

事例ベース NER を用いた中古日本語文のニューラル機械翻訳

Neural Machine Translation of Early Middle Japanese Sentences using Example-Based Named Entity Recognition

中城 裕之[‡] 福原 義久[‡]
Hiroyuki Nakajo Yoshihisa Fukuhara

1. はじめに

一般的に、古文書を読むには歴史学についての専門的な知識が必要とされている。各地の資料館や歴史博物館では実際に古文書が展示されていることも多いが、歴史学を専門的に学んでいない一般層がそれらを詳細に読解することは難しい。古文書のくずし字を光学文字認識で文字起こしするサービスは実用化されており[1][2]、次はそれを翻訳し現代語訳を提供するサービスの実用化が待たれる。また、歴史学分野では近年通説を覆すような新たな発見が次々に行われているが、その一方で歴史学者の不足という課題にも直面している[3]。古文書の翻訳作業は歴史学の研究発展に非常に重要であり、歴史学者不足によりその作業が滞ることも懸念される。古文翻訳システムを実用化し、歴史学者がポストエディットのような形でそのシステムを活用することができれば、翻刻作業の効率化を図ることもできるだろう。

しかしながら、日本語古文のニューラル機械翻訳に関する研究は外国語翻訳と比較して進歩が遅いのが現状である。古文翻訳モデルとして現時点で最も優秀なものは、星野ら[4]が実装した統計的機械翻訳を用いたモデルである。この手法は深層学習を用いない翻訳モデルであり現在では最新の手法ではないが、古文機械翻訳研究の分野では最も高いスコアを残している。また、古文のニューラル機械翻訳としては、臼井ら[5]が実装した大規模な事前学習済みモデルである日本語 T5 の Fine-tuning を用いた研究がある。この手法は外国語翻訳でも使用されている最先端のものだが、古文では統計的機械翻訳を上回ることはできなかった。古文の翻訳システムを実用化するためには、さらに翻訳精度を向上させる必要がある。

高精度なニューラル機械翻訳モデル構築のためには大規模なパラレルコーパスが必要とされているが、現時点で日本語古文のパラレルコーパスは小規模なものしか存在しない。そのため、Fine-tuning 等を駆使しても精度向上には限界があるというのが現状である。一方、パラレルコーパスの規模に依存せずに翻訳モデルの精度を向上させようという研究も存在する。その一つが、文中の固有表現情報を活用した翻訳モデルだ。固有表現は文の中核的な意味を持つ場合が多くより正確に翻訳することが求められるが、翻訳においては学習データ中の出現回数の少なさから未知語となる場合が多い。そのため、原文から固有表現を抽出し、その情報を失わず翻訳文に適用するという研究がニューラル機械翻訳分野では行われており、それらは古文翻訳モデルにも適していると考えられる。

古文のようにデータセットが小規模な分野の固有表現抽出手法として、Ziyadi ら[7]が提案した、事例ベース固有表現抽出(Example-Based Named Entity Recognition)という手法

[‡] 武蔵野大学データサイエンス学部

Musashino University Faculty of Data Science

がある。この手法では、Few-shot 学習を用いることによって 10~100 件程度のデータセットで十分な精度の固有表現抽出を行うことができることが示されている。また、苑ら[8]はこの手法を用いて近代語文の固有表現抽出を行った。この手法は大規模なデータセットを用意できない古文にとって効果的であると言え、より古い時代の文章でも効果を発揮することが期待される。

本論文では、固有表現抽出を活用した中古日本語文の機械翻訳手法を提案する。中古日本語文とは日本語の発展段階の一つであり、平安時代中期頃に用いられたものである[9]。現代の日本語とは語彙などに大きな差があり、固有表現抽出を用いた機械翻訳の対象として適していると考えた。固有表現抽出には前述の事例ベース固有表現抽出手法を用いる。

本論文の構成は以下の通りである。2 節では、本論文の関連研究について示す。3 節では、事例ベース固有表現抽出を用いた中古日本語文のニューラル機械翻訳手法について示す。4 節では、3 節にて示した手法に基づく実験概要について示す。5 節では、実験の結果について示す。6 節では実験結果の考察について述べる。7 節では本論文のまとめを述べる。8 節では本論文の今後の展望を述べる。

