

BioALBERT: A Simple and Effective Pre-trained Language Model for Biomedical Named Entity Recognition

Usman Naseem*, Matloob Khushi† Vinay Reddy‡, Sakthivel Rajendran§, Imran Razzak¶ and Jinman Kim||

*†‡¶||School of Computer Science, University of Sydney, Australia

§School of Information Technology, Deakin University, Australia

Email: *†‡§||FirstName.SecondName@sydney.edu.au, ¶FirstName.SecondName@deakin.edu.au

Abstract—In recent years, with the growing amount of biomedical documents, coupled with advancement in natural language processing algorithms, the research on biomedical named entity recognition (BioNER) has increased exponentially. However, BioNER research is challenging as NER in the biomedical domain are: (i) often restricted due to limited amount of training data, (ii) an entity can refer to multiple types and concepts depending on its context and, (iii) heavy reliance on acronyms that are sub-domain specific. Existing BioNER approaches often neglect these issues and directly adopt the state-of-the-art (SOTA) models trained in general corpora which often yields unsatisfactory results. We propose biomedical ALBERT (A Lite Bidirectional Encoder Representations from Transformers for Biomedical Text Mining) bioALBERT an effective domain-specific language model trained on large-scale biomedical corpora designed to capture biomedical context-dependent NER. We adopted a self-supervised loss used in ALBERT that focuses on modelling inter-sentence coherence to better learn context-dependent representations and incorporated parameter reduction techniques to lower memory consumption and increase the training speed in BioNER. In our experiments, BioALBERT outperformed comparative SOTA BioNER models on eight biomedical NER benchmark datasets with four different entity types. We trained four different variants of BioALBERT models which are available for the research community to be used in future research.

I. INTRODUCTION

The growing volume of the published biomedical literature, such as clinical reports [1] and health literacy [2], are fuelling the advancements in the development of text mining algorithms. Biomedical named entity recognition (BioNER) intends to automatically identify biomedical entities such as diseases, chemicals, genes and proteins, etc., from the biomedical literature. So, a crucial step towards this aim is to build better and effective methods which can automatically recognize and extract biomedical entities. BioNER is an essential building block of many downstream text mining applications such as extracting drug-to-drug interactions [3] and disease-treatment [4] relationships. Traditionally, BioNER relies on feature engineering methods (e.g., lexicon-based, rules-based and statistics-based). However, feature engineering is dependent on domain-specific knowledge which does not perform well on BioNER [5].

Deep learning (DL) with its ability to automatically extract features have become common in BioNER recently [6]. For instance, Long Short-Term Memory (LSTM) is usually em-

ployed to learn vector representations of each word in a sentence, and then as the input to conditional random fields (CRF) have greatly improved the performance in BioNER [6]. Recently state-of-the-art (SOTA) DL based language models such as Embeddings from Language Models (ELMo) [7], Bidirectional Encoder Representations from Transformers (BERT) [8] and (A Lite Bidirectional Encoder Representations from Transformers (ALBERT) [9] obtained SOTA best performance on many NLP tasks.

Although these models show promising results, but applying them on BioNER has multiple challenges and limitations: (i) a limited amount of training data; (ii) an entity could represent multiple entity types depending on its textual context, i.e., a *BRCA1* can be referred as gene name as well as a disease entity depending on its context. Similarly, *heart attack* and *myocardial infarction* refer to the same concept and, (iii) the heavy use of acronyms in biomedical texts makes it challenging to identify concepts, i.e., the abbreviation *RA* may refer to *right atrium*, *rheumatoid arthritis*, or one of several other concepts, where the resolution of the abbreviation is, therefore, context-dependent. Therefore, current models in BioNER rely on various context-independent and transformer-based context-dependent language models which are trained on biomedical corpora [10]–[12].

To overcome the identified limitations, we present Biomedical ALBERT bioALBERT - a context-dependent, fast and effective language model that addresses the shortcomings of recently proposed domain-specific language models. BioALBERT is trained on large biomedical corpora which address the limitation of limited training data. We also innovate in the adoption of cross-layer parameter sharing by learning parameters for the first block and reuse the block in the remaining layers instead of learning unique parameters for each of the layers and sentence-order-prediction (SOP) technique as a measure of coherence loss between sentences. SOP takes two consecutive sentences from training data and creates a random pair from different sentences which helps the model to learn better representations and finally, in BERT based models the size of the embedding was linked to the hidden layer sizes of the transformer blocks. These embeddings are projected directly to the hidden space of the hidden layer whereas in our model we use factorized embedding parameterization

