアにより実現してもよいし、ソフトウェア(OS(オペレーティングシステム)、ミドルウェア、あるいは所定のライブラリとともに実現される場合を含む。)により実現してもよい。さらにソフトウェアおよびハードウェアの混在処理により実現してもよい。

40

【0354】

例えば上記実施形態の各機能部をソフトウェアにより実現する場合、図16に示したハードウェア構成(例えばCPU(GPUであってもよい)、ROM、RAM、入力部、出力部、通信部、記憶部(例えば、HDD、SSD等により実現される記憶部)、外部メディア用ドライブ等をバスBusにより接続したハードウェア構成)を用いて各機能部をソフトウェア処理により実現するようにしてもよい。

【0355】

また上記実施形態の各機能部をソフトウェアにより実現する場合、当該ソフトウェアは、図16に示したハードウェア構成を有する単独のコンピュータを用いて実現されるものであってもよいし、複数のコンピュータを用いて分散処理により実現されるものであってもよい。

50

もよい。

【0356】

また上記実施形態における処理方法の実行順序は、必ずしも上記実施形態の記載に制限されるものではなく、発明の要旨を逸脱しない範囲で、実行順序を入れ替えることができるものである。

【0357】

前述した方法をコンピュータに実行させるコンピュータプログラム、及びそのプログラムを記録したコンピュータ読み取り可能な記録媒体は、本発明の範囲に含まれる。ここでコンピュータ読み取り可能な記録媒体としては、例えば、フレキシブルディスク、ハードディスク、CD-ROM、MO、DVD、DVD-ROM、DVD-RAM、大容量DVD、次世代DVD、半導体メモリを挙げることができる。

10

【0358】

上記コンピュータプログラムは、上記記録媒体に記録されたものに限らず、電気通信回線、無線または有線通信回線、インターネットを代表とするネットワーク等を経由して伝送されるものであってもよい。

【0359】

また、本明細書内の記載、特許請求の範囲の記載において、「最適化」とは、最も良い状態にすることをいい、システム（モデル）を「最適化」するパラメータとは、当該システムの目的関数の値が最適値となるときのパラメータのことをいう。「最適値」は、システムの目的関数の値が大きくなるほど、システムが良い状態となる場合は、最大値であり、システムの目的関数の値が小さくなるほど、システムが良い状態となる場合は、最小値である。また、「最適値」は、極値であってもよい。また、「最適値」は、所定の誤差（測定誤差、量子化誤差等）を許容するものであってもよく、所定の範囲（十分収束したとみなすことができる範囲）に含まれる値であってもよい。

20

【0360】

なお、本発明の具体的な構成は、前述の実施形態に限られるものではなく、発明の要旨を逸脱しない範囲で種々の変更および修正が可能である。

【符号の説明】

【0361】

- 1000 機械翻訳システム
- 100 疑似対訳データ生成装置
- 1 入力データ取得部
- 11 第1入力データ取得処理部
- 12 マスク化処理部
- 2 入力データ埋込部
- 21 トークン埋込部
- 22 位置埋込部
- 23 言語埋込部
- 3 XLM処理部
- 31 XLM用ニューラルネットワークモデル
- 5 機械翻訳処理部
- 51 MT用ニューラルネットワークモデル
- 8 フィルター処理部
- MT1 機械翻訳装置
- DBp(L1-L2) 対訳データ記憶部
- DBm(L1) 第1単言語データ記憶部（他分野の単言語データ用（第1言語））
- DBm(L2) 第2単言語データ記憶部（他分野の単言語データ用（第2言語））
- DBm(R1) 第3単言語データ記憶部（適応先分野の単言語データ用（第1言語））
- DBm(R2) 第4単言語データ記憶部（適応先分野の単言語データ用（第2言語））
- DB1 疑似対訳データ格納部