2. 関連研究

本節では、本論文に関連する研究を紹介する。

2.1 古文の固有表現抽出に関する研究

吉村ら[10]は、古文テキストに対し文字の出現頻度と文字列の出現確率に基づいた単語分割を行い、その結果に対して Support Vector Machine(以下 SVM)による人物表現の抽出を行なった。実験の結果、単語分割を行なった場合は単語分割を用いなかった場合に比べ F1 スコアで約 4%精度が向上するという結果を示した。

永井ら[11]は、吉村ら[10]の手法を基に、単語分割に形態素解析辞書である中古和文 UniDic[12]を用いて江戸時代の書籍に対し人物表現抽出を行なった。提案手法では、まず学習データである古文テキストに対して単語分割を行い、それぞれの文字に単語分割情報や前後の文字情報などの素性の付与を行う。次に SVM を用いて学習データから分類モデルを構築し、テキストデータの文字をそれぞれの素性に分類する。最後に分類したラベルから人物表現の推定・抽出を行なった。実験の結果、F1 スコアで 0.91 という結果を示した。

臼井ら[13]は、深層学習モデルである BiLSTM-CRF を用いた近代語文の固有表現抽出を行った。BiLSTM-CRF は、双方向 LSTM と条件付き確率場から構成されるモデルであり、固有表現抽出を含む系列ラベリングタスクにおいて一定規模の学習データが存在する場合に高い精度を発揮するとされている。実験では、時代や日時といった固有表現タ

イブに対して F1 スコアで 0.9 以上という結果を示した。特に時代や日付、年齢の固有表現タイプに対しては 0.95 以上という高いスコアを示している。

苑ら[8]は、事例ベース固有表現抽出手法を用いて近代語文の固有表現抽出を行い、F1 スコアで最大 0.699 という結果を得た。この結果は十分な量の学習データを使用したモデルには劣るが、学習データの少ない古文に対してもこの手法であれば一定の精度で固有表現抽出を行うことができることが示された。また、実験によって得られた固有表現情報を活用し、可視化手法の比較と分析も行なった。比較対象として NetworkX[14], Multiplex[15], pyvis[16]を選び、それぞれのツールで抽出した人物情報を可視化した。また、可視化ツールとして最も効果が高かった pyvis を用いて人間関係の社会ネットワークを分析した。さらに GraphGPT[17]を用いて人物解説文の可視化を行い、大規模生成モデルによる言語情報可視化の可能性を示した。

2.2 古文の機械翻訳に関する研究

星野ら[4]は、段落単位で対応づけられた対照コーパスをルールベースのスコア関数を用いて分割することでパラレルコーパスに変換し、統計的機械翻訳を行なった。この手法は、文中の句点数が一致しないセグメントでもパラレルコーパスとして活用することができ、データ数確保の点で既存の手法よりも優れている。実験の結果、提案手法の BLEU スコアは対照コーパスをそのまま学習データとして使用した場合や句点による分割を行なった場合と比較して 2.5 ポイント程度上昇した。提案手法は、人工頭脳プロジェクト「ロボットは東大に入れるか」における国語古文問題解答にも利用されている[18]。

高久ら[19]は、古代から近代までの文が含まれる通時的パラレルコーパスでは時代ごとの言語変化に対応できないという問題に着目し、現代文コーパスで事前学習された単語分散表現を時代ごとに分割した古文コーパスで Fine-tuning する手法を提案した。提案手法では、適応元と適応先の単語分散表現の意味のずれが最小になるように近代から時代を遡るように順番に Fine-tuning を行い、得られた単語分散表現を翻訳モデルの初期化に用いた。翻訳モデルは LSTM を用いて構築した。実験の結果、提案手法は Fine-tuning を用いなかった LSTM モデルの BLEU スコアを上回った。なお、この研究が古文のニューラル機械翻訳を実現した最初の事例である。

