

# SimLTE: Simple Contrastive Learning for Long Text Embeddings

Anonymous ACL submission

## Abstract

High-quality text embeddings are essential to various natural language understanding tasks especially for low-resource settings. Existing methods are proposed for general text understanding but representations for long text are less explored. Basically, to have better embeddings for long text, a model should have the ability to recognize key information in a long text. To this end, we present SIMLTE, a simple unsupervised contrastive learning framework for long text embeddings in this paper. Specifically, we pretrain a language model to distinguish if two texts have the same topic without any supervision. The positive pairs are constructed by our key information redundancy assumption for long text. Experimental results on five datasets show that SIMLTE outperforms the state-of-the-art baselines of text embeddings by 3.9% and 12.0% macro-F1 in average respectively under general and few-shot text classification settings.

## 1 Introduction

Learning text embeddings is a fundamental problem in natural language processing (Kiros et al., 2015; Hill et al., 2016; Conneau and Kiela, 2018; Logeswaran and Lee, 2018; Gao et al., 2021; Reimers et al., 2016). Most previous studies (Hill et al., 2016; Logeswaran and Lee, 2018; Gao et al., 2021) focus on sentence-level representations where training data usually contain short text but high-quality long text embeddings are less explored. However, the long text appears frequently in NLP tasks involved in News, scientific articles, and social media. Existing works for long text either design a specific model structure (Pappagari et al., 2019; Grail, 2021; Rohde et al., 2021; Fang et al., 2019; Ainslie et al., 2020) or identify important sentences (Ding et al., 2020; Zaheer et al., 2020; Gidiotis and Tsoumakas, 2020; Huang et al., 2021; Zhong et al., 2021) to improve long text understanding. These methods usually rely on labels

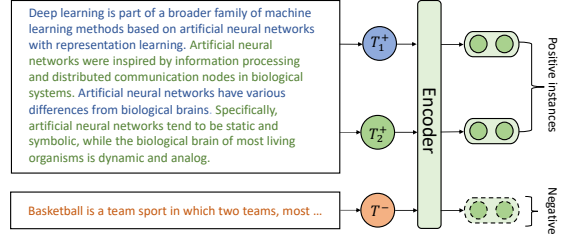


Figure 1: Pre-training SimLTE on a large-scale dataset with positive instances

from human efforts to conduct supervised learning whereas low-resource tasks need high-quality representations from pertaining. Hence, an efficient pretraining method for high-quality long text embeddings is necessary to explore.

To improve long text embeddings, in this paper, we present an unsupervised contrastive learning framework, namely SIMLTE, which can produce superior long text embeddings without any supervision. Instead of model structures, we seek an unsupervised pretraining method for long text embeddings. To conduct contrastive learning, we first need to produce positive pairs. We investigate the information redundancy (details in Appendix A) on five datasets for different lengths of text. We find the information redundancy is larger as the length of the text is increasing. This result indicates long text usually contains repeated information. Based on this observation, we can assume that the model can still learn the main topic of a long text even if we drop some sentences. Therefore, as shown in Figure 1, we randomly divide sentences of a long text into two exclusive subsets and the two subsets work as positive pairs for contrastive learning. The intuition behind this method is that we expect the model will pull representations of two subsets together in the latent space by paying more attention to common keywords so that the model can learn key information from text automatically.

For large-scale pre-training, we follow previous works (Gao et al., 2021; Li et al., 2022) to adopt in-

batch negatives then pretrain BERT (Devlin et al., 2019) and Longformer (Beltagy et al., 2020) on Wikipedia English Corpus<sup>1</sup>. To evaluate the quality of long text embeddings, we conduct general and few-shot text classification on five long text datasets involved in News and scientific articles. The experimental results show that SIMLTE with two kinds of model structures (i.e., BERT and Longformer) can both achieve significant improvements compared to state-of-the-art baselines.

Overall, our contributions are three-fold:

- After analyzing the information redundancy of long text, we propose an unsupervised contrastive learning framework with a specific positive pair construction method for long text.
- For high-quality long text embeddings, we conduct a large-scale pretraining on the Wikipedia dataset with BERT and Longformer.
- Comprehensive experiments on long text classification are conducted and results show that our approach can largely improve the quality of long text embeddings.

