# Transfer Learning based prognostic clinical prediction Framework usingBiomedical NLP

## Sharma Pooja[1], Maheshwari Manish[2]

[1]*Research Scholar, Department of Computer Science & Application , Makhanlal Chaturvedi National University of Journalism and Communication, Bhopal(M.P.)*
[2] *Professor & Head, Department of Computer Science & Application , Makhanlal Chaturvedi National University of Journalism and Communication, Bhopal(M.P.)*

**ABSTRACT:**

In this study, we offer a novel framework for biomedical named entity recognition based on transfer learning using the BioALBERT model. The extraction of a great deal of biomedical knowledge from unstructured texts into organized formats relies heavily on the recognition of biomedical entities in literature, which is a challenging area of study. To implement biological named entity recognition (BioNER), the sequence labeling framework is now the gold standard. The performance of this approach is inconsistent, and it often fails to make full use of the semantic information in the dataset. To have a complete picture of a disease, one must be familiar with its signs and symptoms, medical evaluation, and treatment options. A great deal of medical and scientific work relies on this illness data, including disease diagnosis, consumer health question answering, and the development of medical nomenclature. Instead of approaching the BioNER task as a sequence labeling problem, we present a formulation of the problem as a machine reading comprehension (MRC) issue in this study. Carefully constructed queries can include more prior knowledge into this formulation, and unlike conditional random fields (CRF), no decoding methods are required. Although pre-trained language models like BERT have shown success in extracting syntactic, semantic, and world knowledge from text, we find that they can be further enhanced by specialized information like knowledge about symptoms, diagnosis, and other elements of an illness. Therefore, we combine ALBERT with medical knowledge to enhance BioNER. In specifically, we evaluate a new approach to training that incorporates illness knowledge infusion on BioALBERT. By showing that these models can be enhanced in nearly all cases, the experiments conducted for this task indicate the efficacy of disease knowledge infusion.

## I. INTRODUCTION

The widespread adoption of EHR systems in healthcare has created a wealth of new opportunities for clinical investigation by providing access to a wealth of real-world data sources. Natural language processing (NLP) techniques have been used as an artificial intelligence strategy to extract information from clinical narratives in electronic health records since they include a plethora of important clinical information. However, in free-form texts like electronic health records, a clinical narrative framework still conceals many clinical data. Therefore, biomedical NLP algorithms must be implemented to make full use of EHR data and automatically convert clinical narrative text into structured clinical data. In this way, biomedical NLP applications can be used to improve clinical decision-making, raise awareness of health problems, and successfully delay or prevent disease.[1]

The purpose of BioNER is to automatically detect biological things (e.g., chemicals, diseases, and proteins) in texts. Recognizing biological things accurately is a prerequisite for extracting biomedical knowledge from unstructured texts and turning it into organized formats. As a result, the BioNER problem is very important for the scientific community to tackle. Natural language processing (NLP) and domain-specific expertise are typically leveraged during the feature design phase for developing BioNER methods. Common models used in biomedicine include Lou's joint model [2], TaggerOne [3], DNorm [4], tmChem [5], etc. However, feature engineering relies heavily on human intervention and subject expertise. Furthermore, these characteristics are unique to entities and models alike. In recent years[6], neural networks with independent feature learning capabilities have seen increased usage in NER tasks. In

order to recognize biomedical objects, numerous neural network methods[7-12] have been presented. In these methods, bidirectional long short-term memory (BiLSTM)[13] is utilized to learn vector representations of each word or token in a phrase, which are subsequently fed into conditional random fields (CRF)[14].

Transfer learning has been increasingly prominent in the field of natural language processing as a result of its ability to apply knowledge gained in one context to another. The need for fine-tuning the performance and the dataset decreases. The domain adaptability property is the product of the transfer learning method's word vector mapping between similar words and vectors. To acquire transfer learning abilities and general-purpose knowledge, the entire model is pre-trained on a data-rich task. Later tasks can benefit from this knowledge. The network might be trained using transfer learning and the vast amounts of publicly available text data. For NLP applications, numerous representations have been proposed, including the Pre-trained Transformer (PT), Embeddings from Language Models (ELMo), RoBERTa, and ALBERT. All of these models were shown to be very effective for a certain group of applications, but not for the broader class of applications. Different models have different levels of success on different jobs because they apply different strategies. Therefore, a unified and systematic strategy is required to fully comprehend transfer learning. In addition, the performance of conventional pretrained learning models degrades for biomedical applications because they were trained on generic corpora like Wikipedia and Wordnet. The inefficient performance of these models on biomedical data has been addressed by biomedical natural language processing (BioNLP) researchers by training these models on biological and clinical corpora.

Many NLP tasks have lately been mastered by language models like ELMo [15] and BERT [16], which have attained state-of-the-art (SOTA) performance. The softmax function and BERT that was pre-trained on biomedical corpora were used by Lee et al.[17] to recognize biological things. Their method yielded state-of-the-art performance on a variety of biomedical datasets. Due to their ability to self-learn features, neural network approaches can outperform their feature engineering counterparts. In many existing methods, the BioNER task is recast as a sequence labeling problem, with the goal of using a sequence labeling model to assign a label to each token in a given sequence. Neither the BiLSTM-CRF nor the BioBERT-Softmax models are capable of effectively mining the semantic information present in the sequence labeling framework.. Integrating the semantic knowledge learned at the final BioBERT layer into the sequence labeling framework is difficult. Our work on BioNER (termed BioBERT-MRC) is inspired by the recent movement toward formalizing NLP tasks as machine reading comprehension (MRC) challenges [18]. Responding to a language query that encodes a biological entity type is how the MRC framework identifies them.

BERT and its variants [19] have recently caused a ruckus in the NLP field with their innovative approach to learning and exploiting information. These models learn language and world knowledge in its parameters through self-supervised pre-training over huge volumes of unannotated data, which they then utilize to refine their performance on subsequent tasks. In addition, several BERT models with applications in the medical field are proposed, one of which is BioBERT [20]. A masked language model (MLM) is used to pre-train these models using biomedical corpora, allowing for the prediction of tokens that have been randomly masked based on their context. The purpose of this MLM method is not to learn anything about the sickness itself, but rather to capture the semantic associations that exist between a given environment and a set of randomly masked characters. Since the semantic ties between the linked disease and aspect might not be randomly veiled or stated at all in the disease-descriptive text, MLM is unable to capture them well. Therefore, a novel training method is required to record this health-related data.

