

Prompt Combines Paraphrase: Enhancing Biomedical “Pre-training, Prompt and Predicting” Models by Explaining Rare Biomedical Concepts

Anonymous ACL submission

Abstract

Prompt-based fine-tuning for pre-trained models has proven resultful in general domains for few-shot learning in downstream tasks. As to the biomedical domain, rare biomedical entities, which are quite ubiquitous in health-care contexts, can affect the performance of pre-trained models, especially in low-resource scenarios. We propose a simple yet effective approach to helping models understand rare biomedical words during tuning with prompt. Experiments demonstrate that our method can achieve up to 5% improvement in biomedical tasks without any additional parameters or training steps in few-shot vanilla prompt settings.

1 Introduction

Pre-trained models (PTMs) have achieved a great success in natural language processing (NLP) and become a new paradigm for various tasks (Peters et al., 2018; Devlin et al., 2019; Liu et al., 2019; Qiu et al., 2020). Many studies also have paid attention to PTMs in biomedical NLP tasks (Lee et al., 2020; Lewis et al., 2020; Zhao et al., 2021). However, it is clear that PTMs cannot do very well in biomedical NLP tasks due to its internal characteristics of biomedical texts.

In general, there are two challenges for applying PTMs to biomedical NLP tasks, i.e., 1) **limited data** and 2) **rare biomedical words**. Firstly, it is common that the amount of biomedical labeled data is limited due to data sensitivity (Šuster et al., 2017), high cost and professional requirement for data annotation. PTMs perform poorly with few samples since abundant training samples are essential to optimize task-related parameters. Secondly, biomedical terms are usually low-frequency words and are critical to understanding biomedical texts. As an example of natural language inference (NLI) task in Figure 2 in Appendix A, the model goes wrong when faced with a rare words “*afebrile*¹” in

the premise ,whose meaning is “having no fever”. It’s hard for PTMs to predict the label right without knowing “afebrile”. Thus, PTMs cannot capture the precise semantics of biomedical texts without sufficient information of biomedical rare terms.

With very few annotated samples for a new task, it is hard to fine-tune the PTMs and the new task-specific parameters effectively. Prompt technique has been introduced to smooth the fine-tuning process in few-shot setting by closing the gap between pre-training stage and the downstream task in general domains (Liu et al., 2021), as demonstrated in Figure 1. Similarly, the few-shot setting is also a pervasive challenge in biomedical domain mentioned above. Therefore, it is reasonable to adapt “pre-training, prompt and predicting” framework to biomedical NLP tasks.

Furthermore, the challenge of rare words, which is a critical problem for biomedical PTMs, has not been widely explored. Only a handful of works have studied this issue and they focus on enriching the representation of rare words through pre-training stage (Schick and Schütze, 2020; Yu et al., 2021; Wu et al., 2020). Thus, it’s necessary for them to involve a second-round pre-training to enrich specific rare words upon PTMs, which is highly time-consuming and low-efficiency. Alternatively, we emphasize on tuning stage instead of pre-training, leading to an efficient approach. Specifically, we propose to explain biomedical concepts on the basis of “pre-training, prompt and predicting” framework. The new approach could manage to enhance tuning capability in the aspect of understanding biomedical concepts. Besides, our approach is a plug-in module for specific datasets and model-agnostic, which can be easily transferred to other domains and models².

the BC-RoBERTa-Large and “afebrile” appears only about 100,000 times. For comparison, the frequency of “fever” is 5 times that of “afebrile”.

