

# 使用多种语言模型分类长篇语料句子

李新<sup>1</sup>, 孙润<sup>2</sup>

(1. 广东海洋大学外国语学院语言实训中心, 广东湛江 524088;

2. 广东海洋大学外国语学院, 广东湛江 524088)

**[摘要]**为了进行生态翻译学研究,我们实施了一系列数据提取实验,利用多种语言模型对长篇语料中的句子进行分类研究。该研究旨在探索语言模型在文学作品研究环境下的应用,以进一步提高文学作品研究的质量和效率。通过收集研究语料和训练语料,并运用包括 Bert、ChatGPT、Claude、ChatGLM 等多种语言模型,我们对大规模文本进行了分类和筛选。通过这项实验,我们希望能够识别粗语以及谚语俗语等语义特征,为跨语言翻译提供更深入的理解和支持。在实验结果中,我们观察到利用多种语言模型能够对长篇语料中的句子进行有限的分类,在一定程度上减轻了人工分类的工作量,并且本实验数据能为相关翻译研究提供支持。

**[关键词]**《红高粱家族》; 平行语料; Bert; ChatGPT; ChatGLM; Ernie

**[作者简介]**李新(1987—),男,广东湛江人,工程师,硕士,研究方向:计算机技术与语言处理。孙润(1989—),男,湖北黄冈人,广东海洋大学外国语学院讲师,博士在读,研究方向:普通语言学及文学翻译。

**[基金项目]**本文系广东海洋大学人文社科项目“魔幻现实主义风格下的莫言作品英译研究——虚幻与现实融合下的中国底层人物形象”(项目编号:C22862)。

**[DOI]** <https://doi.org/10.62662/kxwxz0109007>

**[本刊网址]** [www.oacj.net](http://www.oacj.net)

## 一、前期工作回顾

本项目的研究对象为《红高粱家族》的中英文语料,中文为莫言的《红高粱家族》,英文为葛浩文翻译的“Red Sorghum”,相关语料从网络获取。得到语料后,对中英文语料进行清洗、分句和词性标注,使用到的工具包括 nltk、ltp、jieba 等。其中《红高粱家族》句子总数为 7379,“Red Sorghum”句子总数为 8619。

此外,为了研究中英文对应句子的翻译特点,我们还需要对中英文语料进行锚点对齐,将中英文双边连续 3 句高相似度句子作为锚点,经过检查这种策略得到的锚点的准确度较高,并且锚点较多,总共锚点簇数 420、锚点数 768。这意味着可以根据锚点按块对语料进行分块细化对齐。

## 二、利用 ChatGPT 进行长篇双语语料对齐

鉴于 ChatGPT 3.5 近年来火爆全球,且其上下文理解、跨语言分析、自动化处理等通用智能能力让人感觉到其“无所不能”的潜力,且尝试过的一些所谓智能对齐软件,例如 LanguageX,对齐效果也一般,故有必要尝试一下新的手段。当然使用

ChatGPT 进行语料对齐也有“杀鸡用牛刀”的意味。

我们选择一些使用成本较低的、国内能访问的一些 ChatGPT 接入站或者套壳网站来使用,且网站支持 ChatGPT 3.5-16K,它可以支持长文本分析,非常适合用来进行长文本对齐。

每次我们选取 20 句中英文句对进行对齐,并要求 ChatGPT 以句子为单位用 Tuple 元组形式输出对齐结果,然后人工大致分析对齐结果。如果输出不理想,例如出现删除句子、改变句子或者大量调换位置、大量空对齐的现象时,需要重新要求 ChatGPT 输出,或修改提示语或者减少句子数量。我们还发现当中英文结尾句子恰好对应时,ChatGPT 的对齐效果可以很好。另外我们还发现对齐效果可能受时间段影响,某段时间,特别是晚上 9 点,使用人数过多可能造成服务器负担,影响算力的分配。这时我们选择 Claude 模型进行输出也能得到不错的效果。

由于 ChatGPT 存在输出格式错误的情况,例如没有以 Tuple 元组格式输出、转译字符错误,我们还进行了编程,对收集输出后的结果重新划分中英



验。由于测试时有些句子会被识别为民谣和流行语,所以输入提示语时包含了民谣和流行语作为关键字。程序经过2天的运行得到了2300句谚语俗语。由于分类结果较多,且错误率也较高,于是需要进一步进行筛选。我们选择 Ermie 对分类结果进

行二次分类。

使用 Ermie 进行二次分类时需要联网,免费Token有100万个。经过大半天的运行,我们得到了645句分类为谚语俗语的句子,经过人工筛选后得到164句,如图3所示。

Table with 4 columns: 原文, Ermie判断, 人工判断. It shows a list of text segments from a story and their corresponding classification results by the Ermie model and human reviewers.

