

産業翻訳に役立つ自然言語処理技術についての議論の足場

藤田 篤 山田 優 影浦 峽

情報通信研究機構 関西大学 東京大学

1 テーマセッションの開催趣旨

国内の翻訳需要の大半を占めているのは、特許文書や契約文書、製品のマニュアル等、産業活動に関わる文書の翻訳、いわゆる産業翻訳である。種々の機械翻訳(MT)サービスの出現を契機として、産業翻訳の現場におけるMTの活用が活発化しつつあるが、実践に大きく寄与する自然言語処理技術は、MTだけに限られない。例えば、翻訳の初期工程においては対象分野(対象の起点テキスト)の用語集を作成する技術等が、また目標テキストの検査においてはテキスト診断技術等が有用であろう。翻訳作業環境においては、自動校正のほか、統語構造やテキスト間の部分対応等の情報の可視化技術も人間の作業の支援に役立つ。

我々は、産業翻訳への実応用を見据えた自然言語処理技術について、技術を利用する際の条件の見極めや、新たな技術の創出の機会を提供することを目的として、テーマセッション『産業翻訳に役立つ自然言語処理』を企画した。本稿では、本テーマセッションにおける議論の足場とすべく、産業翻訳に関わる概念(2節, 3節), 有用さの基準(4節), 既存の応用例(5節), および自然言語処理技術の産業翻訳応用をめぐる主な論点(6節)を整理しておく。前稿[4, 7]も参考にされたい。

2 用語の整理

翻訳には様々な概念や要素がかかわっている。議論に際し、産業翻訳の関係者と自然言語処理の研究者・技術者の間で齟齬が生じないように、3節で述べる産業翻訳のワークフローに関わるものを中心に、主要な用語を列挙し、各々に関する簡単な説明を表1にまとめた。

まず、「翻訳」が、「原文(起点テキスト)」ではなく、原文(起点テキスト)によって表現された「内容」を扱う¹、という点を確認しておく。Vermeerは、「翻訳するということは、目標文化社会の環境にある目標テキストの受け手と目標テキストの目的のために、目標テキスト受容の状況下でテキストを産出することである」と述べている[12, 1]。これに対して、既存のMT研究のほとんど²において、計算機処理の入力は言語表現であり、かつ単一の文である。

¹岩波国語辞典第7版[10]の語釈文を用いたが、他のいくつかの国語辞典においても同様の記述がなされている。

²文脈[13]や他モジュール[3]を同時に参照する研究も存在する。

次に、表中の「翻訳関連の概念」を確認されたい。既存のMT研究においては、これらは明示的には扱われていない。強いて言えば、モデルの学習用の対訳コーパスにそれらの性質が反映されていると仮定し、対訳コーパスを工夫することで間接的に考慮される場合があるということになる。

以上をふまえると、既存のMTが行っている処理は、翻訳ではなく、X文Y訳[7]であると言えよう。

3 産業翻訳のワークフロー

産業翻訳における大まかなワークフロー³を表2に示す。翻訳の依頼者が翻訳を求めるのには、明確な目的(スコープ)がある。産業翻訳では、その目的を達成するように翻訳戦略が定められ、テキストタイプやレジスタから導かれる翻訳規範、依頼者が求める要件をまとめたブリーフに沿うかたちで翻訳が行われる。

産業翻訳は、出版翻訳よりも極めて短い時間で翻訳を終えることが要求される。このため、複数人からなるプロジェクトで起点テキスト(あるいは複数の起点テキスト群)を翻訳することが一般的である⁴。まずステップ1では、次のような役割のメンバーからなるプロジェクトが編成される。

- リサーチャー(ステップ2a)
- ターミノロジスト(ステップ2b)
- 翻訳メモリ管理者(ステップ2c)
- 翻訳者(ステップ3)
- 修正者(ステップ4)
- レビュアー(ステップ5)

また、その編成・工程管理を司るプロジェクト・マネージャーも配置される。現在の産業翻訳では、表2のようなワークフローに沿った翻訳を支援するために開発された翻訳作業環境⁵が用いられることが多い。近年は、サーバ上でプロジェクトおよび各種アイテムを運用し、翻訳者がWebブラウザを通じて用いる、クラウド型の翻訳作業環境が人気を博している。

