

---

# AN EXPERIMENTAL EVALUATION OF TRANSFORMER-BASED LANGUAGE MODELS IN THE BIOMEDICAL DOMAIN

---

**Paul Grouchy\***    **Shobhit Jain\***    **Michael Liu\***    **Kuhan Wang\***    **Max Tian\***    **Nidhi Arora\***  
 Untether AI    Manulife    Tealbook    CIBC    Adeptmind    Intact

**Hillary Ngai\***    **Faiza Khan Khattak †**    **Elham Dolatabadi ‡**    **Sedef Akinli Kocak §**  
 University of Toronto    Manulife    Vector Institute    Vector Institute

January 1, 2021

## ABSTRACT

With the growing amount of text in health data, there have been rapid advances in large pre-trained models that can be applied to a wide variety of biomedical tasks with minimal task-specific modifications. Emphasizing the cost of these models, which renders technical replication challenging, this paper summarizes experiments conducted in replicating BioBERT and further pre-training and careful fine-tuning in the biomedical domain. We also investigate the effectiveness of domain-specific and domain-agnostic pre-trained models across downstream biomedical NLP tasks. Our finding confirms that pre-trained models can be impactful in some downstream NLP tasks (QA and NER) in the biomedical domain; however, this improvement may not justify the high cost of domain-specific pre-training.

**Keywords** BioBERT · Biomedical data · NLP downstream tasks · Transformer-based models

## 1 Introduction

There have been increased exploration and extension of transformer-based deep learning models for NLP [1, 2, 3, 4] recently. Using transformer-based models, one can pre-train a deep learning model on large datasets and then easily fine-tune it to adapt to downstream NLP tasks. There are three factors that affect the performance of these models: (a) the size of the dataset, (b) the availability of computational resources, and (c) the expressiveness of the model architecture [2, 1]. Because of these factors, the cost and complexity of developing pre-trained models are rising quickly and limit the capability of reproducing results when sufficient resources are not available.

These language models are trained on text corpora of general domains. For example, BERT [3], Bidirectional Encoder Representations from Transformers has been trained on Wikipedia and BooksCorpus. There has also been a rapid growth in NLP for the biomedical domain [5] and many new methods including transformer-based methods are being used for different biomedical tasks involving NLP.

The performance of language models that have been trained on general domains are not yet fully investigated in more specific domains such as biomedical, finance or legal. Therefore, it is worth investigating if large amounts of domain-specific data may help in getting better results, or if similar results may be acquired by using smaller-sized data.

We therefore focus on biomedical domain in order to answer the following questions:

---

\*Equal Contribution

†Corresponding Author: faizakhankhattak@gmail.com

‡Corresponding Author: elham.dolatabadi@vectorinstitute.ai

§Corresponding Author: sedef.kocak@vectorinstitute.ai

1. *Does domain-specific training improve performance compared to baseline models trained on domain-agnostic corpora?*
2. *Is it possible to obtain comparable results from a domain-specific BERT model pre-trained on smaller-sized data?* While it is established fact that with transfer learning, a model can be trained once and be reused for several tasks but there are a few cases when the model needs to be retrained. For example, when the data is dynamic and may change over time due to the data-shift; data can belong to a wide variety of domains, therefore the model may need to be retrained for the new domain data or even sub-domain; data can be confidential to an organization hence the model needs to be retrained for that organization-specific needs. The domain and/or organization specific datasets can be very small. Moreover, not every organization has extensive computational resources required to train large models. In such cases it is helpful to know if a small domain-specific data can be used to get comparable results.

The rest of this paper is organized as follows. We review the existing studies related to our work in Section 2. We explain our pre-training experiments and the results in Section 3.1 and fine-tuning experiments and results in Section 3.2. We complete our paper with discussing and conclusion in Section 4 where we discuss our results and answering our questions.

## 2 Related Work

Large-scale pre-trained language models such as BERT[3], GPT-2 [6], RoBERTa [2] and GPT-3 [7] have shown to outperform state-of-the-art performance in many NLP tasks such as Named Entity Recognition (NER) and Question Answering (QA). Moreover, several studies have used transfer learning and fine-tuning of these models on English NLP tasks (e.g., [8, 9]). These language models are trained on text corpora of general domains but recently there has been a trend of training language models on the domain-specific data. For example, financial version of BERT was introduced by Araci [10] where he only studied sentiment classification task. Ma *et al* [11] fine-tuned BERT on legal documents coming from their proprietary corpus. BioBERT [4] was introduced as an extension of BERT that is further pre-trained on the domain-specific biomedical corpora including PubMed and PubMed Central (PMC).

