

JSTにおける英語・中国語科学技術文献の機械翻訳

Machine translation of English and Chinese science and technology literature in Japan Science and Technology Agency

国立研究開発法人科学技術振興機構（JST）情報企画部主任調査員

岩城 修

1980年東京工業大学大学院総合理工学研究科電子システム専攻修士課程修了。日本電信電話公社（現 日本電信電話株式会社）、株式会社 NTT データ等を経て、2015年7月より現職。機械翻訳システム等の開発・運用業務に従事。博士（工学）。

国立研究開発法人科学技術振興機構（JST）情報企画部主任調査員

松永 務

1988年電気通信大学大学院通信工学専攻修士課程修了。株式会社 NTT データ技術開発本部にてデータマイニング、機械翻訳に関する技術開発に従事。2017年10月より出向し、現職。博士（工学）。

国立研究開発法人科学技術振興機構（JST）情報企画部調査役

堀内 美穂

1988年北海道大学薬学部製薬化学科卒業。同年、日本科学技術情報センター（現、科学技術新興機構）に入職。文献データベースの作成、開発、研究者データベースの運営に従事し、2018年より現職。

1 はじめに

JSTでは、我が国の研究開発活動の効率的実施を促し、科学技術の振興を図ることを目的に、国内外から収集した科学技術や医学・薬学関係の文献に抄録・索引を付与した文献情報データベースを整備し、J-GLOBAL¹⁾やJDream III²⁾などのインターネットアクセス可能な検索サービスを提供している。これらは科学技術に関する国内誌を網羅的に収集している唯一のデータベースであり、外国誌（英語・中国語）については出版者から許諾を得て日本語訳を作成することにより、全て日本語での検索を可能としている。ここに、その日本語への翻訳は永く人手で実施していたが、収録する外国誌の規模拡大と情報提供の迅速性を確保するため、2014年から機械翻訳の導入を図っている。

本稿では、文献情報データベースに蓄積された標題・

抄録の日本語訳に基づく対訳コーパスをはじめ、JSTにおいて構築整備している科学技術用語辞書³⁾を活用して導入した科学技術文献の機械翻訳について述べる。

2 JSTにおける機械翻訳の導入実績

JSTにおける年間記事収録件数の過去5年間の推移と外国誌への機械翻訳導入比率を図1に示す。2014年度からまず中国語文献に機械翻訳を導入し、英語文献には2016年度から導入したことで、外国誌の年間記事収録件数の大幅な増を図から読み取ることができる。2018年度からは外国誌に対して機械翻訳を全面導入し、英語文献で約100万件、中国語文献で約50万件の実績で、2019年度はさらに拡大する計画である。

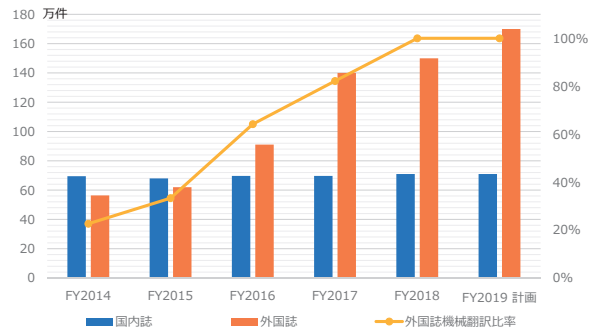


図1 記事収録件数の年次推移と外国誌の機械翻訳導入比率

図2はJDream IIIでの検索結果の回答表示画面の一例で、日本語標題および抄録がそれぞれ英語標題および抄録の機械翻訳により作成されている。機械翻訳の導入にあたり、必要に応じて原文(英文)を参照できるように、図のように併せて表示することとしている。

