

ニューラル機械翻訳による特許機械翻訳システムの開発

Development of Neural Machine Translation System for Patent Translation.

東芝デジタルソリューションズ株式会社

園尾 聡

2009年九州工業大学生命体工学研究科博士課程修了。博士(工学)。同年株式会社東芝入社(2017年より現所属)。機械翻訳、自然言語処理の研究開発に従事。

✉ satoshi.sonoo@toshiba.co.jp

1 はじめに

知財戦略のグローバル化に伴い、特許文献を翻訳する頻度・件数は急増している。一方で、人手による翻訳作業では、時間がかかる、コストが高い、大量の翻訳が難しいなどの課題があり、これらの課題を解決する手段として、機械翻訳の活用が進んでいる。

機械翻訳技術は、従来のルールベース機械翻訳(Rule-Based Machine Translation; RBMT)、統計的機械翻訳(Statistical Machine Translation; SMT)を経て、近年、ニューラルネットワークによるニューラル機械翻訳(Neural Machine Translation; NMT)が主流となってきた。NMTのモデルとして、注意機構(アテンション)を用いたEncoder-Decoderモデル^[1]やTransformerモデル^[2]が提案されており、どのモデルも大規模な対訳データからモデルを学習し、SMTを大きく上回る正確さと流暢さが実現可能であることが報告されている。

このような背景の中、特許庁では、特許情報プラットフォーム(J-PlatPat; JPP)等で利用される機械翻訳システムを刷新し、NMTをベースとした機械翻訳サービスの提供を2019年5月¹に開始した^[3]。この新しい機械翻訳システムでは、国立研究開発法人情報通信研究機構(NICT)が開発したNMTを大規模クラウドプラットフォーム上に構築し、NMT特有の課題に対して

様々な前後処理を適用することで高精度な機械翻訳を実現している^[4]。

本稿では、特許翻訳におけるNMTの実用化を目指して開発された、機械翻訳システムの概要および特徴について紹介する。

2 機械翻訳システム

2.1 システム概要

本機械翻訳システムは、以下の機械翻訳サービスを提供するためのシステムである。

1) 特許情報プラットフォーム(JPP) 審査・審判書類² 情報翻訳

JPPからの翻訳要求に応じて、審査・審判書類の日英翻訳を行い、翻訳結果を返却する。

2) 日本公報情報等翻訳

JPPからの翻訳要求に応じて、日本公報や出願人名情報等の日英翻訳を行い、翻訳結果を返却する。

3) ワン・ポータル・ドシエ(OPD) 審査・審判書類 情報翻訳

OPDからの翻訳要求に応じて、審査・審判書類の日英翻訳を行い、翻訳結果を返却する。

4) 中韓文献翻訳文作成

特許庁内システムからの翻訳要求に応じて、中韓文献の中日翻訳・韓日翻訳を行い、翻訳結果を返却する。翻訳

1 2019年5月、審査・審判書類・日本公報の日英翻訳機能をリリース。2020年度、中国公報の中日翻訳機能、韓国公報の韓日翻訳機能をリリース予定。

2 出願手続において審査・審判官等と出願人・代理人等との間で交わされる文書。

結果は DB に蓄積され、中韓文献の日本語検索サービスに利用される。

翻訳対象となる特許文献は、特許出願に係る審査・審判書類、日本公報、中国公報、韓国公報と多岐にわたり、想定される利用ユーザーも、国内一般利用者、国内審査官、海外審査官と様々である。これに対して、本機械翻訳システムでは、各翻訳要求に応じたインターフェースを提供し、各インターフェースで処理される特許文献の特性を考慮し、NMT をはじめとした複数の翻訳エンジンと翻訳前後処理を組み合わせ、翻訳処理を行う。本機械翻訳システムの詳細について以降で紹介する。

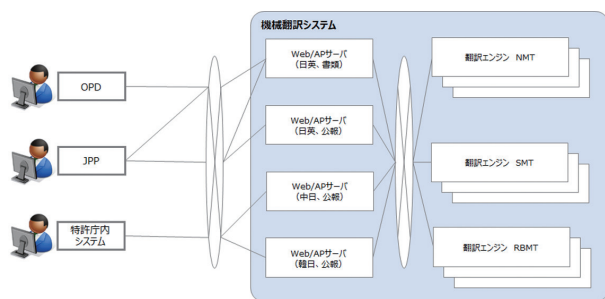


図1 機械翻訳システムの概要

2.2 システム構成

2.2.1 翻訳エンジン

本機械翻訳システムで利用している翻訳エンジンとその主な翻訳対象を表1に示す。NMTは、NICTが世界最大規模の特許対訳データから構築した特許文書向けニューラル機械翻訳エンジンを採用し^[5]、SMTも同様にNICTが開発した特許文書向け統計翻訳エンジンである。NMTは、流暢で高精度であるという特長を最大限に活かし、内容の理解度が求められる特許公報³（要約、