²We plan to release our code at <http://XXX>

¹There are around 4 billion words in pre-training texts of

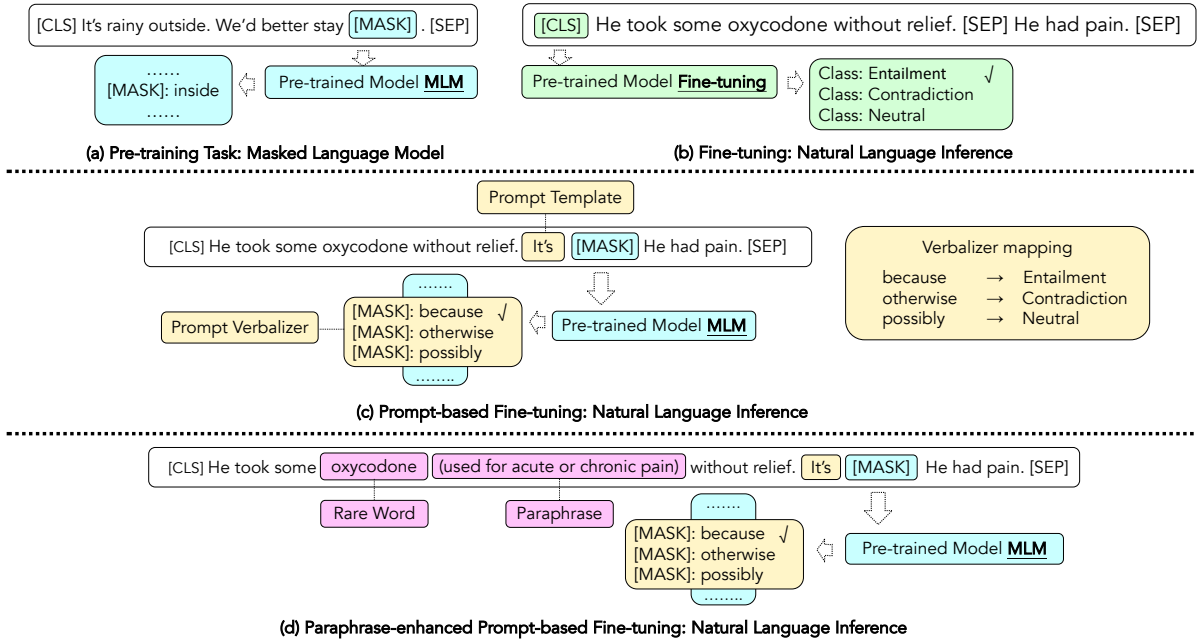


Figure 1: Examples for paradigms of (a) MLM (Masked Language Model) Pre-training. (b) Task-specific fine-tuning. (c) Prompt-based fine-tuning, with same task as pre-training process. (d) Paraphrase-enhanced prompt-based fine-tuning. Best viewed in color.

In summary, our contributions are as follows:

- We investigate a valuable problem of adapting PTMs in scenarios of biomedical text understanding of few samples and rare words, which likely has great impacts on biomedical text mining.
- We propose a novel approach to combine prompt technique and paraphrases of rare words in the PTMs tuning stage to solve the above two challenges.
- We evaluate over two biomedical natural language understanding datasets and our approach can improve the performance by up to 5% in the few-shot setting and 0.6% with a full-size training dataset. Moreover, we discuss how the paraphrases help with the PTMs and provide a perspective about task-related rare words.

2 Related Work

Word frequencies in PTMs Words in the vocabulary list follow a Zipf distribution (Zipf, 2016) by and large. Several previous works have discussed that the word representation space of PTMs is anisotropic and high-frequency words dominates

the representation of a sentence inducing a semantic bias (Gao et al., 2019; Li et al., 2020; Yan et al., 2021). Meanwhile, it has been also proven that rare words hamper the PTMs to perform well in which the uncommon words play a decisive role in the sentence understanding (Schick and Schütze, 2020; Wu et al., 2020; Yu et al., 2021). Schick and Schütze (2020) introduces one-token approximation to infer the embedding of arbitrary rare word by a single token. Wu et al. (2020) proposes taking notes on the fly to maintain a note dictionary for rare words to save the contextual information which helps enhance the representation of pre-training.