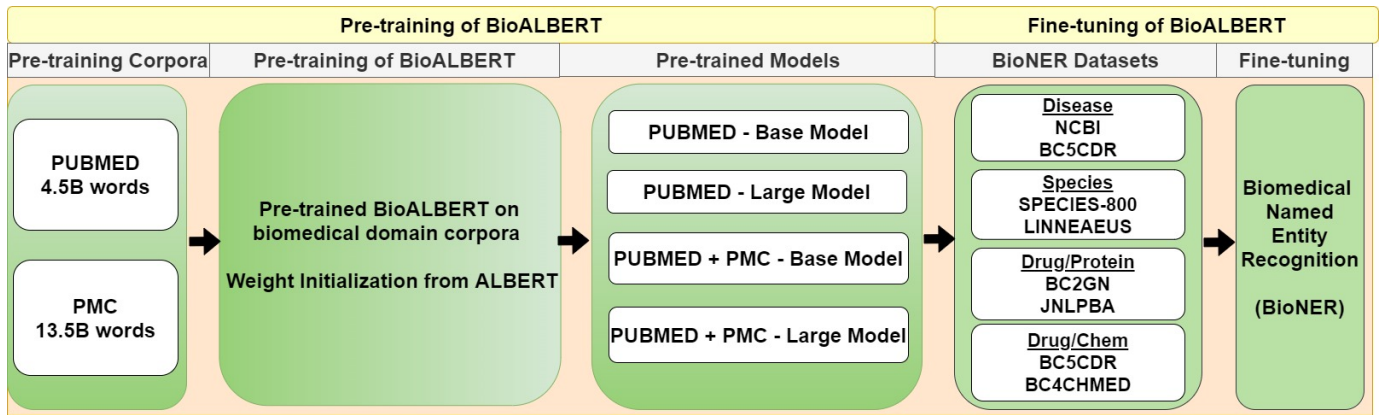


Fig. 1: Overview of the pre-training and fine-tuning of BioALBERT on NER

which decomposes embedding matrix into two small matrices which separate the size of the hidden layers from the size of vocabulary embeddings. This allows for increasing the hidden size without significantly increasing the parameter size of the vocabulary embeddings. BioALBERT is simple and efficient when fine-tuned for BioNER task as compared to other SOTA models. We evaluate our model on eight biomedical named entity recognition benchmark datasets. Our pre-trained BioNER models, along with the source code, will be publicly available.

II. RELATED WORK

A. Language Model

In biomedical text mining research, there is a long history of using shared language representations to capture the semantics of the text. One established trend is a form of word-embeddings [13] that represent syntactic and semantic meaning, and map words into low-dimensional vectors. Similar methods also have been derived for improving embeddings of word sequences by introducing sentence embeddings [14]. These context-independent word embeddings approach such as word2vec [13] were trained on biomedical corpora that contain terms and expressions that are usually not included in a general domain corpus [10]. These methods always require complex neural networks to be effectively used and model context-independent representations.

Another common trend, particularly in recent years, is text representation based on context [7]–[9]. Unlike traditional word-embeddings, this enables a word to change its meaning depending on the context in which it occurs. Several other works have investigated the usefulness of contextual models in clinical and biomedical domains. Several researchers trained ELMo on biomedical corpora and presented BioELMo and found that BioELMo beats ELMo on BioNER tasks [11], [16]. A pre-trained BioELMo model was released along with their work, allowing more clinical research in the area of BioNER. Beltagy et al. [15] released Scientific BERT (SciBERT), where BERT was trained on the scientific texts. BERT has typically been superior and better than ELMo to non-contextual embeddings.

Si et al. [17], trained the BERT on clinical notes corpora, using complex task-specific models to improve both traditional embedding and ELMo embedding on the i2b2 2010 and 2012 BioNER. Similarly in another study, a new domain-specific language model, BioBERT [12], trained a BERT model on a corpus of biomedical articles from PMC abstracts¹ as well as full texts sourced from PubMed² which gave rise to enhanced performance on BioNER. Peng et al [18] introduced Biomedical Language Understanding Evaluation (BLUE), a collection of resources for evaluating and analysing biomedical natural language representation models. Their study also confirmed that the BERT models pre-trained on PubMed abstracts and clinical notes see better performance than most state-of-the-art models.

Despite this success, BERT has some limitations such as BERT has a huge number of parameters which is the cause for problems like degraded pre-training time, memory management issues and model degradation etc [9]. These issues are very well addressed in ALBERT, which is modified based on the architecture of BERT and proposed by Lan et al. [9]. ALBERT incorporates two-parameter reduction techniques that lift the significant obstacles in scaling pre-trained models. (i) Factorized embedding parameterization - decomposes the large vocabulary embedding matrix into two small matrices, (ii) replaces the NSP loss by SOP loss; and (iii) Cross-layer parameter sharing- prevents the parameter from growing with the depth of the network. These techniques significantly reduce the number of parameters used when compared with BERT without significantly affecting the performance of the model, thus improving parameter-efficiency. An ALBERT configuration similar to BERT-large has 18x fewer parameters and can be trained about 1.7x faster.