30

40

50

【図面】
【図 1】

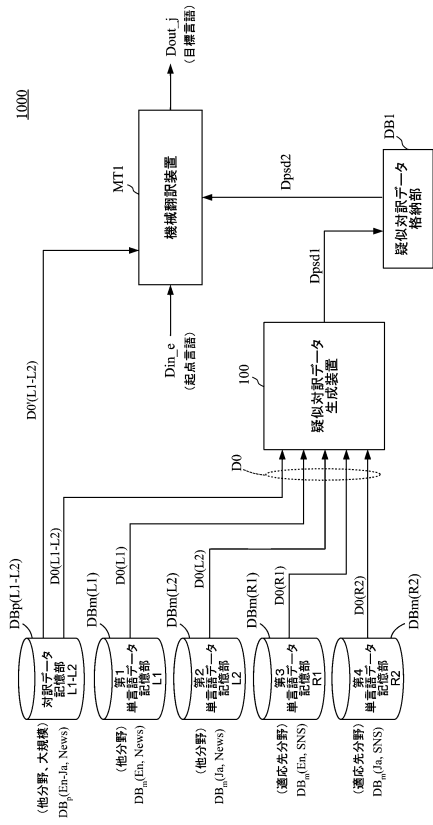


FIG. 1

【図 2】

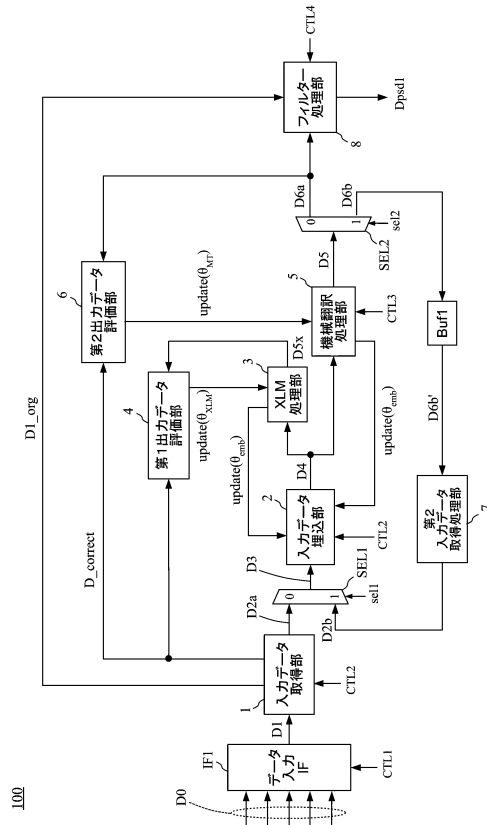


FIG. 2

【図 3】

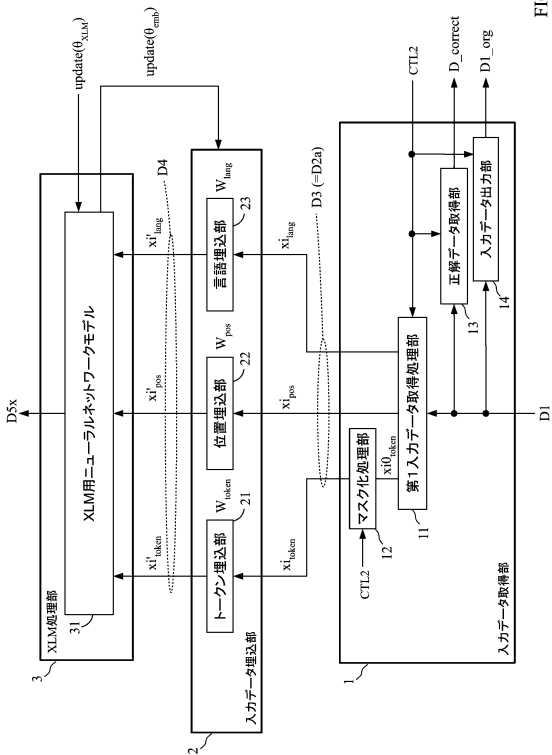


FIG. 3

【図 4】