3 本稿では、公開特許公報を含む公報全般を指す。

請求項、実施形態等）が主な翻訳対象となる。さらに、国内審査官が審査した審査結果を海外審査官に正確に伝えるために、特許NMTモデル（公報版）に対して審査結果の対訳データでカスタマイズ（ドメイン適応）した特許NMTモデル（審査結果版）を審査結果文の翻訳に適用する。また、言語類似性から従来から高い翻訳精度を実現していた韓日方向および、原文が比較的短文かつ名詞句としての訳出が必要となる発明の名称許の翻訳については、特許SMTモデルを利用する。

一方で、NMTは、語彙数の制約から出願人名や住所などの固有名詞を網羅的にカバーすることは難しく、また学習時の対訳データとドメインが異なる文書（特許公報以外の公報）については翻訳精度が劣化し、訳抜けや湧き出しといったNMT特有の誤訳が生じる恐れがある。このため、書誌情報および意匠・商標・審決公報の翻訳に関しては、大規模な固有名詞、専門用語辞書が利用可能で、特許翻訳において長年の実績のあるRBMT（The 翻訳エンタープライズTM^[6]）を採用している。

2.2.2 翻訳フロー

本機械翻訳システムでは、翻訳対象となる特許文献に対して、公報フロー、審査結果フロー、固定表現フローといった複数の翻訳フローが用意されている。翻訳フローとは、対象となる原文（段落）と、その翻訳を実行する翻訳エンジンおよび翻訳前後処理の組み合わせを定義したものである。例えば、公報フローでは、特許NMT（公報版）をベースとし、後述する翻訳前後処理を適用するとともに、NMTで誤訳を検出した場合はRBMTで訳出するという枠組みとなっている。

本機械翻訳システムの特徴として、1つの特許文献（審査・審判書類）に対して複数の翻訳フローを用いて翻訳

表1 機械翻訳システムで利用する翻訳エンジンとその主な翻訳対象

エンジン名称	方式	翻訳方向	主な翻訳対象
特許NMT（公報版）	NMT	日英	特許公報
特許NMT（審査結果版）	NMT	日英	審査結果
特許SMT	SMT	日英	特許公報（発明の名称）
RBMT	RBMT	日英	書誌、意匠・商標・審決公報
特許NMT（中日）	NMT	中日	特許公報
特許SMT（中日）	SMT	中日	特許公報（発明の名称）
特許SMT（韓日）	SMT	韓日	特許公報

することが挙げられる。例えば、重要書類の1つである拒絶理由通知書は、図2に例示されるように出願番号や出願人名等、および定型文が記述される書類ヘッダー部、審査結果文が記述される審査結果部、および注意事項等の定型文と審査官名が記述される書類フッター部から構成されている。これに対し、書類ヘッダー部および書類フッター部は、RBMTを主エンジンとする固定表現フローへ入力し、審査結果部は、特許NMT（審査結果版）をベースとしている審査結果フローへと入力する。各構成要素の切り分けは、キーワードによる分割ルールを適用することで実現している。これにより、原文の特性に適した翻訳エンジンが選択され、種々の情報を含む複雑な書類であっても高精度に翻訳することが可能となる。

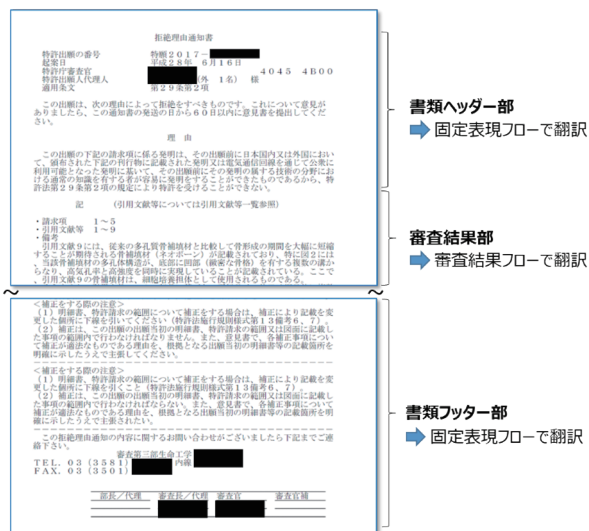


図2 審査・審判書類（拒絶理由通知書）の構成例

2.3 翻訳前後処理

大規模な特許対訳データから学習したNMTは、従来の翻訳エンジンに比べて非常に高い精度を実現しているが、実用化に向けて以下のような課題も挙げられる。

1) 化合物や数式の正確な訳出

特許文書に出現する化合物表現、DNA配列、数式等の記号列の正確な訳出を担保することが困難。

2) 訳抜け・湧き出し

原文にある単語情報が抜け落ちてしまう訳抜け (Under-translation) や、原文にない単語情報が意図せず訳出されてしまう湧き出し (Over-translation)

