

# ニューラル機械翻訳における 大規模語彙および訳抜けへの対応の併用\*

木村 龍一郎<sup>†</sup> 龍 梓<sup>†</sup> 飯田 頌平<sup>‡</sup> 宇津呂 武仁<sup>†</sup> 三橋 朋晴<sup>§</sup> 山本 幹雄<sup>†</sup>  
筑波大学大学院 システム情報工学研究科 <sup>†</sup>東京電機大学 工学部 <sup>§</sup>日本特許情報機構

## 1 はじめに

ニューラル機械翻訳 (Neural Machine Translation; NMT) はひとつの大きなニューラルネットワークで翻訳モデルを構成した機械翻訳である。翻訳モデルは対訳コーパスを用い、正解文を生成する条件付き確率を最大化するように訓練される。NMT は従来の機械翻訳手法と比較して流暢さで優れるものの、目的言語文出力の計算時間が使用する目的言語の単語数に依存するために、大規模語彙を含む翻訳に対応できない。この問題に対応する研究はいくつか存在するが [6, 8], これらの手法は未知語を単語単位で処理しているために、複合名詞の一部として出現する場合に想定通りに翻訳できない問題がある。この問題は専門用語を多く含む特許文の翻訳において特に顕著である。

このような背景から、低頻度語を含みやすい複合名詞をトークンに置き換えて翻訳モデルの訓練を行う手法が提案された [5]。トークンに置き換えられた複合名詞は統計的機械翻訳 (Statistical Machine Translation; SMT) によって翻訳され、MT 出力に代入される。しかしこの手法は、入力文を意味的にすべて翻訳することを保証できずに内容が訳抜けするという、NMT のもう一つの大きな課題に対応できていない。訳抜けに対応するために、出力文が入力文を生成する逆翻訳確率を訳抜けした内容の検出に用いる手法が提案された [3]。そこで、本研究では大規模語彙への対応手法と訳抜けの削減手法を組み合わせることで、翻訳結果の訳抜けの度合いと翻訳精度を評価した。

## 2 訓練・評価文

日英対訳特許文として、NTCIR-7 ワークショップで配布された 180 万文の対訳対のうち単語数が 40 語未満の

表 1: 訓練・評価文

訓練文	開発文	評価文
1,167,198	1,000	1,000

もののみで構成された 110 万文を使用した (NTCIR-7 特許翻訳タスク [2] における日英対訳特許文。詳細は [5] 参照)。日英対訳特許文のうち、1,000 文を評価文、1,000 文を開発文として抽出し、残りを訓練文として使用した。訓練文からは 2,785,108 個の対訳フレーズ対が抽出された。抽出されたフレーズ対の種類は 704,346 種で、日本語フレーズは 511,633 種、英語フレーズは 422,269 種だった。訓練文からは 2,539 個、2,171 種の対訳フレーズ対が抽出された。

## 3 NMT モデル・SMT モデルの訓練条件

ワードアラインメントとフレーズ翻訳テーブルを作成するための SMT モデル作成には、Moses Toolkit を使用した [4]。チューニングには開発文を使用した。NMT モデルを訓練する際のパラメータは [1] で使用されたものを用いた。エンコーダは前向き・後ろ向きの 3 層 LSTM, デコーダは前向き 3 層 LSTM で、各層の時限はいずれも 256 次元とした。入力単語の分散表現は 256 次元とした。原言語と目的言語の語彙は頻度上位 40,000 語として、その他の語は NMT モデルの語彙外の未知語とした。

## 4 文単位の翻訳性能の評価

### 4.1 自動評価

表 2 に BLEU による自動評価の結果を示す。ベースラインの SMT と比較すると、大規模語彙対応とリランキングの併用手法 (5 節) は BLEU が約 5.7 向上したが、ベースラインの NMT と比較して BLEU の向上は見られなかった。

\*NMT Model with an Integrated Framework of Large Vocabulary Translation and Reducing Untranslated Content

