

令和 4 年 6 月 4 日現在

機関番号：34310

研究種目：若手研究

研究期間：2020～2021

課題番号：20K19864

研究課題名（和文）同期式構文解析に基づくニューラル機械翻訳に関する研究

研究課題名（英文）Neural machine translation based on synchronous syntactic structure

研究代表者

田村 晃裕（Tamura, Akihiro）

同志社大学・理工学部・准教授

研究者番号：20804165

交付決定額（研究期間全体）：（直接経費） 3,300,000円

研究成果の概要（和文）：本研究では、翻訳元の言語と翻訳先の言語で対応させた文構造を活用することで、ニューラルネットワークに基づく機械翻訳（NMT）の性能改善を実現した。実施期間中に、同期された文構造を活用するNMTとして、NMTモデル内で文構造を同期させる方法と、既存の同期式構文解析結果をNMTモデルで活用する方法の二つを試行した。そして評価実験を通じて、NMTモデル内で文構造を同期させることにより、日英翻訳、英独翻訳（英語からドイツ語への翻訳）及び英羅翻訳（英語からルーマニア語への翻訳）の性能が改善できることを示した。

研究成果の学術的意義や社会的意義

近年、グローバル化の進展とともに、外国語の利活用を支援する機械翻訳の需要が高まっている。しかし、現在の機械翻訳では構造が異なる言語間の翻訳は難しく、その翻訳性能の改善が大きな課題の一つとなっている。本研究では、その課題を解決するため、翻訳元の言語と翻訳先の言語で対応させた文構造をNMTで活用する初めての試みに取り組んだ。そして、NMTモデル内で文構造を同期させることにより、日英を含む複数の言語対で翻訳性能を改善できることを示し、今後の機械翻訳の研究開発において、同期された文構造を活用する重要性を示唆した。

研究成果の概要（英文）：This research aims to improve the performance of neural network-based machine translation (NMT) by using synchronous syntactic structure, which is a sentence structure aligned across source and target languages. The study has proposed two NMT models based on synchronous syntactic structure: (i) one is to synchronize sentence structures across languages in the NMT model, and (ii) the other is to incorporate synchronous parse trees derived by existing synchronous context free grammar. The evaluations show that the approach (i) improves the performance of Japanese-to-English, English-to-German, and English-to-Romanian translation.

研究分野：情報学 - 人間情報学 知能情報学

キーワード：ニューラル機械翻訳 同期式構文解析 ニューラルネットワーク 機械翻訳 Transformer 同期文構造

科研費による研究は、研究者の自覚と責任において実施するものです。そのため、研究の実施や研究成果の公表等については、国の要請等に基づくものではなく、その研究成果に関する見解や責任は、研究者個人に帰属します。

様式 C - 19、F - 19 - 1、Z - 19 (共通)

1. 研究開始当初の背景

近年、グローバル化の進展とともに外国語に接する機会が増加しており、外国人とのコミュニケーションや外国語で発信された情報の利活用等を補助する高精度な機械翻訳の需要が高まっている。機械翻訳は、自然言語処理の研究分野において古くから研究されており、近年ではニューラルネットワークに基づく機械翻訳 (NMT) が出現したことで、翻訳精度が大きく向上し、世の中への普及が進んでいる。

しかし、NMT を含めた機械翻訳では、語順などの構造が似ている言語間の翻訳に比べて、英語と日本語間などの構造が異なる言語間の翻訳は難しく、その翻訳精度の改善が大きな課題の一つとなっている。そこで構造が異なる言語間の翻訳精度を向上させるため、原言語 (翻訳元の言語) や目的言語 (翻訳先の言語) の文構造を活用する NMT が提案されている。