## 2 Method

In this section, we first introduce background knowledge about contrastive learning and then we present SIMLTE from how to construct positive instances.

### 2.1 Contrastive Learning

Contrastive Learning aims to learn effective representations by pulling semantically close neighbors together and pushing apart non-neighbors in the latent space (Hadsell et al., 2006). It assumes a contrastive instance  $\{x, x^+, x_1^-, \dots, x_{N-1}^-\}$  including one positive and  $N - 1$  negative instances and their representations  $\{\mathbf{h}, \mathbf{h}^+, \mathbf{h}_1^-, \dots, \mathbf{h}_{N-1}^-\}$ , where  $x$  and  $x^+$  are semantically related. we follow the contrastive learning framework (Chen et al., 2020; Li et al., 2022) and take cross-entropy as our objective function:

$$l = -\log \frac{e^{\text{sim}(\mathbf{h}, \mathbf{h}^+)/\tau}}{e^{\text{sim}(\mathbf{h}, \mathbf{h}^+)/\tau} + \sum_{i=1}^{N-1} e^{\text{sim}(\mathbf{h}, \mathbf{h}_i^-)/\tau}} \quad (1)$$

where  $\tau$  is a temperature hyperparameter and  $\text{sim}(\mathbf{h}_1, \mathbf{h}_2)$  is the cosine similarity  $\frac{\mathbf{h}_1^\top \mathbf{h}_2}{\|\mathbf{h}_1\| \cdot \|\mathbf{h}_2\|}$ . In this work, we encode input texts using a pre-trained language model such as BERT (Devlin et al., 2019). Following BERT, we use the first special token

[CLS] as the representation of the input and fine-tune all the parameters using the contrastive learning objective in Equation 1.

### 2.2 SimLTE

The critical problem in contrastive learning is how to construct positive pairs  $(x, x^+)$ . In representation learning for visual tasks (Chen et al., 2020), an effective solution is to take two random transformations of the same image (e.g., flipping, rotation). Similarly, in language representations, previous works (Gao et al., 2021; Karpukhin et al., 2020; Meng et al., 2021; Li et al., 2022) apply augmentation techniques such as dropout, word deletion, reordering, and masking.

In this paper, we propose a new method to construct positive instances for long text. The basic idea of positive instance construction for contrastive learning is adding random noises to the original data for augmentation. The augmented data should have similar representations to the original data. Models trained by contrastive losses on augmented data will have an increased ability to learn important features in the data. To add random noises in long text, we find long text (e.g., paragraphs) usually has higher information redundancy than short text (e.g., sentences) (Table 3 in Appendix). With this observation, we can have an assumption: the semantics of a long text will not be changed even if we drop half of the text. We can construct positive pairs under this assumption easily on any text dataset without supervision. Specifically, for each long text in the dataset, we randomly split sentences in the long text into two subsets and the two sentence sets do not have intersections. In the two subsets, we keep the order of sentences in the original long text to form two new texts. According to our assumption, the two new texts should have the same semantics and hence they are used as a positive pair in contrastive learning.

Consider an example (in Figure 1) to understand our positive instance construction process: Suppose we have a long text  $T = (s_1, s_2, \dots, s_n)$  where  $s_i$  is the  $i$ -th sentence in long text and  $n$  is the number of sentences, each sentence will be sent to anchor set or positive set with the same probability (50%). The sentences in the same set (i.e., anchor or positive) will be concatenated in the same order of  $T$  to form one positive pair  $(T_1^+, T_2^+)$  for contrastive learning. Positive pairs constructed by this method will not contain the same sentence and hence pre-

<sup>1</sup><https://dumps.wikimedia.org/>

vent models from overfitting on recognizing the same sentences. Instead, models are guided to learn keywords appearing in positive instances so as to improve the ability to recognize key information. We split the long text at sentence level instead of word level (e.g., word deletion for augmentation) because the word-level splitting will cause the discrepancy between pretraining and finetuning and then lead to performance decay.

For negative instances, we use in-batch instances following previous contrastive frameworks (Gao et al., 2021; Li et al., 2022).

### 3 Experiments

In this section, we evaluate the effectiveness of our method by conducting text classification tasks. To eliminate the influence of different model structures and focus on the quality of text embeddings. We freeze the parameters of different text encoders and fine-tune only a multi-layer perceptron (MLP) to classify the embeddings of text encoders. We also visualize the attention weights between baselines and SIMLTE.