We describe a novel disease knowledge infusion training strategy in this research that may be used to explicitly enrich BERT-like models with the sickness information. The central idea is to train BERT to infer the related disease and aspect from a disease description text using weakly-supervised signals from Wikipedia. By feeding it text from a section of a Wikipedia page that normally explains a feature of the

disease, BERT is taught to infer the title of the associated article (disease name) and the title of the relevant section (aspect name). BERT is taught to recognize that the diagnostic part of the article "COVID-19" is the source of the sentence "...testing is real-time reverse transcription polymerase chain reaction (rRT-PCR)..." from Table 1. When the condition or perspective is not made clear in the context of a given section, we construct introductory auxiliary sentences including the disease or perspective (for instance, "What is the diagnosis of COVID-19?") and insert them at the beginning of the relevant paragraphs. Then, we use the auxiliary sentence to conceal the disease and feature, and we let BERT-like models infer them from the passage. By doing so, BERT is able to create a semantic connection between a disease-descriptive language and the matching feature and disease. In order to evaluate the efficacy of disease knowledge infusion, we conduct experiments on a BioALBERT, on the tasks of answering CHQs, inferring medical terminology, and recognizing disease names.

The large-scale pretrained model BioALBERT has been fine-tuned for domain-specific adaptation using corpora from the biomedical area. This study details an improved method for extracting EHR system context data for therapeutic use.The proposed work's key contribution is the potential utility to clinicians of extracting context information from the available dataset for a clinical domain. Here, the contextual biological data is used to fine-tune the pretrained model to achieve the desired performance. The improvised ways for sharing parameters provide value to the research effort because they reduce the need for physical memory. Sentence-level tokenization yields word-piece embeddings useful for fine-tuning contextual summary synthesis, which adds novelty to the proposed approach.

Here is how the rest of the paper is structured: The second half of the paper provides a summary of previous presentations by a variety of academics that cover related ground. In Section III, we describe the structural makeup of ALBERT and its biological variants. In Section IV, we demonstrate how to use BioALBERT to retrieve contextual information about the dataset that was provided for biomedical NER. Part V of the experimental inquiry evaluates and discusses the effectiveness of the suggested method, and Part VII wraps up the project.

## RELATED WORK

The exponential rise of the healthcare business has attracted many academics to examine a range of hitherto undiscovered aspects of the medical paradigm It is largely attributed to the digitalization of data and the ability of artificial intelligence to extract valuable information from it,. Natural language processing (NLP) innovations have improved processing standards. In particular, Biomedical NLP has changed the entire medicine development process and the assessment of sequence variants in the medical industry, as well as every other area of data analytics in the medical field. In the biomedical literature, sophisticated LMs did not perform well since they were trained on generic corpora. The BioNLPresearchers resolved the limitation by training the LMs on biological and clinical corpora and verifying the corresponding performance across multiple BioNLP applications.

Most deep learning-based research representing the clinical concept of EHRs use the word2vec model's skip-gram structure [21]. The skipgram approach operates under the premise that the interpretation of a concept is conditional on its immediate setting. So, the skip-gram method uses the sequence of concepts to predict the context of the chosen target thought. Once concepts are represented, researchers can employ analysis to establish the connections among them, such as between diseases or between clinical events and diseases. Similar ideas can be applied by researchers in different tasks to aid doctors in decision making and provide them with useful information.

Attention techniques are also utilized in the examination of electronic medical records. Thanks to the attention mechanism, deep learning models may narrow in on the most relevant details to achieve the desired outcome, rather than wasting time processing irrelevant data. Using the focus mechanisms, it is possible to ascertain if data points validate the model's predictions. Using attention mechanisms and recurrent neural networks, Nigam et al. [22] created a system to learn patient representations from temporal electronic health record data. The model was then used to predict the probability of an

impending hospitalization. Experiment results showed that a deep learning model has the potential to improve prediction accuracy.

Combining several neural networks or methods has been shown in numerous research to improve model accuracy and performance. For instance, [23] used unsupervised deep learning to learn a patient's representation from their EHR. They used a stacked denoising autoencoder with three layers to capture the data's hierarchical relationship and interdependence. The deep learning system Mikolov et al. [24] developed to extract patterns of patient information consists of recurrent neural networks and convolutional neural networks. Using a single-layer decision tree and long short-term memory (LSTM), Rajkomar et al. [25] were able to retrieve data from the dataset. The deep learning model consistently outperforms the gold-standard clinical prediction model.

Biomedical ELMo (BioELMo), a variant of ELMo developed by Jinal. [26] and trained on PubMed abstracts, was first presented. Characteristics of the connections between the various entities in the biomedical data have been retrieved by the authors. Training over scientific content was proposed by Beltagy al.[27], who devised a paradigm they dubbed Scientific BERT (SciBERT). Similarly, Si et al. [28] have used clinical notes as corpora to train BERT-based transfer learning models. They solved the named-entity-recognition (NER) problem for biomedical data and improved upon both traditional non-contextual and contextual word embedding. Peng et al. [15] created a standard called the Biomedical Language Understanding Evaluation (BLUE) score using ten datasets to generate five assignments. Since then, many scientists have used this number to judge the precision of their simulations. The results show that BioNLP performance can be enhanced by training LMs on biomedical corpora. According to Li et al. [29], BioBERT is the most widely used LM in the field of biomedicine. It uses the PubMed and PMC corpus for preliminary training. They've tweaked the proposed model and looked into its potential in three different BioNLP tasks: relational extraction, natural language extraction (NER), and question answering (QA).Using data from PubMed articles, Gu et al. [30] introduced a new LM they called PubMedBERT. They used the LM that had been trained on large corpora to verify the model's efficacy and accuracy and to extract a domain-specific lexicon. KeBioLM[31], a biologically pre-trained LM, utilised data from the Unified Medical Language System (UMLS). Two BioNLP projects were finished using KeBioLM, demonstrating that it was possible to develop a specialized LM for a given domain.

Training these LMs takes a long time and requires a lot of computing power because they all use the BERT design. The extent to which these LMs can be applied outside of the specific BioNLP tasks for which they were validated is also open to debate.ALBERT has been demonstrated to outperform BERT on NLP tasks. Our hypothesis is that a biomedical variant of ALBERT, called BioALBERT, can be effectively trained with biomedical corpora and then fine-tuned for the task of context-aware summary generation, providing a useful tool for clinical analysis in a variety of settings.