ANSWER 1 OF 1 JSTPLUS JST COPYRIGHT JDreamIII転写可能

登録番号: 19A1654306
 和文標題: 成人におけるバンコマイシンの薬物動態に対する好中球減少症の影響 [JST・京大機械翻訳]
 英文標題: The effect of neutropenia on the clinical pharmacokinetics of vancomycin in adults
 著者名: [不明]
 資料名: 日本語標題は英語標題の機械翻訳により作成
 JST資料番号: [不明]
 巻号ページ (発行年月日): Vol.75 No.7 Page.921-928 (2019)
 資料種別: 逐次刊行物(A)
 原書区分: 原著論文(e1)
 発行国: ドイツ(DEU) 言語: 英語(EN)
 抄録: 好中球減少症患者はバンコマイシンのより高い投与量を必要とするという証拠が蓄積している。サブ治療薬理を予防するためには、最初の用量から適切な増量を得ることが最も重要である。バンコマイシンの薬物動態に対する好中球減少症の影響を定量化することを目的とした。(1)血液疾患、(2)固形腫瘍性腫瘍、(3)癌では知られていない患者のマッチした患者コホートからデータを抽出した。薬物動態学的分析はバンコマイシンの分布のクリアランスと半減期に関する二成分共変量として研究された好中球減少症による非線形混合効果薬モデリングを用いて行われた。合計116名の患者が含まれた(39名の血液学的患者、39名の固形腫瘍患者、および38名の癌では知られていない患者)。全部で、742対の時点濃度観察が薬物動態分析に利用された。好中球減少症の存在は有意に(p=0.00157)バンコマイシンのクリアランスを27.7%(95%CI10.2~46.2%)増加させたが、分布容積には影響しなかった(p=0.704)。本研究はバンコマイシンクリアランスが好中球減少症患者で27.7%増加することを示す。したがってバンコマイシン維持量は治療開始時の好中球減少患者において25%増加するべきである。分布容積は影響を受けなかったため、負荷量の調整は必要でなかった。これらの用量調整は、治療薬モニタリングによるさらなる用量個別化の必要性を除外しない。Copyright 2019 The Author(s) Translated from English into Japanese by JST. [JST・京大機械翻訳]
 英文抄録: There is accumulating evidence that neutropenic patients require higher dosages of vancomycin. To prevent sub-therapeutic drug exposure, it is of utmost importance to obtain adequate exposure from the first dose onwards. We aimed to quantify the effect of neutropenia on the pharmacokinetics of vancomycin. Data were extracted from a matched patient cohort of patients known with (1) hematological disease, (2) solid malignancy, and (3) patients not known with cancer. Pharmacokinetic analysis was performed using non-linear mixed effects modeling with neutropenia investigated as a binary covariate on clearance and volume of distribution of vancomycin. A total of 116 patients were included (39 hematologic patients, 39 solid tumor patients, and 38 patients not known with cancer). In total, 742 paired time-concentration observations were available for the pharmacokinetic analysis. Presence of neutropenia showed to significantly (p = 0.00157) increase the clearance of vancomycin by 27.7% (95% CI 10.2-46.2%), whereas it did not impact the volume of distribution (p = 0.704). This study shows that vancomycin clearance is increased in patients with neutropenia by 27.7%. Therefore, the vancomycin maintenance dose should be pragmatically increased by 25% in neutropenic patients at the start of treatment. Since the volume of distribution appeared unaffected, no adjustment in loading dose is required. These dose adjustments do not rule out the necessity of further dose individualization by means of therapeutic drug monitoring. Copyright 2019 The Author(s)
 分類コード: CVG1020D(615.45.03) 生物薬剤学
 GX05030N(615.33.015.1.03) 薬原体に作用する抗生物質の臨床への応用
 シソーラス用語: *好中球減少症, *薬物動力学, 固形腫瘍, 母集団, 用量, 治療法, クリアランス, ヒト, 血液疾患, モデリング, アミノ酸, シクロペプチド, 芳香族複素化合物, オルゴペプチド, 多価フェノール, 二種化合物, ビラシド, 芳香族複素化合物
 漢シソーラス用語: 慢性化, 分布容積, 血液学, 共変量, 悪性腫瘍, [AI/JST], #Vancomycin, #バンコマイシン, #Pharmacokinetics, #薬物動態, #Neutropenia, #好中球減少, #NONMEM, #RMEM, #Hematology, #血液学
 IPC(機械付与) A61K31: 生活必需品>医学または獣医学->医薬用, 薬料用又は化->有効活性成分を含有する医薬品製剤
 物質索引: *バンコマイシン [J24.275K, 1404-90-6]
 DOI情報: doi: 10.1007/s00228-019-02657-6
 リンク情報: [不明]

図2 JDream IIIでの検索結果の回答表示画面の一例(「バンコマイシン」の検索語でヒットした記事候補から選択して得られた表示画面例)

3 ニューラル機械翻訳の導入

3.1 ニューラル機械翻訳エンジンの開発

JSTでは、科学技術文献に対する実用的な翻訳精度を達成することを目標に2013年から2017年まで日中・中日機械翻訳実用化プロジェクト⁴⁾を推進しており、ここで開発されたニューラル機械翻訳エンジン(KyotoNMT)⁵⁾を現在採用している。2016年12月に開催されたアジア言語を対象とした国際的な機械翻訳のワークショップ「WAT 2016」(Workshop on Asian Translation 2016)⁶⁾での評価において、

ニューラル機械翻訳方式を用いて科学技術情報の日中・中日機械翻訳タスクで1位の精度を達成している。

図3はJSTにおける機械翻訳の処理フローの概要である。収録される科学技術文献の記事原文に含まれる特殊文字列の置換処理を経た上で、対訳用語辞書に登録された科学技術用語に対しては、前処理でマスク処理して

