

# Integrating Dictionary Feature into A Deep Learning Model for Disease Named Entity Recognition

Hamada A. Nayel  
Department of Computer Science  
Benha University, Egypt  
hamada.ali@fci.bu.edu.eg

H. L. Shashirekha  
Department of Computer Science  
Mangalore University, India  
hlsrekha@gmail.com

## Abstract

In recent years, Deep Learning (DL) models are becoming important due to their demonstrated success at overcoming complex learning problems. DL models have been applied effectively for different Natural Language Processing (NLP) tasks such as part-of-Speech (PoS) tagging and Machine Translation (MT). Disease Named Entity Recognition (Disease-NER) is a crucial task which aims at extracting disease Named Entities (NEs) from text. In this paper, a DL model for Disease-NER using dictionary information is proposed and evaluated on National Center for Biotechnology Information (NCBI) disease corpus and BC5CDR dataset. Word embeddings trained over general domain texts as well as biomedical texts have been used to represent input to the proposed model. This study also compares two different Segment Representation (SR) schemes, namely IOB2 and IOBES for Disease-NER. The results illustrate that using dictionary information, pre-trained word embeddings, character embeddings and CRF with global score improves the performance of Disease-NER system.

## 1 Introduction

Disease is a principle Biomedical Named Entity (BioNE), which has got attention by biomedical research due to increase in research in health and the impact of disease on public life. Disease-NER is a challenging problem due to multiple challenges such as ambiguity (*same word or phrase refers to different entities*), synonyms (*an entity can be denoted by various names in a synonym relation*), multi-word NEs (most of disease NEs consist of multiple words) and nested NEs (*one NE may occur within a longer NE*). Further, abbreviations which are used frequently in biomedical literature are the main sources of ambiguity. For example, "AS" may refer to "Asperger Syndrome" or "Autism Spectrum" or "Aortic Stenosis" or "Ankylosing Spondylitis" as well as "Angleman Syndrome". In such cases to which entity an abbreviation refers to has to be resolved depending on the context. Figure 1 shows an example of an abstract with disease mentions highlighted. Different approaches have been used for Disease-NER, such as dictionary-based approach, rule-based approach, Machine Learning (ML) approach and hybrid approach [1][2]. Deep Learning (DL) is a big trend in ML, which promises powerful and fast ML algorithms moving closer to the performance of Artificial Intelligence (AI) systems. It is about learning multiple levels of representation and abstraction that help to make sense of any data such as images, sound, and text. Automatically learning features at multiple levels of abstraction allow a system to learn complex functions mapping the input to the output directly from data, without depending on human-crafted features. DL employs the multi-layer Artificial Neural Networks (ANN) for increasingly richer functionality. While the concept of classical ML is characterized as learning a model to make predictions based on past observations, DL approaches are characterized by learning to not only predict but also to correctly represent the data such that it is suitable for prediction. Figure 2 shows the difference between classical ML flow and DL flow.

Identification of APC2, a homologue of the adenomatous polyposis coli tumour suppressor. The adenomatous polyposis coli (APC) tumour-suppressor protein controls the Wnt signalling pathway by forming a complex with glycogen synthase kinase 3beta (GSK-3beta), axin/conductin and betacatenin. Complex formation induces the rapid degradation of betacatenin. In colon carcinoma cells, loss of APC leads to the accumulation of betacatenin in the nucleus, where it binds to and activates the Tcf-4 transcription factor. Here, we report the identification and genomic structure of APC homologues. Mammalian APC2, which closely resembles APC in overall domain structure, was functionally analyzed and shown to contain two SAMP domains, both of which are required for binding to conductin. Like APC, APC2 regulates the formation of active betacatenin-Tcf complexes, as demonstrated using transient transcriptional activation assays in APC -/- colon carcinoma cells. Human APC2 maps to chromosome 19p13.3. APC and APC2 may therefore have comparable functions in development and cancer.

Figure 1: An abstract with disease mentions highlighted

Given a large set of desired input-output mapping, DL approaches work by feeding the data into an ANN that produces consecutive transformations of the input until a final transformation predicts the output. These transformations are learnt from the given input-output mappings, such that each transformation makes it easier to relate the data to the desired label.