In specialized domains like biomedical, recent work has shown that using domain specific data can provide improvement over general-domain language models [12]. In this regard, Wang *et al.* [13] showed that word-embeddings trained on biomedical corpora captured the semantics of medical terms better than those trained on general domain corpora, but may not generalize well to downstream biomedical NLP tasks such as biomedical information retrieval. Zhao *et al* [14] showed that word2vec [15] trained on a smaller and in-domain medical data resulted in better performance than the word2vec trained on a large and general domain dataset. Also, [16] found that performance decreases after 4 million distinct words of training data based on experiments with medical data from PubMed abstracts<sup>5</sup>. In a recent study, Gu *et al* [12] introduced PubMedBERT. They pre-trained the BERT from scratch with PubMed articles and a customized vocabulary (constructed from the PubMed articles). This study indicates that a proper vocabulary helps the performance of downstream tasks in specific domains. However, training the model from scratch is extremely expensive in terms of data and computation. Researchers have specifically built adaptations of BERT that attempt to address different domain related problems but the most effective pretraining process remains an open research problem [12]. We replicated some of the BioBERT original study results, and better tune the training of BioBERT for better understanding domain specific training.

## 3 Methods and Experimental Results

We set BERT<sub>BASE</sub> as our baseline model. We started with BERT<sub>BASE</sub>, then pre-trained it on the PubMed abstracts data (BERT<sub>base+PM</sub>) and leveraged for evaluation of downstream tasks: Name-entity recognition (NER, relational extraction (RE) (Table 1) and question answering (QA) (Table 2). This allows us to conduct a fair comparison between domain-specific pre-training and fine-tuning. PyTorch implementation of BERT<sup>6</sup> [17] was leveraged and the replication experiments were conducted based on the work by McDermott *et. al* [18].

### 3.1 Pre-training language representations in the biomedical domain

The PubMed corpora was used for pre-training consists of paper abstract from millions of samples of biomedical text. While the original BioBERT study considers combined pre-training on PubMed, PMC, and Pubmed+PMC together,

<sup>5</sup><https://www.ncbi.nlm.nih.gov/pubmed/>

<sup>6</sup><https://github.com/huggingface/pytorch-transformers>

our model was pre-trained only on PubMed [19] to check the performance based on smaller data and meet the shared computing resources.

The PubMed data was processed into a format amenable for pre-training. The raw data consisted of approximately 200 million sentences in 30 GBs. The raw sentence data was batch processed into 111 chunks of ready to consume input data for BERT pre-training<sup>7</sup>. Technically, the NSP task was lower-cased sentences with a maximum length of 512 and masked at the sub-token level analogous to the original BERT<sub>BASE</sub>(uncased). The original BERT models were trained on data with maximum sequence lengths of 128 for the initial 90% of the training and 512 for the remainder. The full training loss is shown in Figure 1. Generally, a longer sequence length is preferable if the corpora tends to have longer