白井ら[5]は、大規模事前学習済みモデルである日本語 T5 モデルの Fine-tuning により翻訳モデルを構築し、精度を検証した。T5 の Fine-tuning による翻訳モデル構築は複数の事例があり、小規模パラレルコーパスを用いたものも存在している。実験の結果、提案手法は古文のニューラル機械翻訳モデルの中では最も高い BLEU スコアを示した。また、白井らは逆翻訳による擬似パラレルコーパスを用いた実験も行ったが BLEU スコアは低下し、逆翻訳を用いる手法は有効でないことを示した。

3つの研究の翻訳精度は表 1 に示す通りである。

表 1 古文機械翻訳モデルの BLEU スコア

モデル:	星野	高久	白井	白井(逆翻訳)
スコア:	28.03	19.95	27.73	23.81

2.3 固有表現抽出を用いた翻訳に関する研究

Li ら[20]は、原言語と目的言語間で対応付いた固有表現を別の記号に置き換えた上で翻訳を行う手法を提案した。この手法はニューラル機械翻訳システムにおける未知語が翻訳性能を低下させるという問題への解決策の一つとして提案され、実験結果では固有表現抽出を用いない手法と比べて BLEU スコアで 2.9 ポイントの改善が見られた。

Siekmeier ら[21]は、ニューラル機械翻訳モデルのエンコーダ及びデコーダに固有表現埋め込み層を加え、それらを単語埋め込み層に加算した結合埋め込みにより学習を行う手法を提案した。この手法では、デコーダは目的言語のトークンと固有表現タグ情報を同時に予測する。実験の結果、固有表現抽出を用いない手法と比べて BLEU スコアで 1.61 ポイントの改善が見られた。

南端ら[22]は、学習データの目的言語文に固有表現抽出を適用し、固有表現の前後にタグを付与して翻訳を行う手法を提案した。この手法では、原言語と目的言語間の対応付けに関係なく固有表現を活用することができ、さらに固有表現の構成単語を残したまま翻訳モデルを構築することができるという点において従来手法よりも優れている。実験の結果、目的言語文の固有表現を考慮した翻訳モデルは、考慮しなかった翻訳モデルよりも BLEU スコアが改善し、提案手法が有効であることを示した。

2.4 本論文の位置付け

本論文では、中古日本語文を対象とした固有表現抽出とそれを活用したニューラル機械翻訳を行う。固有表現抽出には Few-shot 学習による手法である事例ベース固有表現抽出を用いる。また、ニューラル機械翻訳では Transformer をモデルとして採用し、固有表現抽出を用いなかった場合、パラレルコーパスの中古日本語文に固有表現抽出を行なった場合、現代語訳文に固有表現抽出を行なった場合、どちらにも固有表現抽出を行なった場合の 4 つを比較する。本論文では、固有表現抽出を用いたニューラル機械翻訳手法が古文の機械翻訳においても有効であることを示す。

3. 事例ベース固有表現抽出を用いた中古日本語文のニューラル機械翻訳

本節では、固有表現抽出および機械翻訳の具体的な手法について示す。

3.1 事例ベース固有表現抽出

本論文では、事例ベース固有表現抽出手法を用いて中古日本語文を対象に固有表現抽出を行う。モデルを図 1 に示す。モデル図は Ziyadi ら[7]が作成したものを参考にしている。この手法は、固有表現の大規模オープンデータで Fine-tuning した BERT をさらに Few-shot 学習することによってシステムを構築している。具体的には、まざクエリ(入力文)とサポート例(Few-shot 学習データ)を 2 つの事前学習済み BERT に別々に送信し、分散表現を獲得する。次に、クエリ内の各トークンの開始位置・終了位置と各サポート例の固有表現の開始位置・終了位置との類似度を計算する。この計算によって各トークンが固有表現の開始位置・終了位置である確率を計算し、固有表現の位置を推定する。最後に、特定された固有表現がどのカテゴリの固有

