

段落単位での対訳データによる大規模言語モデルの翻訳精度向上

近藤海夏斗¹ 宇津呂武仁¹ 永田昌明²

¹ 筑波大学大学院 システム情報工学研究群 ² NTT コミュニケーション科学基礎研究所
s2320743@u.tsukuba.ac.jp utsuro@iit.tsukuba.ac.jp masaaki.nagata@ntt.com

概要

大規模言語モデルは、自然言語処理タスクで優れた性能を発揮しているが、翻訳タスクではモデルサイズや訓練手法による性能差が見られる。本研究では、段落単位の対訳データを活用した継続事前訓練および Supervised fine-tuning が翻訳精度に与える効果を検証した。Llama-3 ベースのモデルを用いて英日・日英翻訳を評価した結果、段落単位の対訳データを継続事前訓練と Supervised Fine-Tuning の両方で利用することで最も高い精度が得られた。また、推論時に段落全体を一括で翻訳する方法が有効であることを確認した。

1 はじめに

GPT-4 や Llama のような大規模言語モデルが、要約や質問応答などの自然言語処理タスクで目覚ましい成果を収めている。機械翻訳においては、GPT-4 や Claude-3 のようなクローズドな大規模言語モデルの翻訳が、既存の翻訳モデルより人手評価が高いことが報告されている [1, 2]。一方で、Llama-2-7B や 13B のような、パラメータ数が 100 億前後の大規模言語モデルでは、既存の翻訳モデルより翻訳精度が大きく劣ることが報告されている [3]。この問題に対し、継続事前訓練のあとに Supervised fine-tuning を行うという訓練手法が提案された [4, 5, 6, 7]。この手法を適用することで、パラメータ数が 100 億前後の大規模言語モデルが、GPT-3.5 や 4 と同等の精度になることが明らかとなった。しかし、この精度は文単位の対訳データをもとにしており、文脈を考慮する必要がある段落単位の翻訳においても有効かは明らかとなっていない。

そこで本研究では、既存の対訳データから 1 文以上 9 文以下の段落単位の対訳データを作成し、段落単位の対訳データを継続事前訓練もしくは Supervised fine-tuning の訓練データに含めて、翻訳精

度の比較を行った。Llama-3 ベースの大規模言語モデルを用いて、英日・日英翻訳について実験した結果、以下のことが明らかとなった。

- 段落単位の対訳データは、継続事前訓練および Supervised fine-tuning の両方に用いた場合が最も精度が高い。
- 段落単位の対訳データを訓練データに含めた場合、1 文ずつ推論するより、段落単位をそのまま入力して推論した場合のほうが、精度が向上する。

2 関連研究

機械翻訳のコンペティションである WMT23 および WMT24 [1, 2] では、GPT-4 のようなクローズドな大規模言語モデルが、既存の翻訳モデルを人手評価で上回る結果が報告されている。一方で、パラメータ数が 100 億前後の大規模言語モデルにおける 8-shot では、既存の翻訳モデルと比較して翻訳精度が大幅に劣ることが指摘されている [3]。これを克服する方法として、継続事前訓練のあとに Supervised Fine-Tuning を行う 2 段階の訓練手法が提案されており [4, 7]、特に Llama-2 のように英語を主に事前訓練された LLM に対しては、単言語データと対訳データの両方を用いた継続事前訓練が翻訳精度向上に有効であることが示されている [5, 6, 8]。

ただし、これらの研究は主に文単位の対訳データを基にしており、段落単位や文書単位のような複数文を対象とした翻訳における有効性については十分に検証されていない。一方で、Karpinska ら [9] は、文学作品から作成した 18 言語対の段落単位の対訳データを用いて、GPT-3.5 と Google 翻訳の翻訳を人手で評価した。その結果、GPT-3.5 の翻訳がより高く評価され、また、段落全体を一括して入力する翻訳手法が最も高い評価を得たことが明らかとなった。

さらに、Llama-2 を用いて段落単位や文書単位の

対訳データを訓練データに取り入れる研究 [10, 11] も行われているが、段落単位の対訳データが翻訳精度に与える具体的な影響については、依然として十分な議論が行われていない。以上を踏まえ、本研究では、継続事前訓練および Supervised fine-tuning の各フェーズにおいて段落単位の対訳データを含める最適な方法を明らかにすることを目的とする。