B. Fine tuning

Fine-tuning has been successfully used to transfer pre-trained weights as initialisation for parameters for various downstream tasks [19]. This improves the efficiency of the

¹<https://www.ncbi.nlm.nih.gov/pmc/>

²<https://www.ncbi.nlm.nih.gov/pubmed/>

TABLE I: Summary of parameters used in the Pre-training

Summary of All parameters used: (Pre-Training)		
Name	BioALBERT 1.0	BioALBERT 1.1
Architecture	ALBERT Base	ALBERT Large
Activation function	GeLU	GeLU
Number of Attention Heads	12	16
Number of layers	12	24
Hidden Size	768	1024
Embedding Size	128	128
Vocab Size	30000	30000
Optimizer	LAMB	LAMB
Train Batch Size	1024	256
Eval Batch Size	16	16
Max Seq Length	512	512
Max Predictions per Seq	20	20
Learning Rate	0.00176	0.00062
Training Steps (PubMed)	200K	200K
Training Steps (PMC)	270K	270K
Warmup Steps	3125	3125

target task when we have limited data and similar tasks [20]. Recently, fine-tuning of pre-trained language models have been widely used in various text mining tasks [21]. Liu et al. [22] pre-trained LSTM with a language model and fine-tuned it, and this has contributed to improved performance for various text classification tasks.

Universal Language Model Fine-tuning (ULMFiT) used general-domain pre-training and fine-tuning techniques to avoid over-fitting and achieved SOTA performance on the datasets with less samples [19]. Similarly, authors of BERT, ALBERT, BioBERT and others tested the performance of their model for a wide range of tasks with minimum fine-tuning efforts and achieved good performance. Recently, multi-task fine-tuning has led to improvements even with many target tasks [23]. In this paper, we fine-tuned our pre-trained language model BioALBERT, which is trained on biomedical corpora for BioNER. BioALBERT offers better performance on BioNER as it addresses the shortcomings of BERT used in BioBERT. However, this is not a trivial task as it requires several optimisations which are discussed in the next section. Weights of our pre-trained model along with complete source code, will be available online.

III. METHODOLOGY

First, we initialized BioALBERT with weights from ALBERT. BioALBERT is pre-trained on biomedical domain corpora (PubMed abstracts and PMC full-text articles). To show the effectiveness of our approach, BioALBERT is fine-tuned on BioNER. We tested various pre-training strategies with different combinations and sizes of biomedical corpora. The overview of our methodology is shown in Fig. 1. Below we present details of each step involved in pre-training and fine-tuning of BioALBERT.

A. Pre-training of BioALBERT

In this section, we present steps involved in the pre-training of BioALBERT, which has the same architecture as ALBERT, which makes it simple. BioALBERT is trained on PubMed abstracts and PMC full-text articles which contains biomedical

TABLE II: BioALBERT Pre-trained Models

Model Version	Model	Trained On	# of words	Machine Used	Batch Size	Steps
BioALBERT 1.0	Base	PubMed	4.3B	GCP TPU v3-8	1024	200K
BioALBERT 1.0	Base	PubMed+PMC	18.8B	GCP TPU v3-8	1024	470K
BioALBERT 1.1	Large	PubMed	4.3B	GCP TPU v3-8	256	200K
BioALBERT 1.1	Large	PubMed+PMC	18.8B	GCP TPU v3-8	256	470K

terms and enables to pre-train the ALBERT model, which was pre-trained on the general text on biomedical corpora. As we cannot use the raw text biomedical corpora (PubMed and PMC) to pre-train the model, so we converted raw text to structured format by converting raw text files into a single sentence in which; (i) all blank lines are removed within a document and converted it into on single paragraph; (ii) Any line less than 20-character length is removed and (iii) there will be a blank line between each document (when combining the multiple files) for training the model. Overall, PubMed contains approximately 4.5 Billion words, and PMC contains about 13.5 Billion words. As an initial step, we have initialized the model weights from ALBERT to create BioALBERT model by pre-training on biomedical corpora and kept the original vocabulary of ALBERT for pre-training.

Sentence embeddings are used for tokenization, and for that, we pre-processed the data as a sentence text. So every line in the input text document is a sentence, and an empty line separates every document. We set the maximum length of each sentence to 512 words. Shorter sentences were padded to make 512 whereas longer sentences were truncated. Learning rate of 0.00176 and number of warm-up steps as 3125 for all the pre-trained versions of the model, except for BioALBERT 1.1 (PubMed +PMC) where the learning rate has been re-adjusted to 0.00062 (by re-scaling the current learning rate with scaled $1/2^{1.5}$ times to avoid the problem of exploding loss. Summary of parameters used in the training process is given in Table I.

We experimented with different settings and found out that both base and large models were successful with larger batch size on V3-8 TPU compute instances. We have used two different embedding sizes, i.e., 128 and 256. The 128 embedding size creates a base model with 12 Million parameters, whereas a large model has 16 Million parameters with 256 embedding size. With these combinations, we have presented a total of four models, given in Table II.