## II. ARCHITECTURAL FRAMEWORK OF ALBERT

ALBERT is a simplified version of BERT that keeps the functionality of the original while making use of a significant reduction in the number of parameters. Integration of the factorized embedding parameterization (FEP) strategy and the following cross-layer parameter (CLP) sharing method is what leads to a reduction in the number of parameters that need to be used. By decomposing the massive vocabulary embedding matrix into two more manageable matrices, FEP makes the dimensionality of the hidden layer less dependent on the vocabulary embedding. This is accomplished by dividing the matrix into two smaller matrices. In addition to this, it enables you to maintain the same amount of parameters while simultaneously raising the total number of hidden levels. The CLP sharing strategy improves parameter efficiency because it limits the rise in parameters even when the network's depth changes. This keeps parameter efficiency at a high level.In addition to this, it demonstrates a more consistent and stable training outcome for the primary BERT model. Despite the fact that ALBERT requires 18 times fewer parameters, it was discovered that the training performance of an ALBERT model is 1.7 times faster than that of a BERT model [32]. This is despite the fact that ALBERT requires less parameters.
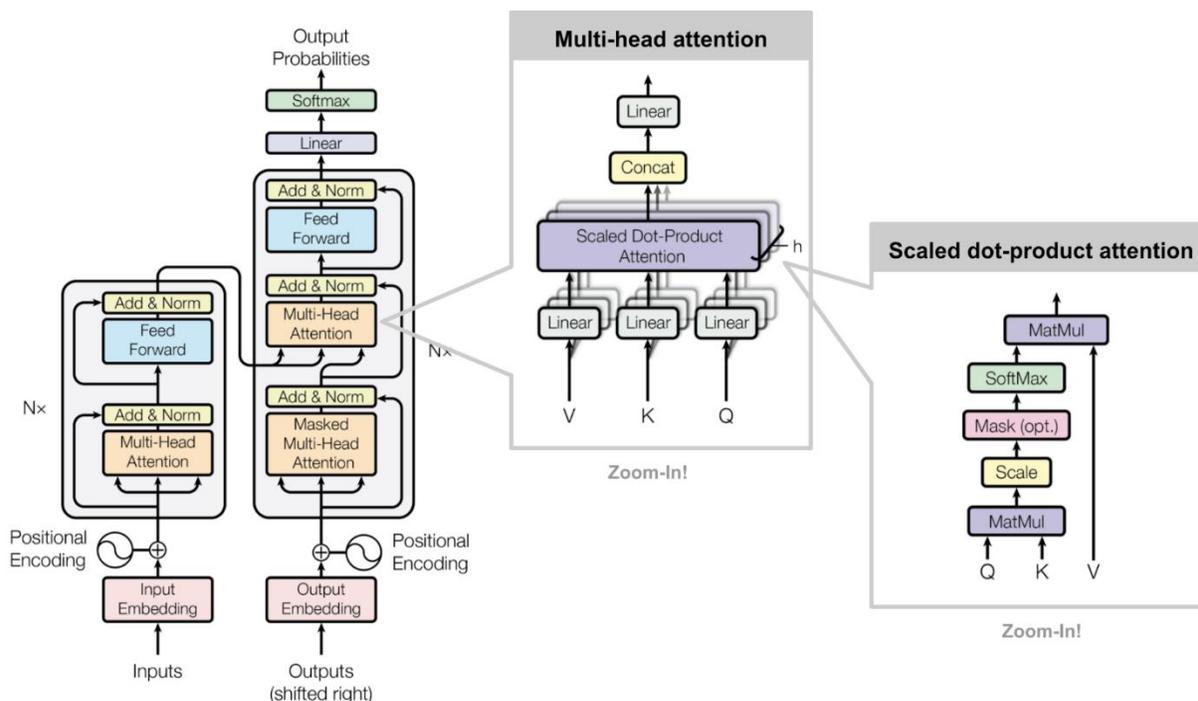
Figure1. Typical Transformer model [32]

Encoder-decoder based models are the foundation upon which ALBERT is constructed, as shown in Figure 1. This model utilizes a transformer architecture with multiple layers. In this strategy, the process of encoding is predicated on the encoder's capacity to concentrate inwardly, whereas the process of decoding is predicated on the encoder's ability to concentrate on its outputs. This structure is made up of a number of blocks stacked one atop the other in ascending order. Each individual unit is made up of a feedforward network as well as a multi-head attention block. It is necessary to have sizeable hidden layer representations in order to accommodate data context and word-level embeddings, which leads to an increase in the total number of parameters. The formula that is typically used to calculate this number is V * H, where V refers to the size of the vocabulary and H refers to the size of the hidden layer.

When E is factored with FEP, the resulting information can be used to determine the size of the embedding. The final parameter order is about V*E + E*H, which is a significant step in the right direction. Stacking the independent layers not only results in a significant increase in the model's level of redundancy but also enhances the capacity of the model to acquire new information. In order to lessen the likelihood of this happening, ALBERT makes advantage of group-wide parameter sharing between layers. In this particular instance, the number of layers and the number of parameters are subject to a give-and-take relationship. Therefore, ALBERT is a simplified version of BERT that is a lot more beneficial than the original.It has the potential to lessen the computational burden of a wide variety of applications and to boost the effectiveness of language comprehension tasks further down the production line.

## III. PROPOSED BIOALBERT BASED NER

This study presents architecture for biomedical NER and disease prediction for prognosis that employs transfer learning for NLP on biological data. This architecture can be used to collect context features in a clinical setting. In this instance, we achieve transfer learning with the help of BioALBERT by refining the massively pretrained LM over biological corpora. The following is a statement that has been made regarding the issue of biomedical NER:

An input sentence $X = \{x_1, x_2, ..., x_N\}$ is given, in which $x_i$ denotes the $i^{th}$ word or token and $N$ is the sentence's length. Sorting every word or token in $X$ and assigning it to a corresponding label $y \in Y$ being a preset list of all potential label types, such as *CHEMICAL*, *DISEASE*, and *PROTEIN*—is the aim of

NER. The labeling-style NER dataset is transformed into a set of (*Context, Query, Answer*) triples by us when we structure the NER task as an MRC task. Context refers to an input sentence (*X*), Query to a query sentence that is constructed based on the sentence (*X*), and Answer refers to the target entity span. We initially created a query $Q_y = \{q_1, q_2, \ldots, q_M\}$ for each sentence for each label type *y*, where *M* is the query's length. After that, we were able to obtain the annotated entities $x_{start,end}$, which are substrings of *X* and *start* $\leq$ *end*, in accordance with the annotated labels *Y*. For instance, the labels "*B O O O O O*" correlate to the statement "Meloxicam - induced liver toxicity" from the BC5CDR-Chem dataset. We may extract the entities and their spans: "Meloxicam"0,0, based on the labels. In the end, we created the triple (*X*, $Q_y$, $x_{start,end}$), which is precisely the triple we require (*Context, Query, Answer*). To create queries, we used biological entities from the target dataset's training and development set.