3 提案手法

本論文では、Llama のような主に英語で事前訓練されている LLM で段落単位の翻訳精度を高める手法として、3 段階の訓練を提案する。3 段階の訓練は、2 段階の継続事前訓練、Supervised Fine-Tuning で構成されている。

3.1 単言語データによる継続事前訓練

主に英語で事前学習された LLM は、英語から英語以外への翻訳精度が低いことが報告されている [4]。そのため、英語以外の言語の生成能力を向上させるために、単言語データを用いた継続事前訓練を行う。継続事前訓練では、大規模言語モデルが出力する全ての次単語予測の損失を計算する。

3.2 文単位・段落単位の対訳データによる継続事前訓練

単言語データによる継続事前訓練を行ったあと、対訳データを用いた継続事前訓練を行う。近藤ら [7] の知見に基づき、`{ 原言語文 }\n{ 目的語文 }` のように、原言語文の後にその翻訳が続く形式のデータを使用する。また、対訳データとして、文単位と段落単位の両方を用いることで、段落単位の翻訳精度の向上を図る。

3.3 文単位・段落単位の対訳データによる Supervised Fine-Tuning

単言語および対訳データによる継続事前訓練のあと、Supervised fine-tuning を行う。なお、Supervised fine-tuning では、プロンプト部分の次単語予測の損失を除外し、目的語文の次単語予測の損失のみを用いる。本研究では、先行研究 [7, 6] を参考に、以下のような原言語で書かれたプロンプトを使用する。

英日翻訳

Translate this from English to Japanese:
English: {source sentence}

Japanese:

日英翻訳

これを日本語から英語に翻訳してください:
日本語: {source sentence}
英語:

4 検証方法

4.1 概要

本論文では、提案手法の 1 段階目の訓練である単言語データを用いた継続事前訓練を行った LLM として、rinna/llama-3-youko-8b¹⁾ (以下 youko-8b という) を利用した。youko-8b は、日本語および英語の単言語データ 220 億トークンを用いて継続事前訓練されたモデルである。本研究では、この youko-8b に対し提案手法の 2 段階目以降の訓練を行い、英日および日英翻訳の評価を行った。

4.2 段落データの作成

対訳データのほとんどは、原言語文と目的語文がそれぞれ 1 文ずつ対応付けされている。一方で、対訳データの中には、文書の境界情報が添付されている場合があり、この境界情報から、1 文書が何行目から何行目までかを知ることができる。本研究では段落データとして、境界情報をもとに復元した文書対がそれぞれ 9 文以下の場合そのまま利用し、10 文以上の場合、先頭から 1~9 文までのランダムな文数で抜き出したものを利用した。また、境界情報が添付されていない対訳データについては、データの先頭から、1 から 9 文までのランダムな数ずつ抜き出したものを疑似段落データとして利用した。

4.3 データセット

継続事前訓練、Supervised fine-tuning、そしてテストデータに利用した対訳データについて述べる。なお、各データの統計情報は付録 A に記載する。

4.3.1 継続事前訓練

訓練データとして、JParaCrawl v3.0 [12]、Kyoto Free Translation Task (KFTT) [13]、TED Talks [14]、News Commentary [1]、Opensubtitles [15]、Globalvoices [16]、日英対訳文対応付けデータ (NICT align) [17]、NTREX-128 [18]、そして Flores-200 [19] の dev データを利用した。JParaCrawl v3.0 については、LEALLA-

1) <https://huggingface.co/rinna/llama-3-youko-8b>

large²⁾ [20] から得た原言語文と目的語文の文埋め込みベクトルのコサイン類似度が 0.4 以上 0.95 未満のものをサンプリングした 2,080 万文対を利用した。また、全ての対訳データは英日・日英それぞれ半分ずつランダムサンプリングして利用した。これにより、継続事前訓練で利用した対訳データは、youko-8b のトークナイザーで約 16 億トークンとなった。開発データは、WMT20 [14] の開発・テストデータ、そして WMT21 [21], WMT22 [22] のテストデータを利用した。