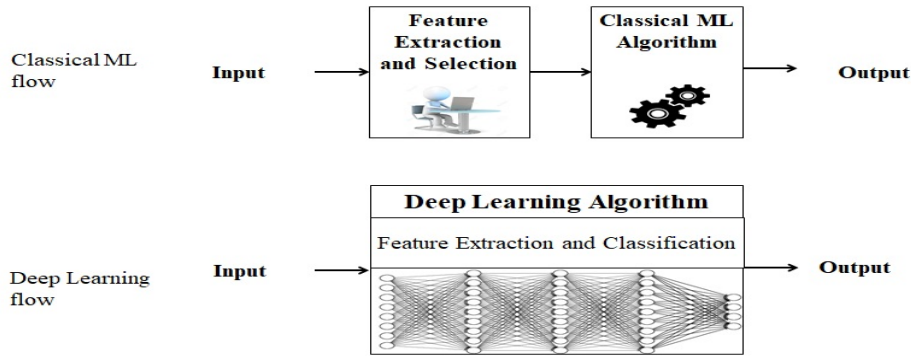


Figure 2: Classical and deep learning flows

Of late, DL algorithms based on ANNs [3, 4] are being used to a larger extent for various NLP tasks such as Biomedical Named Entity Recognition (BioNER) [5] relation extraction for biomedical texts [6] and biomedical event extraction [7]. In this paper, an efficient DL model using ANN for Disease-NER has been developed. NCBI and BC5CDR datasets are used for evaluation of the proposed model and the results are reported in terms of f1-measure.

## 2 Background

### 2.1 Artificial Neural Networks (ANN)

ANN are inspired by the mechanism of brain computation which consists of computational units called neurons. However, connections between ANN and the brain are in fact rather slim. In the metaphor, a neuron has scalar inputs with associated weights and outputs. The neuron multiplies each input by its weight, sums them and transforms to a working output through applying a non linear function called activation function. The structure of a biological neuron and an artificial neuron model with  $n$  inputs and one output is shown in Figure 3. In this example, a neuron receives simultaneous inputs  $X = (x_1, x_2, \dots, x_n)$  associated with weights  $W = (w_1, w_2, \dots, w_n)$ , a bias  $b$  and calculates the output as  $y = f(W \cdot X + b)$  where  $f$  is the activation function. ANN

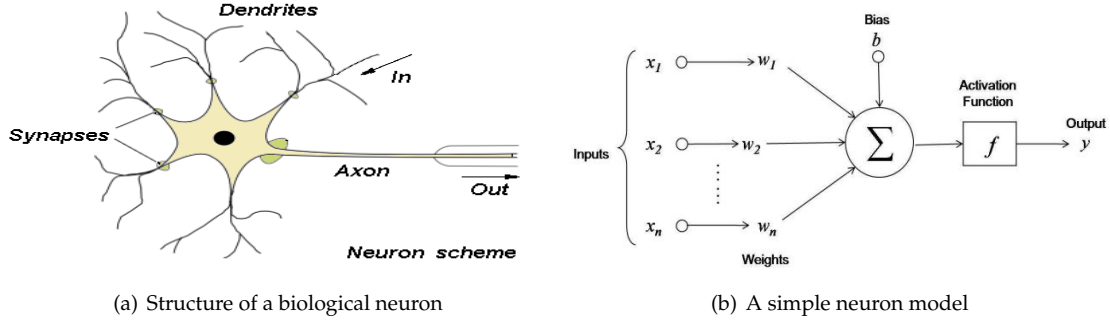


Figure 3: Structure of a biological neuron and a simple neuron model

comprises of a large number of neurons within different layers with each layer having a specific task. An ANN model basically consists of three layers: an input layer, one or more hidden layers and an output layer. Input layer contains a set of neurons called input nodes, which receive raw inputs directly. The hidden layers receive the data from the input nodes and responsible for processing these data by calculating the weights of neurons at each layer. These weights are called connection weights and passed from one node to another. Number of nodes in hidden layers influences the number of connections as well as computational complexity. During learning connection weights are adjusted to be able to predict the correct class label of the input. Using multiple hidden layers helps in detecting more features while learning the model. Output layer receives the processed data and uses its activation function to generate final output. This kind of ANN where information flows in one direction from input layer to output layer through one or more hidden layers is called feed-forward ANN. Figure 4 shows an example of a feed-forward ANN with two hidden layers. An ANN is called fully connected if each node in a layer is connected to all nodes in the subsequent layer.