The entire process is categorized primarily into the following four steps:

a) Data Preprocessing

To test the efficacy of our method, we use the *BC5CDR*[33] datasets that has been preprocessed and made publicly available. To evaluate chemicals and diseases, researchers use the *BC5CDR-Chem* and *BC5CDR-Disease* sub-datasets. Since *BC5CDR-Chem* and *BC5CDR-Disease* were used to evaluate the vast majority of existing methods, we did the same. *BC5CDR-Chem* contains 15,411 annotations and 14,228 sentences related to the chemical/drug entity. Even yet, *BC5CDR*-illness's "disease" object has 14,228 phrases and 12,694 annotations. For the trials, a new training set was created by merging the original training and development sets. The hyper-parameters were then fine-tuned using a validation set sampled from 10% of the new training set. The model was tested with data taken directly from the test set. We followed the standard practice of separating data in this way because it was used by the vast majority of earlier works. Discharge summary features tokenize the text of datasets using the [*CLS*] and [*SEP*] tokens. Tokens at the beginning of a sentence ([*CLS*]) and the end of a sentence ([*SEP*]) indicate their respective locations.

b) Machine reading comprehension

The MRC methods allowed for the extraction of answer spans from the context by use of a targeted query. Predicting where the response spans start and where they terminate is essentially two separate formal classification problems. In recent years, it has become common practice to transform similar NLP tasks into MRC. Using the question-answering architecture, McCann et al.[34] successfully implemented 10 separate NLP tasks, all of which achieved competitive results. There does not appear to be any active research on BioNER for ALBERT inside the MRC framework at the present time, at least none that we are aware of. Our study differs significantly from Li's in that we are only interested in living things. In addition, for the first time, we explore how different parts of the model affect BioNER.

c) Model Training and Testing

Although standard LMs have shown promise for a variety of NLP problems, there are significant risks involved in using them for NER. These issues include contextual dependency, the usage of acronyms, the possibility that one entity could belong to multiple entity types, a lack of training data, and other variables. Consequently, and as was indicated in the introduction, the most sophisticated NER models use transformer-implemented context-dependent language models that are trained on biological corpora. The more sophisticated and rapid LM known as the bioALBERT transformer model was used in this specific experiment. Some of the drawbacks of conventional domain-specific language models can be avoided thanks to the model's ability to adapt to the context in which it is embedded.

First, the context X and the query Qy are concatenated to create the combined sequence {[CLS], X, [SEP], Q, [SEP]}. Following that, BERT is fed the combined sequence, which is described by the following formulas:

$$h_i^0 = W_e t_i + W_b \tag{1}$$

$$h_i^l = Trm(h_i^{l-1}) \tag{2}$$

$$H = [h_1^L; h_2^L; \ldots; h_N^L] \in R^{N \times d} \tag{3}$$

where $t_i$ is the embedding of the $i^{th}$ token. The Transformer block incorporates multi-head attention layers, totally linked layers, and normalizing layers; these characteristics are represented by $W_e$, $W_b$, $L$, and $Trm$. The current layer number is $l(1 \leq l \leq L)$. Furthermore, $N$ is the length of the context, and $H$ is the BERT's output. We simply removed question representations because they are not relevant to the prediction being made by the model. In the MRC framework, there are typically two options for choosing the span. The first predicts the beginning and ending indices using two n-class classifiers, where $n$ is the length of the context. This approach has the limitation that, due to the way the function is constructed, only one span may be generated from every given input sequence. It's also possible to create two binary classifiers. Based on their respective projections, tokens can be categorized as either beginning or ending indexes. Because it allows for the production of a wide variety of start and end indexes for a given sequence, this method can successfully identify all target entities based on $Q_y$. In this inquiry, we followed Option 2. Using the representation matrix $H$ generated by BERT, the model first calculates the probability that each token is a start index. The formula is as follows:

$$L_{start} = linear(HW_{start}) \in R^{N \times 2} \tag{4}$$

Where $W_{start}$ is the weight to learn, and linear is a completely connected layer in this instance. The hidden representation of the index, which is utilized to ascertain the starting position of a target entity for a particular query, is indicated by each row in $L_{start}$. Next, each token's likelihood of being the matching end index is predicted by the model. The equation is:

$$L_{end} = linear\big(HW_{end}; softmax(L_{start})\big) \in R^{N \times 2} \tag{5}$$

where ";" denotes the concatenation operation, linear is a completely connected layer, and $W_{end}$ is the weight to learn.

Lastly, we can get the projected indices, or $I_{start}$ and $I_{end}$, that may represent the start or end positions by using the *argmax* function to each row of $L_{start}$ and $L_{end}$:

$$I_{start} = \{i | argmax(L_{start}^i) = 1, \ i = 1,2, \ldots, N\} \tag{6}$$

$$I_{end} = \{j | argmax(L_{end}^i) = 1, \ j = 1,2, \ldots, N\} \tag{7}$$

where the superscripts $i$ and $j$ designate the $i$ and $j$ rows of the matrix. In setting $X$, there could be several occurrences of entity $Y$. This means it is possible to anticipate multiple *start* and *end* indices from the same input sequence. Since all of the datasets we examined are flat NER datasets, we also used the nearest match concept to synchronize the beginning and ending indices to arrive at our conclusions. When one end/start index corresponds to several start/end indices, only the closest pair will be used.

