

特許請求項翻訳における 単語対応に基づく節分割モデルの有効性

西村 柁人¹ 宇津呂 武仁¹ 永田 昌明²

¹筑波大学大学院 システム情報工学研究群 ²NTT コミュニケーション科学基礎研究所
s2320779@u.tsukuba.ac.jp utsuro@iit.tsukuba.ac.jp
masaaki.nagata@ntt.com

概要

特許請求項は特許範囲を規定する重要部分だが、その長さや独特の書式が原因で NMT モデルでは訳抜けや繰り返しといった誤訳を引き起こしやすい。本論文では、この課題を解決するために、節分割モデルを用いた分割統治翻訳手法を提案する。翻訳元文を節分割モデルを用いて節に分割し、それぞれの節を節翻訳モデルで翻訳、その後並び替え・編集モデルで最終的な翻訳文を生成する。さらに、節分割モデルと節翻訳モデルを訓練するための節単位の対訳コーパスを、単語対応情報をもとに文単位の対訳コーパスから作成する手法を提案する。実験では、提案手法が通常モデルを BLEU で上回り、訳抜けや繰り返しの改善を確認した。

1 はじめに

特許文書における特許請求項は、特許の権利範囲を定義する上で極めて重要な部分である。しかし、その文の長さや独特の記述形式のため、ニューラル機械翻訳 (NMT) モデルを用いた翻訳では、訳抜けや繰り返しといった問題が発生しやすい。長文翻訳におけるこれらの問題に対して加納ら [3] は、英日翻訳において、入力文を構文解析に基づいた節に分割して翻訳後に再構成する「分割統治ニューラル機械翻訳」手法を提案した。ここで、この手法は翻訳精度向上の可能性を示したものの、その一方で、適切な節の単位の選定や、節翻訳後の再構成精度に課題が残っていた。

これに対し石川ら [1] は、文献 [3] で課題とされていた、節分割単位の課題と節分割後の節の翻訳精度の二つの課題に対して、接続詞による節分割の採用と、節翻訳モデルと並び替え・編集モデルの2つに事前学習済みモデルの mBART [5] を用いること

により翻訳精度の向上を試みた。それに加えて、節翻訳モデルを節単位の疑似対訳データでファインチューニングすることによって、節の翻訳精度改善を試みた。実験では、過剰な長さの訳出が大きく減少し、ハルシネーションや繰り返しの抑制を確認している。このように、分割統治翻訳手法を用いることで、長文翻訳において従来の NMT モデルが引き起こしやすい翻訳ミスを抑止できる可能性が示唆されている。しかし、文献 [1] における課題として、節の分割単位が接続詞のため適切に短い節に分割できない長文が存在する点と、節翻訳モデル訓練時に疑似対訳データを用いており、実際に収集された実対訳データを用いていない点が挙げられる。

以上をふまえて、本論文では、特に特許請求項日英翻訳を対象とする分割統治翻訳手法において、先行研究 [3, 1] とは異なるアプローチを提案する。本論文の手法においては、まず、翻訳元の特許請求項文を翻訳モデルが翻訳しやすい単位の節に分割する節分割モデルを提案する。特に、本論文では、対訳文中の単語対応に基づいて設定された節単位が二言語間で整合することを条件とすることにより、最終的な翻訳文で訳抜けや繰り返しといった誤訳の発生を抑え、より正確な翻訳文を生成することを実現した。具体的な手法としては、単語アライメントツールを用いて、元の対訳コーパスから節単位の対訳データを高品質に生成する手法、および、作成した節単位の対訳コーパスを用いて以下の3つのモデルを構築する手法を提案する：

1. 翻訳元の日本語文を節単位に分割する節分割モデル。
2. 節翻訳に特化した節翻訳モデル。
3. 翻訳された節を並び替え・編集し、最終的な翻訳文を生成する並び替え・編集モデル。

実験として、これらのモデルを統合した翻訳手

法を、日英特許対訳コーパス JaParaPat [6] を用いて評価した。その結果、特許請求項の日英翻訳において、提案手法は通常の NMT モデルと比較して、BLEU において統計的に有意な改善を示した。また、通常モデルによる翻訳結果と比較して、提案手法において訳抜けの発生を抑制できていることを確認し、より正確な翻訳結果を得られることが示された。