4.3.2 Supervised Fine-Tuning

Supervised fine-tuning のデータは品質が重要であることが報告されている [4] ため、対訳データの中でも、人手で翻訳された対訳データを利用した。訓練データとして、WMT20 の開発・テストデータ、WMT21 のテストデータ、NTREX-128、そして Flores-200 の dev データを利用した。開発データは、WMT22 のテストデータを利用した。これらのデータは全て、文書情報が添付されているため、文単位もしくは段落単位のデータとして利用可能である。

4.3.3 テストデータ

テストデータとして、WMT23 のテストデータ、そして Flores-200 の devtest データを利用した。両者ともに文書情報が添付されているため、段落単位の対訳データとして利用した。

4.4 比較モデル

(疑似) 段落データを継続事前訓練と Supervised fine-tuning のどちらに用いるのがよいかを検証するため、以下の 4 通りについて比較を行った。なお、JParaCrawl v3.0 はいずれの場合も文単位の対訳データとして利用し、KFTT, TED Talks, News Commentary については文書情報が添付されていないため、疑似段落データとして利用する。

1. (疑似) 段落情報を利用可能なデータも文単位とみなし、すべて文単位の対訳データで継続事前訓練した後に、文単位の対訳データで Supervised fine-tuning
2. (疑似) 段落情報を利用可能なデータも文単位とみなし、すべて文単位の対訳データで継続事前訓練した後に、段落情報を利用可能なデータを使って文・段落単位の対訳データで Supervised

fine-tuning

3. (疑似) 段落情報を利用可能なデータは (疑似) 段落単位、その他は文単位の対訳データで継続事前訓練した後、文単位の対訳データで Supervised fine-tuning
4. (疑似) 段落情報を利用可能なデータは (疑似) 段落単位、その他は文単位の対訳データで継続事前訓練した後、段落情報を利用可能なデータを使って文・段落単位の対訳データで Supervised fine-tuning

また、推論方法と翻訳精度との関係を明らかにするため、以下の 2 通りの推論を行った。

SENT2PARA

文単位で推論し、推論結果を文書情報をもとに結合する。したがって、1 段落にふくまれる文の数だけ推論する。

PARA2PARA

段落をそのまま入力し、一度に翻訳を出力させる。

4.5 ハイパーパラメータ

4.5.1 継続事前訓練

optimizer として AdamW [23] ($\beta_1 = 0.9, \beta_2 = 0.95, \epsilon = 1.0 \times 10^{-6}$) を使用し、weight decay を 0.1, gradient clipping を 1.0 とした。そして、context length を 2,048 とし、1 エポック訓練しながら、全ステップ数の 10% ごとに開発データの誤差を計算した。また、最大の学習率を 2.5×10^{-5} とし、warmup ratio を 1% の cosine scheduler とした。訓練は NVIDIA H100 を 8 台用いて、トータルバッチサイズを 1,024 とした。訓練中には、bfloat16 および deepspeed [24] ZeRO stage 2 を適用した。

4.5.2 Supervised Fine-Tuning

Supervised fine-tuning は、継続事前訓練にて開発データの誤差が最小となるモデルに対して行った。Supervised fine-tuning でも、optimizer として継続事前訓練と同様に AdamW を使用した。Weight decay と gradient clipping についても継続事前訓練と同じ値とした。最大の学習率を 2.5×10^{-6} とし、warmup ratio を 1% の cosine scheduler とした。また、8 台の NVIDIA H100 を用いて、トータルバッチサイズを 64 とし、エポック数は 2 とした。

2) <https://huggingface.co/setu4993/LEALLA-large>

表 1 WMT23 と Flores-200 の COMET スコア. 太字は, 各推論方法およびテストデータで最も高いスコアであることを示す. スコアの末尾にある"*"は, 同じ推論方法かつ, 継続事前訓練と Supervised fine-tuning のデータがどちらも文単位の場合と有意差 ($p < 0.05$) があることを示す.