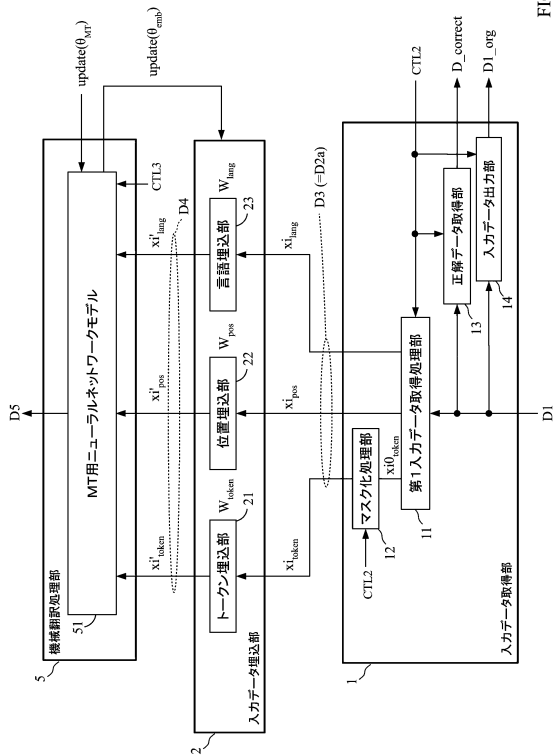


FIG. 4

【図5】

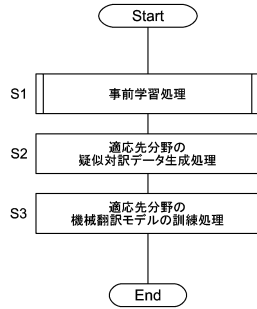


FIG. 5

【図6】

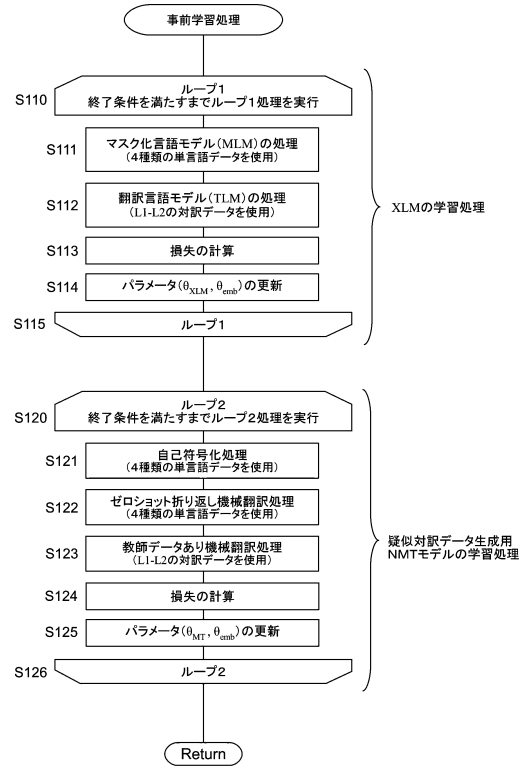


FIG. 6

【図7】

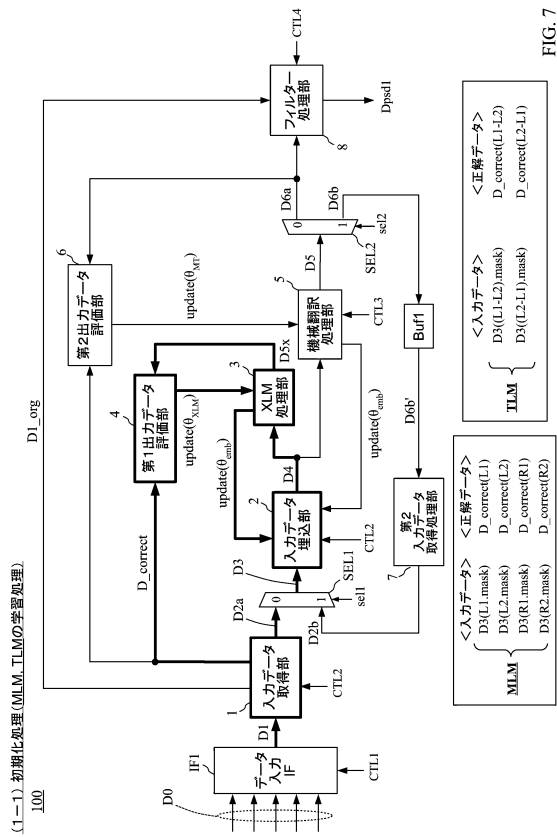