³我々が認識している範囲で蓋然的なものに過ぎないことに注意されたい。すべての翻訳がこのワークフローに沿って行われているという主張ではない。

⁴翻訳サービスに関する国際標準規格[5]でも複数人からなるプロジェクト型の翻訳ワークフローが想定されている。

⁵Trados, Memsource, MemoQ等。プロジェクトの工程管理の他、翻訳メモリ、用語集管理、辞書引き機能、対訳コンコーダンサ、MT訳の表示等の機能を備えている。

ステップ 2a, 2b, 2c では、翻訳に際してプロジェクト内で共有しておくべき情報が収集・整理される。ステップ 2a ではリサーチャーが参照情報源の列挙、情報の真偽の判断、構造化等を担当し、ステップ 2b ではターミノロジストが用語集の作成および管理を担当する。ステップ 2c では翻訳の依頼者から翻訳メモリが提供された場合や 2a のリサーチを通じて既訳が得られた場合に翻訳メモリへの登録を行う。これら 3 つのステップはある程度並行して実施することができるが、互いに情報を補完するとともに、網羅性と一貫性を担保することが求められる。

諸々の準備が整って初めて、ステップ 3 で翻訳を開始できる。翻訳者は下訳の作成にあたり、例えば次に列挙するような操作を組み合わせて用いる。

- セグメント (テキスト全体ではないが、文よりは大きい、通常は段落が中心) 全体の情報の構造を目標言語で表現
- 用語集や翻訳メモリにおける対訳の埋め込み
- 用語以外の語句の辞書引き
- 目標テキストのユーザにとって自明な要素等の省略
- 目標テキストの流暢さ向上のための各種言い換え
- 補足的要素の追加 (文中、但し書き、注釈等)

翻訳者が作成した下訳は、翻訳者とは別の担当者に渡される。ISO の翻訳サービス基準 [5] では、起点テキストも参照して行う修正・校閲 (ステップ 4) と目標テキストのみを参照して行うレビュー (ステップ 5) の 2 つのステップが設けられている。品質評価の観点・基準は様々であるが、表 1 では、広く用いられている 4 つの観点を挙げた。

産業翻訳は、目標テキストを得た時点で終わるわけではない。目標テキストの分量、タグ情報 (強調等) の一貫性、その他のテキスト外の要素 (文書の様式やファイル形式) も含む所定の仕様を満足していることを確認する工程を経て、その成果物を翻訳の依頼者に納品し (ステップ 6)、依頼者からの質問等に回答し⁶、検査に合格してようやく終了となる。

4 技術の有用さの基準

産業翻訳に技術を役立てるということは、

翻訳の品質を人間が担保する前提で、種々のコストに見合う効率化を技術的に実現する。

⁶個々の訳出について理由を説明するコンピテンスが重要視されるようになってきた。表 1 で説明した概念の一部は、そのような説明の道具 (メタ言語) として有用であるが、それらの全体像の把握と整理、およびそれらに基づくコンピテンス形成プロセスの確立は、ともに未解決の重要な課題である。

ということである。より具体的には、次の 3 つの観点で説明できる。

品質 (Quality): 翻訳の依頼者が求める品質⁷を担保することは必須である。技術そのものは翻訳品質の瑕疵に対する責任を取れない。すなわち、十分なコンピテンスを有する翻訳者の介在は不可欠である。

コスト (Cost): 新たな技術の開発には時間と資金を要する。他者による成功例を取り入れる場合は開発のコストをいくらか削減できるが、市場における競争で遅れをとる可能性がある。自ら新たな技術を開発する場合は、成功に至らない場合のリスクも負う。また、導入した技術の運用に際してメンテナンス等のコストを要するし、ユーザ (翻訳者等) が当該技術を十分に活用できるようになるまでには時間を要する。