推論方法	継続事前訓練	SFT	英日翻訳		日英翻訳		平均
			WMT23	Flores-200	WMT23	Flores-200	
SENT2PARA	文単位	文単位	0.8479	0.8989	0.7856	0.8551	0.8469
	文単位	単文・段落単位	0.8511*	0.8990	0.7972*	0.8561	0.8509
	単文・段落単位	文単位	0.8501	0.8992	0.7861*	0.8556	0.8478
	単文・段落単位	単文・段落単位	0.8522*	0.9006*	0.7999*	0.8579*	0.8527
PARA2PARA	文単位	文単位	0.8498	0.8956	0.7891	0.8551	0.8474
	文単位	単文・段落単位	0.8560*	0.9012*	0.8027*	0.8626*	0.8556
	単文・段落単位	文単位	0.8540	0.8999*	0.8003*	0.8591*	0.8533
	単文・段落単位	単文・段落単位	0.8579*	0.9009*	0.8066*	0.8628*	0.8571

4.5.3 推論

Supervised fine-tuning の開発データの誤差が最小となるモデルを用いて, greedy decoding で推論を行った. なお, vllm [25] を用いて高速化した.

4.6 評価指標

評価指標として, COMET [26] を使用した. COMET のモデルは wmt22-comet-da を使用した. なお, 有意差の判定は, 有意水準 5% ($p < 0.05$) で行った.

5 評価結果

5.1 段落単位の対訳データの影響

表 1 より, Flores-200 の英日翻訳を除き, 継続事前訓練と Supervised fine-tuning の両方に段落単位の対訳データを含めた場合が最も高いスコアとなった. したがって, 大規模言語モデルの段落翻訳精度を高めるには, 段落単位の対訳データは, 継続事前訓練と Supervised fine-tuning の両方に段落単位の対訳データを含めるとよいことがわかる. また, 継続事前訓練もしくは Supervised fine-tuning のいずれかのみに段落単位の対訳データを含めた場合でも, 継続事前訓練と Supervised fine-tuning の両方が文単位の場合より, 翻訳精度が有意に改善した. 翻訳の具体例は, 付録 B に記載する.

5.2 推論方法による影響

継続事前訓練と Supervised fine-tuning の両方も文単位の場合, 推論方法による違いは平均で 0.0005

ポイントのみであった. 一方で, 継続事前訓練と Supervised fine-tuning のいずれか, もしくはその両方に段落単位の対訳データが含まれる場合, 1 文ずつ推論すると平均で 0.04 ポイント以上 COMET スコアが減少した. このことから, 継続事前訓練もしくは SFT に段落単位の対訳データを含める場合, 段落をそのまま推論したほうが高い精度となることが明らかとなった.

6 おわりに

本研究では, Llama のような主に英語で事前訓練された大規模言語モデルを用いて, 段落単位の翻訳精度を向上させる手法として, 3 段階の訓練を提案した. 1 段階目に単言語データで継続事前訓練, 2 段階目に文単位・段落単位の対訳データで継続事前訓練, そして 3 段階目に文単位・段落単位の対訳データで Supervised fine-tuning を行う. 提案手法の有効性を検証するため, 段落単位の対訳データを, 継続事前訓練および Supervised fine-tuning に利用する場合と利用しない場合の 4 通りについて実験を行った. その結果, 段落データは, 継続事前訓練と Supervised fine-tuning の両方に用いて, 段落をそのまま推論する場合が最も精度が高いことが明らかとなった.

今後の展望として, オリジナルの Llama-3 を用いて, 第 1 段階目の単言語データによる継続事前訓練に必要なデータ量を明らかにすること, そして Llama-3 以外の大規模言語モデルでも同様の結果が得られるかを明らかにすることが挙げられる.