<sup>†</sup>Ryuichiro Kimura, Zi Long, Takehito Utsuro, Mikio Yamamoto, Graduate School of Systems and Information Engineering, University of Tsukuba

<sup>‡</sup>Shohei Iida, School of Engineering, Tokyo Denki University

<sup>§</sup>Tomoharu Mistuhashi, Japan Patent Information Organization

表 2: 自動評価の結果 (BLEU)

モデル	ja → en
ベースライン SMT [4]	32.3
ベースライン NMT	38.2
大規模語彙に対応した NMT モデル	39.8
大規模語彙に対応した NMT モデル + リランキング	38.0

## 4.2 人手評価

[7]における人手評価尺度である「一対評価」および「JPO 基準に基づく絶対評価」を用い、評価用対訳特許文から無作為に抽出された 100 文を評価対象として、著者が評価を行った。評価対象手法、および、ベースラインとなる手法との間の「一対評価」では、評価対象手法による翻訳精度が、ベースラインとなる手法による翻訳精度を上回った文の数を  $W$ 、評価対象手法による翻訳精度を下回った文の数を  $L$ 、評価対象手法による翻訳精度が、ベースラインとなる手法による翻訳精度と同等となった文の数を  $T$  として、一対評価のスコア (値の範囲は、 $-100 \sim 100$ ) を次式で定義する。

$$\text{score} = 100 \times \frac{W - L}{W + L + T}$$

「JPO 基準に基づく絶対評価」においては、JPO 評価基準<sup>1</sup>に基づき、各翻訳文に対して人手で 1 ~ 5 の値の範囲のスコアを付与し、その平均を「JPO 基準に基づく絶対評価」のスコアとする。「一対評価」結果を表 3 に、「JPO 基準に基づく絶対評価」結果を表 4 に示す。どちらの評価結果も大規模語彙対応とリランキングの併用手法が最も高かった。

表 3: JPO 基準に基づく絶対評価の結果 (スコアの範囲は 1 以上 5 以下)

モデル	ja → en
ベースライン NMT	4.1
大規模語彙に対応した NMT	4.2
大規模語彙に対応した NMT + リランキング	4.5

表 4: ベースライン NMT に対する一対評価の結果 (スコアの範囲は -100 以上 100 以下)

モデル	ja → en
大規模語彙に対応した NMT	18
大規模語彙に対応した NMT + リランキング	37

<sup>1</sup>[https://www.jpo.go.jp/shiryuu/toushin/chousa/pdf/tokkyohonyaku\\_hyouka/01.pdf](https://www.jpo.go.jp/shiryuu/toushin/chousa/pdf/tokkyohonyaku_hyouka/01.pdf)

## 5 訳抜け削減における効果の検証

### 5.1 訳抜け検出

後藤らは、訳抜けした内容を検出する翻訳スコアに基づいてリランキングすることによって、BLEU に関して翻訳評価結果を改善するだけでなく、NMT の訳抜けした内容を検出できることを明らかにした [3]。訳抜けした内容を検出する翻訳スコアの中から、我々は特に影響の大きい逆翻訳確率を採用した。より具体的には、ベースラインの NMT モデルを用いて、大規模語彙に対応した NMT モデルとベースライン NMT の両方の逆翻訳確率を評価した。逆翻訳確率に基づいて訳抜けした内容の多寡を予測した結果、大規模語彙に対応した NMT モデルはベースライン NMT と比較して改善した。

#### 5.1.1 逆翻訳確率

逆翻訳確率は、MT 出力からその入力文へ強制的に翻訳し直すときの確率として定義される。入力単語の内容が MT 出力に欠落している場合、入力単語の逆翻訳確率は小さくなると予想される。この関係から、逆翻訳確率を訳抜けした内容を検出するための手掛かりとして使用した。 $n$ -best MT 出力  $\mathbf{y}^d$  ( $1 \leq d \leq n$ ) の、入力単語  $x_j$  ( $1 \leq j \leq N$ ) についての逆翻訳確率スコア (BT-P)  $b_j^d$  を次のように定義する。

$$b_j^d = -\log p(x_j | x_{<j}, \mathbf{y}^d)$$

BT-P はベースライン NMT モデルと、大規模語彙に対応した NMT モデルの両方について、フレーズトークンのない訓練セットで英語から日本語へ訓練したベースライン NMT モデルを用いて計算する。この BT-P の公式化においては、「翻訳の存在」という以下の仮定を用いる。