納期・スピード (Delivery): 技術が翻訳の効率を改善する保証はない。適切な KPI (例えば起点テキストの分量と最終成果物の作成に要した時間との比) を設定し、定量的に効率を評価する必要がある。

考えられるアプローチは、現在のワークフローの効率化と新たなワークフローの創出という 2 種類に大別できるが、クラウド型の翻訳作業環境の普及を鑑みれば、前者の方が低コストであり、短期的に成果を得られる可能性も高い。ただし、いずれのアプローチでも、重要な意思決定を翻訳者が担うこと、また現在のワークフローでは翻訳の作業 (ステップ 3) に着手する前のステップ 2a, 2b, 2c がその一部を担保するために組み込まれていること、を鑑みれば、新しい技術には、

- 翻訳者の知識不足を補うこと
- 翻訳者の意思決定を支援すること

が求められる。また、技術の精度は意思決定の効率に直結するため、有用であるか否かの判断基準は極めて重要な位置をしめる。身近な例としては、MT システムを用いることが挙げられる。自動評価や人手評価を通じて示される性能がいかに高くても、未知の起点テキストに対する MT 訳は、表 1 の「品質評価の観点」の品質を、(X 文 Y 訳のレベルでも) 保証しえない。また、後述の通りポストエディットと組み合わせて用いるとしても、常に効率を改善できるとは限らない。

5 既存の自然言語処理技術の応用例

次に示すように、表 2 のワークフローの各ステップに対応する自然言語処理技術はいくつかある。

⁷産業翻訳において求められる翻訳品質は、一般人が MT を使って文書を斜め読みする場合に期待する水準とはまったく異なる。

2b. **用語集の作成:** テキストデータからの用語抽出, 対訳データからの対訳辞書の自動構築

2c. **翻訳メモリの設定:** 対訳データの翻訳メモリ利用, 翻訳成果物の動的登録(後述の Adaptive MT [2])

3. **下訳の作成:** MT による訳文候補の生成と提示

4. **修正・校閲** 文法誤り検出(翻訳品質推定)や文法誤り訂正(翻訳の自動後編集)等のテキスト診断技術

これらの中には, 4 節で述べた基準を満たすか否かを確認しながらではあるが, 様々な翻訳作業環境に取り入れられ, 産業翻訳に利用されているものもある。以下では, 国内で比較的広く検討されている技術の例として, MT 訳のポストエディット, Adaptive MT, および我々が経験的に有用であろうと期待している可視化技術について, 簡単に述べる。

MT のポストエディット: MT 訳を人間が修正することによって訳文を得る方法である。一般に MT 訳の品質が不安定であること, ステップ 4 を設けないこと⁸が, 作業者の認知負荷や翻訳成果物の品質に影響する [9]。

Adaptive MT: MT のポストエディットの工程における, 翻訳者の行動や訳出結果に応じて, MT のモデルを動的に更新する技術である [2]。用語や固有表現の翻訳に難のあるニューラル MT を用いる場合でも, モデルの動的更新により用語の訳出の誤りを減らすことができるという報告がある [8]。

各種情報の可視化: 既存の翻訳関連技術に比べて, 各種自然言語解析技術の精度は高い。翻訳者に対して各種解析結果を可視化することは有用であろう。例えば, 次のような用途が考えられる。

- 起点テキストの構造(統語構造等)を可視化すれば, ステップ 3, 4 における起点テキストの内容の理解に役立つ。
- 下訳の産出過程において, 起点テキストの構成要素と訳文中の構成要素との対応を可視化すれば, 訳出済の内容とそうでない内容の区別や情報構造の対応の確認に有用である。

6 論点

本稿では, 産業翻訳を抽象的に捉えて概要を整理したが, 実際には現場によって多様な制約やワークフローが存在する。産業翻訳への技術の導入に際しては, 要求品質の保証を前提として, コストと効率化の多寡の両面からの検討が必要である(4 節)が, そのための方