#### 3.1 Pretraining Details

We use English Wikipedia <sup>2</sup> articles as pretraining data and each article is viewed as one training instance. The total number of training instances is 6,218,825. For pre-training, we start from the pretrained BERT-BASE model (Devlin et al., 2019) and the Longformer (Beltagy et al., 2020) model <sup>3</sup>. We follow previous works (Gao et al., 2021; Li et al., 2022): the masked language model (MLM) loss and the contrastive learning loss are used concurrently with in-batch negatives. Our pretraining learning rate is 5e-5, batch size is 36 and 12 for BERT and Longformer structure respectively. Our model is optimized by AdamW (Kingma and Ba, 2014) in 1 epoch. The temperature  $\tau$  in the contrastive loss is set to 0.05 and the weight of MLM is set to 0.1 following previous work (Gao et al., 2021).

#### 3.2 Datasets

We use the following classic long text datasets to evaluate our method: (1) Fake News Corpus <sup>4</sup>; (2) 20NewsGroups (Lang, 1995); (3) arXiv articles

Datasets	Data Size	Classes	Ave.	Med.
FakeNews	8,558,957	15	467	299
20News	18,846	20	258	153
arXiv	2,162,833	38	138	131
NYT	13,081	5	650	683
BBCNews	2,225	5	133	130

Table 1: Statistics of datasets. Ave. and Med. stand for the average and median number of words respectively in one data instance.

dataset <sup>5</sup>; (4) New York Times Annotated Corpus (NYT) (Sandhaus, 2008); and (5) BBCNews <sup>6</sup>. We do not use semantic textual similarity (STS) tasks (Agirre et al., 2012) because the sentences in these tasks are short which is not suitable to evaluate long text embeddings.

#### 3.3 Baselines

We compare our pre-trained model to the baselines of two groups. (1) BERT based models include BERT (Devlin et al., 2019), SimCSE (Gao et al., 2021), CT-BERT (Carlsson et al., 2021). For a fair comparison, we also train a SimCSE with our pretraining dataset (SimCSE<sub>long</sub>). (2) Transformers specified for long sequences include Longformer (Beltagy et al., 2020) and BigBird (Zaheer et al., 2020). We train two versions of SIMLTE with BERT and Longformer (i.e., SIMLTE<sub>bert</sub> and SIMLTE<sub>long</sub>) for comparison. We do not include RoBERTa (Liu et al., 2019) and IS-BERT (Zhang et al., 2020) as our baselines because SimCSE achieves better results than these methods according to the paper.

#### 3.4 Text Classification

In the general text classification task, we classify text embeddings with the full training set. Training details are in Appendix B.

**Results.** Table 2 shows the evaluation results on different datasets. Overall, we can see that SIMLTE achieves the best performance over the 5 long text datasets and consistently improves the long text embeddings with BERT and Longformer structures. Specifically, methods pretrained with contrastive objectives (i.e., CT-BERT, SimCSE) outperform general language representations (i.e., BERT) which indicates contrastive objectives designed for text embeddings can largely improve the ability of language models to produce high-quality

<sup>2</sup><https://en.wikipedia.org/>

<sup>3</sup>The Longformer checkpoint is pretrained on long documents by MLM task and is available from Huggingface.