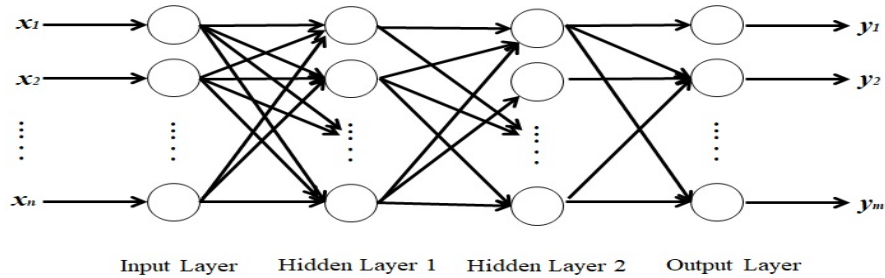


Figure 4: Structure of a simple feed-forward ANN

Recurrent Neural Networks (RNN) is a type of ANN in which hidden layer neurons has self-connections which means output depends not only on the present inputs but also on the previous neuron state. A simple form of RNN which contains an ANN with the previous set of hidden unit activations feeding back into the network along with the inputs is shown in Figure 5. The activations are updated at each time step  $t$  and a delay unit has been introduced to hold activations until they are processed at the next time step. The input vector  $x_0$  at time stamp  $t = 0$  processed using RNN structure is as follows:

$$h_t = f_W(h_{t-1}, x_t) \quad (1)$$

where,  $h_t$  is the output at time stamp  $t$ ,  $h_{t-1}$  is the output at time stamp  $t - 1$ ,  $f_W$  is an activation function with parameter  $W$  and  $x_t$  is the input vector at the time stamp  $t$ .

In case of long sequences, RNNs are biased towards their most recent inputs in the sequence due

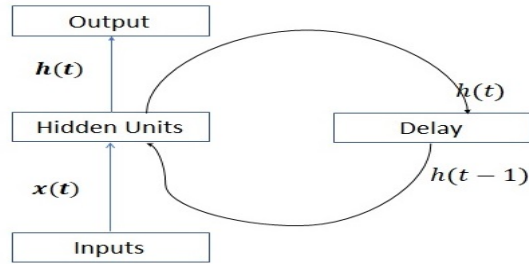


Figure 5: A Simple RNN Structure

to the gradient vanishing problem [8, 9]. While calculating weights in RNN for each time stamp  $t$ , the gradients of error get smaller and smaller as moving backward in the network and gradually vanish. Thus, the neuron in the earlier layers learns very slowly as compared to the neurons in the later layers. Earlier layers in the network are important as they are responsible to learn and detect patterns and are the building blocks of the RNN. Figure 6 shows an example of gradient vanishing problem. In this figure, the dark shade of the node indicates the sensitivity over time of the network nodes to the input at first time stamp. The sensitivity decreases exponentially over time as new inputs overwrite the activation of hidden unit and the network forgets the input at first time stamp. To overcome the gradient vanishing problem, S. Hochreiter and J. Schmidhuber [10] introduced a new RNN architecture called Long Short-Term Memory (LSTM).

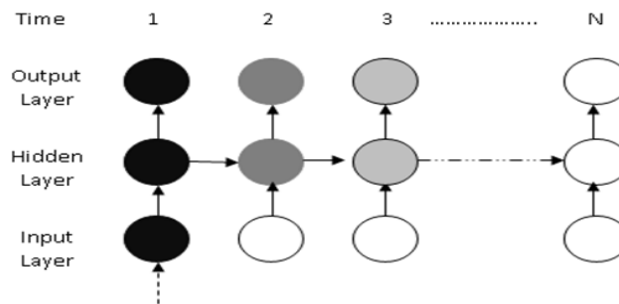


Figure 6: Gradient vanishing problem for RNN

### 2.1.1 Long Short-Term Memory

LSTM [10] is a kind of RNN which handles sequences of arbitrary length and is able to model dependencies between far apart sequence elements as well as consecutive elements. The LSTM architecture consists of a set of RNNs known as memory blocks. Each block contains self-connected memory cell ( $m_t$ ) and three multiplicative units namely input ( $i_t$ ), output ( $o_t$ ) and forget ( $f_t$ ) gates, that provide continuous peers of write, read and reset operations for the cells as shown in Figure 7. These gates regulate the information in memory cell and consists of a sigmoid function. The input gate regulates the proportion of history information that will be kept in memory cell and the output gate regulates the proportion of information stored in the memory cell which will influence other neurons. Forget gate can modify the memory cell by allowing the cell either to remember or forget its previous state. The complete details of LSTM architecture is described by S. Hochreiter and J. Schmidhuber [10].

