# Unlocking the Power of Deep PICO Extraction: Step-wise Medical NER Identification

Tengteng Zhang<sup>1</sup>, Yiqin Yu<sup>2</sup>, Jing Mei<sup>2</sup>, Zefang Tang<sup>2</sup>, Xiang Zhang<sup>1</sup>, Shaochun Li<sup>2</sup>

<sup>1</sup>Southeast University, Nanjing, China; <sup>2</sup>IBM Research, Beijing, China

#### Abstract

The PICO framework (Population, Intervention, Comparison, and Outcome) is usually used to formulate evidence in the medical domain. The major task of PICO extraction is to extract sentences from medical literature and classify them into each class. However, in most circumstances, there will be more than one evidences in an extracted sentence even it has been categorized to a certain class. In order to address this problem, we propose a step-wise disease Named Entity Recognition (DNER) extraction and PICO identification method. With our method, sentences in paper title and abstract are first classified into different classes of PICO, and medical entities are then identified and classified into P and O. Different kinds of deep learning frameworks are used and experimental results show that our method will achieve high performance and fine-grained extraction results comparing with conventional PICO extraction works.

# Introduction

Evidence-based medicine (EBM) is a dominated approach enabling clinical practitioners to utilize the best available information (evidence) in making decision during clinical practice. When seeking evidences in medical literature, the PICO framework (Schardt et al. 2007) is usually used to formulate evidence in the medical domain. PICO represents four elements: Patient/Problem (P), Intervention (I), Comparison (C) and Outcome (O). Unfortunately, the PICO elements are usually not clearly identified in the text contents of medical literature, which need clinical practitioners to read and extract the PICO manually (Xia et al. 2019).

The major task of automatic PICO extraction is to extract sentences from medical literature and classify them into each element (class). Most of the previous methods (Jin and Szolovits 2018) identify PICO elements in abstracts and do classification at the sentence level. Unfortunately in most circumstances, there will be more than one PICO elements in an extracted sentence even it has been categorized to one class with the highest probability. So the fine-grained entity level extraction for PICO elements other than sentence level PICO elements is required.

In this paper, we propose a fine-grained medical Named Entity Recognition (NER) extraction and PICO identification method. With our method, medical NERs are identified and categorized into different classes of PICO. There are a lot of works using natural language processing (NLP) algorithms to extract medical NERs and classify them into different classes. Many deep neural network architectures have been proposed and used, among which the Conditional random field layer combined with the model of bidirectional Long Short-Term Memory (BiLSTM-CRF) produces encouraging results on medical NER extraction (Huang, Xu, and Yu 2015). Conditional Random Field (CRF) is a graph model that is often used to solve sequence labeling problems which focus on sentence level positions. Compared to the softmax classifier, the CRF can capture the context information of the target label. Further, the Convolutional Neural Networks (CNNs) and LSTM networks are applied to encode character-level word embedding. In addition, there are many public knowledge bases in the medical field which are used to enhance the deep learning frameworks. Previous works have compared the performance of these models. However, there are still a lack of validation on classifying the extracted medical NERs into PICO classes.

This paper tries to leverage the dominated deep learning technologies such as Bi-LSTM, CNN and BERT (Devlin et al. 2018; Alsentzer et al. 2019) and properly position them into different steps to streamline the overall entity level PICO identification task. Structured information in medical knowledge graph such as UMLS are also used to enhance the performance. The main contribution of this paper are as follows:

• We propose a step-wise work flow to fulfill the disease entity classification task, which carefully position different deep learning models into different steps. It makes us to improve the performance step-by-step and enhance the explainability of the whole system.

• We introduce mainstream deep learning models for the sentence-level PICO classification task and the disease entity recognition task. Structured information of medical domain knowledge via graph embedding is integrated into the deep learning models. Language rules are also combined with deep methods to provide adjusted entity-level classi-

Copyright © 2020, Disease Computational Modeling (DCM) at the International Joint Conference on Artificial Intelligence. All rights reserved.

fication results.

• A series of experiments for each step are implemented to show step-wise results. A web application is developed to enable the interaction between these results and the end users. Results on a set of risk model papers are also discussed.