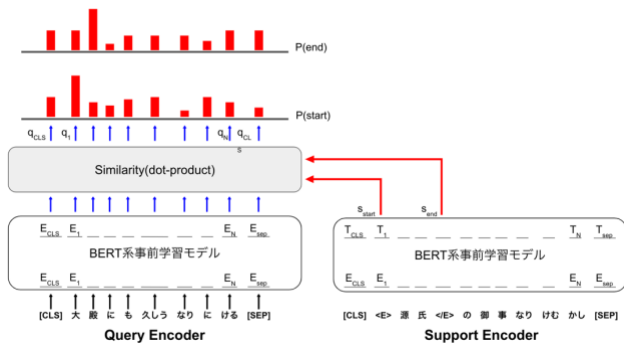


図 1 事例ベース固有表現抽出

表現であるかを識別するというものである。ベースとなる BERT は、苑ら[8]の実験で特に高いスコアを示した deberta-base-japanese-wikipedia[23], roberta-base-japanese-aozora-char[24]の 2 つを用い、それらを日本語現代文の固有表現抽出データセットで Fine-tuning したものを中古日本語文で Few-shot 学習しモデルを構築した。

3.2 固有表現抽出を用いた機械翻訳

本論文では、中古日本語文にのみ固有表現タグを付与したコーパス、現代語訳文にのみ固有表現タグを付与したコーパス、どちらにも固有表現タグを付与したコーパスを用いてそれぞれ翻訳モデルを構築する。翻訳モデルには Transformer を用いた。モデルを図 2 に示す。モデル図は南端ら[22]が作成したものを参考にした。これは中古日本語文に固有表現タグを付与した場合である。まず学習データに対して事例ベース固有表現抽出を適用し、固有表現タグを付与する。その上で Transformer による学習を行い、得られた翻訳モデルによってテストデータの翻訳を行う。翻訳の際、出力に固有表現タグが含まれている場合はこれを除去する。学習及び翻訳タグを記号として設定した。

4. 実験概要

本節では、事例ベース固有表現抽出を用いた中古日本語文のニューラル機械翻訳手法の実験概要について述べる。

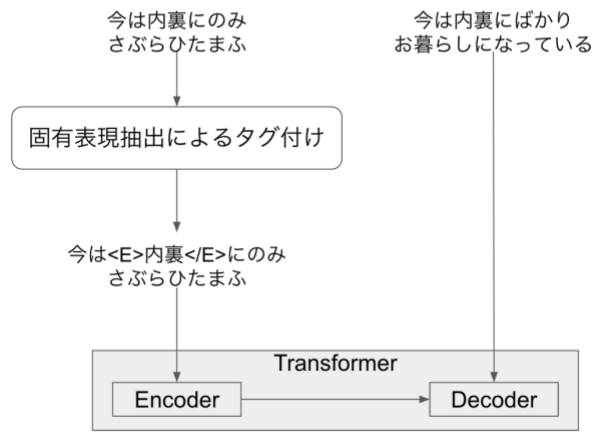
4.1 データセット

本論文では、固有表現抽出及びニューラル機械翻訳の対象として、中古日本語文が使用された平安時代中期の文学作品である源氏物語と紫式部日記を使用した。2 作品のデータは国文学者の渋谷栄一氏が Web 上[25][26]で公開している本文・現代語訳文を使用し、固有表現の正解ラベルは同氏が作成した主要登場人物一覧を参考に行っている。データ例を表 2, 表 3 に示す。また、事例ベース固有表現抽出に用いる BERT の Fine-tuning には、ストックマーク株式会社が作成した Wikipedia を用いた日本語の固有表現抽出データセット[27][28]を使用した。また、日付の固有表現データは前述のデータセットに含まれていなかったため、ChatGPT4.0[29]を用いて作成した。データ例を表 4 に示す。