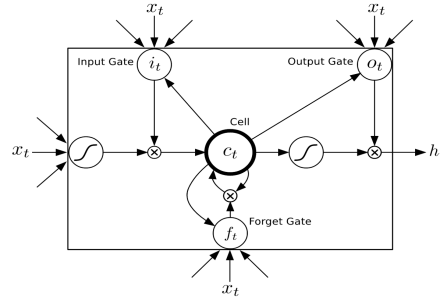


Figure 7: Structure of LSTM unit

In the sequence labeling problem, the typical input is a sequence of input vectors and the output is a sequence of output tags. NER can be considered as a sequence labeling problem where the sentence  $X$  consisting of words  $(w_1, w_2, \dots, w_n)$  is given as input and the required output is the sequence of tags  $T = (t_1, t_2, \dots, t_n)$  that represents the class labels of the words. LSTM is suitable to apply for NER as it can remember up to the first word in the sentence. However, one shortcoming of LSTM is that they process the input only in left context. But, in NER it is beneficial to have access to both left and right contexts as the output tag of a word depends on few previous and few next words (context window). This problem is overcome by a Bidirectional LSTM (BiLSTM) [11] where each sequence is presented in forward and backward direction to two separate hidden states to capture left and right context information respectively. Then the outputs of two hidden states are concatenated to form the final output as shown in Figure 8.

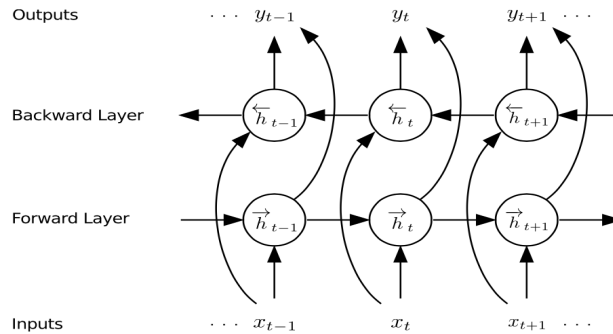


Figure 8: Structure of Bidirectional LSTM

To apply BiLSTM for NER, a sentence  $X = (w_1, w_2, \dots, w_n)$  has to be fed to a BiLSTM and each word has to be replaced by its vector representation. Applying BiLSTM structure on these representations outputs a vector  $Y = (t_1, t_2, \dots, t_n)$  representing the corresponding tags.

ANN in general can be applied for language understanding and NLP tasks such as speech recognition, MT and text summarization using the concept of language modeling (LM). LM is a probabilistic model that is able to predict the next word in the sequence given the words that precede it. The central component in ANN language modelling is the use of an embedding layer which maps discrete symbols (words, characters) to continuous vectors in a relatively low-dimensional space. Vector representation of words and documents has been a leading approach in information retrieval and computational semantics [12]. The primitive method of representing a document as a vector is the bag-of-words model. If  $m$  is the size of the vocabulary, a document is assigned

a  $m$ -dimensional vector with non-zero entries for words corresponding to the entries occurred in the document and zero entries for rest of the words. These are used to learn  $n$ -dimensional real-valued vector representations of words in  $\mathbb{R}^n$ . After the rise in popularity of ANNs, a new approach for representing lexical semantics has become the new default for novel models. These real-valued vectors are called "word embeddings" and have become a primary part of models for various applications, such as information retrieval [13], sentiment analysis [14], automatic summarization [15] and question answering [16].