The method is evaluated on clinical predictive modeling papers (Steverberg and others 2009), which is a major category of medical literature. In these papers, risk models are proposed as an explicit, empirical approach to estimate probabilities of disease(s) or outcome(s). As a powerful tool of disease prevention, risk models are widely used in chronic disease management. Annually, about two-thirds of total deaths around the world are caused by chronic diseases (Oh, Stefani, and Kim 2014). Half of them are cardiovascular diseases. Risk models are essential for physicians to estimate individual patient's risk or discover leading risk factors. However, there are several limitations among risk models, such as low sensitivity, inadequate factors in model(Elosua 2014) or not robust across domains. To automatically extract medical entities and classify them into PICO will be the first step to evaluate and utilize these risk models. To be specific, we address the disease NER extraction and classification in this paper.

#### Method

The popular end-to-end approach for PICO entity extraction is mainly in black-box which reduces the explainability of the model. To address this issue, we divide the complex system into different components. The main work flow of our method is illustrated in Figure 1. Text content of title and abstract of risk model papers are used as inputs, where the papers are retrieved by our previous work (Yu et al. 2019). In a nutshell, we provide a step-wise work flow to extract and classify disease entities in risk model paper title and abstract into PICO classes. Sentences in title and abstract are firstly classified into PICO classes by the PICO Classification Model. Following with the DNER Model, disease entities are then extracted from the sentences. Finally, the Disease Entity Mapping Model adjust and re-classify the disease entities into P and O. There are several advantages for this framework. First, the entity classification task is divided into different tasks, e.g., sentence-level and entitylevel, which makes it possible to leverage different kinds of deep learning models and improve the performance step-bystep. Second, since human interaction can not be avoided, sophistically dividing the task into several steps with meaningful inputs and outputs will improve the explainability of the whole system. In our paper, a web application is also shown to demonstrate how the results of each step are reported to the end user, and how these results can be modified via the user interface.

# **PICO Classification Model**

The PICO classification model is generally a sentence classification model. Its input is the title and abstract of each paper. The output is sentences along with the PICO classes including P, I/C, O, and N. Label N means that this sentence

doesn't belong to any class of PICO. Like many other works (Wallace et al. 2016; Jin and Szolovits 2018), we incorporate the element C into the element I since the "comparison" usually refers to a kind of "intervention" in most literature. The sentence classification example is shown in Figure 2.

Traditional sentence classification models are often based on machine learning techniques, which need complex feature engineering and have poor transferability(Pan and Yang 2009). The popular deep learning models have changed this situation, and excellent models combined with large highquality data have generated amazing results (Kim 2014). One of the major shortcomings of other paper works is the lack of high quality annotated dataset about PICO sentences. The dataset is either not public or coarse-grained labelled. In order to solve this problem, we invite experts to label 500 papers which will be described in experiment part. Two classification models are implemented: CNN (Convolutional Neural Network) (Kim 2014; Yang et al. 2016)and Bi-LSTM (Bidirectional Long Short-Term Memory)(Graves and Schmidhuber 2005). The most classic model for text classification by CNN is the work of kim(Kim 2014). Its input layer is a word vector expression obtained by passing the public corpus through the pre-trained word2vec(Mikolov et al. 2013) model, and the output layer is the sentence classification label probability. The input word sequence can be transformed into a vector by concatenating the corresponding word vectors from the embedding matrix. The word embedding we use is trained under large-scale medical domain texts. It is more suitable for medical related tasks than traditional open world embedding(Chiu et al. 2016). The CNN model has been proved to be effective in many natural language processing tasks(Yin et al. 2017).