4.2 中古日本語文を対象とした事例ベース固有表現抽出

事例ベース固有表現抽出の実装には、この手法を Python

学習時



推論時

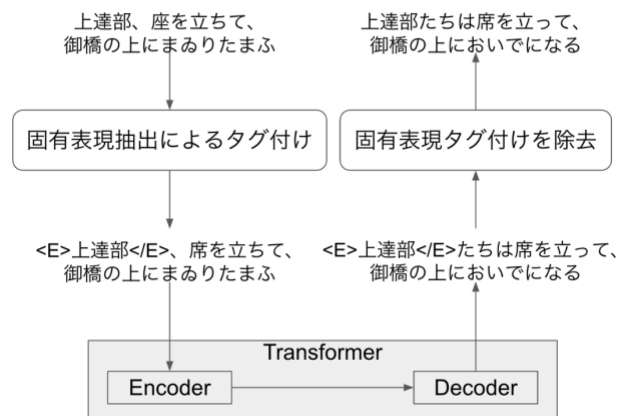


図 2 固有表現抽出を用いたニューラル機械翻訳

で実装した fsner ライブラリ[30]を用いた。Python のバージョンは、3.10.0、fsner のバージョンは 0.1.1 である。BERT は「人名」、「日付」の 2 つの固有表現タイプについて Fine-tuning を行った。Few-shot 学習は学習データを 10 件、20 件、50 件の 3 パターンで実験し、精度の変化を検証した。実験は各パターンにつき 5 回行い、使用するデータはランダムに入れ替えた。なお、Ziyadi ら[7]の実験では Few-shot 学習の学習データ件数についてさらに 100 件、200 件、500 件でも実験を行っているが、学習データ作成に用いた前述の登場人物一覧に掲載されている登場人物が 100 人程度であるため、今回は 50 が限界と判断した。検証に用いるテストデータは 80 文であり、そのうち固有表現を含むのは 52 文、含まないのは 28 文である。テストデータに登場する固有表現は学習データとの重複を避けた。実験結果の検証には F1 スコアを用いた。

4.3 固有表現抽出を用いたニューラル機械翻訳

ニューラル機械翻訳の実装は、Python と Pytorch ライブラリを用いて行う。Python のバージョンは 3.10.0、Pytorch のバージョンは 2.0.1 である。コーパスの単語分割には Mecab v0.996[31]を使用し、古文側には中古和文 UniDic v202308[32]を、現代文側には現代書き言葉 UniDic v2.0[33]を辞書として利用した。学習データは前述の源氏物語及び

表 2 固有表現抽出モデルの Few-shot 学習用データセット

固有表現タイプ	データ例
人名:	[E]源氏[E]の御事なりけむかし。
日付:	[E]弥生の十余日[E]のほどに、平らかに生まれたまひぬ。

表 3 ニューラル機械翻訳モデルの学習用データセット

時代	データ例
中古日本語文:	はじめより我はと思ひ上がりたまへる御方がた、めざましきものにおとしめ嫉みたまふ
現代語訳文:	最初から自分こそはと気位い高くいらっしやった女御方は、失敬な者だと貶んだり嫉んだりなさる
中古日本語文:	げに、御容貌ありさま、あやしきまでぞおぼえたまへる
現代語訳文:	なるほど、ご容貌や姿は不思議なまでによく似ていらっしやった

表 4 固有表現抽出モデルの Fine-tuning 用データセット

固有表現タイプ	データ例
人名:	[E]松友美佐紀[E]は、日本のバドミントン選手。
日付:	次のセミナーは[E]2024年10月14日[E]に予定されています。

紫式部日記であり、その内訳は訓練データが 13844 件、検証データが 4766 件、テストデータが 916 件である。データの分割はランダムに行なった。epoch 数は 100 で実験を行った。また、テストデータ 916 件のうち、固有表現を含む文 485 件を抽出し固有表現を含む文のみでの実験を行った。実験結果の検証には BLEU スコア [34] を用いた。BLEU スコアの導出式は以下に示す通りである [34]。

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-\frac{r}{c})} & \text{if } c \leq r \end{cases} \quad (1)$$

$$w_n = \frac{1}{N} \quad (2)$$

$$p_n = \frac{\sum_i \text{翻訳文}i \text{と正解文}i \text{で一致した } N - \text{gram 数}}{\sum_i \text{翻訳文}i \text{中の全 } N - \text{gram 数}} \quad (3)$$

$$BLEU = BP \cdot \exp\left(\sum_{n=1}^N w_n \log p_n\right) \quad (4)$$

ここで、 c は翻訳文の長さ、 r は正解文の長さである。BP は短い翻訳文が正解文よりも短い場合にペナルティを付与する。また、 N は $n=1$ から N までの整数で、 N -gram の長さを表す。 N -gram は「隣り合う連続した N 文字」を表す。BLEU スコアは、プロ翻訳者による正解文により近い出力が良い翻訳結果であるという定義に基づき、 N -gram を用いて翻訳文と正解文の類似度を提示する。