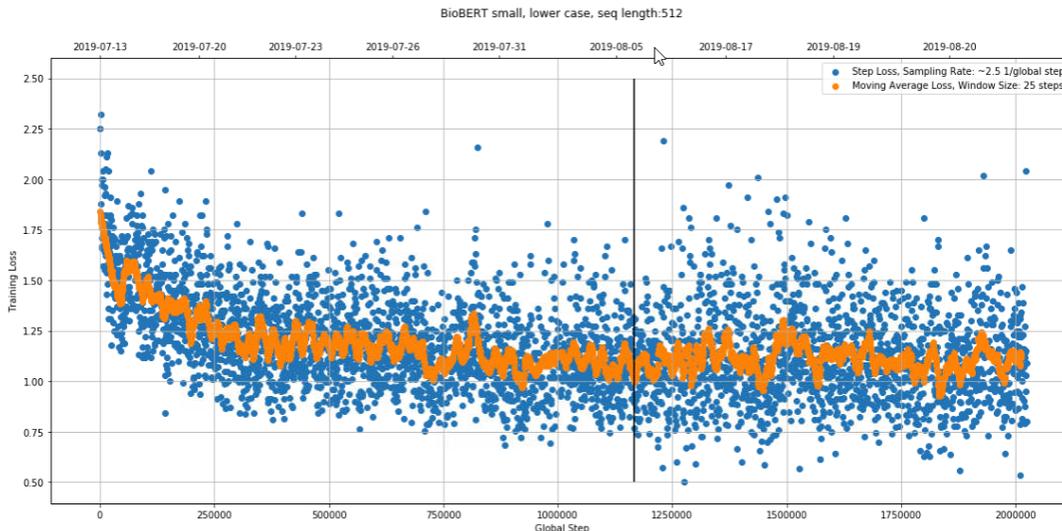


Figure 1: Training loss. The effect of chunking the data into shards is seen in the fluctuation of the training loss. The vertical line indicates separation between batch size 28 and 14 in the training paradigm. The top (bottom) x-axis indicate calendar date (global step).

passages but the amount of data and time to train is also increased. In this experiment, the training was maintained at the maximum length of 512 tokens throughout.

The pre-training experiment was designed to accommodate some realistic computing resource limitations. In a time-shared computing cluster GPU resources were allocated on limited priority basis by Slurm Workload Manager. In order to facilitate stable training and to accommodate a hardware environment where resources may be reallocated to higher priority users at any time and without warning, the following training features were implemented: (1) storing model weights and gradients at regular time intervals during training; (2) querying and automating job submission from system task scheduler; (3) automatically restoring training from data chunk and step as GPU resources became available.

### 3.2 Fine-tuning language representations in the biomedical domain

Following previous work, our domain specific pre-trained model (BERT<sub>base</sub>+PM) as well as BERT<sub>BASE</sub>(uncased) were evaluated on common biomedical NLP tasks as described below.

#### 3.2.1 Named-entity recognition (NER)

The biomedical NER task involves the extraction of biomedical named entities such as genes, proteins, diseases and species from unstructured biomedical text. This is a challenging task because of the unique characteristics of biomedical named entities such as sharing of head nouns, spelling forms per entity, ambiguous abbreviations, descriptive naming convention, and nested names [20]. In the original BioBERT study nine different biomedical NER datasets were tested and their best model (BERT<sub>BASE</sub>+PM+PMC) was able to outperform state-of-the-art on the majority. In this experiment, three NER datasets including NCBI Disease [21], JNLPBA [22], and Species-800 [23] were used. The model was

<sup>7</sup>This implies that the sampling of the second sentence for the Next Sentence Prediction (NSP) training task is not from the entire corpus but from within its respective chunk. As each chunk is still close to 2 million separate sentences, this was judged to be an acceptable compromise.

trained on 50 tokens long sentences with a learning rate of  $5x10^{-5}$ . These experiments were repeated several times and F1 scores corresponding to the epoch with the minimum validation loss are reported in Table 1. The standard deviation in F1 scores across different runs make definitive comparisons with the results reported in BioBERT study, but a surprising outcome is that the scores for the BERT<sub>BASE</sub> starting point models are higher on average than the ones fine-tuned with our BERT<sub>base</sub>+PM. This runs counter to the notion that pre-training on biomedical domains corpora improves downstream NER.

### 3.2.2 Relation Extraction (RE)

The GAD dataset [24] was used for the RE experiment. The goal of RE on the GAD dataset is to correctly classify whether a gene and disease are related, based on a given sentence. Two models were fine-tuned on GAD, one starting with BERT<sub>BASE</sub> and the other starting with BERT<sub>base</sub>+PM. A learning rate of  $5x10^{-5}$  was used for all runs. The best validation F1 score over several runs for each fold is reported. Our BERT<sub>base</sub>+PM outperformed other BERT models but the differences between models are not significant (see Table 1).