RNN (Recurrent Neural Network) is another popular deep learning model and has achieved great success and been widely applied in many NLP (Cho et al. 2014) tasks. The purpose of RNNs is to process sequence data. In the traditional neural network model, the layers are fully connected between the input layer and the hidden layer as well as the hidden layer and the output layer. The nodes between each layer are disconnected. But this common neural network is powerless for many problems. For example, if you want to predict what the next word is in a sentence, you usually need to use the previous word, because the words in a sentence are not independent. The LSTM (Hochreiter and Schmidhuber 1997) is a special RNN model proposed to solve the problem of gradient vanishing. But there is a problem with modeling sentences using LSTM: it is impossible to encode information from back to front. Bi-LSTM(Schuster and Paliwal 1997) solves this problem and can better capture two-way semantic dependencies. Follow this architecture, The dynamic of a LSTM cell is controlled by an input vector  $(x_t)$ , a forget gate  $(f_t)$ , an input gate  $(i_t)$ , an output gate  $(o_t)$ , a cell state  $(c_t)$ , and a hidden state  $(h_t)$ , which are computed as:

$$i_t = \sigma(W_i * [h_{t-1}, x_t] + b_i)$$
 (1)

$$f_t = \sigma(W_f * [h_{t-1}, x_t] + b_f)$$
(2)

$$o_t = \sigma(W_o * [h_{t-1} + b_o]) \tag{3}$$



Figure 1: The system work flow. There are four main modules: PICO classification model, diseases named entity recognition model, mapping model and visualization system.

$$g_t = \tanh(W_g * [h_{t-1}, x_t] + b_g)$$
 (4)

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \tag{5}$$

$$h_t = o_t \odot \tanh\left(c_t\right) \tag{6}$$

where  $c_{t_1}$  and  $h_{t-1}$  are the cell state and hidden state respectively from previous time step,  $\sigma$  is the sigmoid function  $(\frac{1}{1+e^{-x}})$ , tanh is the hyperbolic tangent function  $(\frac{e^x - e^{-x}}{e^x + e^{-x}})$ . The Bi-LSTM neural network model can automatically capture deep semantic features and classify sentences.

We compare the sentence level PICO classification performance of different models. The input to the model is sentences consisting of the title and abstract of a risk model paper. Through the PICO model, the output is a tagged sentence of four labels: P, I/C, O, N, as shown in Figure 2.

## **Diseases Named Entity Recognition Model**

The main task of the Diseases Named Entity Recognition (DNER) model is to identify and extract disease-related entities in the text. Its input are the labeled sentences produced by the PICO classification model, and the output are the disease entities in the sentences. As we observed, most of disease entities are contained in P and O sentences. To simplify the problem, in this paper we only address the DNER task with sentences labeled with P and O. The extraction task can be converted to a sequence labeling problem(Huang, Xu, and Yu 2015) by assigning the annotated entities with appropriate tag representations. In order to be more versatile, we use the standard "BIO" schema, in which each word is assigned to a label as following: B = beginning of an entity, I = inside an entity, and O = outside of an entity, and the training data are the NCBI-BIO dataset (Doğan, Leaman, and Lu 2014). The NCBI disease corpus is fully annotated at the concept level to serve as a research resource for the biomedical natural language processing community. The detailed statistical information is shown in Table 1.

In order to get better performance, our paper proposes a method to improve the model by encoding structured information of medical knowledge bases. Specifically, plain text contains a large amount of unstructured information, and the traditional model has low utilization rate. In the medical domain, there are many public structured medical knowledge bases, such as Unified Medical Language System, also known as UMLS(Bodenreider 2004), DrugBank(Wishart et al. 2006), etc. The knowledge base can be represented into a graph. The nodes in the graph are entities, and the edges between nodes are the relationships among entities. For example, the semantic relationship between different medical entities is included in UMLS, which promotes the performance of the model.

We use the graph representation learning model, Deep-Walk(Perozzi, Al-Rfou, and Skiena 2014), to encode the nodes in the medical knowledge base into low-dimensional space vectors. The DeepWalk draws on the idea of the famous word embedding model word2vec in nature language processing. The basic processing element of word embedding is single word, and corresponding to the representation of network is graph node. Word embedding analyzes the implicit information of word sequence of a sentence, graph embedding captures structured information corresponding to node in the knowledge base. In a similar way, DeepWalk encodes the list of nodes in a random walk of a knowledge base. The so-called random walk can be understood in this way: It repeatedly and randomly selects the node of walking path on the network, and finally forms various of fixedlength node paths. Starting from a random node first, each step of the walk is randomly selected from the edge connected to the current node. Then moving along the selected edge to the next vertex, and the process is repeated. Given a graph G = (V, E), V is the node set and E is the edge set.  $w_i = v_o, v_1, \ldots, v_n$  is a sequence of several nodes, where  $v \in V$ . v is a node and can not be calculated, so a



Figure 2: The PICO sentences classification model. Its input is the sentences of paper and outputs are sentences with labels.