法論が確立されれば, 技術の導入は促進され, 要求品質に対する公正な価格設定も可能になるだろう。

ただしそのためには, 翻訳のワークフローや翻訳品質に関する合意形成という, より大きな課題を解決する必要がある。例えば, 表 1 に示した各種概念, アイテム, 行為等, 表 2 に示した行為やコンピテンス等を整理した規範的なプロセスモデルを確立するというアプローチが考えられる。その上で, 技術の導入とそこで想定されることが人間のコンピテンスにどのような影響を与えるかについても考える必要がある。

翻訳という知的営みを多面的に捉え, 上記のような課題に関する議論の場を提供するため, 今後も本テーマセッションと同様の企画を継続したい。

参考文献

- [1] M. Baker and G. Saldanha. 藤濤文子(監修・編訳). 伊原紀子, 田辺希久子(訳). 翻訳研究のキーワード. 研究社, 2013.
- [2] M. Denkowski, A. Lavie, I. Lacruz, and C. Dyer. Real time adaptive machine translation for post-editing with cdec and TransCenter. In *Proceedings of the Workshop on Humans and Computer-assisted Translation*, pp. 72–77, 2014.
- [3] D. Elliott, S. Frank, and E. Hasler. Multilingual image description with neural sequence models. *CoRR*, abs/1510.04709, 2015.
- [4] 藤田篤, 山田優. 翻訳の品質と効率: 実社会におけるニーズと工学的実現可能性. 言語処理学会第 23 回年次大会発表論文集, pp. 915–918, 2017.
- [5] ISO/TC27. ISO 17100:2015 translation services: Requirements for translation services, 2015.
- [6] ISO/TC27. ISO 18587:2017 translation services: Post-editing of machine translation output: Requirements, 2017.
- [7] 影浦峽. 改めて, 翻訳とは何か: Google NMT が使える時代に. 言語処理学会第 23 回年次大会発表論文集, pp. 931–934, 2017.
- [8] 梶木正紀. 機械と人間の協働 LSP perspective: 学習型機械翻訳から翻訳プラットフォーム. 言語処理学会第 24 回年次大会発表論文集, pp. 750–752, 2018.
- [9] 森口功造. ISO17100 と ISO/DIS18587.2 の要求事項の比較とポストエディット現場への影響. 言語処理学会第 23 回年次大会発表論文集, pp. 1157–1159, 2017.
- [10] 西尾実, 岩淵悦太郎, 水谷静夫(編). 岩波国語辞典第 7 版. 岩波書店, 2009.
- [11] 田中千鶴香. 実務翻訳における翻訳品質評価プロセス. 言語処理学会第 23 回年次大会発表論文集, pp. 923–926, 2017.
- [12] H. J. Vermeer. Skopos and commission in translational action. In L. Venuti, editor, *The Translation Studies Reader*, pp. 227–238. Routledge, 1989/2004.
- [13] E. Voita, P. Serdyukov, R. Sennrich, and I. Titov. Context-aware neural machine translation learns anaphora resolution. In *Proceedings of ACL*, pp. 1264–1274, 2018.

⁸ISO の翻訳サービス基準 [6] では, ポストエディットの作業者と別修正者を立てることを要求していない。

表 1: 本テーマセッションに関連する用語とその説明.