Biomedical PTMs With the booming trend of PTMs in NLP tasks (Peters et al., 2018; Devlin et al., 2019; Liu et al., 2019), various trials have been made in biomedical domain (Peng et al., 2019; Lee et al., 2020; Huang et al., 2019) by pre-training on biomedical texts. And then, Lewis et al. (2020) and Gu et al. (2021) further get the domain-specific vocabulary list to amend representation of biomedical words. Recently, biomedical PTMs are guided with domain knowledge. Zhang et al. (2021) amplifies the biomedical entities with type semantic information of neighbor entities. Michalopoulos et al. (2021) learns clinical term embedding with relevant meaning and semantic type.

PTMs tuning with prompt Many works are dedicated to applying prompt in fine-tuning by adapting the downstream tasks to the paradigm of pre-training tasks. Prompts that have been employed by now fall into two groups: discrete prompt, described by natural language (Schick and Schütze, 2021; Gao et al., 2021); and continuous prompt, based on trainable vectors (Li and Liang, 2021; Shin et al., 2020b).

3 Rare Biomedical Words and Paraphrases

In this section, we introduce how we find the rare biomedical words and the method we adopt to supplement paraphrases to those words with the prompt-based tuning of PTM.

3.1 Selection of Rare Biomedical Words

“Rare” is a relative concept, which is context-relevant in most cases. We use the RoBERTa-Large model proposed by Lewis et al. (2020) that has been pre-trained adequately on biomedical corpora (details in Appendix B). We download the entire corpora above and loop them through to obtain the frequency of each word in the pre-training phase. In place of involving all rare words, we opt for rare words in biomedical domain for two reasons:

- 1) Word distribution in general domains differs from that in biomedical domain (Lee et al., 2020).
- 2) Biomedical rare words can be worth more than general rare words to biomedical tasks. To obtain rare words, we set a threshold on the word frequency in the pre-training corpora empirically as a hyper-parameter similar to Yu et al. (2021). Afterwards, with the help of an online dictionary - Wiktionary³, we can retrieve the paraphrases of rare words along with the category labels. Only rare words with health-related labels are reserved as rare biomedical words. Full list of selected labels is available in Appendix C.

3.2 Selection of Paraphrases

To avoid introducing noise information from paraphrases, rare biomedical words with more than one paraphrase are eliminated. Also, there should be no additional rare words in the paraphrases. Therefore, we filter out the paraphrases in which frequency of any word is lower than the same threshold mentioned before.

³<https://en.wiktionary.org/>

3.3 Prompt-based Fine-Tuning with Paraphrases

When we read and come across new words, we will consult a proper dictionary for their definitions. Analogously, when the biomedical PTM deal with downstream tasks, we provide the model with paraphrases of biomedical rare words surrounded by brackets attached to the rare words, as Figure 1(d). In this case, given a PTM, paraphrases of biomedical rare words can be considered as a portable plug-in module and generated for any datasets instantly before prompt-based fine-tuning.

4 Experiments

4.1 Setup

Model We use a Biomedical-Clinical-RoBERTa-Large model mentioned in Section 3.1 as a strong baseline to verify our approach.

Datasets Note that rare words hinder the PTMs more in natural language understanding (NLU) (Schick and Schütze, 2020) than in other NLP tasks. However, most biomedical and clinical NLP tasks fall in the category of information extraction (Shin et al., 2020a; Gu et al., 2021). Thus, we evaluate our method over two NLU-relevant biomedical tasks - MedNLI (Romanov and Shivade, 2018) and MedSTS⁴ (Wang et al., 2020). Respectively, MedNLI is an NLI dataset in which premises are made up with clinical notes in MIMIC-III and MedSTS is a semantic textual similarity dataset gathered from a clinical corpus at Mayo Clinic. Semantic Textual Similarity is a regression task and we adapt the task following Gao et al. (2021). Statistics of datasets can be found in Appendix D. We sample from 16 up to 256 samples from the original training sets as training and development sets with 10 different random seeds and use full-size testing sets.