5. 実験結果

本節では、提案手法の実験結果について示す。

5.1 中古日本語文の固有表現抽出

表 5 にモデル別の F1 スコアを示す。DeBERTa は debert-base-japanese-wikipedia を、RoBERTa は roberta-base-japanese-aozora-char を表している。実験の結果から、今回は DeBERTa が RoBERTa を大きく上回った。また、DeBERTa ではサポート数を増やすとスコアが改善したが、

RoBERTa ではサポート数 20 の場合のスコアを 50 の場合のスコアが上回った。

5.2 固有表現抽出を用いたニューラル機械翻訳

表 6 に翻訳結果の BLEU スコアを示す。Baseline は固有表現抽出を行わずに学習させたモデルの翻訳結果である。source は学習データの原言語文である中古日本語文に対して固有表現抽出を適用したモデル、target は目的言語文である現代文に対して固有表現抽出を適用したモデル、all はどちらにも固有表現抽出を適用したモデルである。実験の結果、Baseline のスコアを source のスコアが、source のスコアを target のスコアがやや上回った。また、固有表現を含む文のみで構成されるテストデータに対しての翻訳結果を表 7 に示す。こちらの場合でも、Baseline のスコアを source のスコアが、source のスコアを target のスコアが上回った。また、どちらの場合においても、all のスコアは他のモデルを大きく下回る結果となった。

6. 考察

本節では、実験結果の考察を述べる。

6.1 中古日本語文の固有表現抽出

実験の結果、中古日本語文を対象とした固有表現抽出は、苑ら [8] が行った近代語文を対象とした固有表現抽出実験の

表 5 固有表現抽出実験結果

サポート数	DeBERTa	RoBERTa
10	0.595	0.421
20	0.606	0.480
50	0.651	0.472

表 6 ニューラル機械翻訳実験結果

モデル	Baseline	source	target	All
BLEU スコア	15.40	15.79	16.03	8.9

表 7 ニューラル機械翻訳実験結果-固有表現を含む文のみ

モデル	Baseline	source	target	All
BLEU スコア	17.36	17.65	18.09	10.46

F1 スコアをやや下回る結果となった。中古日本語文固有表現抽出のスコアが近代語文固有表現抽出のスコアを下回った理由としては、中古日本語文がより現代語から離れた文法であること、Few-shot 学習の学習データセット規模に差があったことが考えられる。Few-shot 学習の学習データセットを拡張するためには、中古日本語文の充実した固有表現データが必要だろう。

また、DeBERTa モデルのスコアを RoBERTa モデルのスコアが上回った理由としては、deberta-base-japanese-wikipedia の語彙サイズが 32000 トークンであるのに対して roberta-base-japanese-aozora-char の語彙サイズは 9415 トークンであり[35]、語彙サイズがより大きい DeBERTa モデルの方が古文の語彙に適応できたのではないかと考えられる。

6.2 固有表現抽出を用いたニューラル機械翻訳

実験の結果、固有表現抽出を用いたニューラル機械翻訳手法は中古日本語文の翻訳にある程度有効であることが示された。中古日本語文へ固有表現抽出を適用したモデルのスコアを現代文へ固有表現抽出を適用したモデルのスコアが上回ったのは、現代文側はベースとなる BERT モデル、Fine-tuning に用いた固有表現抽出データセット、Few-shot 学習に用いたデータセットが全て同時代の文章であったため固有表現抽出がより正確に行われたことが影響していると考えられる。また、学習データの中には誤った固有表現のタグ付け例も見られ、固有表現抽出の精度が向上すれば翻訳のスコアも改善することが考えられる。また、原言語文と目的言語文どちらにも固有表現抽出を適用したモデルのスコアが大きく低下したのは、原言語文と目的言語文の間で固有表現タグが付与された箇所には差があったためではないかと考えられる。南端ら[22]の研究ではこのモデルが最も高いスコアを示しており、今回の結果は固有表現抽出の精度が強く影響した結果であると言える。