## 2.2 Word Embeddings

Word embeddings is a distributed representation of words in a vector space that captures semantic and syntactic information for words [17]. The basic idea behind word embeddings is to use distributional similarity based representations by representing a word by means of its neighbors. Distributed representations of words help to enhance the performance of learning algorithms in various NLP tasks by grouping similar words. Mikolov et al. [18] introduced the skip-gram model and Continuous Bag of Words (CBOW) model for learning word embeddings from huge amounts of text data. An ANN structure has been used for learning word embeddings which encode many linguistic regularities and patterns explicitly. Some of these patterns can be represented as linear operations, e.g. the result of vector calculation:  $\vec{king} - \vec{man} + \vec{woman}$  is closer to  $\vec{queen}$  than to any other word vector as shown in Figure 9. CBOW model trains the word embed-

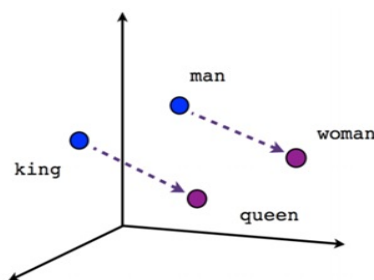


Figure 9: Example of word vector calculations

dings by predicting a word in a sentence using its surrounding words, while Skip-gram model trains the word embeddings by predicting the surrounding words for a given word in the input layer. The word embeddings  $e_t$  representing a word located at the  $i$ -th position in a sentence is calculated by maximizing the average log probability as follows:-

- In CBOW model

$$\frac{1}{T} \sum_{t=1}^T \log p(e_t | e_{t-\frac{n}{2}}, \dots, e_{t-1}, e_{t+1}, \dots, e_{t+\frac{n}{2}})$$

- In Skip-gram model

$$\frac{1}{T} \sum_{t=1}^T \log p(e_{t-\frac{n}{2}}, \dots, e_{t-1}, e_{t+1}, \dots, e_{t+\frac{n}{2}} | e_t)$$

where  $e_{t-\frac{n}{2}}, \dots, e_{t-1}$  are vectors for  $\frac{n}{2}$  preceding words,  $e_{t+1}, \dots, e_{t+\frac{n}{2}}$  are vectors for  $\frac{n}{2}$  subsequent words, and  $T$  is the number of tokens in the sentence.

**Character-level Embeddings** Word embeddings maintain semantic and syntactic information of words, but does not capture orthographical information at character level such as capitalization, numeric characters, special characters and hyphenation. However, such information plays

an important role in Disease-NER as disease names can be characterized by the existence of combination of these information. Character-level word representation has been introduced to represent orthographical and morphological information [19].

### 3 Related Works

Researchers have explored many approaches for Disease-NER. Robert Leaman et al. [20] have addressed Disease-NER by learning the similarities between the disease mentions and concept names. They evaluated their approach using NCBI dataset and achieved a f1-measure of 80.9%. This was the first work to use pairwise learning to rank Disease NER approach but the performance was not high. A multi-task learning approach applied to BioNER using Neural Network architecture by Gamal C. et al. [21] has achieved a f1-measure of 80.73% for Disease-NER using NCBI dataset. In addition to NCBI dataset, 15 biomedical corpora have been used to train the system. These corpora have joint articles with NCBI test set and this affects the model evaluation. Leaman and Lu [22] used semi-Markov models to build TaggerOne, a tool for BioNER which reported a high competitive f1-measure of 82.9% for NCBI dataset. Sunil and Ashish [23] designed a RNN model for Disease-NER using Convolution Neural Network for representing character embeddings and bidirectional LSTM for word embeddings. A pre-trained word embeddings trained over a corpus of PubMed articles have been used for training the model and it is not enough to represent words in general domain. They evaluated their system on NCBI dataset and achieved f1-measure of 79.13%. Wei et al. [24] developed an ensemble-based system for Disease-NER. At the base level, they built CRF with a rule-based post-processing system and Bi-RNN based system. At the top level, they used SVM classifier for combining the results of the base systems. The proposed system uses manually extracted features for SVM and CRF as well as hand-crafted rules which depends on human experts. BC5CDR corpus has been used for evaluation and the model reported a f1-measure of 78.04%.