Table 1: NCBI dataset statistics.							
	number of sentences	avg sentence length	number of tokens	number of unique tokens	number of annotations		
Train-NCBI	5,576	23	132,584	9,805	2,911		
Valid-NCBI	918	25	23,456	3,580	487		
Test-NCBI	941	25	24,019	3,679	535		

mapping function  $\Phi$  is introduced. The mapping function  $\Phi$ maps each node in the network into a d-dimensional vector.  $\Phi$  is actually a matrix with a total of  $|V| \times d$  hyperparameters which need to be updated and optimized. So the original optimization goal changed from formula 7 to formula 8. This objective function can be trained and optimized with the skip-gram model.

$$P_r(v_n | v_0, v_1, \dots, v_{n-1}) \tag{7}$$

$$P_r(v_n | \Phi(v_0), \Phi(v_1), \dots, \Phi(v_{n-1}))$$
 (8)

UMLS has nearly 3.6 million nodes and 80 million edges. In order to save training time, we filter out other languages and only keep the English version. The parameters of the training are updated appropriately on the basis of the original DeepWalk paper. The length of randomly walk list is 32, the number of nodes sequence is 10 and the word2vec window size is 5. The training takes 2 days on a computer with 128G RAM and 32 core CPU. We use the downstream task, NER, to evaluate the quality of the entity embedding, which in turn tunes the training of the DeepWalk. After that we get the medical knowledge base entity embedding and we use the knowledge graph information the same as the way of word embedding. For comparison, we implement different models, such as Bi-LSTM, Bi-LSTM + CRF(Huang, Xu, and Yu 2015) and conduct full experiments to prove the effectiveness of the methods. For the credibility of the experimental results, the hyper-parameters are consistent with relevant papers.

### Mapping Model

The mapping model evaluates and adjusts the negative results from the previous steps. Intuitively the following issues are highlighted: First, DNER based on P-labeled or Olabeled sentences would get empty result. The solution of this issue is to look for the results in the I/C-labeled and Nlabeled sentences. Second, there is an intersection between the recognition results of P and O which we should keep only one if an entity appear in both of P and O at the same time. The solution is to leverage the "soft" label of the PICO sentences classification model which are the probabilities of corresponding classes on each sentence. As we mentioned previously, a sentence might contain more than two PICO elements. For example, the first half of sentence belongs to P label while the latter part is belonging to O label. The PICO classification model cannot solve this problem. In order to solve this issue, we incorporate the method of linguistic rules. The final score calculation formula for disease entity is as follows:

$$Score = \lambda f(d_i) + (1 - \lambda)g(d_i) \tag{9}$$

where  $d_i$  represents a disease entity which has two scores  $d_i = (s_1, s_2), i$  is the index of disease entity.  $s_1$  represents the probability that the entity belongs to P label, and  $s_2$  represents the probability that the entity belongs to O label.  $\lambda$  is an adjustable specific coefficient.

This probability is generated from the aforementioned PICO classification model.  $g(d_i)$  stands for the regular expression rule model. We manually formulate some typical linguistic rules, which are deterministic for the segmentation of sentence, e.g., if they belong to P or O. If the disease entity  $d_i$  is covered by the rule, then the value of  $g(d_i)$  is set to 1, otherwise 0. The results of g(s) are (1,0), (0,1) and (0,0), respectively. Re-dividing the entity associated label set based on the score can improve the overall performance and reduce the effects caused by the error accumulation of the previous models. In order to verify the method, we invite several medical experts to help us label 100 articles. We compare the results of adding rules with no rules and calculate the recall rate.

# Results

In order to fully demonstrate the effectiveness of our method, we conducted complete experimental evaluation for each module, and obtained convincing experimental results. Precision, Recall and F1 score are used as indicators which are the most commonly used measurements in the field of machine learning and deep learning.