従来の文構造に基づく NMT では、依存構造や句構造という文構造が使われているが、どのような文構造が NMT にとって最適なものは明らかになっていない。また、従来の NMT で使われる原言語や目的言語の文構造は、その言語の構文解析器を用いて、翻訳相手となる文の構造とは独立に解析される。一方で、NMT 以前の統計的機械翻訳 (SMT) においては、同期式文脈自由文法に基づく文構造など、対訳文 (対訳関係にある原言語と目的言語の文対) の間で同期をとった文構造が有効であることが示されている。これらのことから、従来の文構造に基づく NMT で活用されている単一言語で解析された文構造は、NMT おける翻訳の手がかりとして最適であるとは限らない可能性があり、文構造に基づく NMT は改善の余地を残している。

2. 研究の目的

本研究では、文構造に基づく NMT の翻訳精度を向上させることを目的として、原言語の文構造と目的言語の文構造の間で対応をもたせた文構造を活用する NMT の実現を目指す。NMT 以前の SMT においては、同期式文脈自由文法に基づく文構造などの原言語と目的言語間で同期された文構造の有効性が示されていることから、NMT においても同期された文構造が翻訳性能の改善に寄与することが期待される。

これまで、同期された文構造を NMT で活用する試みは行われていない。そして、NMT において同期された文構造を有効に活用する方法は自明ではない。そこで本研究では、NMT において同期された文構造を活用する方法として、「NMT モデル内で文構造を同期させる方法」と「既存の同期式構文解析結果を NMT モデルで活用する方法」の二つを実現して翻訳性能を比較する。

構造が異なる言語間の翻訳実験として日英翻訳実験を行い、同期された文構造を NMT で活用することで、構造が異なる言語間の翻訳精度が向上するかどうかを明らかにする。また、NMT モデル内で文構造を同期させる方法と既存の同期式構文解析結果を NMT モデルで活用する方法のどちらが有効であるかを実験的に確認する。

3. 研究の方法

(1) 同期された文構造を活用する NMT の開発

本研究では、提案手法のベースラインモデルとして、現在様々な翻訳タスクで最高精度を達成している Transformer に基づく NMT [1] を採用し、このモデルを拡張することで 2 種類の提案手法を開発する。各手法の実装は Python の機械学習ライブラリ PyTorch を使って行う。

(2) 日英翻訳実験による評価

日英翻訳実験は、評価型ワークショップ Workshop on Asian Translation 2017 の Scientific paper Subtasks の日英翻訳タスクで行う。実験では、開発した 2 種類の提案手法の翻訳性能評価とともに、ベースラインモデルとして、文構造を利用しない通常の Transformer NMT [1] と、原言語と目的言語で独立に解析した文構造を用いる従来の文構造に基づく Transformer NMT [2] の翻訳性能も評価する。そして、それらの性能を比較することにより、同期された文構造を活用することで、構造が異なる言語間の翻訳精度が改善するかどうかを検証する。

4. 研究成果

(1) NMT モデル内で文構造を同期させる方法を創出した。提案手法は、参考文献[2]の文構造に基づく Transformer NMT モデルにおいて、言語間注意が捉える原言語と目的言語の対応関係を用いて、エンコーダ及びデコーダの自己注意が捉える原言語と目的言語の文構造が、言語間で整合性を持つように解析を行う。具体的には、エンコーダの自己注意 (原言語の文構造) を言語間注意 (原言語と目的言語の対応関係) に基づいて目的言語側の確率空間に変換し、目的言語側の空間に変換したエンコーダの自己注意とデコーダの自己注意 (目的言語の文構造) との差を NMT モデルの学習時の損失に加えて学習を行うことで実現した。提案手法の概要を図 1 に示す。この提案手法の具体的なアルゴリズムや方法は ACL/IJCNLP SRW 2021 にて発表した。

(2) 研究成果(1)に記載の提案手法の日英翻訳性能を評価し、文構造を利用しない場合及び原言語と目的言語で独立に解析した文構造を用いる場合の性能と比較した。実験結果を表 1 に示す。表 1 より、提案手法の BLEU (翻訳性能の評価指標) は、原言語と目的言語で独立に解析した文構造を用いる場合より 0.27 ポイント、文構造を用いない場合より 0.9 ポイント高いことが分かる。この実験結果により、Transformer NMT モデル内で文構造を同期させることで日英翻訳性能