が発生する。

3) 繰り返し出力

湧き出しの一部であり、意図しない単語 (列) が繰り返されて訳出される現象。リカレントな構造を持つモデルに起因する誤訳であり、抜本的な解決が難しい。

特に、1) および2) に関しては、訳文自体は非常に自然な文となっているので誤訳となっていることに気づきにくく、間違った情報伝達に繋がり、特許翻訳において大きな課題となる。これらの課題を解決するために本機械翻訳システムで開発した翻訳前後処理について以降で説明する。

2.3.1 固有表現退避

NMTでは、図3に示すように、局所的に数字や記号が抜け落ちる (または湧き出す) 誤訳が発生する。これを回避するために固有表現待避処理を行う。

固有表現 (化合物表現) 退避処理の概要を図4に示す。まず、原文中の化合物表現や数式をパターンマッチによって抽出し、抽出された固有表現を特殊タグへと置き換える。この際、化合物表現であれば、化学物質名は対象言語に逐次翻訳 (リバビリン → ribavirin、リボフラノシル → ribofuranosyl、トリアゾール → triazole、カルボキサミド → carboxamide) し、記号列であれば全角半角変換をすることで対象言語の文字列へと正規化する。置き換え後の原文をNMTによって翻訳したのち、訳文に現れる特殊タグを正規化された文字列で復元させる。このような処理により、化合物や数式を訳抜け・湧き出しを軽減し、正確な訳出を実現する。

[原文]リバビリン (1-B-D-リボフランシル-1-H-1,2,4-トリアゾール-3-カルボキサミド) は、インソシン5'-リン酸ヒドロゲナーゼ (IMPDH) の阻害剤であり、HCVの治療においてIFN-αの有効性を増強する。

[NMT訳例] Ribavirin (1-B-D-ribofuranosyl-2H-1-1-triazole-2-carboxamide) is an inhibitor of Inosine 5'-monophosphate dehydrogenase (IMPDH) and enhances the efficacy of IFN-α in the treatment of HCV. 3. 1. 2. 4.

図3 化合物表現に対する誤訳例

た。

3.2 人手評価

本機械翻訳システムの運用において、自動評価に加えて、人手評価を定期的実施する。人手評価の対象は、実際の翻訳要求からサンプリングによって選定し、特許文献機械翻訳の品質評価基準^[11]に基づいた以下の観点で評価を行う。

- 1) 情報の過不足がない適切な内容伝達
- 2) 技術用語・法律用語の最適な翻訳
- 3) 翻訳文の流暢さ
- 4) 文法・構文の適切な翻訳
- 5) その他（発明の名称や請求項特有の表現、化合物や数式の表現、箇条書き等）

人手評価結果をフィードバックし、リリース後も継続して翻訳精度の維持・向上を実施していく予定である。

4 おわりに

本稿では、特許翻訳における NMT の実用化を目指して開発した、機械翻訳システムについて紹介した。本機械翻訳システムは、特許翻訳に特化した NMT をベースとし、複数の翻訳エンジンおよび前後処理を組み合わせることで、特許文献の高精度な機械翻訳システムを実現した。

2019 年 5 月より本機械翻訳システムの日英翻訳機能がリリースされ、これまで⁸⁾に 8 千万件以上の特許文献が翻訳されている。2020 年度には、中日翻訳機能、韓日翻訳機能もリリース予定である。

今後も、利便性の高い機械翻訳サービスが提供できるよう、翻訳精度の維持・向上の改善プロセスを継続していく予定である。

参考文献

- [1] Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." arXiv preprint arXiv:1409.0473 (2014).
- [2] Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems. 2017.
- [3] 古田敦浩, 特許情報普及活用施策に関する最近の取組と今後の展開, Japio YEAR BOOK 2019 (2019) : 66-71.
- [4] 特許庁から受注した機械翻訳システムの稼働開始について,
<https://www.toshiba-sol.co.jp/news/detail/20190531.htm>
- [5] 特許庁“次期機械翻訳サービス”の中核技術として NICT の技術が採用,
<https://www.nict.go.jp/press/2018/07/10-1.html>
- [6] https://www.toshiba-sol.co.jp/pro/hon_yaku/seihin/server/index_j.htm
- [7] Mi, Haitao, et al. "Coverage embedding models for neural machine translation." arXiv preprint arXiv:1605.03148 (2016) .
- [8] 後藤功雄, and 田中英輝. "ニューラル機械翻訳での訳抜けした内容の検出." 自然言語処理 25.5 (2018) : 577-597.
- [9] Papineni, Kishore, et al. "BLEU: a method for automatic evaluation of machine translation." Proceedings of the 40th annual meeting on association for computational linguistics. Association for Computational Linguistics, 2002.
- [10] 平尾努, et al. "語順の相関に基づく機械翻訳の自動評価法." 自然言語処理 21.3 (2014) : 421-444.
- [11] 特許文献機械翻訳の品質評価手順について,
https://www.jpo.go.jp/system/laws/sesaku/kikaihonyaku/tokkyohonyaku_hyouka.html

8 2019 年 8 月末現在



4

機械翻訳技術の向上