<sup>4</sup><https://github.com/several27/FakeNewsCorpus>

<sup>5</sup><https://www.kaggle.com/datasets/Cornell-University/arxiv>

<sup>6</sup><http://mlg.ucd.ie/datasets/bbc.html>

Datasets	FakeNews		20News		arXiv		NYT		BBCNews	
Metrics	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
<i>Text Classification</i>										
BERT	54.98	42.17	62.34	54.19	68.52	20.46	95.11	92.65	91.06	90.34
CT-BERT	55.19	42.53	65.76	63.37	71.61	26.09	95.69	91.59	90.32	88.87
SimCSE	58.48	47.46	74.02	72.57	74.46	30.01	97.17	94.69	94.12	93.86
SimCSE <sub>long</sub>	58.37	47.56	73.51	72.05	73.16	29.41	97.25	93.83	94.22	94.30
SIMLTE <sub>bert</sub>	<b>60.04</b>	<b>50.14</b>	<b>76.89</b>	<b>74.85</b>	<b>76.66</b>	<b>32.24</b>	<b>98.20</b>	<b>96.05</b>	<b>95.56</b>	<b>95.58</b>
LongFormer	65.72	57.66	73.69	72.47	71.66	25.92	94.36	88.39	96.33	94.75
BigBird	57.44	47.87	70.35	68.91	71.58	25.05	97.13	94.33	94.11	94.62
SIMLTE <sub>long</sub>	<b>71.60</b>	<b>61.66</b>	<b>75.44</b>	<b>74.38</b>	<b>77.68</b>	<b>33.26</b>	<b>97.90</b>	<b>95.43</b>	<b>96.67</b>	<b>95.91</b>
<i>Few-shot Text Classification</i>										
BERT	23.96	23.73	19.94	18.71	24.08	10.14	51.85	43.90	54.22	52.73
CT-BERT	23.71	23.06	24.11	23.53	27.02	13.53	47.23	36.83	59.56	58.95
SimCSE	25.04	22.68	42.63	41.42	32.61	17.19	86.51	78.41	83.56	83.75
SimCSE <sub>long</sub>	26.39	23.26	48.65	47.81	23.42	12.66	85.36	75.90	84.44	83.96
SIMLTE <sub>bert</sub>	<b>27.79</b>	<b>24.65</b>	<b>55.79</b>	<b>55.43</b>	<b>35.79</b>	<b>18.52</b>	<b>90.52</b>	<b>83.71</b>	<b>86.86</b>	<b>86.31</b>
LongFormer	26.56	25.12	44.42	42.41	25.04	13.36	73.06	54.87	84.89	85.47
BigBird	25.36	23.28	39.14	39.06	23.62	10.18	86.66	78.96	79.11	76.63
SIMLTE <sub>long</sub>	<b>29.17</b>	<b>27.13</b>	<b>51.18</b>	<b>50.96</b>	<b>34.33</b>	<b>18.80</b>	<b>89.78</b>	<b>82.88</b>	<b>86.78</b>	<b>86.66</b>

Table 2: For all performance measures, larger numbers are better. Our pre-trained model achieves the best results in all cases.

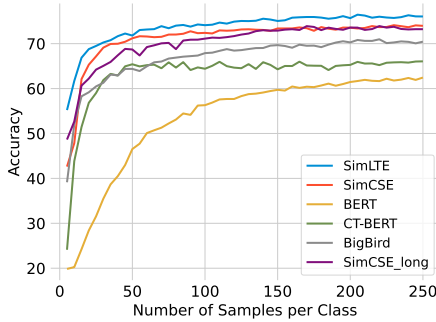


Figure 2: Performance of different models with different numbers of instances per class under few-shot setting.

text embeddings. SimCSE pretrained with our long text data (i.e., SimCSE<sub>long</sub>) has similar results as the original SimCSE which indicates simply increasing the length of pretraining text cannot improve long text embeddings. Compared to SimCSE and Longformer, our model achieves 3.9% and 9.4% average macro-F1 improvements with BERT and Longformer structures respectively. Hence, our contrastive learning method is effective for long text embeddings.

### 3.5 Few-shot Text Classification

To show the performance of different text embeddings under low-resource settings, we evaluate our model with few-shot training instances. Training details are in Appendix B.

**Results.** Table 2 shows the results of few-shot text classification on these five datasets. We can see that SIMLTE (i.e., SIMLTE<sub>bert</sub> and SIMLTE<sub>long</sub>) achieves 12.0% and 24.3% macro-F1 improve-

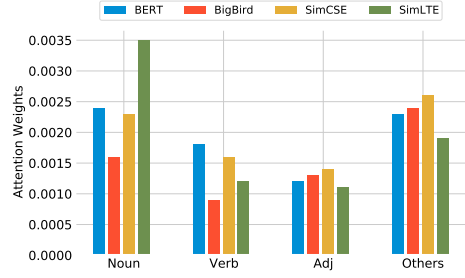


Figure 3: Attention weights from different models on the NYT dataset.

ments compared to SimCSE and Longformer respectively. These improvements are higher than general text classification. Besides, we also compare the performance of different baselines and SIMLTE<sub>bert</sub> with different numbers of training instances on 20News. The results in Figure 2 show the improvements from our method become larger as the number of training instances decreases indicating the importance of high-quality long text embeddings for low-resource settings. Furthermore, our method achieves the best results under different numbers of training instances.