Tasks	Datasets	Metrics	BioBERT [4]		Our Experiments	
			BERT <sub>base</sub>	+PubMed	BERT <sub>base</sub>	BERT <sub>base</sub> +PM
NER	NCBI Disease	F1	85.63	<b>89.71</b>	84.08±2.07	80.33±2.40
	JNLPBA	F1	74.94	<b>77.49</b>	76.43±2.29	75.16±3.63
	Species 800	F1	71.63	74.06	<b>78.38±7.87</b>	74.59±5.28
RE	GAD	F1	79.30	79.83	79.60	<b>81.50</b>

Table 1: NER and RE Performance Results. Our experiments including fine-tuning of the BERT<sub>BASE</sub>(uncased) model and the BERT<sub>BASE</sub>(uncased) model further pre-trained on PubMed abstracts (+PM). For comparison, the results of BioBERT study using BERT<sub>BASE</sub> (i.e. Wiki+Books) and the BioBERT version of the PubMed trained model (+PubMed) are also shown. Best scores are shown in bold.

### 3.2.3 Question Answering (QA)

The goal of QA tasks is to automatically find the answer to a question posed in human language, usually from a context paragraph. In this study, we explored the performance of BERT and BioBERT on two common biomedical QA tasks including BioASQ [25] and PubmedQA [26]. For this task, three versions of BERT model (BERT<sub>large</sub>, BERT<sub>base</sub>, and our BERT<sub>base</sub>+PM) were fine-tuned on the QA datasets.

**BioASQ:** As was suggested in the BioBERT study [4], all BERT models were initially fine-tuned on the SQuAD [27] dataset (with intermediate evaluations), and then on the BioASQ training set before finally evaluating on the BioASQ test sets (Table 2). It has been also shown in [28] that pre-training BERT on SQuAD 1.1 generated better results when fine-tuned on BioASQ in comparison to the model pre-trained on SQuAD 2.0. For tuning the model on SQuAD 1.1 and SQuAD 2.0, we took inspiration from the training schemes outlined in [2] and [29] to adjust the hyperparameters, namely the learning rate,  $\beta_2$ , learning rate schedule, batch size, and number of training epochs. We found that training with a cosine learning rate schedule with no warm-up steps with  $\beta_2 = 0.98$  consistently resulted in the best performance, and thus carried over the heuristics generated from fine-tuning on SQuAD to further fine-tune our models on BioASQ. *Training on domain-specific data:* Overall, there is some evidence (Table 2) that pre-training on biomedical domain corpus improves performance on the downstream BioASQ QA task. However, the improvement is not so large as to be entirely convincing and carefully fine-tuned BERT models (such as BERT<sub>large</sub> in our case) can perform comparably to BioBERT.

*Zero-shot setting:* We also conducted an experiment in order to evaluate the performance of the BERT models on BioASQ datasets in a zero-shot setting where BERT<sub>BASE</sub> was fine-tuned on SQuAD and *not* on BioASQ training data. We averaged the scores over the 5 test sets within Strict Accuracy (S), Lenient Accuracy (L), and Mean Reciprocal Rank (M) metrics and the results were 31.37, 46.8, and 37.16, respectively. As we can see the zero-shot evaluation results on test sets are worse than the results indicated in Table 2 which emphasizes the effects of fine-tuning on BioASQ.

*Number of epochs:* Additionally, we found that fine-tuning BERT<sub>large</sub> model on BioASQ for more epochs than the typically recommended 1-3 epochs resulted in much better results. Our best results across the three evaluation metrics on BioASQ came from fine-tuning BERT<sub>large</sub> for 20 epochs and BERT<sub>base</sub> for 4 epochs with initial learning rate of  $5x10^{-6}$ .

**PubMedQA:** All three versions of the BERT models were fine-tuned on PQA-L i.e., 1k expert-annotated generated QA instances of the PubMedQA dataset with the results shown in Table 2. Since PubMedQA hasn't been experimented

Datasets	Metrics	BioBERT [4]		Our Experiments		
		BERT <sub>base</sub>	+PubMed	BERT <sub>large</sub> *	BERT <sub>base</sub>	BERT <sub>base</sub> +PM
BioASQ (4b)	S	27.33	27.95	<b>31.8</b>	31.2	31.54
	L	44.72	44.10	<b>51.8</b>	44.58	48.36
	M	33.77	34.72	<b>40.0</b>	36.25	38.39
BioASQ (5b)	S	39.33	46.00	43.0	41.61	<b>46.05</b>
	L	52.67	<b>60.00</b>	55.8	56.8	57.94
	M	44.27	<b>51.64</b>	48.2	47.86	50.54
BioASQ (6b)	S	33.54	<b>42.86</b>	35.8	36.55	37.57
	L	51.55	57.77	54.4	<b>59.3</b>	58.87
	M	40.88	<b>48.43</b>	43.4	47.59	46.5
Average over 4b,5b, and 6b	S	33.4	<b>38.93</b>	36.86	37.48	38.39
	L	49.65	52.68	<b>55.12</b>	53.56	55.06
	M	39.6	<b>44.93</b>	43.86	43.9	45.14
PubMedQA	K-Fold Acc	—	<b>57.28</b>	56.52	55.20	56.20
	K-Fold F1	—	<b>28.70</b>	26.14	23.71	23.98