BioALBERT was trained using a massive dataset constructed from biomedical corpora with the intention of improving its learning performance via CLP sharing. Layers beyond the first can learn the model's parameters with the help of the information learned in the first block thanks to CLP sharing. This eliminates the need to independently acquire the model's parameterization on a layer-by-layer basis. By assessing the coherence loss of each sentence and making an effort to minimize it during training, sentence-order-prediction (SOP) is an extra method that supplements contextual information. This method improves the model's representation and understanding by leveraging the training corpus to generate random word pairs. In order to reduce the hidden layer's reliance on the vocabulary embedding, the model has been FEP-enhanced by splitting the huge vocabulary embedding matrix into two smaller matrices. This is because, according to the transformer model, the size of the embedding is directly proportional to the level of the hidden layer, thus we can gradually increase the number of hidden layers while maintaining the same number of word embedding settings. CLP sharing, on the other hand, limits the increase in parameters with increasing network depth. Without compromising BERT's efficacy, it boosts the efficiency of the parameters and the consistency and dependability of the training results. Vocabulary embedding parameter values can be increased independently of the number of hidden layers used. The next section delves deeper into the initial training and fine-tuning of BioALBERT.

A. BioALBERT's pre-training

To prepare the model for pre-training, the text from the datasets needs to be further processed in a more systematic way. This was required in order to make use of the material. Several common processing operations are applied to this raw text, such as removing blank lines, creating paragraphs, removing sentences with a length of fewer than twenty characters, and inserting blank lines between each document. The data was initially handled as a sentence text because tokenization relies on sentence embeddings. Tokenization needed to function properly, thus this was done. There is a blank line between each document because every sentence in the input text document is represented by a line in the document. The maximum length for each declaration was 512 characters. Using this updated dataset, the BioALBERT model was pretrained so that it wouldn't require the ALBERT model's established vocabulary.

B. BioALBERT fine-tuning

In this section, we will examine potential improvements that could be made to the contextual summary work that BioALBERT does. In this approach, the individual words that make up a sentence are isolated and processed separately from one another. In order to further enhance the model that was proposed, the discharge information summaries that were gathered from the datasets were utilized. The goal of this exercise is to familiarize yourself with some of the most important terms and phrases in the field of biomedicine so that you can make intelligent estimates about them. The purpose of this project is for you to create a clinically meaningful summary out of the symptom information that has been provided to you. When compared to pre-training, fine-tuning is a simpler process that demands a lower level of user attentiveness.

N-tuple representations of the start and end indexes of the biological entity, $Y_{start}$ and $Y_{end}$, respectively, that correspond to the ground-truth labels of each token $x_i$ are provided by $Y_{start}$ and $Y_{end}$. These representations are referred to as $Y_{start}$ and $Y_{end}$. The specification of the loss function looks like this:

$$Loss_{start} = CE(L_{start}, Y_{start}) \tag{8}$$
$$Loss_{end} = CE(L_{end}, Y_{end}) \tag{9}$$
$$Loss = (Loss_{start} + Loss_{end})/2 \tag{10}$$

In this context, CE refers to the cross-entropy loss function. The test's beginning and finishing indexes are selected ad hoc based on the values of two variables, $I_{start}$ and $I_{end}$. Then, using the nearest match technique, the answers are derived by comparing the beginning and ending indexes.

The proposed approach makes use of informal methods for exchanging parameters and makes only modest use of physical memory. Tokenizing sentences allows for the production of word embeddings, which in turn can be used to generate more accurate contextual summaries. The one-of-a-kind dataset is used to generate model-specific fine-tuning exercises. The model's weights were set using the previously constructed pre-trained model. Before performing the tuning operation, we trained with lowercase texts using a batch size of 32, a learning rate of $1 \times 10^{-5}$, and a total of 5336 steps. For the first round of BioALBERT model training, we relied on the Tensorflow TPUv3-8 processing unit. The majority of the hyper-parameters from the traditional ALBERT model have been preserved in their default settings. The model was tested at regular intervals, and its performance was rated using the evaluation checkpoint. The ultimate forecast was calculated using the top-performing model and the stubborn development dataset. *PyTorch*-transformers with *XLNet* and BERT implementations were employed in the studies. The entire training procedure is depicted in Figure 2, from start to finish. Model updates were performed without interfering with the layering process. Initially, Adam [35] was employed as the optimizer, but the training time was significantly reduced when the researchers switched to the Layer wise Adaptive Large Batch (LAMB) optimizer [36]. The overall performance of the development set was taken into account while deciding which hyperparameters to use and how to optimize them. The Transformer model appears to be highly sensitive to the settings it is provided; most attempts to apply it with varying learning rates resulted in subpar performance. The most effective training speeds were $7 \times 10^4$ and $6 \times 10^6$. The codes were evaluated depending on their level of confidence using the binary cross entropy logits.
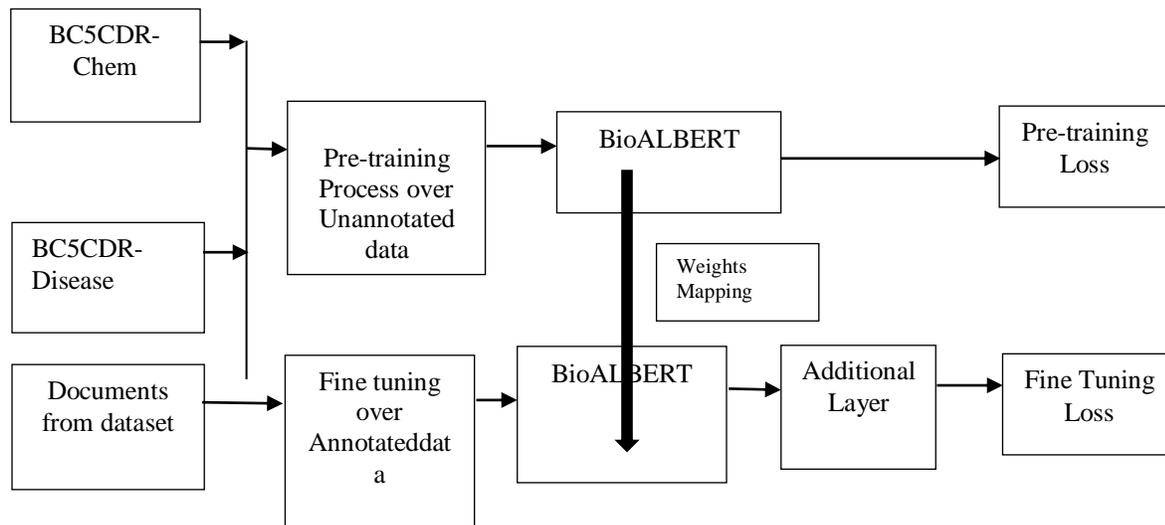
Figure 2. End-to-end training procedure