2 関連研究

長文翻訳の課題に対しては、従来からさまざまなアプローチが検討されてきた。Sudoh ら [11] は、長い文を翻訳するため、統計的機械翻訳に分割統制的翻訳手法を用いた。入力文を構文解析に基づいた節の単位に区切って翻訳し、その結果を節の階層構造に基づいて並べ替えることで、精度向上を実現させた。

NMT において、長文を節に分割し、分割された節ごとに翻訳した後に、前から順に結合し直す手法として Pouget-Abadie ら [10] の自動分割の手法がある。この手法は RNN を用いて、長文を分割して翻訳する際の最適な位置を予測し、モデルが翻訳し易いように入力文を分割する手法である。しかし、この手法は英仏翻訳を対象としており、本論文で行う日英翻訳のような語順が大きく異なる言語対においては、翻訳後の再構成において、不自然な語順が生じやすいという課題が残されていた。この課題に対し、加納ら [3] は、英日翻訳において、長文を分割して翻訳した後にそれぞれの節を適切な順序に並び替えるニューラルネットワークモデルを作成した。さらに石川ら [2] は、等位接続詞を基準とした新しい英文の節の分割単位及び、分割後の節翻訳において文内コンテキストを参照して翻訳を行うモデルの学習手法を提案し、長文の英日翻訳において翻訳精度向上を達成した。

3 提案手法

図 1 に、提案手法全体の構造を示す。翻訳元の日本語特許請求項文は、まず節分割モデルによって複数の節に分割される。その後、各節が節翻訳モデルによって翻訳され、最後に並び替え・編集モデルによって統合され、翻訳後の英語特許請求項が生成される。本手法により、従来の NMT モデルで発生しやすい訳抜け・繰り返しを抑制することを目指す。

3.1 節単位の対訳コーパス

本論文では、Zhang ら [15] の長文の対訳データから対訳部分文を生成する手法を参考に、節単位の対訳コーパスを自動的に生成する手法を提案する。この手法では、文単位の対訳コーパスに対して、単語アライメントツールの WSPAlign [14] を用いて単語対応情報を取得する。単語対応情報を基に、文中の対応する節を抽出し、節単位の対訳データを生成する。節単位の対訳コーパスは以下のような手順で作成する。本論文では文献 [15] の報告を基に、単語含有割合の閾値を 0.5 と設定した。

1. 特許対訳データの対訳文に対して WSPAlign を使用することで単語対応を取得する。
2. 日英の対訳文それぞれに対して、「、」「,」「。」「.」「;」「:」などの区切り記号の位置で複数の節に分割する。
3. 単語対応情報から、対訳節ごとの単語含有割合を計算し、割合が 0.5 を超えた場合、対象の節は対応関係があると判定する。
4. 節対応が 1 対多、多対多となる場合は、1 対 1 になるように複数節側を 1 つにまとめる。

上記の手順を特許対訳コーパスから抽出した対訳文に適用することで節対訳コーパスを作成する

3.2 節分割モデル

本論文では、Wicks と Post [13] が提案した文分割モデル ERSATZ を節分割モデルの基盤とし、このモデルを節単位の分割に適用するように訓練を行うことによって節分割モデルを作成した。ERSATZ は、文分割をバイナリ分類タスクとして定式化しており、文分割候補となる句点(「。」や「.」など)に対して「文中」か「文末」かを予測する。本論文では、この文分割を節分割に拡張するために、節分割候補となる読点(「、」や「,」など)を用いて節分割モデルを作成する¹⁾。モデルの訓練には 3.1 節で提案した節単位の対訳コーパスを用いた。

訓練に使用するデータは、節対訳コーパス内の日本語節を文単位で抽出し、節の分割位置の句読点に文末ラベルを付与することで正解データを作成した。これにより、日本語特許請求項文に対して単語対応情報をもとにした節分割が行えるモデルを作成