Prompt settings We combine the prompt settings from Schick and Schütze (2021) and Gao et al. (2021) for the NLI and STS tasks without further adaption (details in Appendix E) to explore the effectiveness of paraphrases of biomedical rare words rather than the prompt paradigm.

4.2 Main Results and Analysis

We report average accuracy for MedNLI and pearson correlation coefficient for MedSTS along with

⁴We use ClinicalSTS-2018 and 2019 which are sub-datasets of MedSTS provided by the maintainers of MedSTS project.

MedNLI					
#Samples	16	32	64	128	256
Model					
BC-RoBERTa-Large	51.3 (5.9)	60.6 (6.7)	71.0 (3.7)	80.6 (1.3)	83.1 (1.3)
+ paraphrase	56.6 (5.0)	62.3 (6.0)	74.5 (3.0)	81.1 (1.5)	83.6 (1.0)

MedSTS: Clinical-2019					
#Samples	16	32	64	128	256
Model					
BC-RoBERTa-Large	41.1 (11.8)	53.9 (6.7)	67.9 (7.4)	73.1 (5.0)	80.4 (3.1)
+ paraphrase	45.2 (9.3)	57.3 (6.8)	67.2 (7.5)	74.5 (3.7)	79.6 (2.6)

MedSTS: Clinical-2018					
#Samples	16	32	64	128	256
Model					
BC-RoBERTa-Large	54.2 (8.1)	63.9 (9.2)	73.3 (3.8)	77.4 (2.7)	81.5 (1.5)
+ paraphrase	53.0 (7.4)	67.2 (6.6)	74.5 (2.7)	79.1 (1.6)	81.8 (1.2)

Table 1: Our main results on three dataset: MedNLI, MedSTS: Clinical-2018 and Clinical-2019, using BC-RoBERTa-Large (Biomedical and Clinical RoBERTa-Large) (Lewis et al., 2020) with different size of training sets. We report average (and standard deviation) performance (accuracy for MedNLI and pearson correlation coefficient for MedSTS) over 10 different random seeds. + paraphrase: with paraphrases of rare biomedical words.

standard deviation. Table 1 shows results for biomedical natural language inference and semantic textual similarity tasks. Model with paraphrases for rare biomedical words can outperform the baseline in most cases. Paraphrases bring about up to 5% improvement on average for few-shot learning with 16 training samples as to MedNLI task and 0.5% increment with 256 training samples. We can see that PTMs tend to learn more about rare words with more training samples but paraphrases still act well. As to MedSTS, appended paraphrases are also shown as an effective strategy for most cases. In addition, tuning with paraphrases also generally improves model stability and reduces the variance of model prediction in few-shot scenarios.

5 Discussion

Train with more samples Apart from infusing dictionary paraphrases in few-shot scenarios, we also attempt with more training samples, even with full-size training dataset. Table 4 in Appendix F demonstrates that with larger amount of training samples, our method still advances the PTMs for majority cases, implying that paraphrases of rare biomedical words are not only impactful in few-sample situations.

Which to look up? By far, experiment results have attested that paraphrases of rare biomedical words help with PTMs in training with few or more samples. Nevertheless, it may not always work

well. We scrutinize the cases that model predicts differently after paraphrases being appended and display several cases in Table 5 in Appendix G. Table 5 shows that paraphrases of rare words which are task-related and decisive in understanding the whole sentence can be beneficial to PTMs. Otherwise, paraphrases can involve more confusion than certainty. When human reads, we probably won't look up a new word until it blocks our understanding. Similarly, it is worthwhile to explore how to attach *helpful* paraphrases or utilize knowledge selectively in future research.