7. まとめ

本論文では、事例ベース固有表現抽出を用いた中古日本語文のニューラル機械翻訳手法について述べた。本論文では、Few-shot 学習に中古日本語文を用いた固有表現抽出モデルをパラレルコーパスに適用し、固有表現情報を活用した中古日本語文の翻訳モデルを構築した。実験の結果、事例ベース固有表現抽出手法が中古日本語文の固有表現抽出に有効であり、また中古日本語文のニューラル機械翻訳に固有表現抽出を活用することでスコアが向上することが確認された。

8. 今後の展望

今後は、固有表現抽出の精度を向上させることでニューラル機械翻訳の精度を改善を目指していく。固有表現抽出の精度向上には、中古日本語文の固有表現データセットの構築、他の BERT モデルの実験などが考えられる。

本手法の応用などにより古文のニューラル機械翻訳精度を向上することができれば、実用的な古文翻訳システムを開発することができるだろう。

謝辞

本論文では、源氏物語の本文を国文学者の渋谷栄一氏が現代語に翻訳したデータを使用しました。渋谷栄一氏、並びにホームページの再構築を行なった宮脇文経氏に御礼申し上げます。

参考文献

- [1] TOPPAN, “古文書解読とくずし字資料の利活用サービス「ふみのは」”, <https://www.toppan.com/ja/joho/fuminoha/>, (参照: 2024-06-02).
- [2] 山本純子, 大澤留次郎, “古典籍翻刻の省力化: くずし字を含む新方式 OCR 技術の開発”, 情報管理誌, 58 巻, 11 号, pp819-827, 2016.
- [3] 東京大学基金, 研究者インタビュー, 山本浩司教授, “若手歴史研究者を育成し「歴史的思考法」をひろく日本社会と共有したい”, <https://utf.u-tokyo.ac.jp/project/person/623>, (参照: 2024-06-02).
- [4] 星野翔, 宮尾祐介, 大橋駿介, 相澤彰子, 横野光, “対照コーパスを用いた古文の現代語機械翻訳”, 言語処理学会第 20 回年次大会発表論文集, pp816-819, 2014
- [5] 白井久生, 古宮嘉那子, “T5 を用いた古文から現代文への翻訳”, 言語処理学会第 29 回年次大会発表論文集, pp2522-2527, 2023
- [6] Google Cloud “モデルの評価 | AutoML Translation Documentation”, <https://cloud.google.com/translate/automl/docs/evaluate?hl=ja>, (参照: 2024-06-08)
- [7] Morteza Ziyadi, Yuting Sun, Abhishek Goswami, Jade Huang and Weizhu Chen, Example-Based Named Entity Recognition, arXiv:2008.10570v1 [cs.CL], 2020.
- [8] 苑広媛, 李康穎, 後藤真, 木村文則, 前田亮, “事例ベース固有表現抽出を用いた古文からの歴史人物情報の抽出および活用”, DEIM2023 最終論文集, 2023
- [9] 「中古語」, 小学館国語辞典編集部, 『精選版 日本国語大辞典』, 小学館, 2005-12
- [10] 吉村衛, 木村文則, 前田亮, “古文テキストからの人物表現抽出”, 人文科学とコンピュータシンポジウム論文集, pp97-102, 2013
- [11] 永井規善, 前田亮, 木村文則, 赤松亮, “役者評判記からの人物表現抽出手法の提案”, 人文科学とコンピュータシンポジウム論文集, pp145-150, 2014
- [12] 小木曾智信, 小椋秀樹, 田中牧郎, 近藤明日子, 伝康晴, “中古和文を対象とした形態素解析辞書の開発”, 情報処理学会研究報告 人文科学とコンピュータ, Vol.2010-CH-85, No.4, pp1-8, 2010
- [13] 白井圭佑, 森信介, 後藤真, “人名辞典からの知識抽出”, 情報処理学会論文誌, 63 巻, 2 号, pp293-301, 2020-02-15
- [14] NetworkX developers, “NetworkX”, 2024, <https://networkx.org/>, (参照: 2024-06-13)
- [15] Nicholas Mamo, “multiplex-plot”, GitHub, 2022-06-18, <https://github.com/NicholasMamo/multiplex-plot>, (参照: 2024-06-13)
- [16] West Health Institute, “pyvis”, 2018, <https://pyvis.readthedocs.io/en/latest/index.html>, (参照: 2024-06-13)
- [17] Varun Shenoy, “GraphGPT”, <https://graphgpt.vercel.app/>, (参照: 2024-06-13)
- [18] 横野光, “人工頭脳プロジェクト「ロボットは東大に入れるか」国語試験古文問題解答に向けて”コーパスと日本語史研究, pp149-166, 2015-10
- [19] 高久雅史, 平澤虎庄, 小町守, 古宮嘉那子, “通時的な領域適応を行った単語分散表現を利用した古文から現代文へのニューラル機械翻訳”, 言語処理学会第 26 回年次大会発表論文集, 26 巻, pp3-9, 2020