### 3.6 Attention Weights

To explore the difference between SIMLTE and other models, we analyze the attention weights of Transformers in different models on the NYT dataset (details in Appendix C). The average weights of different kinds of words are shown in Figure 3. We can see that our model has more than 40% higher attention weights on nouns compared to BERT and SimCSE. Martin and Johnson (2015) shows nouns are more informative than other words in the document understanding. Hence, our pre-training method increases the attention weights of models on nouns which results in higher performance on long text classification.

## 4 Conclusion

In this work, we propose an unsupervised contrastive learning framework for long text embeddings. Our method provides a new method for long text data augmentation without any supervision and language models can get large-scale pretraining on any long text. We conduct extensive experiments on text classification tasks under fully supervised and few-shot settings. Results show that our pre-trained model greatly outperforms state-of-the-art text embeddings, especially when the training data is limited.



## References

- Eneko Agirre, Daniel Matthew Cer, Mona T. Diab, and Aitor Gonzalez-Agirre. 2012. Semeval-2012 task 6: A pilot on semantic textual similarity. In *\*SEM-EVAL*.
- Joshua Ainslie, Santiago Ontañón, Chris Alberti, Václav Cvicek, Zachary Kenneth Fisher, Philip Pham, Anirudh Ravula, Sumit K. Sanghai, Qifan Wang, and Li Yang. 2020. Etc: Encoding long and structured inputs in transformers. In *Conference on Empirical Methods in Natural Language Processing*.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *ArXiv*, abs/2004.05150.
- Fredrik Carlsson, Amaru Cuba Gyllensten, Evangelia Gogoulou, Erik Ylipää Hellqvist, and Magnus Sahlgren. 2021. Semantic re-tuning with contrastive tension. In *ICLR*.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. 2020. A simple framework for contrastive learning of visual representations. *ArXiv*, abs/2002.05709.
- Alexis Conneau and Douwe Kiela. 2018. Senteval: An evaluation toolkit for universal sentence representations. *ArXiv*, abs/1803.05449.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.
- Ming Ding, Chang Zhou, Hongxia Yang, and Jie Tang. 2020. Cogltx: Applying bert to long texts. In *NeurIPS*.
- Yuwei Fang, S. Sun, Zhe Gan, Rohit Radhakrishna Pillai, Shuohang Wang, and Jingjing Liu. 2019. Hierarchical graph network for multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing*.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. *ArXiv*, abs/2104.08821.
- Alexios Gidiotis and Grigorios Tsoumakas. 2020. A divide-and-conquer approach to the summarization of long documents. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:3029–3040.
- Quentin Grail. 2021. Globalizing bert-based transformer architectures for long document summarization. In *Conference of the European Chapter of the Association for Computational Linguistics*.
- Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. Dimensionality reduction by learning an invariant mapping. *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, 2:1735–1742.
- Felix Hill, Kyunghyun Cho, and Anna Korhonen. 2016. Learning distributed representations of sentences from unlabelled data. In *NAACL*.
- Luyang Huang, Shuyang Cao, Nikolaus Nova Parulian, Heng Ji, and Lu Wang. 2021. Efficient attentions for long document summarization. In *North American Chapter of the Association for Computational Linguistics*.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Yu Wu, Sergey Edunov, Danqi Chen, and Wen tau Yih. 2020. Dense passage retrieval for open-domain question answering. *ArXiv*, abs/2004.04906.
- Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980.
- Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-thought vectors. *ArXiv*, abs/1506.06726.
- Ken Lang. 1995. Newsweeder: Learning to filter net-news. In *ICML*.
- Jiacheng Li, Jingbo Shang, and Julian McAuley. 2022. [UCTopic: Unsupervised contrastive learning for phrase representations and topic mining](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6159–6169, Dublin, Ireland. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *ArXiv*, abs/1907.11692.
- Lajanugen Logeswaran and Honglak Lee. 2018. An efficient framework for learning sentence representations. *ArXiv*, abs/1803.02893.
- Fiona Martin and Mark Johnson. 2015. [More efficient topic modelling through a noun only approach](#). In *Proceedings of the Australasian Language Technology Association Workshop 2015*, pages 111–115, Parramatta, Australia.
- Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia Song. 2021. Coco-lm: Correcting and contrasting text sequences for language model pretraining. *ArXiv*, abs/2102.08473.
- R. Pappagari, Piotr Żelasko, Jesús Villalba, Yishay Carmiel, and Najim Dehak. 2019. Hierarchical transformers for long document classification. *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 838–844.
- Nils Reimers, Philip Beyer, and Iryna Gurevych. 2016. Task-oriented intrinsic evaluation of semantic textual similarity. In *COLING*.