B. Fine-tuning of BioALBERT

Fine-tuning of BioALBERT on BioNER task is presented in this section. BioNER involves annotating words in a sentence as named-entities. The labelled datasets that were used for this task include four categories representing Disease, Species, Drug/Proteins, and Drugs/Chemicals. The objective is to train and make a prediction on the labels, which are the proper nouns within the domain area. More formally, given an input sentence $S = \{x_1, x_2, \dots, x_z\}$, where x_i is the i -th word and z represents the length of the sentence. The goal of BioNER is to classify each word in S and assign it to a corresponding label $y \in Y$, where Y is a predefined list of all possible label.

Fine-tuning is simple as compared to the pre-training, and the computational requirements are also not that significant.

TABLE III: Summary of parameters used in fine-tuning

Summary of All parameters used: (Fine Tuning)	
Name	BioNER
Optimizer	adamw
Train Batch Size	32
Eval Batch Size	16
Save Checkpoint	200
Max Seq Length	512
Learning Rate	1.00E-05
Training Steps	5336
Warmup Steps	320

BioALBERT which takes less physical memory and improvised parameter sharing techniques. The BioNER fine-tuning is trained to learn the word embeddings using the sentence piece tokenization while the BioBERT model was based on the word piece embeddings. For each of the pre-trained models, we constructed a fine-tuning task by using the specific dataset.

The model setup uses the weights of the pre-trained model that are created previously. We have used $1e-5$ learning rate, batch size as 32 and lower case texts and finally fine-tuned for 5336 steps. Pre-trained BioALBERT models were pre-trained using TPUv3-8. All of the hyper-parameters used are same as default ALBERT unless stated otherwise. All the tested datasets contain a list of words along with a label B, I, and O where B denotes the beginning of an entity, I stands for inside and is used for all words comprising the entity except the first one, and O means the absence of an entity. For our experiments, we used these datasets as-is and passed to it our pre-trained models for the downstream task. The adamw optimiser was used with evaluation checkpoint so that the model will be evaluated at different time intervals using the holdout development dataset to identify the best model for final predictions. The predictions were performed on the test datasets, and the performance is compared with baseline models that were established previously by calculating the F1 Score, Precision and Recall. The summary of all parameters used in fine-tuning is given in Table III.

IV. EXPERIMENTAL ANALYSIS

In this section, we present the dataset used, baselines and results to show the effectiveness of our model.

A. Datasets

Our model is evaluated on eight biomedical NER benchmark datasets which contain four types of entities and provided by Lee et al. [12]. Table IV shows the statistics of the datasets used. Below we briefly explained each dataset:

- **BC5CDR**: The Bio-creative community challenge for the chemical-disease relation extraction task (BC5CDR) the corpus was made available in a Bio-creative workshop [24]. The two sub-tasks of BC5CDR are identifying; (i) **chemical** and (ii) **disease** entities from Medline abstracts.
- **BC4CHEMD**: This dataset is provided by Bio-Creative community challenge IV for the development and evaluation of tools for Chemical NER [25]. BC4CHEMDNER

TABLE IV: Statistics of the BioNER datasets

Entity Type	Dataset	# of Annotations
Disease	NCBI Disease	6,881
	BC5CDR	12,694
Drug/Chem	BC5CDR	15,411
	BC4CHEMD	79,842
Drug/Protein	BC2GN	20,703
	JNLPBA	35,460
SPECIES	LINNAEUS	4,077
	Species-800	3,708

was used for the recognition of chemical compounds and drugs from Pubmed abstracts.

- **NCBI Disease**: To promote disease NER-system research, American National Institutes of Health released the NCBI disease corpus for disease NER-research. The NCBI disease corpus is large-scale and high-quality; it is based on the corpus released by Leaman et al. [26].
- **JNLPBA**: We also used the JNLPBA corpus in our experiments which are provided by Kim et al. [27]. This corpus contains five entity types, including DNA, RNA, Cell Type, Cell Line and Protein.
- **BC2GM**: BC2GM is provided by Ando [28], the state-of-the-art system in the Bio-Creative II gene mention recognition task is a semi-supervised learning method using alternating structure optimization.
- **LINNAEUS**: The LINNAEUS corpus was provided by Gerner et al. [29] which consists of 100 full-text documents from the PMC Open access document set which were randomly selected. All mentions of species terms were manually annotated and normalized to the NCBI taxonomy IDs of the intended species.
- **Species-800**: Species-800, which is also known as S800 [30], is a novel abstract-based manually annotated corpus. S800 comprises 800 PubMed abstracts in which organism mentions were identified and mapped to the corresponding NCBI taxonomy identifiers.