分類と用語	説明
翻訳関連の活動・機械処理	
翻訳 (translation)	ある言語で表現された文章の内容を、原文に即して他の言語に移しかえること [10].
通訳 (interpretation)	互いに言語が異なって話が通じない人々の間に立って、双方の言うことを翻訳して話を通じさせること [10]. 著者補足: 一定のまとまりをもった講演や談話のセッション全体を基本単位としてなされる.
産業翻訳 (industrial translation)	産業活動の国際化等で必要となる翻訳. 知的財産関連文書 (特許文書等), 製品等のマニュアル文書等を扱う.
機械翻訳 (machine translation)	ある言語表現を、計算機を用いて別の言語に変換すること. 自動翻訳とも.
音声翻訳 (speech translation)	通訳と同様に音声発話を扱う機械翻訳の下位区分. ただし、セッションを基本単位とする通訳とは異なり、主に発話ごとの変換を想定する.
翻訳関連の概念	
スコ-pos (skopos)	翻訳の目的. これにより翻訳成果物, 人間の行為とその下位範疇としての翻訳が決定される [12].
翻訳戦略 (strategy)	スコ-posを達成するための戦略で, 種々の判断の基準となる概念.
翻訳規範 (norm)	翻訳成果物が使用される目標文化社会に受容される表現や言い回しに関する判断基準と慣例. 翻訳者の翻訳行為・翻訳成果物は翻訳規範に支配される. 翻訳成果物のユーザは期待規範を有する.
テキストタイプ (text type)	テキストの様式. メール, 報告書, 特許出願書類等.
レジスタ (register)	言語使用域. 特定の社会的場面で用いられる言語変種 [1, p.27].
スタイル (style)	送り仮名, 括弧, 記号等の使い方, 文体, 表現の難度等.
翻訳関連のアイテム	
ブリーフ (brief)	テキストタイプ, レジスタ, スタイル等に基づいて翻訳における種々の選択点に対する制約を与えるための指示書.
参照情報源 (reference)	起点テキストの翻訳に際して, 翻訳規範を確認するために参照されるもの. 起点テキストと同じテキストタイプ, レジスタのテキスト等, および内容の確認のために参照される関連する内容のテキストやレファレンス・ツール等.
用語集 (terminology)	起点テキストが属する分野における用語およびその対訳を列挙したリスト.
翻訳メモリ (translation memory)	起点テキストが属する分野における既訳のリスト. 文単位が中心だが, 用語集のエントリを含みうる.
成果物およびそれをめぐる行動	
下訳 (translation draft)	大まかな訳を行うこと. またその成果物. 仮訳とも.
修正・校閲 (revision)	目標テキストにおける内容や表現の誤り等を, 起点テキストも参照して修正すること.
レビュー (review)	目標テキストのみを参照して内容や表現等の適否を確認すること.
リライト (rewrite)	誤りではないが流暢さを向上させるために目標テキストを編集すること.
品質評価の観点	
適切さ (adequacy)	起点テキストの内容が目標テキストに正しく反映されているかどうか. 正確さ (accuracy) [11] とも.
流暢さ (fluency)	目標言語を母語とするネイティブ話者が目標テキストを読んだときに, スムーズに読めるかどうか [11].
結束性 (cohesion)	テキストを構成する要素間のつながりのよさ.
首尾一貫性 (coherence)	スタイル, 論理, 用語の訳出等の一貫性.

表 2: 産業翻訳のワークフロー.

テキストタイプ, レジスタ, スタイル, およびスコ-posはすべてのステップで共有される. [ST]: 起点テキスト, [規]: 翻訳規範, [ブ]: ブリーフ, [参]: 参照情報源, [語]: 用語集, [メ]: 翻訳メモリ, [TT]: 目標テキスト. “(✓)” は該当するアイテムが存在する場合にのみ参照される.

ステップ	関連する概念・アイテム							担当者に求められるコンピテンス
	[ST]	[規]	[ブ]	[参]	[語]	[メ]	[TT]	
1. 依頼内容の確認, プロジェクトの設置	✓	-	(✓)	(✓)	(✓)	(✓)	-	人的資源等の管理, 受注の判断
2a. 調査一般	✓	✓	-	✓	(✓)	(✓)	-	起点テキストの翻訳に関連する情報の収集, 信頼性の検証, 情報の構造化
2b. 用語集の作成	✓	✓	✓	(✓)	✓	(✓)	-	用語の認定基準, 既訳の踏襲
2c. 翻訳メモリの設定	✓	✓	✓	(✓)	(✓)	✓	-	翻訳メモリへの既訳の登録
3. 下訳の作成	✓	✓	✓	✓	✓	✓	✓	所与の時間内に一定水準以上の品質の訳を産出すること, 個々の訳出について理由を説明すること
4. 修正・校閲	✓	✓	✓	✓	✓	✓	✓	同上
5. レビュー	-	✓	✓	✓	✓	-	✓	目標言語のテキストの品質評価
6. 最終確認, 納品後対応	✓	✓	✓	✓	✓	-	✓	テキスト外の要素も含む仕様の満足