6 Conclusion

Biomedical terms, which are pervasive in biomedical texts, are sometimes rare in the whole corpora and domain-specific rare words understanding remains as a tough challenge for pre-trained models. In this paper, we present a simple yet effective method to help biomedical pre-trained models grasp the semantics of rare biomedical words, that is attaching paraphrases to rare biomedical words as a plug-in approach in the prompt-tuning datasets without additional parameters to train during pre-training and downstream task-related tuning. Experiments show that our method can substantially boost the performance of biomedical pre-trained model in few-shot setting and bring about plausible enhancement with more training data, even full-size of training set.

278	Ethical Considerations		
279	In this work, we propose an approach to explaining		
280	rare biomedical words for biomedical PTMs to		
281	help understand sentences with rare biomedical		
282	words. We conduct our experiments on the public		
283	biomedical datasets MedNLI and MedSTS with		
284	the authorization from the respective maintainers		
285	of the datasets. All biomedical data involved have		
286	been de-identified by dataset providers and only		
287	used for research.		
288	References		
289	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and		
290	Kristina Toutanova. 2019. Bert: Pre-training of deep		
291	bidirectional transformers for language understand-		
292	ing. In <i>NAACL-HLT (1)</i> .		
293	Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and		
294	Tieyan Liu. 2019. Representation degeneration prob-		
295	lem in training natural language generation models.		
296	In <i>International Conference on Learning Representa-</i>		
297	<i>tions</i> .		
298	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021.		
299	Making pre-trained language models better few-shot		
300	learners . In <i>Proceedings of the 59th Annual Meet-</i>		
301	<i>ing of the Association for Computational Linguistics</i>		
302	<i>and the 11th International Joint Conference on Natu-</i>		
303	<i>ral Language Processing (Volume 1: Long Papers)</i> ,		
304	pages 3816–3830, Online. Association for Computa-		
305	tional Linguistics.		
306	Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto		
307	Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng		
308	Gao, and Hoifung Poon. 2021. Domain-specific lan-		
309	guage model pretraining for biomedical natural lan-		
310	guage processing. <i>ACM Transactions on Computing</i>		
311	<i>for Healthcare (HEALTH)</i> , 3(1):1–23.		
312	Kexin Huang, Jaan Altosaar, and Rajesh Ranganath.		
313	2019. Clinicalbert: Modeling clinical notes and		
314	predicting hospital readmission. <i>arXiv preprint</i>		
315	<i>arXiv:1904.05342</i> .		
316	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon		
317	Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang.		
318	2020. Biobert: a pre-trained biomedical language		
319	representation model for biomedical text mining.		
320	<i>Bioinformatics</i> , 36(4):1234–1240.		
321	Patrick Lewis, Myle Ott, Jingfei Du, and Veselin Stoy-		
322	anov. 2020. Pretrained language models for biomed-		
323	ical and clinical tasks: Understanding and extending		
324	the state-of-the-art. In <i>Proceedings of the 3rd Clini-</i>		
325	<i>cal Natural Language Processing Workshop</i> , pages		
326	146–157.		
327	Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang,		
328	Yiming Yang, and Lei Li. 2020. On the sentence		
329	embeddings from pre-trained language models . In		
	<i>Proceedings of the 2020 Conference on Empirical</i>		330
	<i>Methods in Natural Language Processing (EMNLP)</i> ,		331
	pages 9119–9130, Online. Association for Computa-		332
	tional Linguistics.		333
	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:		334
	Optimizing continuous prompts for generation . In		335
	<i>Proceedings of the 59th Annual Meeting of the Asso-</i>		336
	<i>ciation for Computational Linguistics and the 11th</i>		337
	<i>International Joint Conference on Natural Language</i>		338
	<i>Processing (Volume 1: Long Papers)</i> , pages 4582–		339
	4597, Online. Association for Computational Lin-		340
	guistics.		341
	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,		342
	Hiroaki Hayashi, and Graham Neubig. 2021. Pre-		343
	train, prompt, and predict: A systematic survey of		344
	prompting methods in natural language processing.		345
	<i>arXiv preprint arXiv:2107.13586</i> .		346
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-		347
	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,		348
	Luke Zettlemoyer, and Veselin Stoyanov. 2019.		349
	Roberta: A robustly optimized bert pretraining ap-		350
	proach. <i>arXiv preprint arXiv:1907.11692</i> .		351
	George Michalopoulos, Yuanxin Wang, Hussam Kaka,		352
	Helen Chen, and Alexander Wong. 2021. Umlsbert:		353
	Clinical domain knowledge augmentation of contex-		354
	tual embeddings using the unified medical language		355
	system metathesaurus. In <i>Proceedings of the 2021</i>		356
	<i>Conference of the North American Chapter of the</i>		357
	<i>Association for Computational Linguistics: Human</i>		358
	<i>Language Technologies</i> , pages 1744–1753.		359
	Yifan Peng, Shankai Yan, and Zhiyong Lu. 2019. Tran-		360
	sfer learning in biomedical natural language process-		361
	ing: An evaluation of bert and elmo on ten bench-		362
	marking datasets. In <i>Proceedings of the 18th BioNLP</i>		363
	<i>Workshop and Shared Task</i> , pages 58–65.		364
	Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt		365
	Gardner, Christopher Clark, Kenton Lee, and Luke		366
	Zettlemoyer. 2018. Deep contextualized word rep-		367
	resentations. In <i>Proceedings of NAACL-HLT</i> , pages		368
	2227–2237.		369
	Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao,		370
	Ning Dai, and Xuanjing Huang. 2020. Pre-trained		371
	models for natural language processing: A survey.		372
	<i>Science China Technological Sciences</i> , pages 1–26.		373
	Alexey Romanov and Chaitanya Shivade. 2018.		374
	Lessons from natural language inference in the clini-		375
	cal domain. In <i>Proceedings of the 2018 Conference</i>		376
	<i>on Empirical Methods in Natural Language Process-</i>		377
	<i>ing</i> , pages 1586–1596.		378
	Timo Schick and Hinrich Schütze. 2020. Rare words:		379
	A major problem for contextualized embeddings and		380
	how to fix it by attentive mimicking. In <i>Proceedings</i>		381
	<i>of the AAAI Conference on Artificial Intelligence</i> ,		382
	volume 34, pages 8766–8774.		383