FIG. 7

【図8】

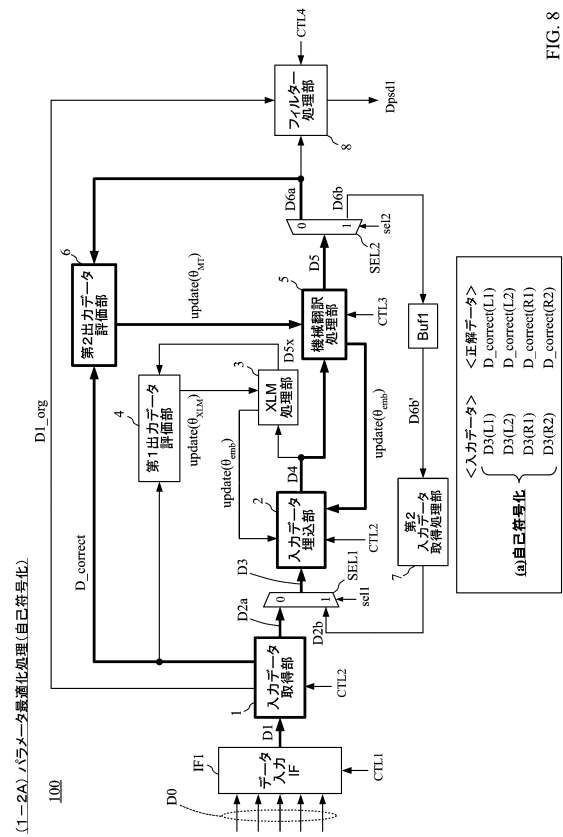


FIG. 8

10

20

30

40

50

【図 9】

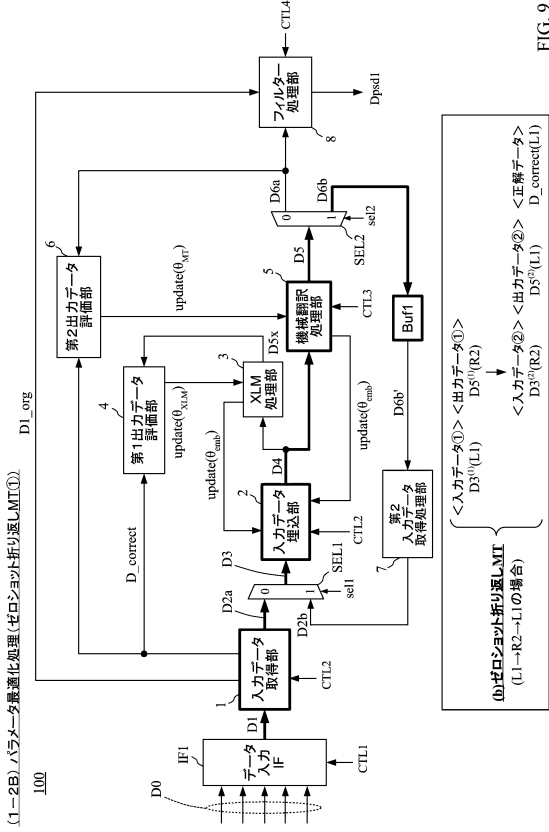


FIG. 9

【図 10】

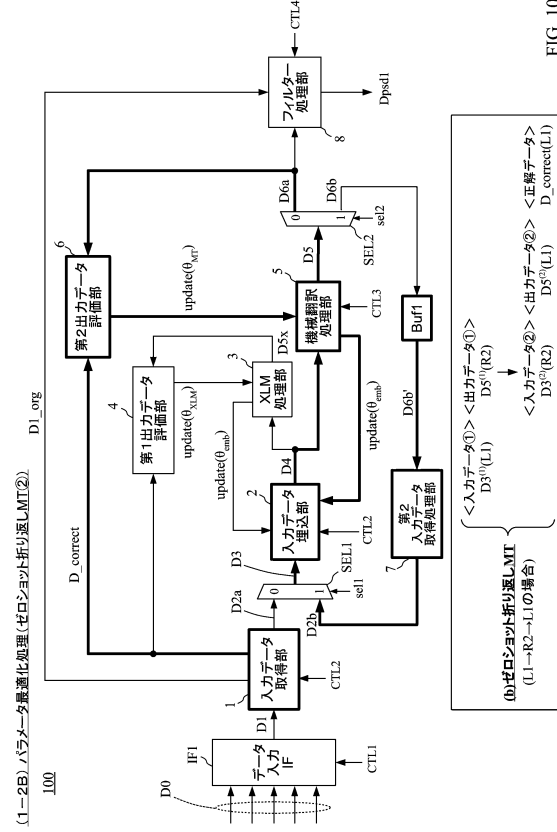