In the area of information retrieval, true positive (TP) indicates the number of labeled instances the method find. False positive (FP) indicates the number of unlabeled instances that the method find. And false negative (FN) indicates the number of labeled instances that the method doesn't find. Then precision, recall and F1 score are calculated as following:

$$Precision = \frac{TP}{TP + FP} \tag{10}$$

$$Recall = \frac{TP}{TP + FN} \tag{11}$$

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall}$$
(12)

## **Results of PICO Classification Model**

The 500 papers for the PICO classification model are downloaded from PubMed sorted by time. Two experts are invited to label the same batch of papers at the same time independently. If there is a disagreement, it will be confirmed by the third expert. We finally get our dataset and divide it into train set and validation set with a 7:3 ratio. In order to ensure the adequacy of the experiment, two classification models are implemented: CNN (Kim 2014) and Bi-LSTM (Graves and Schmidhuber 2005). However, the scores of the Bi-LSTM model are much higher than the CNN model's, so here we only show the result of the Bi-LSTM model as in Table 2. We compare our classification results with Jin et al.(Jin and Szolovits 2018), which extracted 489k medical abstracts from MEDLINE(https://www.medline.com/). Yuan et al. (Xia et al. 2019) proposed a soft-margin SVM classification model which merges other features that combines 1-2 gram analysis with TF-IDF(Pan and Yang 2009) method. We also compare with them.

### **Results of DNER Model**

The performance of the DNER model has a major impact on the final results. We compare the performance of different models and the same model under different hyperparameters. Most of the hyper-parameters follow Reimers (Reimers and Gurevych 2017): dropout=0.25, lstm hidden layer size = 100, optimizer=adam, char cnn embedding size=30, char LSTM embedding size=25, and mini batch size=32. In addition, in order to improve performance, we introduce structured information in the medical knowledge base via graph representation model training and prove its effectiveness from experimental results as shown in Table 3. We also use the BERT model which is the dominated model in various fields for medical entity recognition tasks. The experimental code refers to BERTpytorch (https://github.com/codertimo/BERT-pytorch). Two different versions of BERT are implemented: BERT-base and BERT-large. With fewer parameters, BERT-base has faster training speed. but BERT-large achieves the best performance over all the other models with precision 0.8507, recall 0.8844 and F1 score 0.8672.

## **Results of Mapping Model**

The mapping model is used to perform further processing of the DNER model and obtain the disease entity classification results, so that the clinical practitioners can quickly get the key evidence with less time. Several clinical researchers read series of medical papers of risk models and summarize a set of linguistic rules such as "... risk of < outcome >". If a sentence matches this rule, the entity that appears in the later part is extracted and classified as "O". 100 papers are labeled by the clinical researchers following the same protocol with the PICO classification labeling work. The results are shown in Table 4.

## Web Application

A web application is developed to enable the method and provide interface (see Figure 3) for user to interact with the results. When the user uploads a paper, the back-end system will call the service API on the server to parse the text in title and abstract, and display the PICO classification result and the corresponding disease entity recognition result. At the same time, to further improve the accuracy of the system, we also provide labeling and modification buttons. If there is an error in the recognition result, the user can manually modify the result and store it in the back-end database. If the system collects very few samples, it will use the language rules and remember it to avoid mistakes next time. When samples are accumulated to a certain extent, the back-end system will retrain the model to improve the performance of the model. As in Figure 3, the title and abstract of the risk model paper are displayed on the left side of the web page, the panel in the middle shows the results of the sentencelevel PICO classification. The right panel shows the final disease entity mapping results in P and O respectively.

## Discussion

The methods mentioned earlier provide a step-wise framework for automatically extracting disease entities from risk model papers based on the PICO framework. The following sections analyze and discuss some details of the experimental results for each module.

	Р		I/C			0			
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Jin2018(Jin and Szolovits 2018)	0.885	0.828	0.856	0.749	0.815	0.781	0.845	0.832	0.838
Yuan2019 (Xia et al. 2019)	0.925	0.838	0.879	0.842	0.789	0.814	0.886	0.897	0.891
Our Model (Bi-LSTM)	0.947	0.867	0.906	0.773	0.850	0.810	0.867	0.918	0.891

Table 2: Comparison of PICO classification models.