- 384 Timo Schick and Hinrich Schütze. 2021. Exploiting
385 cloze-questions for few-shot text classification and
386 natural language inference. In *Proceedings of the*
387 *16th Conference of the European Chapter of the Asso-*
388 *ciation for Computational Linguistics: Main Volume,*
389 pages 255–269.
- 390 Hoo-Chang Shin, Yang Zhang, Evelina Bakhturina,
391 Raul Puri, Mostofa Patwary, Mohammad Shoeybi,
392 and Raghav Mani. 2020a. Bio-megatron: Larger
393 biomedical domain language model. In *Proceed-*
394 *ings of the 2020 Conference on Empirical Methods*
395 *in Natural Language Processing (EMNLP)*, pages
396 4700–4706.
- 397 Taylor Shin, Yasaman Razeghi, Robert L Logan IV,
398 Eric Wallace, and Sameer Singh. 2020b. Eliciting
399 knowledge from language models using automati-
400 cally generated prompts. In *Proceedings of the 2020*
401 *Conference on Empirical Methods in Natural Lan-*
402 *guage Processing (EMNLP)*, pages 4222–4235.
- 403 Simon Šuster, Stéphan Tulkens, and Walter Daelemans.
404 2017. A short review of ethical challenges in clin-
405 ical natural language processing. *arXiv preprint*
406 *arXiv:1703.10090*.
- 407 Yanshan Wang, Naveed Afzal, Sunyang Fu, Liwei
408 Wang, Feichen Shen, Majid Rastegar-Mojarad, and
409 Hongfang Liu. 2020. Medsts: a resource for clinical
410 semantic textual similarity. *Language Resources and*
411 *Evaluation*, 54(1):57–72.
- 412 Qiyu Wu, Chen Xing, Yatao Li, Guolin Ke, Di He, and
413 Tie-Yan Liu. 2020. Taking notes on the fly helps bert
414 pre-training. *arXiv preprint arXiv:2008.01466*.
- 415 Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang,
416 Wei Wu, and Weiran Xu. 2021. **ConSERT: A con-**
417 **trastive framework for self-supervised sentence repre-**
418 **sentation transfer**. In *Proceedings of the 59th Annual*
419 *Meeting of the Association for Computational Lin-*
420 *guistics and the 11th International Joint Conference*
421 *on Natural Language Processing (Volume 1: Long*
422 *Papers)*, pages 5065–5075, Online. Association for
423 Computational Linguistics.
- 424 Wenhao Yu, Chenguang Zhu, Yuwei Fang, Donghan Yu,
425 Shuohang Wang, Yichong Xu, Michael Zeng, and
426 Meng Jiang. 2021. Dict-bert: Enhancing language
427 model pre-training with dictionary. *arXiv preprint*
428 *arXiv:2110.06490*.
- 429 Taolin Zhang, Zerui Cai, Chengyu Wang, Minghui
430 Qiu, Bite Yang, and Xiaofeng He. 2021. Smedbert:
431 A knowledge-enhanced pre-trained language model
432 with structured semantics for medical text mining.
433 In *Proceedings of the 59th Annual Meeting of the*
434 *Association for Computational Linguistics and the*
435 *11th International Joint Conference on Natural Lan-*
436 *guage Processing (Volume 1: Long Papers)*, pages
437 5882–5893.
- 438 Sendong Zhao, Chang Su, Zhiyong Lu, and Fei Wang.
439 2021. Recent advances in biomedical literature min-
440 ing. *Briefings in Bioinformatics*, 22(3):bbaa057.
- George Kingsley Zipf. 2016. *Human behavior and the
principle of least effort: An introduction to human
ecology*. Ravenio Books.

