

MATCHA：専門家が平易化した記事を用いたやさしい日本語パラレルコーパス

宮田 莉奈¹ 惟高 日向¹ 山内 洋輝¹ 柳本 大輝¹ 梶原 智之¹ 二宮 崇¹ 西脇 靖紘²¹ 愛媛大学大学院理工学研究科 ² 株式会社 MATCHA

{miyata, koretaka, yamauchi, yanamoto}@ai.cs.ehime-u.ac.jp

{kajiwara@cs., ninomiya.takashi.mk@ehime-u.ac.jp nishiwaki@matcha-jp.com}

掲載号の情報

31巻2号 pp. 590-609.

doi: <https://doi.org/10.5715/jnlp.31.590>

概要

本研究では、テキスト平易化のための日本語パラレルコーパスを構築し、公開¹⁾した。テキスト平易化とは、文の意味を保持しつつ、難解な文を平易に言い換えるタスクのことである。既存の日本語コーパスとしては、文を非専門家が平易化したSNOW [1, 2] や、専門家が平易化したJADES [3] がある。SNOWの平易化は、作者らによって定義された基礎語彙2,000語²⁾に基づくため、この語彙で網羅できない表現については、却ってわかりにくく言い回しが見られる。また、JADESは小規模である。既存の日本語テキスト平易化パラレルコーパスには、高品質かつ大規模なコーパスは存在しない。

そこで本研究では、高品質かつ大規模な日本語のテキスト平易化パラレルコーパスを構築するため、訪日観光客向けメディアMATCHA³⁾の日本語記事とその専門家による平易化版の記事⁴⁾の対から人手で文アライメントをとった。ただし、記事単位で平易化されており、全ての難解文に対して意味的に完全に対応する平易文が存在するわけではないため、完全一致と部分一致の2種類の文アライメントをとった。文対の抽出後、著者らが人手で誤字脱字を訂正し、括弧の用法などのスタイルを統一した。我々のMATCHAコーパスは、完全一致11,000文対と部分一致5,000文対の計16,000文対からなる。

1) <https://github.com/EhimeNLP/matcha>2) <https://www.jnlp.org/GengoHouse/list/>語彙3) <https://matcha-jp.com/>

4) 日本語能力試験新4級相当（旧3級相当）の文法や単語を用いて、通常の日本語記事を書き換えて作成されている。

表1 日本語テキスト平易化コーパスの特徴

	専門家による 言い換え	言い換えの 多様性	言い換えの 品質	規模
SNOW	-	★	★★	85,000
JADES	✓	★★★	★	3,907
MATCHA	✓	★★★	★★★	16,000

コーパスを人手評価した結果を表1に示す。専門家が平易化したパラレルコーパスは非専門家によるものと比べて多様な平易化操作を含んでいることが明らかになった。また、SNOWの平易化操作が語句の置換に集中している一方で、JADESやMATCHAには語句の挿入や削除、並び替えや文分割など多様な変換が含まれることもわかった。さらに、本コーパスは他のコーパスよりも流暢かつ意味を保持した平易化が行われていることを確認した。

SNOWおよびMATCHAを用いて事前訓練済み系列変換モデルをファインチューニングし、テキスト平易化モデルを構築した。実験の結果、MATCHAで訓練したモデルはSNOWで訓練したモデルに比べて、流暢かつ意味を保持した平易化ができるのを確認した。MATCHA上での訓練においては、完全一致の文対のみで訓練したモデルの方が流暢性および同義性に優れる一方、平易性については部分一致の文対も訓練に含めた方が高い性能を発揮した。

参考文献

- [1] Takumi Maruyama and Kazuhide Yamamoto. Simplified Corpus with Core Vocabulary. In **Proc. of LREC**, pp. 1153–1160, 2018.
- [2] Akihiro Katsuta and Kazuhide Yamamoto. Crowdsourced Corpus of Sentence Simplification with Core Vocabulary. In **Proc. of LREC**, pp. 461–466, 2018.
- [3] Akio Hayakawa, Tomoyuki Kajiwara, Hiroki Ouchi, and Taro Watanabe. JADES: New Text Simplification Dataset in Japanese Targeted at Non-native Speakers. In **Proc. of TSAR**, pp. 179–187, 2022.