1) 実際には、補足説明を表す () 中の文分割位置において節分割を行うために、読点(「、」や「,」など)だけでなく句点(「。」や「.」など)もあわせて用いる。

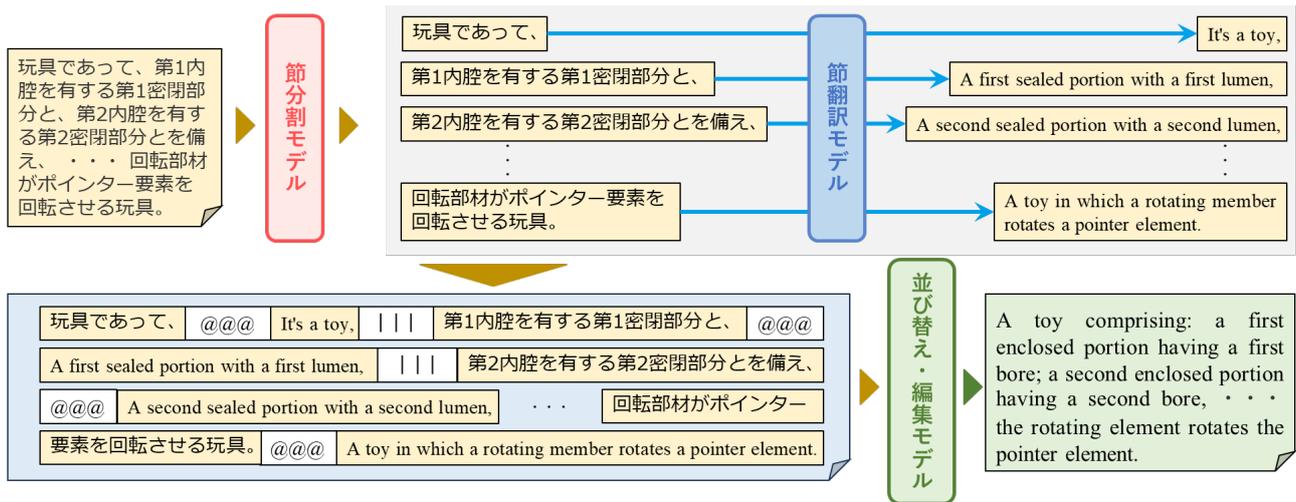


図1 提案手法の流れ

した。

3.3 節翻訳モデル

本論文では、節翻訳に特化した節翻訳モデルを作成するために、事前に特許対訳コーパスから作成した日英翻訳モデルを、3.1節の手法により作成された節単位の対訳コーパスでファインチューニングした。この節翻訳モデルは、分割された節を翻訳する際に分割によって失われた文脈を勝手に補おうとする事象を抑制し、より正確な翻訳を実現することを目的としている。

3.4 並び替え・編集モデル

並び替え・編集モデルを用いる理由としては、節翻訳モデルによって翻訳した複数の節を目的言語文である英語の一つの文に戻すためである。日本語と英語の語順は大きく異なることから、日本語文を節ごとに分割して接続するのみでは、日本語と英語間での語順の変化に対応することができない。そのため、節翻訳モデルによって得られた、節に分割された翻訳節を接続する際に並び替え・編集モデルを用いることで、適切な語順に並び替えることを期待している。

並び替え・編集モデルの訓練データは、コーパスの日本語文を節分割モデルで分割し、日本語の節とその翻訳結果の英語節を繋げたものを入力データ、コーパスの英語文を教師データとして利用する。モデルの語彙の中に特殊トークンとして“@@@”と“|||”を登録しておく。“@@@”は日本語の節とその翻訳結果の英語節を繋げるために用い、“|||”はそれらの節のペアを接続するために用いる。

入力データをこのような形にしている理由としては、翻訳結果の英語節だけを入力データとして用いると、翻訳された各節どうしがどのような関係性をもつかという情報が欠落してしまうため、日本語文の情報を付与するためである。

並べ替え・編集モデルの入力は日英両方の単語が含まれた文となるため、日英両方の理解力が必要となる。そのため並び替え・編集モデルは日英双方向の翻訳モデルを上記の手順で作成した訓練データでファインチューニングすることによって作成する。この事前訓練がモデルの性能に与える影響については、付録Bに示す。

