

多言語翻訳ソリューションに向けた翻訳モデルのカスタマイズによる翻訳精度の向上

Translation Accuracy Improvement by Customization of a Machine Translation Model for Multilingual Translation Solutions

生田 久敏* 荒木 昭一**
Hisatoshi Ikuta Shoichi Araki

訪日外国人との多言語による業務コミュニケーションを支援するため、接客業務向け会話コーパスを強化したニューラル機械翻訳(NMT:Neural Machine Translation)モデルを開発した。従来の統計的機械翻訳(SMT:Statistical Machine Translation)に基づく翻訳モデルに対して、ホテルや旅館、公共交通機関などでの業務会話を想定した性能評価において、日本語と英語、日本語と中国語のそれぞれ双方向の翻訳精度が10%~23%向上した。

To support multilingual operational communications, we developed Neural Machine Translation (NMT) model with enhancement of the service conversation corpus. The efficiency became higher by 10 to 23 percent compared to conventional Statistical Machine Translation (SMT) by translation accuracy evaluation assuming the use of operational conversation in situations at Hotels, Japanese Inns, Mass Transport Systems, and so on, for both directions for Japanese-English and Japanese-Chinese respectively.

1. 翻訳モデルの業務会話向けカスタマイズ

東京2020オリンピック・パラリンピック競技大会の開催を1つの重要な契機として、訪日外国人数は急増しており、大会後も引き続き増加することが期待されている。このように急増する訪日外国人とのコミュニケーション支援を目的として、機械翻訳技術の近年の急速な進歩を背景に、多言語翻訳ソリューションの実用化が進んでいる[1]。当社は2017年11月より、ホテルや旅館などでの対面接客業務を支援する多言語音声翻訳サービス「対面ホンヤク^(注1)」を提供している[2]。「対面ホンヤク」では、国立研究開発法人情報通信研究機構(NICT)と株式会社みらい翻訳(以下、みらい翻訳)が共同研究で開発した翻訳エンジン[3][4]を活用している。サービス開始時は、統計的機械翻訳(SMT:Statistical Machine Translation)エンジン、2018年度からは深層学習に基づくNMT(Neural Machine Translation)エンジンに順次移行し、翻訳精度の向上に取り組んできた。

翻訳精度の向上には、翻訳エンジンの進化とともに、サービス対象の業種やお客様それぞれに特有の業務会話で十分な翻訳精度が得られるように翻訳モデルをカスタマイズすることが重要である。翻訳モデルのカスタマイズでは、基本的な会話に関する対訳コーパス(“こんにちは”⇔“Hello”など、翻訳する言語間での対訳データの集合)に、業務に特有の対訳コーパスを追加して翻訳モデルを学習させることが有効である。当社では、NICTから利用許諾を得た対訳

(注1) 当社の日本国内における登録商標。

* コネクティッドソリューションズ社 イノベーションセンター
Innovation Center, Connected Solutions Company

** 要素技術開発センター

Core Element Technology Development Center

コーパスに、当社が提供するサービス関連(接客会話など)の対訳コーパスおよびみらい翻訳の対訳コーパスを追加した翻訳モデルの学習を、みらい翻訳のカスタムモデル学習サービスを用いて実行し、対象業務での翻訳精度の向上を図っている。以下、SMTおよびNMTモデルのカスタマイズによる翻訳精度向上について解説する。

2. SMTカスタムモデルによる翻訳精度向上

2016年には、ホテルや旅館でのサービス提供を想定し、お客様との実証実験を含めて、接客時の多様な言い回しに対応した独自の対訳コーパスの追加による翻訳精度向上の可能性を検証した。具体的には、「〇時に△△駅から出発するツアーがあります。」などの基本的な接客会話のコーパスに対して、例えば、「△△駅から〇時に出発するツアーがございませう」など、語順や丁寧表現などの言い回しが異なるコーパスを追加して翻訳モデルを学習させた。第1表は、2016年当時のNICTのSMTの翻訳エンジンに対して、ホテルや旅館などでの接客会話の独自テストセット424文を用いた翻訳精度の主観評価結果である。翻訳結果の主観評価は、「Good(4):意味が完全に通じる」、「So-So(3):ほぼ意味が通じる」、「Not-good(2):誤りがあり意味が通じない」、「Bad(1):全く意味が通じない」の4段階とし、「Good」「So-So」と判定された文の割合[%]を翻訳精度とした。日英翻訳(日本語から英語への翻訳)では、NICTやみらい翻訳の既存コーパスで学習された従来モデルの翻訳精度は71%~74%であるが、NICTおよびみらい翻訳の既存コーパスに、当社独自の接客会話コーパスを追加したコーパス結合モデルでは86%に向上することを確認した。同様に、日中翻訳(日本語から中国語への翻訳)では、従来モ

第1表 コーパス結合による翻訳精度の向上

Table 1 Improvement of translation accuracy by combination of the plurality of corpora

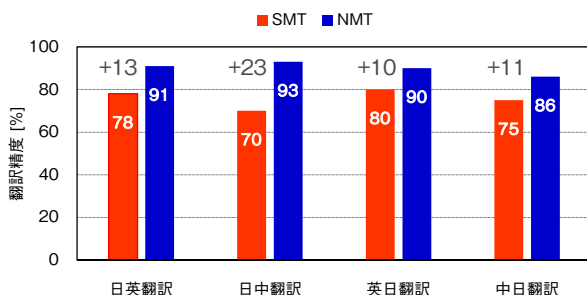
	日英翻訳	日中翻訳
従来モデル [%]	71~74	64~66
コーパス結合モデル [%]	86	87