**仮定: 翻訳の存在** 任意の入力単語の翻訳  $x_j$  ( $1 \leq j \leq N$ ) は、目的言語側に対訳が全く存在しない場合を除いて、 $n$ -best MT 出力  $\mathbf{y}^d$  ( $1 \leq d \leq n$ ) のどこかに存在する。

したがって、 $n$ -best MT 出力のうち最小のスコアをもつ出力は  $x_j$  の内容を最も含んでいる確率が高いと考えられ、 $\min_{1 \leq d' \leq n} b_j^{d'}$  は  $x_j$  の内容を含む出力のスコアと考えることができる。

次に、 $\mathbf{y}^d$  から  $x_j$  の内容が欠落したスコアとして、逆翻訳確率の比に基づく BT-R スコア  $q_j^d$  を考え、次のように定義する。

$$q_j^d = b_j^d - \min_{1 \leq d' \leq n} b_j^{d'}$$

これは、各 MT 出力のスコアと、最も  $x_j$  の内容を含む確率が高いもののスコアである  $n$ -best MT 出力の最小スコアとの差である。

最後に、このスコアを入力文のすべての入力単語  $\mathbf{x} = (x_1, \dots, x_N)$  について足し合わせることで、入力文  $\mathbf{x}$  に対する MT 出力  $\mathbf{y}^d$  の逆翻訳確率の比に基づく BT-R スコアは次のように得られる。

$$\text{BT-R}(\mathbf{x}, \mathbf{y}^d) = \sum_j q_j^d$$

本研究では BT-R スコアを  $n$ -best MT 出力それぞれで計算し、もっとも値の小さい出力を最終的な翻訳文としてランキング結果として採用した。

### 5.1.2 訳抜け検出結果

評価文についての逆翻訳確率の比に基づく BT-R スコアを測定した。ベースライン NMT モデルと、大規模語彙に対応した NMT モデルの、評価文における BT-R スコアの平均をそれぞれ表 5(a) に示す。大規模語彙に対応した NMT モデルは、ベースライン NMT モデルよりも低い BT-R スコアを達成した。この結果から、提案された NMT モデルによる MT 出力は、ベースライン NMT モデルによる MT 出力よりも訳抜けした内容が少なく考えられる。次に、各評価文  $\mathbf{x}$  に対して、大規模語彙に対応した NMT モデルとベースライン NMT モデルとの BT-R スコアの差を測定する。BT-R( $\mathbf{x}, \mathbf{y}^d$ ) を大規模語彙に対応した NMT モデルによる MT 出力の BT-R スコア、BT-R( $\mathbf{x}, \mathbf{y}'^d$ ) をベースライン NMT モデルによる MT 出力の BT-R と定義すると、BT-R スコアの差は以下のように定義される。

$$\text{BT-R}(\mathbf{x}, \mathbf{y}^d) - \text{BT-R}(\mathbf{x}, \mathbf{y}'^d)$$

評価文に対する BT-R スコアの差を表 5(b) に示す。評価文の 57.6% について、大規模語彙に対応した NMT モデルの BT-R スコアはベースライン NMT と比べて小さかった。また、ベースライン NMT の BT-R スコアが大規模語彙に対応した NMT モデルと比べて 5 以上小さかったものは 14.7% であったが、大規模語彙に対応した NMT モデルの BT-R スコアがベースライン NMT と比べて 5 以上小さかったものは 25.6% で、10.9% 大きかった。

表 5: 評価文における訳抜け予測の評価

(a) 評価文における文ごとの BT-R スコアの平均

モデル	ja → en
ベースライン NMT	16.3
大規模語彙に対応した NMT モデル	14.0
大規模語彙に対応した NMT モデル + リランキング	6.9