B. Baselines

To assess the performance of the proposed method, an exhaustive comparison is performed with several advanced SOTA methodologies along with their published results³. Our model is compared with the following methods.

Yoon et al. [35] introduced CollaboNet, which consists of multiple BiLSTM-CRF models, for BioNER. While existing models were only able to handle datasets with a single entity type, CollaboNet leverages multiple datasets and achieves the highest F1 scores. CollaboNet is built upon multiple single-task NER models (STMs) that send information to each other for more accurate predictions.

Lou et al. [33] proposed a transition-based model for joint disease entity-recognition and normalization, based on transition-based structured prediction framework using structured perceptron with early-update training and beam-search

³The published results were acquired from respective original publication

TABLE V: Comparison of performance in biomedical named entity recognition (BioNER) task.

Type	Datasets	Metrics	SOTA	BioBERT v1.0	BioBERT v1.0	BioBERT v1.0	BioBERTv1.1	BioALBERT 1.0	BioALBERT 1.0	BioALBERT 1.1	BioALBERT 1.1
				(PubMed)	(PMC)	(PubMed+PMC)	(PubMed)	(PubMed)	(PubMed+PMC)	(PubMed)	(PubMed+PMC)
				Base	Base	Base	Base	Base	Base	Large	Large
Disease	NCBI Disease	p	88.30	86.76	86.16	<u>89.04</u>	88.22	97.45	96.84	97.18	97.38
		R	89.00	88.02	89.48	89.69	<u>91.25</u>	94.39	94.40	97.18	94.37
		F	88.60	87.38	87.79	89.36	<u>89.71</u>	95.89	95.61	97.18	95.85
	BC5CDR	p	89.61	85.80	84.67	85.86	86.47	99.69	99.11	99.27	99.39
		R	83.09	86.60	85.87	87.27	87.84	95.72	96.17	96.33	95.85
		F	86.23	86.2	85.27	86.56	87.15	97.66	97.62	97.78	97.61
Drug/Chem	BC5CDR	p	94.26	92.52	92.46	93.27	93.68	99.99	99.99	99.99	99.99
		R	92.38	92.76	92.63	<u>93.61</u>	93.26	95.89	96.24	95.62	95.68
		F	93.31	92.64	92.54	93.44	<u>93.47</u>	97.9	98.08	97.76	97.79
	BC4CHEMD	p	92.29	91.77	91.65	92.23	<u>92.80</u>	97.76	97.71	97.71	97.88
		R	90.01	90.77	90.30	90.61	<u>91.92</u>	94.22	94.83	94.83	94.63
		F	91.14	91.26	90.97	91.41	<u>92.36</u>	95.96	96.25	96.25	96.23
Drug/Protein	BC2GM	p	81.81	81.72	82.86	<u>85.16</u>	84.32	97.86	97.84	98.26	98.02
		R	81.57	83.38	84.21	<u>83.65</u>	<u>85.12</u>	94.87	94.27	95.72	94.70
		F	81.69	82.54	83.53	84.40	<u>84.72</u>	96.34	96.02	96.97	96.33
	JNLPBA	p	<u>74.43</u>	71.11	71.17	72.68	72.24	85.14	85.60	86.23	85.56
		R	83.22	83.11	82.76	83.21	<u>83.56</u>	80.43	80.98	81.90	81.49
		F	78.58	76.65	76.53	<u>77.59</u>	77.49	82.72	83.22	84.01	83.53
Species	LINNAEUS	p	92.80	91.83	91.62	<u>93.84</u>	90.77	99.95	99.98	99.98	99.92
		R	94.29	84.72	85.48	86.11	85.83	99.47	99.44	99.48	99.55
		F	<u>93.54</u>	88.13	88.45	89.81	88.24	99.71	99.72	99.73	99.73
	Species-800	p	<u>74.34</u>	70.60	71.54	72.84	72.80	99.17	98.93	99.10	98.75
		R	75.96	75.75	74.71	<u>77.97</u>	75.36	98.34	98.04	98.95	98.69
		F	74.98	73.08	73.09	<u>75.31</u>	74.06	98.76	98.49	99.02	98.72

Notes: Precision (P), Recall (R) and F1 (F) scores on each dataset are reported. The best scores are in **bold**, and the second best scores are underlined. We list the scores of the state-of-the-art (SOTA) models on different datasets as follows: scores of Xu et al. [31] on NCBI Disease, scores of Sachan et al. [32] on BC2GM, scores of Lou et al. [33] on BC5CDR-disease, scores of Luo et al. [34] on BC4CHEMD, scores of Yoon et al. [35] on BC5CDR-chemical and JNLPBA and scores of Giorgi and Bader [36] on LINNAEUS and Species-800

decoding. In another study, Lou et al. [34] proposed a neural network approach, i.e. attention-based bidirectional Long Short-Term Memory with a conditional random field layer (Att-BiLSTM-CRF), to document level chemical NER. The method leverages document-level global information obtained by attention mechanism to enforce tagging consistency across multiple instances of the same token in a document. It achieves better performances with little feature engineering.