デルでは64 %~66 %であるが、コーパス結合モデルでは87 %に向上した。

2017年には、総務省のグローバルコミュニケーション計画に基づいてNICTが中心となって開発した多言語パラレルコーパス[5]（日英中韓はおのおの200万文規模）の提供も受け、「対面ホンヤク」のサービス開始時には、ホテルや旅館に加え、公共交通機関や商業施設など、複数の業種、複数のお客様での接客業務への対応を想定して翻訳モデルを学習した。想定用途に対して独自に設計した1000文の評価データで、翻訳精度が日英翻訳82 %、英日翻訳83 %、日中翻訳71 %、中日翻訳78 %となった。

3. NMTカスタムモデルによる翻訳精度向上

SMTで約20 %~30 %生じていた誤翻訳を低減し、かつ、より複雑な文脈の会話にも対応するため、SMTの開発と並行して、2016年度から当社独自のNMTカスタムモデルの開発を開始した。評価用データとしては、2章でのSMT評価用データ、お客様の実使用ログデータ、NMT評価専用データを合わせて数千文規模のデータセットを構築した。NMT評価専用データでは、例えば「駐車場はまだ空いています。」などSMTでも翻訳可能な比較的短い基本的な業務会話文から、「別館の駐車場はまだ空いていますが、それ以外の駐車場には空きがありません。」など、より複雑な文脈の会話文を設計した。第1図は、SMTとNMTのコーパス結合によるカスタムモデルの翻訳精度の比較結果である。SMTでは翻訳精度が70 %~80 %であるが、NMTでは、翻訳方式の変更に加え、誤翻訳要因であった複数種のコーパスでの語彙の表記ゆれを学習時に統一するなどの対策により86 %



第1図 翻訳精度 (SMT vs. NMT)

Fig. 1 Translation accuracy (SMT vs. NMT)

~93 %に向上し、2章でのSMTでは得られなかった精度に改善した。

第2表に翻訳結果の例を示す。例1、例2では、SMTはフレーズ内文脈しか扱えず、文中で離れた位置にある単語の関係を学習できないため、「観光案内所の近くに」や「別館に空きがない」と誤翻訳されているが、NMTは入力文全体の情報を文脈として扱えるため正しく翻訳されている。

例3では、原文に含まれていた語彙が翻訳結果に現れない「訳抜け」と呼ばれるNMT特有の現象が、比較的短い文でも確認されており、今後改善すべき課題の1つである[6]。

第2表 翻訳結果の例

Table 2 Examples of translation results

例1	原文	ラーメン屋さんの近くに観光案内所があります。
	SMT	There is a ramen shop near the Tourist information center.
	NMT	There is a tourist information center near the ramen shop.
例2	原文	別館の駐車場はまだ空いていますが、それ以外の駐車場には空きがありません。(商業施設など)
	SMT	The parking lot is still available but other than that there is no vacancy in the parking lot of the Annex.
	NMT	The parking lot in Annex is still available, but there is no vacancy for other parking lots.
例3	原文	オペレーター9番へ申し込んでください。
	SMT	Dial nine and ask the operator.
	NMT	Operator 9, please.

4. 今後の展望

NMTは現在も進化の途上であり、訳抜けなどの特有の課題への対応を含め、より文脈が複雑で多岐にわたる業務会話文に対応させることが重要である。また、改善モデルについては、当社の提供サービスの「対面ホンヤク」だけでなく、さまざまなシステムへ多言語翻訳ソリューションを組み込むことで、多岐にわたる応用を進めていく。

参考文献

- [1] 国立研究開発法人情報通信研究機構 (NICT), “国立研究開発法人情報通信研究機構 (NICT) の多言語音声翻訳技術を活用した民間の製品・サービス事例,” http://gcp.nict.go.jp/news/products_and_services_GCP.pdf, 参照 Oct. 24, 2019.
- [2] パナソニック (株), “多言語音声翻訳サービス「対面ホンヤク」,” <https://panasonic.biz/cns/invc/taimenhonyaku/>, 参照 Oct. 24, 2019.
- [3] 国立研究開発法人情報通信研究機構 (NICT), “ニューラル機械翻訳で音声翻訳アプリVoiceTraが更なる高精度化を実現,” <https://www.nict.go.jp/press/2017/06/28-1.html>, 参照 Oct. 24, 2019.
- [4] 株式会社みらい翻訳, “TOEIC900点以上の英作文能力を持つ深層学習による機械翻訳エンジンをリリース,” <https://miraitranslate.com/news/313>, 参照 Oct. 24, 2019.
- [5] 今村賢治 他, “グローバルコミュニケーション計画のための

多言語パラレルコーパス,” 言語処理学会 第24回年次大会 講演論文集, pp.512-514, 2018.

- [6] I. Goto et al., “Detecting Untranslated Content for Neural Machine Translation,” Proc. the First Workshop on Neural Machine Translation, pp.47-55, 2017.