FIG. 10

【図 11】

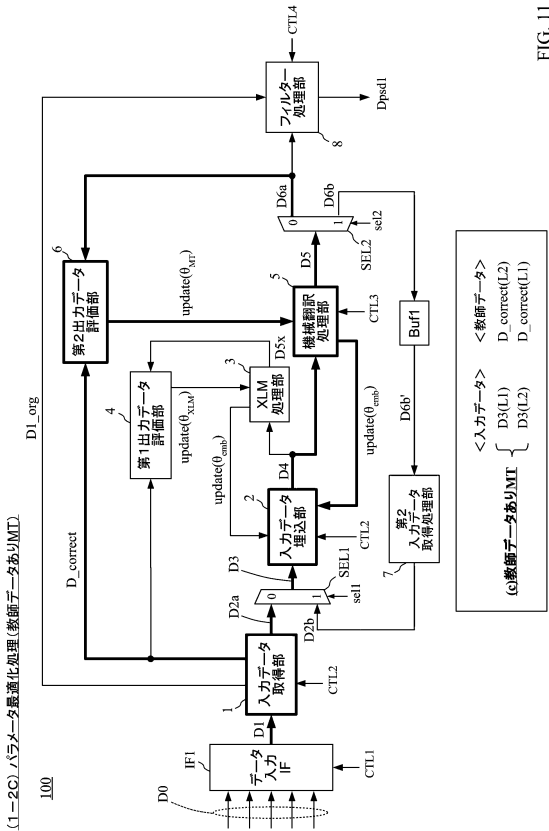


FIG. 11

【図 12】

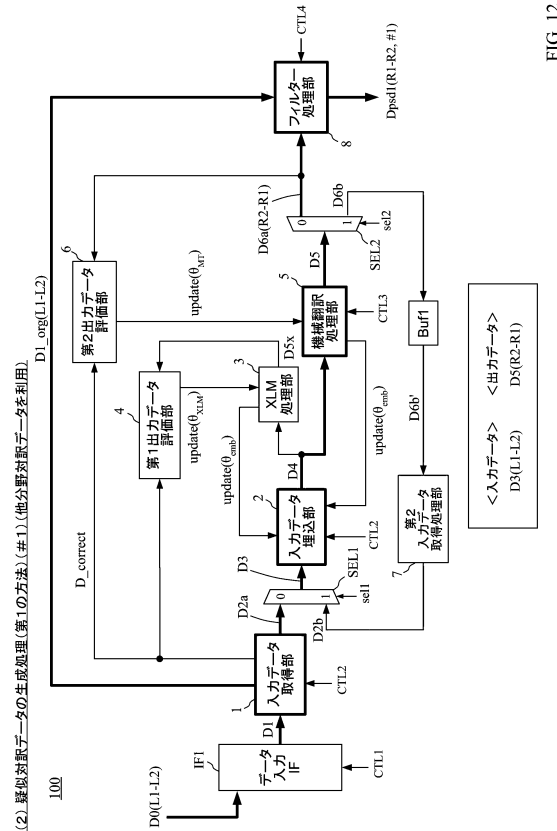


FIG. 12

(1-2C) パラメータ最適化処理 (教師データありMIT)

(2) 類似データ生成処理 (第1の方法)(#1)(他分野対照データを利用)