Xu et al. [31] proposed a novel dictionary-based and document-level attention mechanism with a deep neural network NER method, named as DABLC. DABLC tags the consistency of multiple instances in a document at the document level and combines an external disease dictionary that is constructed with five disease resources containing a rich collection of disease entities. The authors adopted the efficient exact string matching method for dictionary matching; this method can effectively and accurately match the disease names.

Lee et al. [12] introduced BioBERT, which is a pre-trained language model for biomedical text mining. Authors showed that pre-training BERT on biomedical corpora is crucial in applying it to the biomedical domain. BioBERT outperforms previous models on biomedical text mining tasks such as NER, RE and QA. We compare BioALBERT with both BioBERT (v1.0 and V1.1) models and other SOTA models used in BioNER task. We selected those methods because they are the SOTA, and based on the conducted meta-analysis exhibit the highest performance among the techniques so far developed.

C. Results

Table V presents the performance of all the variants of BioALBERT and contrasts them to baseline methodologies.

Our model outperforms all other methods on all eight datasets and entity types. For, (i) Disease-type datasets, BioALBERT improved the performance by 7.47% and 10.63% for NCBI-Disease and BC5CDR-Disease datasets respectively; (ii) Drug/Chem type datasets increase in performance by 4.61% and 3.89% for BC5CDR-Chem and BC4CHEMD datasets respectively; (iii) Gene/Protein type datasets increase in performance by 12.25% and 6.42% for BC2GM and JNLPBA datasets respectively and; (iv) Species type datasets increase in performance by 6.19% and 23.71% respectively is observed.

We have performed multiple comparisons to analyse the effectiveness of BioALBERT. Fig. 2 presents, the performance comparison of the same versions (trained on same corpora and for the same number of steps) of both BioALBERT and BioBERT. We can see that in Fig. (2a), we compared BioBERT v1.0-base model which is trained on PubMed with BioALBERT 1.0-base model trained on Pubmed and similarly in Fig. (2b), we compared BioBERT v1.0-base model trained on Pubmed and PMC with BioALBERT 1.0-base model trained on PubMed and PMC biomedical corpora. In both cases, BioALBERT outperformed BioBERT on all eight datasets. We also compared the training time of BioALBERT with BioBERT. We found that all models BioALBERT beat BioBERT with a considerable margin. The run time statistics of both pre-trained models are given in Fig. 3.

D. Discussions

BioALBERT gives better performance and addresses the previously mentioned challenges in the biomedical domain. We attribute this to the BioALBERT built on top of the

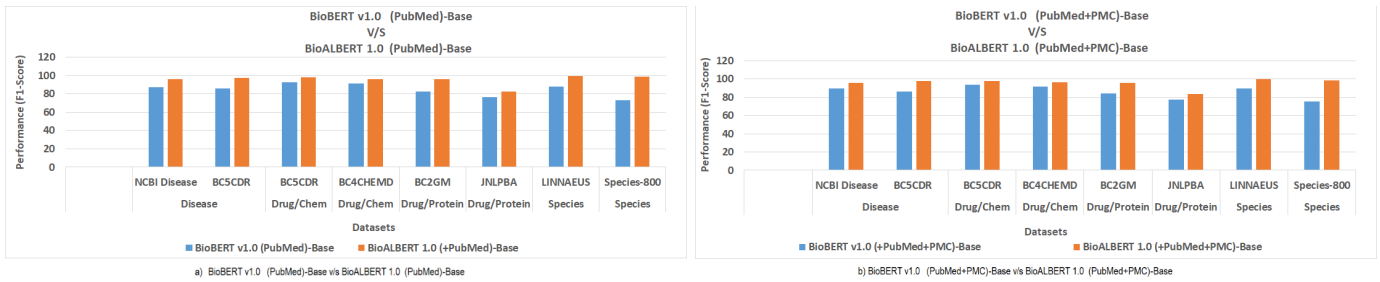


Fig. 2: Comparison of a) BioBERT v1.0 (Pubmed)-Base v/s BioALBERT 1.0 (Pubmed)-Base ; b) BioBERT v1.0(Pubmed+PMC)-Base v/s BioALBERT 1.0(Pubmed+PMC)-Base