Appendix

A Wrong Case of Biomedical PTM

Task:	Medical Natural Language Inference
Premise:	Lactate only 1.3 and pt afebrile .
Hypothesis:	Temperature was within normal range.
Gold label:	Entailment
Model Prediction:	Neutral
Paraphrase:	afebrile - having no fever

Figure 2: Wrong case of BC-RoBERTa-Large model fully pre-trained on Biomedical and Clinical texts (Lewis et al., 2020) and fine-tuned on MedNLI tasks caused by not understanding biomedical rare word - “afebrile”.

B Pre-trained Corpora of the Model

We use a biomedical and clinical RoBERTa-Large (Lewis et al., 2020) trained on biomedical corpora, including PubMed abstract⁵, PubMed Central⁶ (PMC) full-text and MIMIC-III dataset⁷.

C Word Labels for Rare Biomedical Words

We focus on the rare words which have been tagged with labels that contain any of following medicine-related strings:

[‘medical’, ‘medicine’, ‘disease’, ‘symptom’, ‘pharma’]

D Dataset

We conduct our experiments on MedNLI and MedSTS datasets. Specifically, we use the available sub-datasets ClinicalSTS-2018 and ClinicalSTS-2019 for MedSTS provided by the maintainer of MedSTS project. The statistics of involved datasets can be found in Table 2. Note that there is no development set split in MedSTS. Therefore, we sample the development set for MedSTS from its training set with the same quantity as sampled few-shot training set and make sure there is no overlap between training and development set.