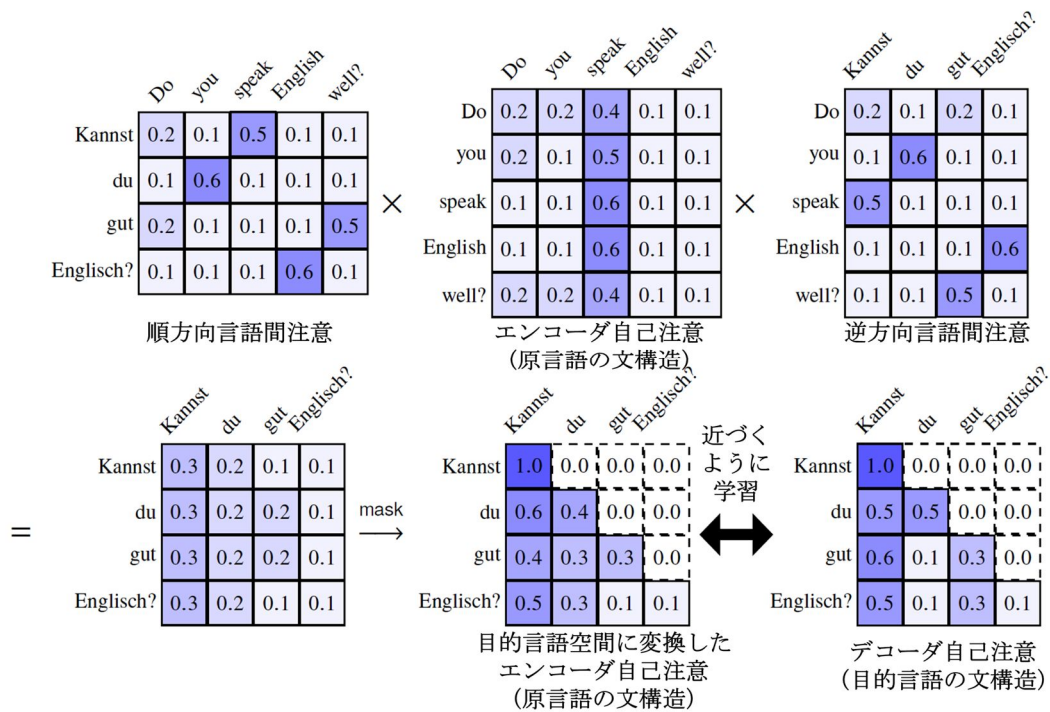


図1：NMTモデル内で文構造を同期させる方法の概要図

表1：NMTモデル内で文構造を同期させる方法の有効性（数値：BLEU (%)）

手法	日英	英独	英羅
文構造を用いないNMT [1]	28.94	27.23	23.83
従来の文構造を用いるNMT [2]	29.57	27.31	24.13
提案NMT	29.84	27.69	24.33

が改善することを実験的に示した。

また、提案手法の汎用性を確かめるために、WMT 2014の英独翻訳（英語からドイツ語への翻訳）タスク及びWMT 2016の英羅翻訳（英語からルーマニア語への翻訳）タスクでの性能比較も行った。その結果も表1に示す。表1より、英独翻訳においては、提案手法のBLEUは、原言語と目的言語で独立に解析した文構造を用いる場合より0.38ポイント、文構造を用いない場合より0.46ポイント高いことが分かる。また、英羅翻訳においては、提案手法のBLEUは、原言語と目的言語で独立に解析した文構造を用いる場合より0.2ポイント、文構造を用いない場合より0.5ポイント高いことが分かる。これらより、日英に限らず複数の言語対において文構造を同期させることで翻訳性能が改善することを実験的に示し、提案手法の有効性が汎用的であることを示した。