- 415 Tobias Rohde, Xiaoxia Wu, and Yinhan Liu. 2021. Hi-  
416 erarchical learning for generation with long source  
417 sequences. *ArXiv*, abs/2104.07545.
- 418 Evan Sandhaus. 2008. The new york times annotated  
419 corpus. *Linguistic Data Consortium, Philadelphia*,  
420 6(12):e26752.
- 421 Manzil Zaheer, Guru Guruganesh, Kumar Avinava  
422 Dubey, Joshua Ainslie, Chris Alberti, Santiago  
423 Ontañón, Philip Pham, Anirudh Ravula, Qifan  
424 Wang, Li Yang, and Amr Ahmed. 2020. Big  
425 bird: Transformers for longer sequences. *ArXiv*,  
426 abs/2007.14062.
- 427 Yan Zhang, Ruidan He, Zuozhu Liu, Kwan Hui Lim,  
428 and Lidong Bing. 2020. An unsupervised sentence  
429 embedding method by mutual information maxi-  
430 mization. In *Conference on Empirical Methods in*  
431 *Natural Language Processing*.
- 432 Ming Zhong, Yang Liu, Yichong Xu, Chenguang Zhu,  
433 and Michael Zeng. 2021. Dialoglm: Pre-trained  
434 model for long dialogue understanding and summa-  
435 rization. In *AAAI Conference on Artificial Intelli-*  
436 *gence*.

## A Redundancy

Length	(1)	(2)	(3)	(4)	(5)	All
FakeNews	1.06	1.21	29	1.35	1.52	1.37
20News	1.12	1.18	1.24	1.31	1.50	1.28
arXiv	1.12	1.25	1.36	1.49	1.62	1.34
NYT	1.00	1.14	1.21	1.31	1.48	1.45
BBCNews	1.05	1.14	1.20	1.29	1.46	1.19

Table 3: Information redundancies for different lengths (i.e., word numbers) of text: (1) 0-50 (2) 51-100 (3) 101-200 (4) 201-300 (5) more than 300.

We evaluate the redundancy of the text by counting the repeated verbs and nouns in the text. Specifically, we first use SpaCy<sup>7</sup> to find verbs and nouns and get their lemmatizations. Intuitively, if the redundancy of a document is high, nouns and verbs will be repeated frequently to express the same topic. Hence, redundancies  $R$  in our paper are computed as:

$$R = \frac{N_{\text{nouns,verbs}}}{D_{\text{nouns,verbs}}} \quad (2)$$

where  $N_{\text{nouns,verbs}}$  denotes the number of nouns and verbs in a document and  $D_{\text{nouns,verbs}}$  is the number of distinct nouns and verbs.

## B Training Details

For text classification, the learning rate for fine-tuning is  $3e-4$ ; the batch size is 8; the maximum sequence length is 512 tokens. We fine-tune the last MLP layer on these five datasets and evaluate the classification performance with accuracy and macro-F1 scores. For few-shot text classification, we sample 10 data instances per class for the FakeNewsCorpus dataset and the arXiv dataset and 5 data instances per class for the other three datasets. Other settings are the same as the general text classification. Since there is randomness in sampling, we repeat every experiment 10 times and take the average value of metrics.

## C Attention Weights

We compute the attention weights for Transformers as follows: (1) we first extract the attention weights between [CLS] token and all the other tokens; (2) we compute the averaged weights along different heads in multi-head attention; (3) the attention weights of the last layer in Transformers are

used as the weights for words. Averaged values are computed for nouns, verbs, adjectives, and other words.

<sup>7</sup><https://spacy.io/>