4 実験

4.1 実験設定

本論文では、日英特許対訳コーパスのJaParaPat [6]を基に、日英翻訳の実験を行った。使用したデータは、2016年から2020年までの特許全文の対訳データを訓練データ、2021の前半までの特許請求項対訳データをテストデータとした。

機械翻訳ソフトウェアはfairseq [7]を使用し、ベースラインモデル、節翻訳モデル、並び替え・編集モデルのモデルアーキテクチャにTransformer big [12]を使用した。トークナイズにはsentencepiece [4]を使用した。特許対訳データから10M文対をランダムにサンプリングしてモデルの訓練を行った。語彙数は日英ともに32Kとした。また、節分割モデルの訓練にはERSATZ²⁾を利用した。

ベースラインモデルと提案手法の3種類のモデ

2) <https://github.com/rewicks/ersatz>

表 1 提案手法およびベースラインモデルに使用したデータの概要

モデル	使用データ	データ数
ベースラインモデル	JaParaPat2016~2020	61,364,685 文対
節分割モデル	節対訳コーパス (claims)	200,462 文
節翻訳モデル	JaParaPat2016~2019	49,474,547 文対
	節対訳コーパス	5,480,682 節対
並び替え・編集モデル	JaParaPat2016~2020(双方向)	109,028,682 文対
	JaParaPat2016~2020(claims)	2,613,107 文対

表 2 各評価対象の文数および BLEU

評価対象	文数	ベースライン	提案手法
全体	238,902	55.5	56.6**
長文	33,959	50.1	51.6**

ルの訓練に使用したデータの概要を表 1 に示す。節単位の対訳コーパスは、2020 年のデータの半数 (5,976,295 文対) に対して WSPAlign³⁾ を使用して単語対応情報を取得し、3.1 節の手法に沿って作成した。これによって 5,480,682 対の節対訳データを持つ節単位の対訳コーパスが作成された。節分割モデルの訓練には節単位の対訳コーパスから、特許請求項文を分割して作成した日本語節データを使用した。節翻訳モデルは、2016 年から 2019 年の特許全文の対訳データで事前訓練し、節単位の対訳コーパスの全データを使用してファインチューニングすることで作成した。並び替え・編集モデルは、2016 年から 2020 年の特許全文の双方向の対訳データで事前訓練し、2016 年から 2020 年の特許請求項の日本語文を、節分割モデルと節翻訳モデルを用いて 3.4 節の手法で作成した訓練データでファインチューニングして作成した。

評価指標として、BLEU [8] を使用した。BLEU は sacreBLEU⁴⁾ [9] を用いて計測した。特許翻訳では専門用語を正しく訳出できることが重視されるため、本論文では BLUE を主たる評価尺度とした。

4.2 結果

本論文では、2021 年の特許データから特許請求項のみを抽出したテストデータ (238,902 文) を用いて提案手法の性能を評価した。その結果、提案手法は BLEU において 56.6 を記録し、ベースラインモデルの 55.5 を統計的に有意に上回った ($p < 0.01$)。これにより、提案手法が全体的な翻訳精度を向上させることが確認された。

3) <https://github.com/qiyuw/WSPAlign>

4) <https://github.com/mjpost/sacrebleu>

また、テストセットの中で翻訳元日本語文のサブワードトークンが 100 トークン以上の長文のみを抽出したサブセットに対しても性能を評価した。結果としては、提案手法は 51.6 となり、ベースラインモデルの 50.1 を統計的に有意に上回った。また、テストセット全体では BLEU の上昇幅が 1.1 ポイントであったが、長文においては 1.5 ポイントの上昇が見られ、提案手法は特に長文において顕著な改善を示した。

さらに、提案手法が訳抜けや繰り返しといった翻訳ミスを抑制できるかを分析するため、ベースラインモデルおよび提案モデルの訳出文と参照訳文の文長比に着目し、その傾向を観察した。テストセット全体の分析結果では、訳抜けが発生する可能性が高い文の数は、提案手法で 515 文であるのに対し、ベースラインモデルでは 900 文であった。この結果から、提案手法はベースラインモデルに比べて訳抜けを効果的に抑制できていると考えられる。一方で、繰り返しに関しては、提案手法で発生した文が 376 文であるのに対し、ベースラインモデルでは 273 文であった。このことは、ベースラインモデルによる翻訳の方が繰り返しの発生が少ない可能性を示唆している。繰り返しの多くは、並び替え・編集モデルで発生している可能性があり、並び替え・編集モデルにおける繰り返し抑制を行う必要がある。これらの分類結果および詳細な分析については、付録 C に記載する。