(3) 既存の同期式構文解析結果をNMTモデルで活用する方法を創出した。提案手法は、まず、NMTの教師データである対訳文から同期式文脈自由文法を学習し、その同期式文脈自由文法に基づき各教師データに対して同期された構文木を導出する。同期式文脈自由文法の学習と同期された構文木の導出は、文献[3]の方法を用いる。その後、同期された構文木の情報を取り込みながらNMTモデルを学習する。構文木の情報を取り込むNMTの学習方法は文献[4](文献[4]のTable 1中の12番目の設定)を用いる。具体的には、同期された構文木の原言語側の木及び目的言語側の木、それぞれに基づき、Neural Syntactic Distance(NSD)を算出し、NSDに基づくTransformer NMTを学習する。NSDに基づくTransformer NMTでは、原言語の文構造の情報は、原言語側の木のNSDをTransformerエンコーダの埋め込み層及び位置エンコーディングに組み込むことで考慮する。一方、目的言語の文構造の情報は、目的言語側の木のNSD系列を推定するようにTransformer NMTモデルを学習することで組み込む。

推論時には、まず、翻訳対象の文に対して、学習した同期式文脈自由文法を用いて構文木を導出し、導出した木のNSDを求める。そして、翻訳対象の文と導出した構文木のNSDを、NSDに基づくTransformer NMTに入力して翻訳を行うことで実現する。

(4) 研究成果(2)に記載の提案手法の日英翻訳性能を評価し、文構造を利用しない場の性能と比

較した。実験結果を表 2 に示す。表 2 より、提案手法の BLEU は、文構造を用いない場合より 0.07 ポイント高いことが分かる。しかし、この性能差をブートストラップによる検定手法 [5] により有意差水準 10% で検定した結果、有意ではなかった。このことから、研究成果 (2) に記載の既存の同期式構文解析結果を用いる NMT モデルは翻訳性能を悪化させるわけではないが、本実験設定の日英翻訳においては有効ではないことが分かった。

表 2: 既存の同期式構文解析結果を NMT モデルで活用する方法の有効性 (数値: BLEU (%))

手法	日英
文構造を用いない NMT [1]	28.94
提案 NMT	29.01

(5) NMT において同期された文構造を有効に活用する方法は自明ではない状況で、同期された文構造を活用する NMT として、NMT モデル内で文構造を同期させる方法と既存の同期式構文解析結果を NMT モデルで活用する方法を日英翻訳において比較した。表 1 と表 2 の結果より、既存の同期式構文解析結果を NMT モデルで活用する方法よりも、NMT モデル内で文構造を同期させる方法の方が有効であることを確認した。これは、既存の同期式構文解析結果を NMT モデルで活用する方法では同期された文構造が機械翻訳モデルとは独立に導出されるのに対して、NMT モデル内で文構造を同期させる方法では文構造と翻訳を同時に最適化するため、NMT により有用な文構造を活用できるからであると考えられる。

参考文献:

[1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention Is All You Need", In Proc. of NIPS 2017, pp. 5998-6008, 2017.

[2] Hiroyuki Deguchi, Akihiro Tamura, and Takashi Ninomiya. "Dependency-Based Self-Attention for Transformer NMT", In Proc. of RANLP 2019, pp. 239-246, 2019.

[3] Hao Zhang, Liang Huang, Daniel Gildea, and Kevin Knight. "Synchronous Binarization for Machine Translation", In Proc. of HLT-NAACL 2006, pp. 256-263, 2006.

[4] Chunpeng Ma, Akihiro Tamura, Masao Utiyama, Eiichiro Sumita, Tiejun Zhao. "Improving Neural Machine Translation with Neural Syntactic Distance", In Proc. of NAACL-HLT 2019, pp. 2032-2037, 2019.

[5] Philipp Koehn. "Statistical Significance Tests for Machine Translation Evaluation", In Proc. of EMNLP 2004, pp. 388-395, 2004.