# Table 3: Result of different DNER models.

model	Precision	Recall	F1
LSTM+SOFTMAX	0.7910	0.8125	0.8016
LSTM+CRF	0.8179	0.8375	0.8275
LSTM+CRF+CNN CHAR	0.8298	0.8333	0.8316
LSTM+CRF+LSTM CHAR	0.8283	0.8697	0.8485
LSTM+CRF+CNN CHAR+GRAPH EMBEDDING	0.8433	0.8583	0.8508
LSTM+CRF+LSTM CHAR+GRAPH EMBEDDING	0.8348	0.8635	0.8489
BERT BASE	0.8200	0.8729	0.8456
BERT LARGE	0.8507	0.8844	0.8672
Zhai, Zenan 2018	0.8213	0.8366	0.8289
TaggerOne (Leaman and Lu, 2016)	0.852	0.802	0.826
Dnorm (Leaman et al., 2013)	0.820	0.795	0.807

AUTO PICO

My Accounts 💙 Contact Updates

Hi Jack! 06/16/2019		
t mark PICO components		
. Upload your text file:	2. PICO classification	3. DNER
File address: /Users/Jack/ner_27504006.csv	Population/Patient: Ornostatic hypotension (NH) is associated with hypertension and diabetes mellitus.	Population/Patient: Push to Outcome
Upload		Orthostatic hypotension OH hypertension diabetes
File content: Orthostatic hypotension (OH) is associated with hypotension and dialsets mellitus, isoverus, in mellitus, its prevaions, the effect of intensive versus standard systelic blood pressure (DP) targets on incident participants in the ACCOPID Action to Control Cardiovascular Risk in Diabetes) BP trial, seated BP was measured 3-% followed by reading very minute for a standard systelic pand disatolic BV was assessed and standard systelic pand disatolic BV was assessed and between OH and clinical successful and tertificants in the distribution of the advection of the tertificants and the distribution of the advection of the tertificant and the distribution of the advection of the tertificant of the advection of the advection of the tertificant of the distribution of the advection of the tertificant of the advection of the advection of the advection of the tertificant of the advection of the advection of the advection of the tertificant of the advection of the advection of the advection of the tertificant of the advection of the advecti	Outcome: - The relationship between 0x4 and clinical obtomse Infraction, nonfratal artosk, heart fallur hospitalization or death and the primary composite outcome of nonfratal participationscular deathy was assessed unique gropperinonal hazarda analysis. - associated with increased risk of the non- balantial and the second risk of the second second second participation of the second second second second second and the second second second second second second and second	Outcome: Push to Patient cardiovascular death death bospitalization heart failure myocardial infarction
Ihe primary composite outcome of nonfratal myocardial infraction, nonfratal stroke, and confronzeouslar death) Consensus OH, defined by orthostatic decline in systolic BP ±20 mm 169 of taskinoli BP ±10 mm 149, occurred task BP ±10 mm 149 of taskinoli BP ±10 mm 149, occurred total tandard, ±104 mm 149 user anisate to OH incidence. Over a median follow up of 4.6 mmonthy. OH was tandard, ±104 confidence interval. 11.12.30 and heart failure death/hospitalization (hazard ratio, 18, 95% outcome) or thom prespective outcomes. In patients with who may be the strong task of t	Intervention/Comparison: - Orthostatic BP change, calculated as the minimum tanding minus the mean seated systel (DP and disatolic BP was assessed at baseline, 12 months, and - Nother age nor systel(DB P treatment target (interactive, -120 mm Hg versus standard, +140 mm Hg) was related to OH incidence.	Save the modification
type 2 diabetes meltitus and hypertension, OH was common, not associated with intensive versus atmanded and heart failure events.	None: - Orthoastic Hypotension in the ACCORD (Action to Control Cardiovascular Risk in Diabetes) Blood Pressure That Prevalence, inclusions, and Porgonatics Biginflance. Phypertension, OH was common, not associated with Intensive versus atandard BP treatment polas. Add - In 4366 participants in the ACCORD (Action to Control Cardiovascular Risk in Diabetering IPT trait action glack diabetering - In 4366 participants in the ACCORD (Action to Control Cardiovascular Risk in Diabetering IPT trait general general participants) - In 4366 participants in the ACCORD (Action to Control Cardiovascular Risk in Diabetering IPT trait general participants)	

Figure 3: The visualization System.