## IV. RESULTS AND DISCUSSION

To test the efficacy of our method, we use the *BC5CDR* datasets that has preprocessed and made publicly available. To evaluate chemicals and diseases, researchers use the BC5CDR-Chem and BC5CDR-Disease sub-datasets. Since *BC5CDR-Chem* and *BC5CDR-Disease* were used to evaluate the vast majority of existing methods, we did the same. *BC5CDR-Chem* contains 15,411 annotations and 14,228 sentences related to the "*chemical/drug*" entity. Even yet, BC5CDR-illness's "*disease*" object has 14,228 phrases and 12,694 annotations. For the trials, a new training set was created by merging the original training and development sets. The hyper-parameters were then fine-tuned using a validation set sampled from 10% of the new training set. The model was tested with data taken directly from the test set. Because of its low processing complexity, ALBERT is the basis for many recent efforts. To facilitate comparisons with previous investigations, the BERT models in this work are all based on the ALBERT framework.

The F1-score (F1) is used to evaluate efficacy; it gives equal weight to recall (*R*) and precision (*P*). First, find F1 using the formula: *2PR/(P + R)*. For each experiment performed five times, we report the highest possible F1-score (denoted by max), the average F1-score (denoted by mean), and the standard deviation (denoted by *std*). We also utilize T-TEST to conduct statistical significance tests and offer the confidence interval. Our experiments show that BioALBERT is most effective on one- or three-epoch versions of the BC5CDR datasets. This has two root causes: There are two notable features of these data sets: The first is their massive scale, and the second is the powerful feature learning capabilities provided by BioALBERT. Another thing to keep in mind is that the input of BioALBERT is "[*CLS*] Context [*SEP*] Query [SEP], not "[CLS] Query [SEP] Context [SEP]". Using "[*CLS*] Query [*SEP*] Context [*SEP*]" as an input was also investigated, but this did not result in any noticeable gains in model performance.

The effects of ALBERT inside the sequence labeling framework and the MRC framework were extensively explored in this work. Our first step was to evaluate ALBERT's BioNER capabilities in both the MRC and sequence labeling environments. Training batches for the BioALBERT base model were originally set at 1,024, but this number was reduced to 256 due to constraints on available computing capacity. Table 1 details the hyperparameters used in the training procedure.

Table 1. Pre-training Parameters

| Parameters | Description/Values | | Parameters | Description/Values |
|---|---|---|---|---|
| Baseline model | ALBERT | | Number of layers | 12 |
| Optimization technique | LAMB | | Hidden layers | 768 |
| Kernel | GeLU | | Size of vocabulary | 30000 |
| Maximum length of sentence | 512 | | Warm-up steps | 3125 |
| Number of attention heads | 12 | | Size of evaluation batch | 16 |
| Size of training batch | 1024 | | Size of embedding | 128 |

BioALBERT's processing requirements were not much higher than those of other baseline models, while fine-tuning only required less computation than pre-training. BioALBERT's spontaneous parameter sharing techniques and lower physical memory requirement significantly improve its performance. Tokenizing sentences into sentences allows the training process to learn word embeddings more quickly. When fine-tuning the model, the pre-trained BioALBERT's weights are employed.Here, the *AdamW* optimizer is utilized with a batch size of 32 and a learning rate of 0.00001. We capitalized every word and set a maximum phrase length of 128 characters for every job. In the end, we used 512 warm-up steps to refine our pre-trained models for 10,000 training steps during the fine-tuning process.

The best tuning of the suggested model is presented in figure 5, which also displays the training performance in terms of cost function. The estimation error variation is depicted in Figure 6, and the error's convergent nature reflects the learning characteristics.
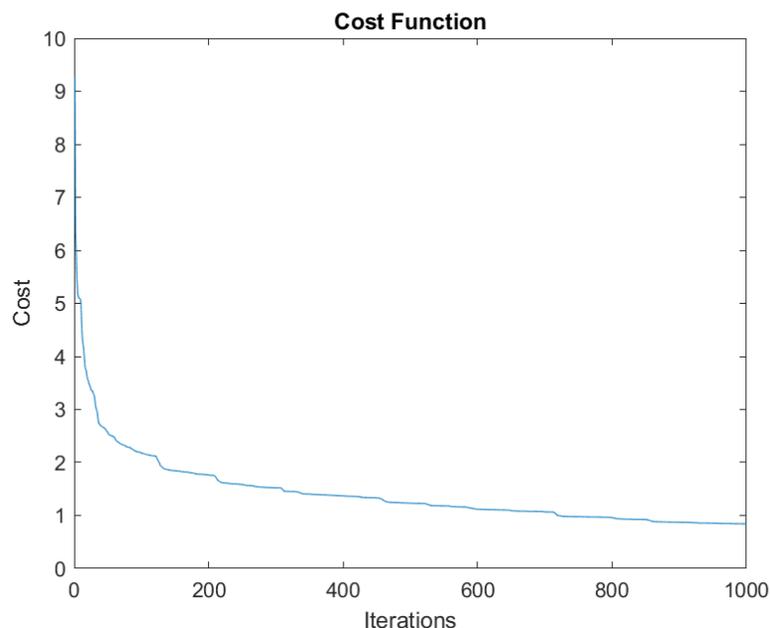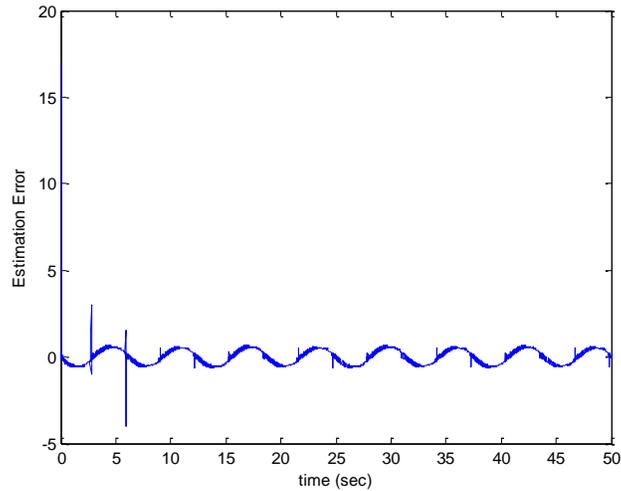