5. 主な発表論文等

〔雑誌論文〕 計2件（うち査読付論文 2件／うち国際共著 0件／うちオープンアクセス 2件）

1. 著者名 Nishihara Tetsuro, Tamura Akihiro, Ninomiya Takashi, Omote Yutaro, Nakayama Hideki	4. 巻 28
2. 論文標題 Supervised Visual Attention for Multimodal Neural Machine Translation	5. 発行年 2021年
3. 雑誌名 Journal of Natural Language Processing	6. 最初と最後の頁 554 ~ 572
掲載論文のDOI (デジタルオブジェクト識別子) 10.5715/jnlp.28.554	査読の有無 有
オープンアクセス オープンアクセスとしている (また、その予定である)	国際共著 -

1. 著者名 Deguchi Hiroyuki, Utiyama Masao, Tamura Akihiro, Ninomiya Takashi, Sumita Eiichiro	4. 巻 28
2. 論文標題 Bilingual Subword Segmentation for Neural Machine Translation	5. 発行年 2021年
3. 雑誌名 Journal of Natural Language Processing	6. 最初と最後の頁 632 ~ 650
掲載論文のDOI (デジタルオブジェクト識別子) 10.5715/jnlp.28.632	査読の有無 有
オープンアクセス オープンアクセスとしている (また、その予定である)	国際共著 -

〔学会発表〕 計8件（うち招待講演 0件／うち国際学会 3件）

1. 発表者名 Tetsuro Nishihara, Akihiro Tamura, Takashi Ninomiya, Yutaro Omote, Hideki Nakayama
2. 発表標題 Supervised Visual Attention for Multimodal Neural Machine Translation
3. 学会等名 The 28th International Conference on Computational Linguistics (国際学会)
4. 発表年 2020年

1. 発表者名 Hiroyuki Deguchi, Masao Utiyama, Akihiro Tamura, Takashi Ninomiya, Eiichiro Sumita
2. 発表標題 Bilingual Subword Segmentation for Neural Machine Translation
3. 学会等名 The 28th International Conference on Computational Linguistics (国際学会)
4. 発表年 2020年

1. 発表者名 出口 祥之, 内山 将夫, 田村 晃裕, 二宮 崇, 隅田 英一郎
2. 発表標題 ニューラル機械翻訳のためのバイリンガルなサブワード分割
3. 学会等名 情報処理学会 第246回自然言語処理研究会
4. 発表年 2020年

1. 発表者名 岩本 裕司, 田村 晃裕, 二宮 崇
2. 発表標題 画像生成による疑似教師データを用いたマルチモーダル機械翻訳
3. 学会等名 情報処理学会 第246回自然言語処理研究会
4. 発表年 2020年

1. 発表者名 出口 祥之, 田村 晃裕, 二宮 崇
2. 発表標題 同期注意制約を与えた依存構造に基づくTransformer NMT
3. 学会等名 言語処理学会 第27回年次大会
4. 発表年 2021年

1. 発表者名 張 瀟廬, 二宮 崇, 田村 晃裕
2. 発表標題 ニューラル機械翻訳のためのアテンション確率のスムージングとゲーティング学習
3. 学会等名 言語処理学会 第27回年次大会
4. 発表年 2021年

1. 発表者名 岩本 裕司, 田村 晃裕, 二宮 崇
2. 発表標題 画像生成による疑似教師データを用いたマルチモーダルニューラル機械翻訳
3. 学会等名 言語処理学会 第27回年次大会
4. 発表年 2021年

1. 発表者名 Hiroyuki Deguchi, Akihiro Tamura, Takashi Ninomiya
2. 発表標題 Synchronous Syntactic Attention for Transformer Neural Machine Translation
3. 学会等名 The 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop (国際学会)
4. 発表年 2021年

〔図書〕 計0件

〔産業財産権〕

〔その他〕

-

6. 研究組織

氏名 (ローマ字氏名) (研究者番号)	所属研究機関・部局・職 (機関番号)	備考
---------------------------	-----------------------	----

7. 科研費を使用して開催した国際研究集会

〔国際研究集会〕 計0件

8. 本研究に関連して実施した国際共同研究の実施状況

共同研究相手国	相手方研究機関
---------	---------