また、実際の提案手法による訳抜け・繰り返しの翻訳改善例を表 3 に示す。図中の例のようにベースラインモデルで訳抜け・繰り返しが起きている文に対して、提案手法ではそれらの翻訳ミスが改善され、より正確な翻訳が行えることが確認できた。

5 おわりに

本論文では、特許請求項の日英翻訳における課題である「訳抜け」や「繰り返し」といった翻訳ミスを解決するため、特許請求項が長文や独特な構造を持つことに着目し、特許請求項をより翻訳しやすい単位に分割する節分割モデルを用いた翻訳手法の提案を行った。実験の結果、提案手法は BLEU においてベースラインモデルを統計的に有意に上回る性能を示し、訳抜けの発生を抑制できていることを確認した。これにより、提案手法が特許請求項の翻訳においてより正確な翻訳が可能であることが示された。

参考文献

- [1] 石川隆太, 加納保昌, 須藤克仁, 中村哲. 事前学習モデルによる分割統治ニューラル機械翻訳. 言語処理学会第 29 回年次大会論文集, pp. 1451–1456, 2023.
- [2] 石川隆太, 加納保昌, 須藤克仁, 中村哲. 文内コンテキストを利用した分割統治ニューラル機械翻訳. 言語処理学会第 30 回年次大会論文集, pp. 2342–2347, 2024.
- [3] 加納保昌, 須藤克仁, 中村哲. 分割統治的ニューラル機械翻訳. 言語処理学会第 27 回年次大会論文集, pp. 148–153, 2021.
- [4] T. Kudo and J. Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In **Proc. EMNLP**, pp. 66–71, 2018.
- [5] Y. Liu, et al. Multilingual denoising pre-training for neural machine translation. **TACL**, Vol. 8, pp. 726–742, 2020.
- [6] M. Nagata, M. Morishita, K. Chousa, and N. Yasuda. JaParaPat: A large-scale Japanese-English parallel patent application corpus. In **Proc. LREC-COLING**, pp. 9452–9462, 2024.
- [7] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli. fairseq: A fast, extensible toolkit for sequence modeling. In **Proc. NAACL**, pp. 48–53, 2019.
- [8] K. Papineni, S. Roukos, T. Ward, and W. Zhu. Bleu: a method for automatic evaluation of machine translation. In **Proc. 40th ACL**, pp. 311–318, 2002.
- [9] M. Post. A call for clarity in reporting BLEU scores. In **Proc. 3rd WMT**, pp. 186–191, 2018.
- [10] J. Pouget-Abadie, D. Bahdanau, B. van Merriën-Boer, K. Cho, and Y. Bengio. Overcoming the curse of sentence length for neural machine translation using automatic segmentation. In **Proc. 8th SSST**, pp. 78–85, 2014.
- [11] K. Sudoh, K. Duh, H. Tsukada, T. Hirao, and M. Nagata. Divide and translate: Improving long distance reordering in statistical machine translation. In **Proc. 5th WMT**, pp. 418–427, 2010.
- [12] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In **Proc. 30th NIPS**, pp. 5998–6008, 2017.
- [13] R. Wicks and M. Post. A unified approach to sentence segmentation of punctuated text in many languages. In **Proc. 59th ACL**, pp. 3995–4007, 2021.
- [14] Q. Wu, M. Nagata, and Y. Tsuruoka. WSPAlign: Word alignment pre-training via large-scale weakly supervised span prediction. In **Proc. 61st ACL**, pp. 11084–11099, 2023.
- [15] J. Zhang and T. Matsumoto. Corpus Augmentation for Neural Machine Translation with Chinese-Japanese Parallel Corpora. **Applied Sciences**, Vol. 9, No. 10, 2036, 2019.