Figure 5 Cost function v/s time

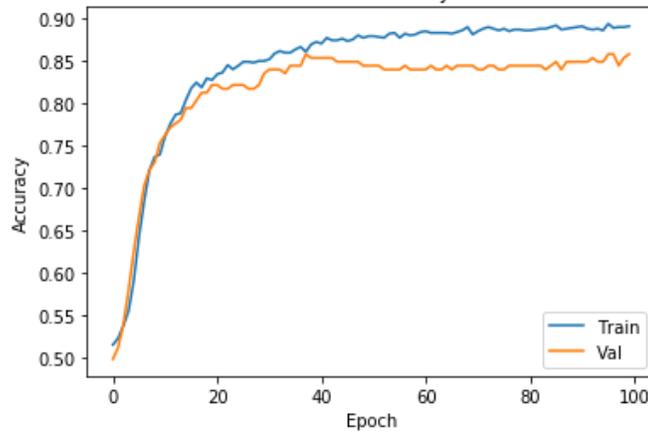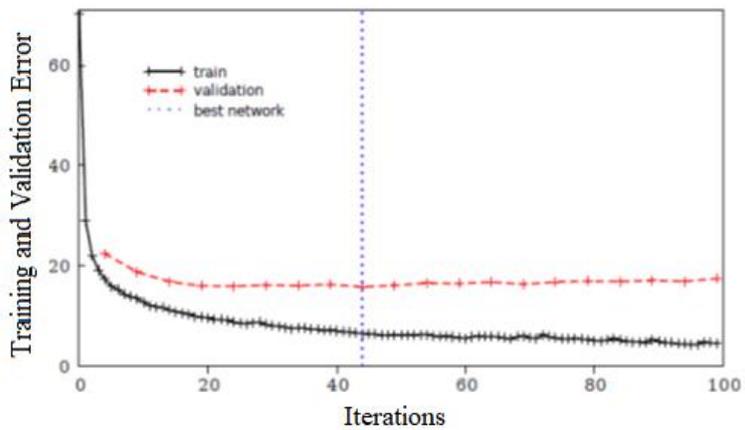Figure 6 Estimation error v/s time



Figure 7 Modeling Accuracy v/s number of iterations



b

Figure 8 Training and Validation Error wrt number of iterations

Figure 7 and 8 show the modelling accuracy and error during the training and validation over time for the training run on the online database.Each word in the input sequence is labeled with a Begin, Inside, and Outside (BIO) tag using these methods, treating BioNER as a token-level sequence labeling problem. The goal is to predict the BIO label for each token in the output sequence. The softmax method is used by BioALBERT-Softmax to perform the BIO classification on the BioALBERT output for each token in the sequence. Taking into account both the query $Q_y$ and the input sequence $X$, BioALBERT predicts the response spans $x_{start,end}$, which differs from the other three methods. It has been found that BioALBERT-Softmax has an average F1-score of 92.99% on the dataset, which is higher than both *BioBERT-CRF*[37](91.52%) and *BioBERT-BiLSTM-CRF*[37] (91.39%). One possible explanation is that CRF and LSTM aren't particularly good at finding long-distance hidden representations in this dataset because of the prevalence of really long words. However, BioBERT-MRC shows the most improvement over the other three methods. BioBERT-MRC can greatly improve the performance of BioNER tasks compared to sequence labeling strategies, regardless of the size of the corpus or the length of the sentences. These results show that BERT performs better in the MRC framework than in the sequence labeling framework when it comes to identifying biomedical items.

## V. CONCLUSION

This research uses the BioALBERT model to propose a novel framework for biomedical named entity recognition through transfer learning. A difficult area of study is the recognition of biomedical entities in literature, which serves as the basis for extracting a significant amount of biomedical knowledge from unstructured texts into organized formats. Currently, the standard approach for implementing biological named entity recognition (BioNER) uses the sequence labeling framework. However, this method's performance is not always good, and it frequently fails to fully utilize the semantic information present in the dataset. Understanding an illness entails knowing about its many facets, including its symptoms, diagnosis, and course of therapy. Many scientific and health-related tasks, such as diagnosing diseases, answering consumer health questions, and deriving medical terminology, depend on this illness information. In this paper, we formulate the BioNER job as a machine reading comprehension (MRC) problem, rather than considering it as a sequence labeling problem. With carefully crafted queries, this formulation can incorporate more prior knowledge and does not require decoding procedures like conditional random fields (CRF). We find that pre-trained language models, such as BERT, can be further enhanced by specialized information such as knowledge about symptoms, diagnosis, treatments, and other elements of an illness, even though they have demonstrated success in extracting syntactic, semantic, and world knowledge from text. Therefore, in order to improve BioNER, we integrate ALBERT with illness knowledge. In particular, we assess a novel illness knowledge infusion training method on BioALBERT. Experiments conducted on this job demonstrate the viability of illness knowledge infusion, as these models may be improved in almost all cases.

## REFERENCES

[1]    S. A. Hasan and O. Farri, ``Clinical natural language processing with deep learning,'' in Data Science for Healthcare. Springer, 2019, pp. 147-171.

[2]    Y. Lou, Y. Zhang, T. Qian, F. Li, S. Xiong, D. Ji, A transition-based joint model for disease named entity recognition and normalization, Bioinformatics 33 (15) (2017) 2363–2371.

[3]    R. Leaman, Z. Lu, TaggerOne: joint named entity recognition and normalization with semi-Markov Models, Bioinformatics 32 (18) (2016) 2839–2846.

[4]    R. Leaman, R. Islamaj Doˇgan, Z. Lu, Dnorm: disease name normalization with pairwise learning to rank, Bioinformatics 29 (22) (2013) 2909–2917.

[5]    R. Leaman, C. Wei, Z. Lu, tmChem: a high performance approach for chemical named entity recognition and normalization, J. Cheminformat. 7 (1) (2015) 1–10.

[6]    G. Lample, M. Ballesteros, S. Subramanian, K. Kawakami, C. Dyer, Neural Architectures for Named Entity Recognition, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics, 2016, pp. 260–270.

[7]     A. Jagannatha, H. Yu, Structured prediction models for RNN based sequence labeling in clinical text, in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, 2016, pp. 856–865.