(b) 評価文におけるベースライン NMT と大規模語彙に対応した NMT モデルの間の BT-R スコアの差 (BT-R( $\mathbf{x}, \mathbf{y}^d$ ) - BT-R( $\mathbf{x}, \mathbf{y}'^d$ )) (%) の分布

< 0					> 0				
< -20	-20 ~ -10	-10 ~ -5	-5 ~ -1	-1 ~ 0	0 ~ 1	1 ~ 5	5 ~ 10	10 ~ 20	> 20
4.9	8.4	12.3	19.1	12.9	12.9	14.8	8.0	4.4	2.3
57.6					42.4				

(c) 評価文におけるベースライン NMT とリランキングした大規模語彙に対応した NMT モデルの出力の間の BT-R スコアの差 (BT-R( $\mathbf{x}, \mathbf{y}^d$ ) - BT-R( $\mathbf{x}, \mathbf{y}'^d$ )) (%) の分布

< 0					> 0				
< -20	-20 ~ -10	-10 ~ -5	-5 ~ -1	-1 ~ 0	0 ~ 1	1 ~ 5	5 ~ 10	10 ~ 20	> 20
11.8	20.3	23.0	27.7	10.5	3.1	2.9	0.4	0.3	0.0
93.3					6.7				

表 6: 評価文 1,000 文における、入力日本語文中の単語訳抜け数の人手評価

(a) 入力日本語文中の単語訳抜け数

モデル	ja → en
ベースライン NMT	73
大規模語彙に対応した NMT モデル	51
大規模語彙に対応した NMT モデル + リランキング	31

(b) 単語訳抜け数の分布 (%)

モデル	単語訳抜け数										
	0	1	2	3	4	5	6	7	8	9	≥ 10
ベースライン NMT	64	24	6	2	0	1	1	0	1	0	1
大規模語彙に対応した NMT モデル	74	14	4	4	3	1	0	0	0	0	0
大規模語彙に対応した NMT モデル + リランキング	83	10	3	3	0	0	1	0	0	0	0

## 5.2 入力日本語文中の単語訳抜け数の人手評価

無作為に選んだ 100 文の評価文について、日本語から英語への翻訳タスクにおいて、英語に翻訳されない入力日本語文中の単語の数を数えた。表 5(a) に示すように、大規模語彙に対応した NMT モデルを逆翻訳確率でリランキングすることで、翻訳されなかった単語の数はベースライン NMT と比較して約 40 % 程度に減少した。表 6(b) に翻訳されなかった単語の数の分布を示す。大規模語彙に対応した NMT モデルは、MT 出

表 7: 大規模語彙に対応した NMT モデルによって、訳抜けした内容が削減された具体例 (下線部は訳抜けした内容の部分を示す)

(a) 英語参照訳

英語参照訳	in an air bag device for a driver , it is favorable for a driver that the air bag 16 quickly and largely extends vertically and horizontally ( rightward and leftward in fig . 2 ) .
-------	--

(b) ベースライン NMT モデルによる翻訳結果

入力日本語文	運転席用エアバッグ装置においては、エアバッグ 16 が運転手にとって上下及び左右方向 (第 2 図の左右方向) に大きく且つ素早く展開することが好ましい。
MT 出力	in the air bag apparatus for the driver seat , it is preferable that the air bag 16 is in a large and left direction ( left and right in fig . 2 ) , and is rapidly deployed and developed .
BT-R スコア	34.7

(c) 大規模語彙に対応した NMT モデルによる翻訳結果

入力日本語文	運転席用エアバッグ装置においては、エアバッグ 16 が運転手にとって上下及び左右方向 (第 2 図の左右方向) に大きく且つ素早く展開することが好ましい。
フレーズトークンを含む MT 出力	in the $T_1^t$ , it is preferable that the $T_2^t$ 16 is larger and quickly in the $T_3^t$ and $T_4^t$ ( the $T_5^t$ of fig . 2 ) .
MT 出力	in the driver 's seat airbag device , it is preferable that the airbag 16 is larger and quickly in the driver and lateral direction ( the lateral direction of fig . 2 ) .
BT-R スコア	30.1