参考文献

- [1] T. Kocmi, et al. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In **Proc. of 8th WMT**, pp. 1–42, 2023.
- [2] T. Kocmi, et al. Findings of the WMT24 general machine translation shared task: The LLM era is here but MT is not solved yet. In **Proc. 9th WMT**, pp. 1–46, 2024.
- [3] W. Zhu, H. Liu, Q. Dong, J. Xu, S. Huang, L. Kong, J. Chen, and L. Li. Multilingual machine translation with large language models: Empirical results and analysis. In **Findings of NAACL**, pp. 2765–2781, 2024.
- [4] H. Xu, Y. Kim, A. Sharaf, and H. Awadalla. A paradigm shift in machine translation: Boosting translation performance of large language models. In **Proc. 12th ICLR**, 2024.
- [5] D. Alves, et al. Tower: An open multilingual large language model for translation-related tasks. In **Proc. 1st CoLM**, 2024.
- [6] J. Guo, H. Yang, Z. Li, D. Wei, H. Shang, and X. Chen. A novel paradigm boosting translation capabilities of large language models. In **Findings of NAACL**, pp. 639–649, 2024.
- [7] M. Kondo, T. Utsuro, and M. Nagata. Enhancing translation accuracy of large language models through continual pre-training on parallel data. In **Proc. 21st IWSLT**, pp. 203–220, 2024.
- [8] H. Xu, K. Murray, P. Koehn, H. Hoang, A. Eriguchi, and H. Khayrallah. X-alma: Plug & play modules and adaptive rejection for quality translation at scale. **arXiv:2410.03115**, 2024.
- [9] M. Karpinska and M Iyyer. Large language models effectively leverage document-level context for literary translation, but critical errors persist. In **Proc. 8th WMT**, pp. 419–451, 2023.
- [10] M. Kondo, R. Fukuda, X. Wang, K. Chousa, M. Nishimura, K. Buma, T. Kano, and T. Utsuro. NTTSU at WMT2024 general translation task. In **Proc. 9th WMT**, pp. 270–279, 2024.
- [11] R. Rei, et al. Tower v2: Unbabel-IST 2024 submission for the general MT shared task. In **Proc. 9th WMT**, pp. 185–204, 2024.
- [12] M. Morishita, K. Chousa, J. Suzuki, and M. Nagata. JParaCrawl v3.0: A large-scale English-Japanese parallel corpus. In **Proc. 13th LREC**, pp. 6704–6710, 2022.
- [13] G. Neubig. The Kyoto free translation task. <http://www.phontron.com/kfft>, 2011.
- [14] L. Barrault, et al. Findings of the 2020 conference on machine translation (WMT20). In **Proc. 5th WMT**, pp. 1–55, 2020.
- [15] P. Lison, J. Tiedemann, and M. Kouylekov. OpenSubtitles2018: Statistical rescoring of sentence alignments in large, noisy parallel corpora. In **Proc. 11th LREC**, 2018.
- [16] P. Prokopidis, V. Papavassiliou, and S. Piperidis. Parallel Global Voices: a collection of multilingual corpora with citizen media stories. In **Proc. 10th LREC**, pp. 900–905, 2016.
- [17] U. Masao and Mayumi T. English-japanese translation alignment data, 2003.
- [18] C. Federmann, T. Kocmi, and Y. Xin. NTREX-128 – news test references for MT evaluation of 128 languages. In **Proc. 1st SUMEval**, pp. 21–24, 2022.
- [19] NLLB Team, et al. No language left behind: Scaling human-centered machine translation. **arXiv:2207.04672**, 2022.
- [20] Z. Mao and T. Nakagawa. LEALLA: Learning lightweight language-agnostic sentence embeddings with knowledge distillation. In **Proc. 17th EAACL**, pp. 1886–1894, 2023.
- [21] F. Akhbardeh, et al. Findings of the 2021 conference on machine translation (WMT21). In **Proc. of the 6th WMT**, pp. 1–88, 2021.
- [22] T. Kocmi, et al. Findings of the 2022 conference on machine translation (WMT22). In **Proc. 7th WMT**, pp. 1–45, 2022.
- [23] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In **Proc. 7th ICLR**, 2019.
- [24] J. Rasley, S. Rajbhandari, O. Ruwase, and Y. He. Deep-speed: System optimizations enable training deep learning models with over 100 billion parameters. In **Proc. 26th ACM SIGKDD**, pp. 3505–3506, 2020.
- [25] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In **Proc. 29th ACM SIGOPS**, 2023.
- [26] R. Rei, J. C. de Souza, D. Alves, C. Zerva, A. Farinha, T. Glushkova, A. Lavie, L. Coheur, and A. Martins. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In **Proc. 7th WMT**, pp. 578–585, 2022.