图3 机器分类和人工筛查的结果

最后我们再次使用 ChatGLM3-6B 对粗语再次分类和人工筛选,并和 Bert 模型的结果合并去重后得到155句粗语,如图4所示。

Table with 2 columns: 原文, 译文. It shows a list of text segments and their corresponding machine-translated versions.

图4 粗语进行机器分类和人工筛查的结果

五、总结

我们通过利用多种语言模型对《红高粱家族》及其译文按句子进行对齐,以及对粗语和谚语俗语进行分类,最后得到155句粗语的平行语料以及164句谚语俗语的平行语料。该实验为下一步生态

翻译学研究奠定了基础,为相关研究提供借鉴作用。

通过此项研究,我们感觉到语言模型对于特殊领域的语义特征的识别在性能上依然不足,同时缺少能直接使用的面向专业领域的模型,期待未来的语言工具能在某一垂直领域做得更智能化。

另外我们总结了此次研究工作中的不足和未来可以改进的地方:1. 我们使用 ChatGPT 进行语料对齐时,由于人工检查的不到位,语料有某些部分存在错误,例如句子重复、中英文位置错误等;2. 训练的文本数据量较少,可以尝试进行数据增强,利用已有的大模型进行句子改写和生成;3. 我们从中文的角度出发分类粗语和谚语俗语,或许从英文的角度出发进行分析可能得到不一样的数据;4. 尝试多种提示语,或者多种提示语结合,或者将句子分割进行多段分析综合判断。

#### 参考文献:

[1] J. Devlin, M. W. Chang, K. Lee, K. Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding[A]. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computa-

tional Linguistics: Human Language Technologies [C]. Kerrville: Association for Computational Linguistics, 2019.

[2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin. Attention Is All You Need[A]. NIPS' 17: Proceedings of the 31st International Conference on Neural Information Processing Systems [C]. New York: Curran Associates Inc., 2017.

[3] Jian Zhu, Zuoyu Tian, Sandra Kübler. UM-IU@ LING at SemEval-2019 Task 6: Identifying Offensive Tweets Using BERT and SVMs [EB/OL]. <https://arxiv.org/abs/1904.03450>, 2019-4-6.

[4] Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, Hua Wu. ERNIE: Enhanced Representation through Knowledge Integration [EB/OL]. <https://arxiv.org/abs/1904.09223>, 2019-4-19.

### Classifying Sentences of Long Textual Corpus by Using Various Language Models

LI Xin<sup>1</sup>, SUN Run<sup>2</sup>

- (1. Language Training Center, College of Foreign Languages, Guangdong Ocean University, Zhanjiang Guangdong 524088;
2. College of Foreign Languages, Guangdong Ocean University, Zhanjiang Guangdong 524088, China)

**Abstract:** In order to conduct research in eco-translation studies, we carried out a series of data extraction experiments that involved using multiple language models to classify sentences in lengthy corpora. This study aims to explore the application of language models in the context of literary works research to enhance the quality and efficiency of literary studies. By collecting research and training corpora, and employing various language models including Bert, ChatGPT, Claude, and ChatGLM, we performed classification and filtering on a large-scale text corpus. Through this experiment, we hoped to identify semantic features such as profanity and idioms, and provide a deeper understanding and support for cross-language translation. In the experimental results, we observed that utilizing multiple language models allowed for limited classification of sentences in lengthy corpora, thus alleviating the workload of manual classification to some extent. Additionally, this experimental data can provide support for related translation research.

**Key words:** “Red Sorghum”; parallel corpus; Bert; ChatGPT; ChatGLM; Ernie