In table 2, we compare the PICO classification results of different models based on Bi-LSTM architecture. Our model achieves the best performance on the P label sentences classification comparing with the state-of-the-art models as well as some indicators of I/C and O labels, which proves the effectiveness of our model. In addition to the performance comparison, we also analyze the negative samples. Statistics show that in some cases, due to the frequent occurrence

Table 4: Result of Mapping Model with different method.

Model	Recall		
WIOdel		0	
Bi-LSTM+CRF+CNN CHAR+GRAPH EMBEDDING	0.7262	0.60	
BERT LARGE	0.7588	0.6293	
BERT LARGE+ LINGUISTIC RULE	0.7911	0.7474	

of some typical words, the sentence of P label is difficult to distinguish from I/C, the same as O and I/C. In addition, in some sentences, the part of one sentence belongs to P label, however, the other part belongs to O label. This increases the difficulty of classification. We compare with the results of two recent related papers, and we can clearly see that our results in most indicators are better.

In the DNER module, we experiment with leading deep models on NER task (Table 3). For models with Bi-LSTM, our model Bi-LSTM + CRF + CNN char + graph embedding improves the baseline of F1 score from 80.16% to 85.08% and precision score from 79.10% to 84.33%. Compared to current state-of-the-art models, we achieve significant improvements.

The BERT model outperforms all other models on all the measurements (the precision score is 85.07%, the recall score is 88.44%, the F1 score is 86.72%). What we can not ignore is that the model is too complex and the calculate process is very slow, so each time the parameter adjustment takes a lot of time.

For the mapping model, since some sentences may have key elements of P and O at the same time, classification is difficult and has a great impact on the results of DNER. The mapping model uses a set of statistical rules and language rules to alleviate this problem and improve the overall performance. In table 4, we only give results about recall because we care more about how many entities are found which we need. The precision actually is low due to following issues: (1) there are still errors from the PICO classification model and the DNER model which have not been well addressed, and (2) there's lack of public standard labeling strategies. Note that the manual work of labeling disease entities with PICO elements is a non-trivial work for most of the clinical researchers even we restrict the papers with risk model type. During the experiment, we observe that a good deal of labels done by the first two researchers are mismatched. That also increase the work of the third researchers to further do a decision. By considering this situation and less number of labeled samples, we finally decide to use rule-based method in the mapping model. Specifically for risk model disease entity extraction and classification, this is the first try to do such experiments and can serve as a benchmark for future research.

One advantage of step-wise approach is that the results of each step can be output and presented to the user. In Figure 3, we show an example with the medical paper (Fleg et al. 2016) which is randomly retrieved from PubMed. The middle of the web page is the classification result of PICO elements and the DNER recognition results are shown on the right panel. It can be found that the classification result of PICO sentences is reasonable, but the results of DNER have some errors. For example, the item "hospitalization" should not belong to the disease entity caused by the misclassification of the DNER model. At this point, the user can continuously improve the performance of the model by deleting this unreasonable word and saving the result to the back-end system.

# Conclusion

In this paper, we propose a step-wise method for extracting medical entities based on the PICO framework. Main steps include PICO sentence classification, disease entity recognition and disease mapping. Mainstream deep neural networks such as CNN and Bi-LSTM are used to classify sentences into PICO elements. The disease entity recognition is based on dominated deep learning frameworks. Structural medical knowledge is embedded with probabilistic knowledge mapping model. Experimental results show that our method achieves reasonable performance on precision, recall and F1 score. An online web system is also developed to facilitate clinical researchers to do medical entity based PICO extraction.

### References

Alsentzer, E.; Murphy, J.; Boag, W.; Weng, W.-H.; Jindi, D.; Naumann, T.; and McDermott, M. 2019. Publicly available clinical BERT embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, 72–78. Minneapolis, Minnesota, USA: Association for Computational Linguistics.

Bodenreider, O. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research* 32(suppl\_1):D267–D270.