[8]     M. Habibi, L. Weber, M.L. Neves, D.L. Wiegandt, U. Leser, Deep learning with word embeddings improves biomedical named entity recognition, Bioinformatics 33 (14) (2017) i37–i48.

[9]     T.H. Dang, H.-Q. Le, T.M. Nguyen, S.T. Vu, D3NER: biomedical named entity recognition using CRF-biLSTM improved with fine-tuned embeddings of various linguistic information, Bioinformatics 34 (20) (2018) 3539–3546.

[10]    L. Luo, Z. Yang, P. Yang, Y. Zhang, L. Wang, H. Lin, J. Wang, An attention-based BiLSTM-CRF approach to document-level chemical named entity recognition, Bioinformatics 34 (8) (2018) 1381–1388.

[11]    D.S. Sachan, P. Xie, M. Sachan, E.P. Xing, Effective Use of Bidirectional Language Modeling for Transfer Learning in Biomedical Named Entity Recognition, Proc. Mach. Learn. Res. (2018) 383–402.

[12]    X. Wang, Y. Zhang, X. Ren, Y. Zhang, M. Zitnik, J. Shang, C. Langlotz, J. Han, Cross-type biomedical named entity recognition with deep multi-task learning, Bioinformatics 35 (10) (2018) 1745–1752.

[13]    W. Yoon, C.H. So, J. Lee, J. Kang, CollaboNet: collaboration of deep neural networks for biomedical named entity recognition, BMC Bioinformatics 20 (10) (2019) 55–65.

[14]    S. Hochreiter, J. Schmidhuber, Long Short-Term Memory, Neural Comput. 9 (8) (1997) 1735–1780.

[15]    Y. Peng, S. Yan and Z. Lu,"Transfer learning in biomedical natural language processing: An evaluation of BERT and ELMo on ten benchmarking datasets", In Proceedings of the 18th BioNLP Workshop and Shared Task, pp. 58–65, 2019.

[16]    J.Devlin, M. W. Chang, K. Lee and K. 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, vol. 1, pp. 4171–4186, 2019.

[17]    J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C.H. So, J. Kang, BioBERT: a pre-trained biomedical language representation model for biomedical text mining, Bioinformatics 36 (4) (2019) 1234–1240.

[18]    O. Levy, M. Seo, E. Choi, L. Zettlemoyer, Zero-Shot Relation Extraction via Reading Comprehension, in: Proceedings of the 21st Conference on Computational Natural Language Learning, 2017, pp. 333–342.

[19]    C. Sun, X.Qiu, Y. Xu, X. Xuang,"How to fine-tune BERT for text classification?", In: China National Conference on Chinese Computational Linguistics, pp. 194-206, Springer, 2019.

[20]    J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "BioBERT: a pre-trained biomedical language representation model for biomedical text mining, 2019.

[21]    J. Huang, C. Osorio and L. Wicent Sy,"An empirical evaluation of deep learning for ICD-9 code assignment using MIMIC-III clinical notes", Computer Methods and Programs in Biomedicine, vol. 177, pp.141–153, 2019.

[22]    P. Nigam,"Applying deep learning to ICD-9 multi-label classification from medical records", Technical report, Stanford University, 2016.

[23]    F. Li,W. Liu and H. Yu,"Extraction of information related to adverse drug events from electronic health record notes: design of an end-to-end model based on deep learning", JMIR medical informatics, vol. 6, no. 4, e12159, 2018.

[24]    T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, in: International Conference on Learning Representations, 2013.

[25]   A. Rajkomar, E. Oren, K. Chen, A. M. Dai, N. Hajaj, M. Hardt, P. J. Liu, X. Liu, J. Marcus, and M. Sun, ``Scalable and accurate deep learning with electronic health records,'' NPJ Digit. Med., vol. 1, no. 1, p. 18, 2018.

[26]   Jin Q, Dhingra B, Cohen WW, Lu X. Probing biomedical embeddings from language models. 2019. arXiv:1904.02181.

[27]   I. Beltagy, K. Lo, and A. Cohan, "SCIBERT: A pretrained language model for scientific text", In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 3606–3611, 2019.

[28]   Y. Si, J. Wang, H. Xu and K. Roberts,"Enhancing clinical concept extraction with contextual embeddings", Journal of American Medical Informatics Association,vol.  26, pp. 1297–1304, 2019.

[29]   Yichong Xu, Xiaodong Liu, Chunyuan Li, Hoifung Poon,  and  Jianfeng Gao. 2019. Doubletransfer at mediqa 2019: Multi-source transfer learning for natural language understanding in the medical domain. arXiv preprint arXiv:1906.04382.

[30]   Y. Gu, et al.,"Domain-specific language model pretraining for biomedical natural language processing",arXiv preprint arXiv:2007.15779, 2020.

[31]   Z. Yuan, Y. Liu, C. Tan, S. Huang and F. Huang, "Improving biomedical pretrained language models with knowledge",arXiv preprint arXiv:2104.10344, 2021.

[32]   Z.Lan, et al.,"Albert: A lite BERT for self-supervised learning of language representations", In International Conference on Learning Representations, 2019.

[33]   J. Li, Y. Sun, R.J. Johnson, D. Sciaky, C.-H. Wei, R. Leaman, A.P. Davis, C. J. Mattingly, T.C. Wiegers, Z. Lu, BioCreative V CDR task corpus: a resource for chemical disease relation extraction, Database 2016 (2016).

[34]   B. Mccann, N.S. Keskar, C. Xiong, R. Socher, The Natural Language Decathlon: Multitask Learning as Question Answering, arXiv: Computation and Language, 2018.

[35]   D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization", arXiv preprint arXiv:1412.6980, 2014.

[36]   You, Yang and Li, Jing and Reddi, Sashank and Hseu, Jonathan and Kumar, Sanjiv and Bhojanapalli, Srinadh and Song, Xiaodan and Demmel, James and Keutzer, Kurt and Hsieh, Cho-Jui, "Large Batch Optimization for Deep Learning: Training BERT in 76 minutes", 2019.

[37]   Cong Sun, Zhihao Yang, Lei Wang, Yin Zhang, Hongfei Lin, Jian Wang, Biomedical named entity recognition using BERT in the machine reading comprehension framework, Journal of Biomedical Informatics, Volume 118,2021,