(d) 大規模語彙に対応した NMT モデルとランキングの併用による翻訳結果

入力日本語文	運転席用エアバッグ装置においては、エアバッグ 16 が運転手にとって上下及び左右方向 (第 2 図の左右方向) に大きく且つ素早く展開することが好ましい。
フレーズトークンを含む MT 出力	in the $T_1^t$ , it is preferable that the $T_2^t$ 16 is larger and quickly in the $T_3^t$ and $T_4^t$ ( the $T_5^t$ of fig . 2 ) .
MT 出力	in the driver 's seat airbag device , it is preferable that the airbag 16 is largely and quickly developed to the upper and lower lateral direction ( the lateral direction of fig . 2 ) for the driver .
BT-R スコア	22.4

力内の翻訳されなかった単語を減少させた。これは、翻訳されなかった単語の多くは語彙外で、ベースライン NMT によっては翻訳できないためと考えられる。

提案されたシステムは、フレーズの一部としてそれらの語彙外単語を抽出し、NMT のデコーディングの前にそれらのフレーズをトークンに置き換え、抽出したフレーズは SMT によって翻訳して MT 出力に挿入するため、語彙外の単語が確実に翻訳される。

表 7 に大規模語彙に対応した NMT モデルを逆翻訳確率でランキングした翻訳例とベースライン NMT の翻訳例の比較を示す。大規模語彙に対応した NMT モデルは、ベースライン NMT モデルと比較して、翻訳されなかった単語を減少させる。

## 6 おわりに

本論文では、大規模語彙に対応した NMT モデル [5] に基づき、訳抜けした内容の削減という観点から、提案された NMT モデルの効果をさらに検証した。日本語から英語への翻訳タスクにおいて、訳抜けした内容検出の結果を示し、英語に翻訳されない入力日本語文中の単語の数を調査した。大規模語彙に対応した NMT モデルは、ベースライン NMT モデルと比較して訳抜けした内容が減少していることを確認した。また、[3] で報告された逆翻訳確率に基づくランキングを大規模語彙に対応した NMT モデルに適用することで、BLEU において改善は見られないものの、訳抜けした内容がより減少することを確認した。今後は、提案された NMT モデルをサブワード単位に基づくもの [8] と比較する。

## 参考文献

- [1] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *Proc. 3rd ICLR*, 2015.
- [2] A. Fujii, M. Utiyama, M. Yamamoto, and T. Utsuro. Toward the evaluation of machine translation using patent information. In *Proc. 8th AMTA*, pp. 97–106, 2008.
- [3] I. Goto and H. Tanaka. Detecting untranslated content for neural machine translation. In *Proc. 1st NMT*, pp. 47–55, 2017.
- [4] P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst. Moses: Open source toolkit for statistical machine translation. In *Proc. 45th ACL, Companion Volume*, pp. 177–180, 2007.
- [5] Z. Long, T. Utsuro, T. Mitsuhashi, and M. Yamamoto. Translation of patent sentences with a large vocabulary of technical terms using neural machine translation. In *Proc. 3rd WAT*, pp. 47–57, 2016.
- [6] M. Luong, I. Sutskever, O. Vinyals, Q. V. Le, and W. Zaremba. Addressing the rare word problem in neural machine translation. In *Proc. 53rd ACL*, pp. 11–19, 2015.
- [7] T. Nakazawa, H. Mino, I. Goto, G. Neubig, S. Kurohashi, and E. Sumita. Overview of the 2nd workshop on Asian translation. In *Proc. 2nd WAT*, pp. 1–28, 2015.
- [8] R. Sennrich, B. Haddow, and A. Birch. Neural machine translation of rare words with subword units. In *Proc. 54th ACL*, pp. 1715–1725, 2016.