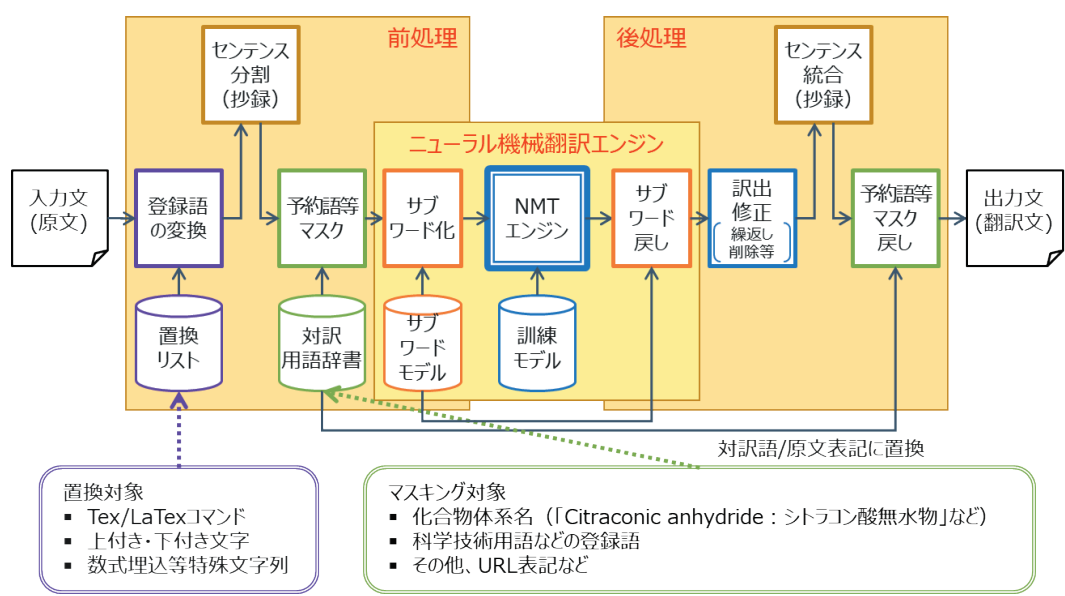


図3 ニューラル機械翻訳の処理フロー概要

後処理で対応する語に置換して翻訳文を得る予約語化処理を行っている。(4章で事例を詳説)

3.2 対訳コーパスの構築と整備

JSTでは、医学・薬学分野を含めて広く科学技術文献を扱う中で、英語標題・抄録の人手による日本語訳作成から英日対訳を蓄積保有し、また記事への索引付けを目的に科学技術用語辞書を構築整備しており、これらの言語データ資産を活用して機械翻訳を開発している。

対訳コーパスの構築において、複数文からなる抄録に関しては文アライメント処理⁷⁾により文対形式にして機械翻訳に用いる対訳コーパスとしている。中国語・日本語の対訳コーパスについては日中・中日機械翻訳実用化プロジェクトにおいて中国科学技術情報研究所(ISTIC)と連携して整備した。現行の対訳コーパスの規模は、英日対訳の標題約2,300万件、抄録約1,300万件で、中日対訳は標題と抄録で約500万件である。

一般に、対訳コーパスの量が多いほど翻訳精度が向上するとされているが、上記文アライメント処理時に科学技術用語含有で高スコアのものから選択するなど、科学技術文献の機械翻訳に向け、対訳コーパスと用語辞書の整備を併せて行っている。

4 科学技術用語の翻訳

機械翻訳の導入にあたり、特に科学技術用語の正しい訳出に注力しており、保有する科学技術用語辞書³⁾を活用して、用語辞書に登録された語に対する予約語化処理(図3参照)を導入している。以下、代表的な科学技術用語として化合物名を取り上げ、その効果について述べる。

図4aおよびbは化合物(有機低分子化合物)体系名を含む英日機械翻訳例であり、化合物名を予約語化対象に適用した場合としない場合で翻訳出力を示している。

原英文	<u>Citraconic anhydride</u> -modified ovalbumin (Cit-OVA), as model antigen, was incorporated into PHMs via electrostatic interaction, giving antigen-loaded micelles of around 150nm in size.
機械翻訳出力和文(予約語化処理有)	モデル抗原としての <u>シトラコン酸無水物</u> 修飾卵白アルブミン(CIT-OVA)を静電相互作用によりPHMに組み込み、サイズ約150nmの抗原負荷ミセルを得た。
機械翻訳出力和文(予約語化処理無)	モデル抗原としての <u>クエン酸無水物</u> 修飾卵白アルブミン(CIT-OVA)を静電相互作用によりPHMに組み込み、サイズ約150nmの抗原負荷ミセルを得た。

図4a 化合物名を含む英日機械翻訳例
(予約語化処理が無いと「シトラコン酸無水物」を「クエン酸無水物」に誤る)