Table 2: The performance results of three versions of BERT model (BERT<sub>large</sub>, BERT<sub>base</sub> and BERT<sub>base</sub>+PM) on two distinct QA datasets. For comparison, the results of BioBERT study using BERT<sub>BASE</sub> (i.e. Wiki+Books) and the BioBERT version of the PubMed trained model (+PubMed) are also shown. Since, we already explored latest versions of BERT<sub>large</sub> and BERT<sub>base</sub>, we didn’t examine the BERT<sub>base</sub> version used in the BioBERT study. Best scores are shown in bold. \*Used BERT<sub>large</sub> as other versions were not providing good results.

in the BioBERT study, in order to make a fair comparison between our study and the BioBERT study [4], we ran an additional experiment where we fine-tuned PubMed trained model (+PubMed) from the BioBERT study on the PubMedQA. 10-fold cross-validation was performed with only 450 training instances in each fold of validation. As seen in Table 2, the BioBERT version of the PubMed trained model (+PubMed) has the highest accuracy 57.28 and F1 Score 28.27. This suggests that pre-training on biomedical corpus improves the performance of the downstream PubMedQA task. However, the performance improvement compared to BERT<sub>large</sub>, BERT<sub>base</sub>, and BERT<sub>base</sub>+PM is minimal.

### 3.2.4 Text Summarization

Text summarization refers to automatic generation of summary of a given text. Extractive summarization is done by extracting the most important sentences from the document that summarize the whole document. Abstractive summarization refers to condensing the document into shorter versions while preserving its meaning [30].

# docs (train/val/test)	CNN/DailyMail			XSum			BioASQ 7b (Our experiments)		
	ROUGE-1	ROUGE-2	ROUGE-3	ROUGE-1	ROUGE-2	ROUGE-3	ROUGE-1	ROUGE-2	ROUGE-3
	196,961/12,148/10,397			204,045/11,332/11,334			22,462/4,814/4,800		
<b>Abs</b>	<b>41.72</b>	19.39	<b>38.76</b>	38.81	16.33	31.15	38.20	<b>27.12</b>	34.15
<b>ExtAbs</b>	<b>42.13</b>	19.6	<b>39.18</b>	38.76	16.50	31.27	38.89	<b>29.29</b>	36.01
<b>Ext</b>	<b>43.25</b>	20.24	<b>39.63</b>	†	†	†	33.30	<b>26.50</b>	32.20

Table 3: Text summarization F1 score. † Results were poor hence were not reported by the authors [31]

BERTSum [32, 31] is a simple variant of BERT, for text summarization that produces abstractive and extractive summaries. BERTSumAbs is based on an encoder-decoder architecture where the encoder is pre-trained whereas decoder is a randomly initialized transformer which is trained from scratch. There are three versions of BERTSum. BERTSumAbs produces abstractive summary, BERTSumExt outputs extractive summary, and BERTSumExtAbs also produces abstractive summary but is based on two-stage approach where the encoder is fine-tuned twice, first with an extractive objective followed by an abstractive one.

We explored the performance of these methods on BioASQ 7b<sup>8</sup> dataset consisting of links to research papers as well as summaries are embedded in the json file itself. BERT was trained from scratch on BioASQ 7b. The results are shared in Table 3. For CNN/DailyMail dataset all three methods outperform on BioASQ dataset according to ROUGE-2 score,

<sup>8</sup>[http://bioasq.org/participate/challenges\\_year\\_7](http://bioasq.org/participate/challenges_year_7)

while produce comparable results in all other cases. For XSum dataset [33], we have comparable results according to ROUGE-1, while BioASQ 7b performs better according to ROUGE-2 and ROUGE-3 (for Ext summary results were not reported by the authors [31]). This maybe due to the fact that BioASQ dataset is much smaller in size as compared to the other two datasets. This also shows that smaller-sized dataset can be used for text-summarization.