⁵<https://pubmed.ncbi.nlm.nih.gov>

⁶<https://www.ncbi.nlm.nih.gov/pmc>

⁷<https://physionet.org/content/mimiciii/1.4/>

Dataset	Train	Dev	Test
MedNLI	11232	1395	1422
MedSTS: ClinicalSTS-2019	1642	/	412
MedSTS: ClinicalSTS-2018	750	/	318

Table 2: Statistics of datasets MedNLI and MedSTS

E Prompt Settings

We adopt the prompt settings empirically from Schick and Schütze (2021) and Gao et al. (2021) for the natural language inference and semantic textual similarity tasks shown in Table 3 since the prompt paradigm is not the core of this work and our method is prompt-agnostic.

Task	Prompt Template	Prompt Verbalizers
MedNLI	<Sent1>. [MASK]. <Sent2>	Yes/No/maybe
MedSTS	<Sent1>. [MASK]. <Sent2>	Yes/No

Table 3: Prompt settings for MedNLI and MedSTS

F Train with More Samples

Besides few-shot scenarios, we also train with more samples for MedNLI since it has 11,232 training samples. Experiment results are shown in Table 4.

G Case Analysis

We display several cases in which model predicts differently with or without paraphrases of rare biomedical words from MedNLI in Table 5. From the cases, we can see that paraphrases of rare biomedical words that are determinant in sentence understanding can be helpful to pre-trained model while paraphrases of those irrelevant rare biomedical words may confuse the model.

MedNLI						
#Samples	512	1024	2048	4096	8192	full-size
Model						
BC-RoBERTa-Large	85.0(0.8)	85.6(0.8)	86.6(0.6)	86.7(0.7)	86.2(0.6)	86.1(0.6)
+ paraphrases	85.2(0.7)	86.3(0.9)	86.4(0.7)	86.4(0.5)	86.7(0.6)	86.7(0.7)

Table 4: Test results on MedNLI dataset with larger size of training sets. We report average (and standard deviation) accuracy over 10 different random seeds. BC-RoBERTa-Large: **B**iomedical and **C**linical RoBERTa-Large (Lewis et al., 2020). + paraphrase: with paraphrases of rare biomedical words.

Sentence Pairs	w/o paraphrases	w/ paraphrases
P: She was found to have BRBPR (<i>bright red blood per rectum</i>) on rectal exam. H: the patient had bright read blood per rectum	Neutral	Entailment (right answer)
P: Antenatal history - pregnancy complicated by chronic hypertension with increased gestational hypertension leading to admission 3 days prior to delivery followed by cesarean section. H: The patient had proteinuria (<i>The presence of protein in the urine</i>) during pregnancy	Entailment	Neutral (right answer)
P: Following this rehab admission she was sent to a different OSH on [**2725-10-26**], for acute CHF (<i>congestive heart failure</i>) and at least one PEA arrest. H: The patient has a poorly functioning heart.	Contradiction	Entailment (right answer)
P: The patient was sent to the HD unit prior to coming to the floor for workup (<i>A general medical examination to assess a persons health and fitness</i>) of fever. H: The patient has an infection	Neutral (right answer)	Contradiction
P: - COPD (<i>chronic obstructive pulmonary disease</i>) - obesity - unspecified hypoxemia - CNS lymphoma c/b CVAs x3 (posterior circulation) and seizure d/o - history of SAH while on coumadin - diastolic heart failure - coronary artery disease - atrial fibrillation - hypertension - hyperlipidemia - severe OSA (did not tolerate CPAP in the past) - primary hyperparathyroidism/25-vit D deficiency c/b nephrolithiasis - toxic multinodular goiter with subclinical (<i>Less than is needed for clinical reasons</i>) hyperthyroidism - neovascular glaucoma c/b right eye blindness H: Patient has a history of malignancy	Neutral (right answer)	Entailment

Table 5: Cases that BC-RoBERTa-Large predicts differently after the supplement of paraphrases for rare biomedical words in MedNLI. “P” for Premise and “H” for Hypothesis. Words in **bold** are rare biomedical words and expressions in *italic* inside the brackets are the paraphrases of rare words.