- [20] Xiaoqing Li, Jiajun Zhang, Chengqing Zong, “Neural Name Translation Improves Neural Machine Translation”, Proc.of CWMT 2018, Springer, pp93-100, 2018
- [21] Aren Siekmeier, WonKee Lee, Hongseok Kwon, Jong-Hyeok Lee, Tag Assisted Neural Machine Translation of Film Subtitles, Proc.of IWSLT 2021, pp255-262, 2021
- [22] 南端尚樹, 田村晃裕, 加藤恒夫, “目的言語文の固有表現タグ付与に基づく Transformer ニューラル機械翻訳”, 言語処理学会第28回年次大会発表論文集, pp937-941, 2022
- [23] KoichiYasuoka, “deberta-base-japanese-wikipedia”, Hugging Face, <https://huggingface.co/KoichiYasuoka/deberta-base-japanese-wikipedia>, (参照: 2024-05-21)
- [24] KoichiYasuoka, “roberta-base-japanese-aozora-char”, Hugging Face, <https://huggingface.co/KoichiYasuoka/roberta-base-japanese-aozora-char>, (参照: 2024-05-21)
- [25] 渋谷栄一, “源氏物語の世界”, 2024-06-05, <http://www.sainet.or.jp/~eshibuya/index.html>, (参照: 2024-06-06)
- [26] 宮脇文経, “源氏物語の世界 再編集版”, 2018-05-04, <http://www.genji-monogatari.net/>, (参照: 2024-06-06)
- [27] ストックマーク株式会社, “Wikipedia を用いた日本語の固有表現抽出データセット”, GitHub, 2023-09-02, <https://github.com/stockmarkteam/ner-wikipedia-dataset>, (参照: 2024-05-22),
- [28] 近江崇宏, Wikipedia を用いた日本語の固有表現抽出のデータセットの構築, 言語処理学会第27回年次大会発表論文集, pp350-352, 2021
- [29] OpenAI, “Hello ChatGPT4o”, 2024-05-13, <https://openai.com/index/hello-gpt-4o/>, (参照: 2024-05-22)
- [30] Sayef, “fsner”, GitHub, <https://github.com/sayef/fsner>, 2022-03-30, (参照: 2024-05-20)
- [31] 工藤拓, MeCab: Yet Another Part-of-Speech and Morphological Analyzer, <https://taku910.github.io/mecab/>, (参照: 2024-05-28)
- [32] 古文用 Unidic, “中古和文 Unidic”, Unidic, 2023-08, https://clrd.ninjal.ac.jp/unidic/download_all.html#unidic_wabun, (参照: 2024-05-15)
- [33] “現代書き言葉 Unidic”, Unidic, https://clrd.ninjal.ac.jp/unidic/download.html#unidic_bccwj, (参照: 2024-05-15)
- [34] Kishore Papineni, Salim Roukos, Todd Ward, Wei-Jing Zhu, Bleu: a Method for Automatic Evaluation of Machine Translation, ACL '02: Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, pp311-318, 2002-07-06
- [35] 安岡孝一, 青空文庫 DeBERTa モデルによる国語研長単位係り受け解析, 「東洋学へのコンピュータ利用 第35回研究セミナー」予稿集, 35巻, pp29-43, 2022-07