## 4 Discussion and Conclusion

In this study, we present our findings for experiments conducted in conjunction with further pre-training of BERT model in the biomedical domain as well as evaluating both the domain specific and domain agnostic pre-trained models across downstream biomedical NLP tasks. Our experiments also included considering data and optimization related factors. Further pre-training was conducted on PubMed abstracts only to evaluate the performance of the models trained on smaller-size datasets and also the effects of learning rate was carefully evaluated for fine-tuning tasks.

(1) *Does domain-specific training improves performance?* This study confirms that unsupervised pre-training in general could improve the performance on fine-tuning tasks. However, the effectiveness of *domain specific* pre-training as a way of *further* improving the performance of supervised downstream tasks does not significantly outperform the effectiveness of domain agnostic pre-trained models considering the high cost of domain-specific pre-training which makes it challenging for most of researchers and NLP developers. In the biomedical domain, however, this conclusion may not be wholly substantiated owing to a lack of consistent evidence particularly in downstream NER, QA, and Text Summarization tasks. For the SQuAD task, it has been shown that fine-tuning results depend on the size and duration of training [34] [35]. However, given the small size of the biomedical QA (BioASQ) datasets, it is not possible to run the same experiments, as the current pre-trained models easily overfit. Therefore, it should be beneficial for the biomedical community to curate and expand the magnitude of benchmark datasets. Understandably it is difficult and expensive, but with the ever increasing size of the deep-learning models and significant advances in the development of pre-training language representations, it is necessary in order to facilitate reproducible research for health.

(2) *Is it possible to obtain comparable results using BERT model pre-trained on smaller-sized data?* We present results of the experiments that were conducted using the models pre-trained on the PubMed abstracts only. These results were comparable to the results produced by the model trained on the PubMed, PMC, and Pubmed+PMC together (please check [4]). Although it requires further investigation but empirically it shows that small-sized datasets may be used as a surrogate. Using small-sized dataset can be especially useful when the models need to retrained instead of using pre-trained publicly available models due to (a) data-shift, (b) wide variety of data-domains, (c) confidential data not publicly available to train the model on, (d) small-size of the domain specific data available, and (e) a lack of computing resources.

We position these experiments as complementary to existing literature on the applications of transformer-based models for biomedical NLP. Our results provide some evidence for the validity and the limitations of existing language representations for pre-training & fine-tuning in biomedical domain.

We would like to expand our experiments on other variants of BERT, with more domain specific datasets and a variety of downstream tasks. This would also allow us to use our experiments for other medical corpora<sup>9</sup> for external validation and generalization of the results.

## Acknowledgements

We want to thank Vector Institute industry sponsors, researchers and technical staff who participated in the Vector Institute’s NLP Project (<https://vectorinstitute.ai/wp-content/uploads/2020/12/nlp-report-final.pdf>).

## References

- [1] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019.
- [2] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