【図 1 3】

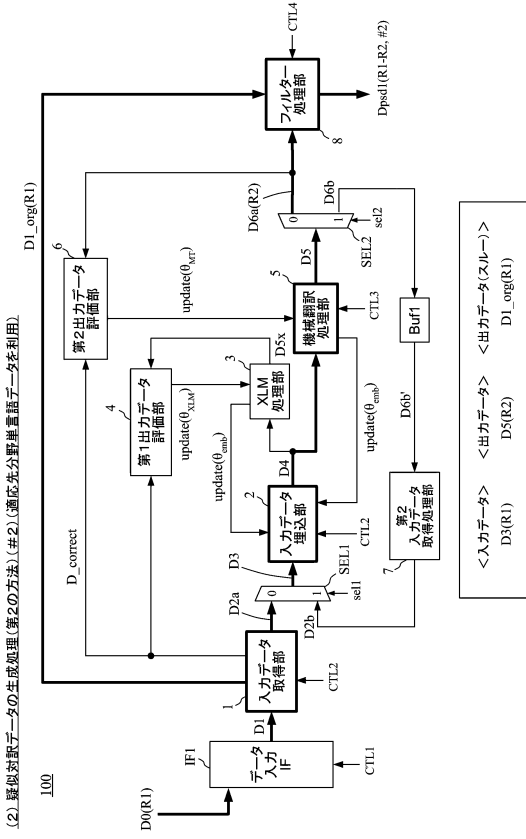


FIG. 13

【図 1 4】

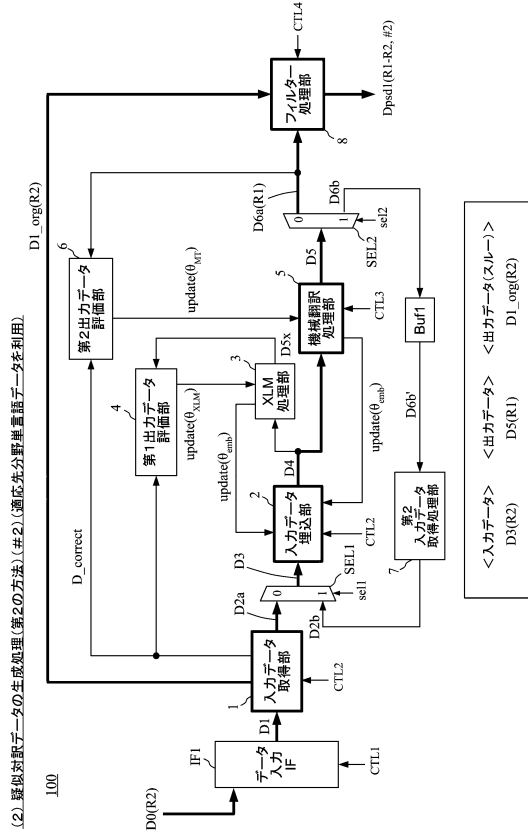


FIG. 14

【図 1 5】

原文

#COVID19 in NYC, plz dont go out and stay @home!!

分野適応なし

#NYCのCOVID19、PLZは外出せず@home!

本発明

#COVID19 NYC、外出しない、家についてください!!!

FIG. 15

【図 1 6】

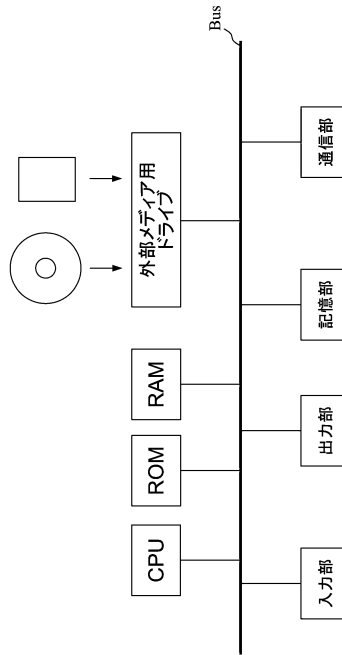


FIG. 16

フロントページの続き

東京都小金井市貫井北町4 - 2 - 1 国立研究開発法人情報通信研究機構内

審査官 成瀬 博之

(56)参考文献 特開2020 - 112915 (JP, A)

特開2018 - 116324 (JP, A)

森田知熙 他2名, 双方向ニューラル機械翻訳の反復的な教師なし適応の検討, 言語処理学会
第25回年次大会 発表論文集, 言語処理学会, 2019年03月04日, 1451-1454頁

(58)調査した分野 (Int.Cl., DB名)

G06F 40/20 - 40/58