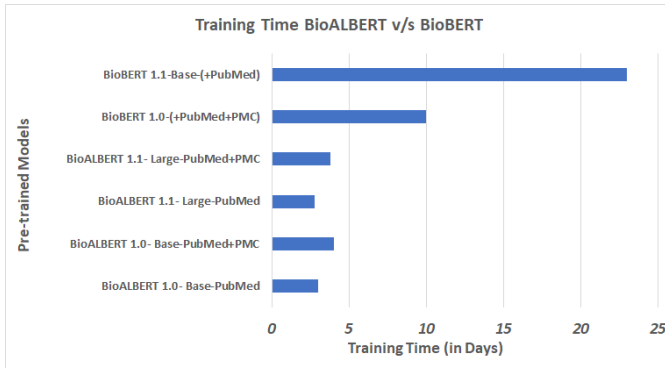


Fig. 3: Comparison of run-time statistics of BioALBERT v/s BioBERT

transformer-based language model that learns contextual relations between words (or subwords) in a text. As opposed to directional models, which read the text input sequentially (left-to-right or right-to-left), the transformer encoder reads the entire sequence of words at once. This characteristic allows the model to learn the context of a word based on all of its surroundings (left and right of the word) and address the issue of contextual representation.

Our model addresses the shortcomings of BERT based biomedical models. At first, BioALBERT uses cross-layer parameter sharing and reduces 110 million parameters of 12-layer BERT-base model to 31 million parameters while keeping the same number of layers and hidden units by learning parameters for the first block and reuse the block in the remaining 11 layers. Secondly, our model uses the SOP, which takes two segments from the training corpus that appear consecutively and constructs a random pair of segments from different documents. This enables the model to learn about discourse-level coherence characteristics from a finer-grained distinction and leads to better learning representation in downstream tasks. Thirdly, our model uses factorized embedding parameterization in which a smaller size layer vocabulary and hidden layer to decompose the embedding matrix into two small matrices which reduce the number of parameters between vocabulary and the first hidden layer whereas in BERT based biomedical models embedding size is equal to the size of the hidden layer. Furthermore, finally, our model is trained on massive biomedical corpora to be effective on BioNER to address the issue of the shift of word distribution

from general domain corpora to biomedical corpora. All these, when combined, address all the issues associated with BioNER earlier. As our model offers a consistent improvement over all other methods for all tested datasets, we can conclude that it is a robust solution for BioNER.

To extend our analysis, we analysed the performance of different pre-trained models of BioALBERT. We found out that the performance of all BioALBERT models are almost equally good, but BioALBERT 1.1-large trained on PubMed works better than others (shown in Fig. 4). BioALBERT 1.1-large model, which is trained on PubMed with dup-factor as five performs better on most of the datasets. This shows that the relevance of duplication data in NLP tasks.

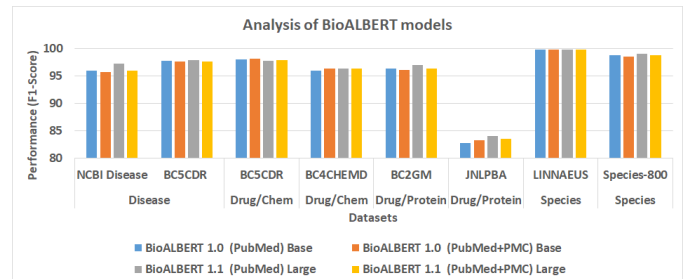


Fig. 4: Comparison of different variants BioALBERT models

V. CONCLUSION

In this study, we presented BioALBERT, which is a pre-trained language model for biomedical named entity recognition. We presented four different variants of BioALBERT models which are trained on huge biomedical corpora for a different number of steps. We showed that training ALBERT on biomedical corpora is a crucial step in applying it to BioNER. As future works, we plan to pre-trained other versions which include hybrid of general and biomedical corpora of ALBERT on biomedical corpora with more training steps and fine-tune on biomedical text mining task. We also plan to fine-tune BioALBERT on other text mining tasks to show the effectiveness of our model.

REFERENCES

- [1] Stéphane M Meystre, Guergana K Savova, Karin C Kipper-Schuler, and John F Hurdle. Extracting information from textual documents in the electronic health record: a review of recent research. *Yearbook of medical informatics*, 17(01):128–144, 2008.