原英文	In this study, gel-like liquid silk (LS) obtained from the middle division of the middle silk glands of fully grown larvae of the domesticated <i>Bombyx mori</i> silkworm was dissolved into <u>1,1,1,3,3,3-hexafluoro-2-propanol</u> (HFIP) and electrospinning was performed, and then the crystal modification of silk fibroin (SF) in the nanofibers by the water vapor or ethanol treatment at room temperature was investigated.
機械翻訳出力和文(予約語化処理有)	本研究では、栽培された <i>Bombyx mori</i> カイコの完全に成長した幼虫の中絹糸腺の中央部から得られたゲル様液体絹(LS)を <u>1,1,1,3,3,3-ヘキサフルオロ-2-プロパノール</u> (HFIP)に溶解し、電気紡糸を行い、次に、室温での水蒸気またはエタノール処理によるナノ繊維中の絹フィブロイン(SF)の結晶改質を研究した。
機械翻訳出力和文(予約語化処理無)	本研究では、栽培された <i>Bombyx mori</i> カイコの完全に成長した幼虫の中絹糸腺の中央部から得られたゲル様液体絹(LS)を <u>1,1,1,3,3,3,3,3-ヘキサフルオロ-2-プロパノール</u> (HFIP)に溶解し、次に、室温での水蒸気またはエタノール処理による絹フィブロイン(SF)の結晶改質を研究した。

図4b 化合物名を含む英日機械翻訳例
(予約語化処理が無いと「3,3」の過剰な出力がみられる)

翻訳出力の比較から、予約語化対象に適用しない場合には、図 4a では「シトラコン酸無水物」を「クエン酸無水物」に誤る翻訳文、図 4b では「3, 3」の過剰な出力がみられる。低頻度語や記号列に対して不正確となるニューラル機械翻訳の問題⁹⁾との関わりが考えられるが、対応する化合物名を予約語化処理することで正しい訳出となることが図から確認できる。

なお化合物体系名については、日本化学物質辞書⁹⁾を用いることにより IUPAC 命名規則に準拠して予約語化の処理を行っている。

新語への対応をはじめ、訳質維持改善に向けて対訳コーパスの拡充に取り組んでいるが、ここで取り上げた化合物名については、日本化学物質辞書において新規登録されたもの（2018 年で例えば「1,5-dibutoxynaphthalene - 2,6 - dicarbaldehyde」や「(2E) -3- (4-chlorophenyl) -N- (4-oxocyclohexyl) prop-2-enamide」) を基に、それらを含むようなコーパスから効率的に収集するなど、用語辞書の更新と併せた対応に取り組んでいる。

5 おわりに

文献情報データベースは、これまでの書誌、抄録、索引情報から全文データ活用の方向性にある。また、科学技術戦略や研究テーマの策定、研究パートナーの探索、研究成果の評価など、科学技術文献情報の利用ニーズはこれまでの文献検索から付加価値利用に拡大している。

外国語文献を幅広く大規模に扱う重要性が増して機械翻訳の役割がますます高まる状況にある中で、構築整備している科学技術用語辞書の活用展開を基に、機械翻訳の効果的導入の取り組みを進めているところである。

参考文献

- 1) J-GLOBAL : <http://jglobal.jst.go.jp/>
- 2) JDream III : <http://jdream3.com/>, 提供は株式会社ジー・サーチ
- 3) 川村隆浩, 渡邊勝太郎, 松邑勝治, 榎田達矢, 古崎晃司 : JST 科学技術用語シソーラスの Linked Data 化 : 科学技術情報をリンクする知識インフラの構築に向けて, 情報管理, 59 巻, 12 号 (2017) pp.839-848
- 4) 日中・中日機械翻訳実用化プロジェクト : https://jipsti.jst.go.jp/jazh_zhja_mt/
- 5) KyotoNMT : <https://github.com/fabiencro/knmt/>
- 6) WAT 2016 : <http://lotus.kuee.kyoto-u.ac.jp/WAT/WAT2016/>, Proceedings of the 3rd Workshop on Asian Translation (WAT 2016), December 11-16, 2016, Osaka, Japan
- 7) Antoine BURLON, Chenhui Chu, Toshiaki Nakazawa and Sadao Kurohashi: Simultaneous Sentence Boundary Detection and Alignment with Pivot-based Machine Translation Generated Lexicons, Proceedings of the 10th Conference on International Language Resources and Evaluation (LREC2016), Portoroz, Slovenia (2016.5)
- 8) 本間奨 : 機械翻訳の近未来 第 6 回 NMT の今後, 日本翻訳ジャーナル, No.289 (2017 年 5 月 /6 月号)
- 9) 木村考宏, 榎田達矢 : 日化辞 RDF データの公開と化合物情報の統合, 情報管理, 58 巻, 3 号 (2015) pp.204-212