表2 対訳データの文章の数および段落の数

データ	文章の数	段落の数
JParaCrawl v3.0	20,813,885	-
KFTT	440,286	85,956
TED Talks	241,738	47,193
News Commentary	1,887	368
Globalvoices	56,014	12,185
NICT align	110,908	22,227
Opensubtitles	2,083,600	418,449
NTREX-128	1,997	460
Flores-200 dev	997	282
WMT20 dev・test	5,989	1,333
WMT21 test	2,005	422
WMT22 test	4,045	893
WMT23 test En-Ja	2,074	429
WMT23 test Ja-En	1,992	437
Flores-200 devtest	1,012	279

表3 1段落あたりの平均文数, 平均単語数(英語), 平均文字数(日本語)

データ	平均文数	平均単語数(En)	平均文字数(Ja)
KFTT	5.0	40.2	101.2
TED Talks	5.0	44.2	74.5
News Commentary	5.0	106.7	241.9
Globalvoices	4.6	84.9	232.0
NICT align	5.0	62.1	163.8
Opensubtitles	5.0	32.4	64.8
NTREX-128	4.3	91.4	242.7
Flores-200 dev	3.5	74.3	193.8
WMT20 dev・test	4.4	102.9	242.1
WMT21 test	4.8	113.9	261.3
WMT22 test	4.5	70.7	171.8
WMT23 test En-Ja	4.8	81.8	217.4
WMT23 test Ja-En	4.6	76.8	153.8
Flores-200 devtest	3.6	78.5	204.1

A 対訳データの統計情報

表2は, 利用した対訳データの文章の数および段落の数である。また, 表3は, 1段落あたりの平均文数, 英語の平均単語数, そして日本語の平均文字数をまとめたものであり, 英語の単語数は半角スペースごとに文を区切り算出した。なお, KFTT, TED Talks, そして News Commentary には文書情報が添付されていないため, 段落の数は疑似段落であることを注意である。

表4 WMT23の PARA2PARA 推論の英日翻訳の具体例。CPTは継続事前訓練, SFTはSupervised Fine-Tuningを表す。

原言語文	It is time to fight back. Go out and protest now. Get some bumper stickers that are anti-Trump and plaster them all over. Yell back at those who would try to talk over us. This is not normal. This is not OK...
参照文	反撃の時だ。外に出て抗議しよう。反トランプのバンパーステッカーを買って、そこら中に貼ろう。議論をふっかけて来るやつには大声で怒鳴り返そう。こんなの普通じゃない。こんな的大丈夫じゃない...
文単位 CPT+文単位 SFT	これは異常だ。これは許されない。
文単位 CPT+文・段落単位 SFT	これは反撃する時だ。外に出て今すぐ抗議しよう。トランプに反対するステッカーを手に入れて、至る所に貼ろう。私たちが押しのけようとする人たちに大声で叫ぼう。これは普通じゃない。これは許されない...
文・段落単位 CPT+文単位 SFT	これは反撃する時だ。外に出て今すぐ抗議しよう。トランプに反対するステッカーを買って貼りまくろう。私たちが押しのけようとする人たちに叫ぼう。これは普通じゃない。これは許されない...
文・段落単位 CPT+文・段落単位 SFT	これからは反撃の時だ。外に出て今すぐ抗議しよう。トランプに反対するステッカーを買って、貼りまくろう。私たちが押しのけようとする人たちに大声で叫ぼう。これは普通じゃない。これは許されない...

B 翻訳の具体例

表4は, WMT23の英日翻訳の具体例である。継続事前訓練およびSFTの両方が文単位の場合, 1文目と最後の文のみ翻訳され, 中間の翻訳が抜けている。一方, 継続事前訓練もしくはSFTのいずれかに段落単位の対訳データが含まれる場合, 中間の翻訳も出力されている。したがって, 段落単位の対訳データを含めることで, 訳抜けしづらい翻訳モデルになることが示唆される。