- [2] Lena Mårtensson and Gunnel Hensing. Health literacy—a heterogeneous phenomenon: a literature review. *Scandinavian journal of caring sciences*, 26(1):151–160, 2012.
- [3] Sangrak Lim, Kyubum Lee, and Jaewoo Kang. Drug drug interaction extraction from the literature using a recursive neural network. *PLOS ONE*, 13:e0190926, 01 2018.
- [4] Barbara Rosario and Marti Hearst. Classifying semantic relations in bioscience texts. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*, pages 430–437, Barcelona, Spain, July 2004.
- [5] Aaron M Cohen and William R Hersh. A survey of current work in biomedical text mining. *Briefings in bioinformatics*, 6(1):57–71, 2005.
- [6] Vikas Yadav and Steven Bethard. A survey on recent advances in named entity recognition from deep learning models, 2019.
- [7] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *CoRR*, abs/1802.05365, 2018.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [9] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations, 2019.
- [10] Sampo Pyysalo, Filip Ginter, Hans Moen, Tapio Salakoski, and Sophia Ananiadou. Distributional semantics resources for biomedical text processing. 2013.
- [11] Qiao Jin, Bhuwan Dhingra, William W. Cohen, and Xinghua Lu. Probing biomedical embeddings from language models, 2019.
- [12] 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, 2019.
- [13] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–9, 2013.
- [14] Qingyu Chen, Yifan Peng, and Zhiyong Lu. Biosentvec: creating sentence embeddings for biomedical texts. *2019 IEEE International Conference on Healthcare Informatics (ICHI)*, pages 1–5, 2018.
- [15] Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text, 2019.
- [16] Henghui Zhu, Ioannis C. Paschalidis, and Amir M. Tahmasebi. Clinical concept extraction with contextual word embedding. *NIPS Machine Learning for Health Workshop*, Dec 2018.
- [17] Yuqi Si, Jingqi Wang, Hua Xu, and Kirk Roberts. Enhancing clinical concept extraction with contextual embeddings. *Journal of the American Medical Informatics Association*, 26(11):1297–1304, Jul 2019.
- [18] Yifan Peng, Shankai Yan, and Zhiyong Lu. Transfer learning in biomedical natural language processing: An evaluation of bert and elmo on ten benchmarking datasets, 2019.
- [19] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [20] Dani Yogatama, Cyprien de Masson d’Autume, Jerome Connor, Tomas Kocisky, Mike Chrzanowski, Lingpeng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu, Chris Dyer, et al. Learning and evaluating general linguistic intelligence. *arXiv preprint arXiv:1901.11373*, 2019.
- [21] Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. Pre-trained models for natural language processing: A survey. *arXiv preprint arXiv:2003.08271*, 2020.
- [22] Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. Recurrent neural network for text classification with multi-task learning. *arXiv preprint arXiv:1605.05101*, 2016.
- [23] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Multi-task deep neural networks for natural language understanding. *arXiv preprint arXiv:1901.11504*, 2019.
- [24] Jiao Li, Yueping Sun, Robin J. Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J. Mattingly, Thomas C. Wieggers, and Zhiyong Lu. Biocreative v cdr task corpus: a resource for chemical disease relation extraction. *Database : the journal of biological databases and curation*, 2016, 2016.
- [25] Martin Krallinger, Obdulia Rabal, Florian Leitner, Miguel Vazquez, David Salgado, Zhiyong lu, Robert Leaman, Yanan Lu, Donghong Ji, Daniel Lowe, Roger Sayle, Riza Batista-Navarro, Rafal Rak, Torsten Huber, Tim Rocktschel, Srgio Matos, David Campos, Buzhou Tang, Hua Xu, and Alfonso Valencia. The chemdner corpus of chemicals and drugs and its annotation principles. *Journal of Cheminformatics*, 7:S2, 03 2015.
- [26] Rezarta Islamaj Doundefinedan, Robert Leaman, and Zhiyong Lu. Ncbi disease corpus. *J. of Biomedical Informatics*, 47(C):110, February 2014.
- [27] 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, JNLPBA 04*, page 7075, USA, 2004. Association for Computational Linguistics.
- [28] Rie Kubota Ando. Biocreative ii gene mention tagging system at ibm watson. 2007.
- [29] Martin Gerner, Goran Nenadic, and Casey M Bergman. Linnaeus: a species name identification system for biomedical literature. *BMC bioinformatics*, 11(1):85, 2010.
- [30] 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):1–6, 06 2013.
- [31] Kai Xu, Zhenguo Yang, Peipei Kang, Qi Wang, and Wenyin Liu. Document-level attention-based bilstm-crf incorporating disease dictionary for disease named entity recognition. *Computers in biology and medicine*, 108:122–132, 2019.
- [32] Devendra Singh Sachan, Pengtao Xie, Mrinmaya Sachan, and Eric P Xing. Effective use of bidirectional language modeling for transfer learning in biomedical named entity recognition, 2017.
- [33] Yinxia Lou, Yue Zhang, Tao Qian, Fei Li, Shufeng Xiong, and Donghong Ji. A transition-based joint model for disease named entity recognition and normalization. *Bioinformatics (Oxford, England)*, 33, 03 2017.
- [34] Ling Luo, Zhihao Yang, Pei Yang, Yin Zhang, Lei Wang, Hongfei Lin, and Jian Wang. An attention based bilstm crf approach to document level chemical named entity recognition. *Bioinformatics*, 34:13811388, 2018.
- [35] Wonjin Yoon, Chan Ho So, Jinhyuk Lee, and Jaewoo Kang. Collabonet: collaboration of deep neural networks for biomedical named entity recognition. *BMC Bioinformatics*, 20(S10), May 2019.
- [36] John M Giorgi and Gary D Bader. Transfer learning for biomedical named entity recognition with neural networks. *Bioinformatics*, 34(23):4087–4094, 2018.