Chiu, B.; Crichton, G.; Korhonen, A.; and Pyysalo, S. 2016. How to train good word embeddings for biomedical NLP. In *Proceedings of the 15th Workshop on Biomedical Natural Language Processing*, 166–174. Association for Computational Linguistics.

Cho, K.; van Merrienboer, B.; Gulcehre, C.; Bougares, F.; Schwenk, H.; and Bengio, Y. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)*.

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Doğan, R. I.; Leaman, R.; and Lu, Z. 2014. Ncbi disease

corpus: a resource for disease name recognition and concept normalization. *Journal of biomedical informatics* 47:1–10.

Elosua, R. 2014. Cardiovascular risk functions: usefulness and limitations. *Revista Española de Cardiología (English Edition)* 67(2):77–79.

Fleg, J. L.; Evans, G. W.; Margolis, K. L.; Barzilay, J.; Basile, J. N.; Bigger, J. T.; Cutler, J. A.; Grimm, R.; Pedley, C.; Peterson, K.; et al. 2016. Orthostatic hypotension in the accord (action to control cardiovascular risk in diabetes) blood pressure trial: prevalence, incidence, and prognostic significance. *Hypertension* 68(4):888–895.

Graves, A., and Schmidhuber, J. 2005. Framewise phoneme classification with bidirectional LSTM networks. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 4, 2047–2052. IEEE.

Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.

Huang, Z.; Xu, W.; and Yu, K. 2015. Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*.

Jin, D., and Szolovits, P. 2018. PICO element detection in medical text via long short-term memory neural networks. In *Proceedings of the BioNLP 2018 workshop*, 67–75. Association for Computational Linguistics.

Kim, Y. 2014. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1746–1751. Association for Computational Linguistics.

Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G. S.; and Dean, J. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, 3111–3119.

Oh, S. M.; Stefani, K. M.; and Kim, H. C. 2014. Development and application of chronic disease risk prediction models. *Yonsei medical journal* 55(4):853–860.

Pan, S. J., and Yang, Q. 2009. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* 22(10):1345–1359.

Perozzi, B.; Al-Rfou, R.; and Skiena, S. 2014. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 701–710. ACM.

Reimers, N., and Gurevych, I. 2017. Reporting score distributions makes a difference: Performance study of lstm-networks for sequence tagging. *arXiv preprint arXiv:1707.09861*.

Schardt, C.; Adams, M. B.; Owens, T.; Keitz, S.; and Fontelo, P. 2007. Utilization of the pico framework to improve searching pubmed for clinical questions. *BMC medical informatics and decision making* 7(1):16.

Schuster, M., and Paliwal, K. K. 1997. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing* 45(11):2673–2681.

Steyerberg, E. W., et al. 2009. *Clinical prediction models*, volume 381. Springer.

Wallace, B. C.; Kuiper, J.; Sharma, A.; Zhu, M.; and Marshall, I. J. 2016. Extracting pico sentences from clinical trial reports using supervised distant supervision. *The Journal of Machine Learning Research* 17(1):4572–4596.

Wishart, D. S.; Knox, C.; Guo, A. C.; Shrivastava, S.; Hassanali, M.; Stothard, P.; Chang, Z.; and Woolsey, J. 2006. Drugbank: a comprehensive resource for in silico drug discovery and exploration. *Nucleic acids research* 34(suppl\_1):D668–D672.

Xia, Y.; Liao, X.; Li, S.; Shi, Q.; Wu, J.; and Li, K. 2019. Extracting PICO elements from RCT abstracts using 1-2gram analysis and multitask classification. *CoRR* abs/1901.08351.

Yang, Z.; Yang, D.; Dyer, C.; He, X.; Smola, A.; and Hovy, E. 2016. Hierarchical Attention Networks for Document Classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 1480– 1489. San Diego, California: Association for Computational Linguistics.

Yin, W.; Kann, K.; Yu, M.; and Schütze, H. 2017. Comparative study of cnn and rnn for natural language processing. abs/1901.08351.

Yu, Y.; Xu, E.; Xia, E.; Huang, H.; Hao, B.; and Zhang, S. 2019. A method to accelerate and visualize iterative clinical paper searching. *MedInfo* 2019.