---

<sup>9</sup><http://www.nactem.ac.uk/resources.php>

- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [4] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.
- [5] Faiza Khan Khattak, Serena Jeblee, Chloé Pou-Prom, Mohamed Abdalla, Christopher Meaney, and Frank Rudzicz. A survey of word embeddings for clinical text. *Journal of Biomedical Informatics: X*, page 100057, 2019.
- [6] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [7] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [8] Matthew E Peters, Sebastian Ruder, and Noah A Smith. To tune or not to tune? adapting pretrained representations to diverse tasks. *arXiv preprint arXiv:1903.05987*, 2019.
- [9] Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. End-to-end open-domain question answering with bertserini. *arXiv preprint arXiv:1902.01718*, 2019.
- [10] Dogu Araci. Finbert: Financial sentiment analysis with pre-trained language models. *arXiv preprint arXiv:1908.10063*, 2019.
- [11] Xiaofei Ma, Peng Xu, Zhiguo Wang, Ramesh Nallapati, and Bing Xiang. Domain adaptation with bert-based domain classification and data selection. In *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, pages 76–83, 2019.
- [12] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *arXiv preprint arXiv:2007.15779v3*, 2020.
- [13] Yanshan Wang, Sijia Liu, Naveed Afzal, Majid Rastegar-Mojarad, Liwei Wang, Feichen Shen, Paul Kingsbury, and Hongfang Liu. A comparison of word embeddings for the biomedical natural language processing. *Journal of biomedical informatics*, 87:12–20, 2018.
- [14] Mengnan Zhao, Aaron J Masino, and Christopher C Yang. A framework for developing and evaluating word embeddings of drug-named entity. In *Proceedings of the BioNLP 2018 workshop*, pages 156–160, 2018.
- [15] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [16] Yongjun Zhu, Erjia Yan, and Fei Wang. Semantic relatedness and similarity of biomedical terms: examining the effects of recency, size, and section of biomedical publications on the performance of word2vec. *BMC medical informatics and decision making*, 17(1):95, 2017.
- [17] Hugging Face. A library of state-of-the-art pretrained models for natural language processing (nlp). <https://github.com/huggingface/pytorch-transformers>, 2019.
- [18] Matthew McDermott, Shirly Wang, Nikki Marinsek, Rajesh Ranganath, Marzyeh Ghassemi, and Luca Foschini. Reproducibility in machine learning for health. *arXiv preprint arXiv:1907.01463*, 2019.
- [19] US National Library of Medicine National Institute of Health. Pubmed, pubmed central, (accessed June, 2019). <https://www.ncbi.nlm.nih.gov/pubmed/>.
- [20] Amy Neustein, S Sagar Imambi, Mário Rodrigues, António Teixeira, and Liliana Ferreira. Application of text mining to biomedical knowledge extraction: analyzing clinical narratives and medical literature. *Text mining of web-based medical content. De Gruyter, Berlin*, 50, 2014.
- [21] Rezarta Islamaj Doğan, Robert Leaman, and Zhiyong Lu. Ncbi disease corpus: a resource for disease name recognition and concept normalization. *Journal of biomedical informatics*, 47:1–10, 2014.
- [22] Jin-Dong Kim, Tomoko Ohta, Yoshimasa Tsuruoka, Yuka Tateisi, and Nigel Collier. Introduction to the bio-entity recognition task at jnlpba. In *Proceedings of the international joint workshop on natural language processing in biomedicine and its applications*, pages 70–75. Citeseer, 2004.
- [23] Evangelos Pafilis, Sune P Frankild, Lucia Fanini, Sarah Faulwetter, Christina Pavloudi, Aikaterini Vasileiadou, Christos Arvanitidis, and Lars Juhl Jensen. The species and organisms resources for fast and accurate identification of taxonomic names in text. *PLoS One*, 8(6):e65390, 2013.

- [24] Kevin G Becker, Kathleen C Barnes, Tiffani J Bright, and S Alex Wang. The genetic association database. *Nature genetics*, 36(5):431, 2004.
- [25] George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R Alvers, Dirk Weissenborn, Anastasia Krithara, Sergios Petridis, Dimitris Polychronopoulos, et al. An overview of the bioasq large-scale biomedical semantic indexing and question answering competition. *BMC bioinformatics*, 16(1):138, 2015.
- [26] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W. Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146v1*, 2019.
- [27] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- [28] Wonjin Yoon, Jinhyuk Lee, Donghyeon Kim, Minbyul Jeong, and Jaewoo Kang. Pre-trained language model for biomedical question answering. *arXiv preprint arXiv:1909.08229*, 2019.
- [29] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [30] Vishal Gupta and Gurpreet Singh Lehal. A survey of text summarization extractive techniques. *Journal of emerging technologies in web intelligence*, 2(3):258–268, 2010.
- [31] Yang Liu and Mirella Lapata. Text summarization with pretrained encoders. *arXiv preprint arXiv:1908.08345*, 2019.
- [32] Yang Liu. Fine-tune bert for extractive summarization. *arXiv preprint arXiv:1903.10318*, 2019.
- [33] Shashi Narayan, Shay B Cohen, and Mirella Lapata. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *arXiv preprint arXiv:1808.08745*, 2018.
- [34] Alon Talmor and Jonathan Berant. Multitqa: An empirical investigation of generalization and transfer in reading comprehension. *CoRR*, abs/1905.13453, 2019.
- [35] Georg Wiese, Dirk Weissenborn, and Mariana L. Neves. Neural domain adaptation for biomedical question answering. *CoRR*, abs/1706.03610, 2017.