BANNER [25], a tool which implements CRF algorithm for BioNER using general features such as orthographical, linguistic and syntactic dependency features has reported a f1-measure of 81.8% on NCBI dataset. Xu et al. [26] designed a system (CD-REST) for chemical-induced disease relations from biomedical texts. A CRF-based module for NER as first step is designed using character level, word level features, context features and distributed word representation features learned from external un-annotated corpus. They evaluated the system using BC5CDR dataset and reported a f1-measure of 84.43%. Zhao et al. [27] developed a system for Disease-NER based on convolutional neural network. The proposed system integrated with dictionary information and a post-processing module has been used for performance enhancements. NCBI and BC5CDR datasets have been used for evaluation of the proposed system and reported a f1-measure of 85.17% and 87.83% respectively. Haodi Li et al. [28] designed an end-to-end system used for chemical-disease relation extraction. The proposed system uses an ensemble approach using SVM and CRF as base classifiers with SVM as a meta-classifier to build a module for Disease-NER. BC5CDR dataset has been used to evaluate the system and reported a f1-measure of 86.93%. Hsin-Chun Lee et al. [29] presented an enhanced CRF-based system for Disease-NER. Rich feature set in addition to dictionary based features extracted from different lexicons have been used to train CRF. BC5CDR dataset has been used to evaluate has been used for system evaluation and reported f1-measure of 86.46%.

In this paper, a DL model has been designed for Disease-NER using BiLSTM for learning the model and CRF for decoding the results with the following objectives:

- Using dictionary information for each token
- Using pre-trained word embeddings trained over a huge corpus of texts from biomedical domain as well as generic domain

- Learning a BiLSTM model for character-level word representations instead of using hand engineered features

## 4 Methodology

General structure of the proposed model is given in Figure 10. Proposed model accepts a sequence of words, for example, the sentence “The risk of colorectal cancer was significantly high.” and the associated tags “(O, O, O, B-Disease, I-Disease, O, O, O, O)” as input and gives a vector representation for each word containing information about the word itself and the neighbouring words within a sentence, denoted as contextual representation. In addition to word embeddings and character-level embeddings, dictionary information for each token has been extracted from Merged Disease Vocabulary (MEDIC) [30]. MEDIC is a comprehensive and publicly available dictionary for disease entities which provides information including disease names, concept identifiers, definitions of diseases and synonyms. Dictionary information for each word is represented as a binary vector containing information about the existence of the token in the dictionary either solo or as a part of multi-word disease name or abbreviation or a synonym of a disease. Charac-

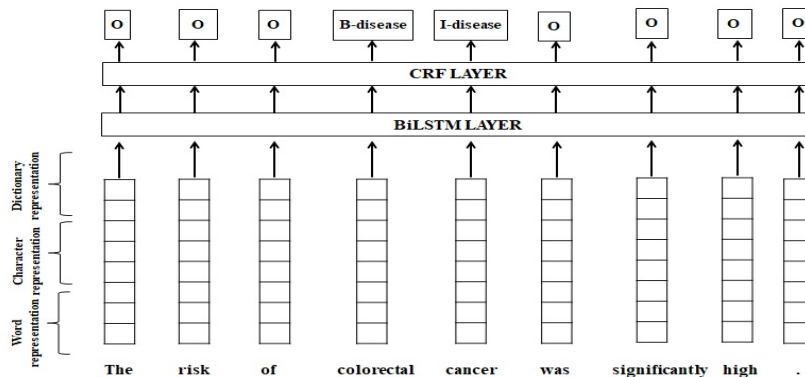


Figure 10: General structure of the proposed model

ter level representation is concatenated with word embeddings and the dictionary information. At this level, every word is represented as a vector comprising of character level information, word level information and dictionary information. Feeding the vector representation of word sequence of a sentence in direct and reverse order to a BiLSTM network will output a contextual representation for each word as shown in Figure 11.

### 4.1 Decoding

Decoding is the final step, which converts the contextual representations of the tokens into corresponding tags. For sequence labeling problem, it is beneficial to consider the correlations between labels of the tokens in the neighbourhoods and jointly decode the best chain of labels for a given input sentence. CRF [31] model is used for decoding as CRF considers the contextual information for decoding the label for each token. There are two approaches for calculating the scores of output tags; local scores and global scores. Local scores use scores represented in the final contextual representations for each word, while global scores use transition scores as well as scores represented in the final contextual representations for each word.

A fully connected neural network has been used to convert the contextual representations of the tokens to a vector where each entry corresponds to a score for each output tag. For an input sentence,  $\mathbf{X} = (x_1, x_2, \dots, x_m)$ ,  $\mathbf{Y} = (y_1, y_2, \dots, y_m)$  is the sequence of output tags,  $S \in R^{T \times m}$  is the



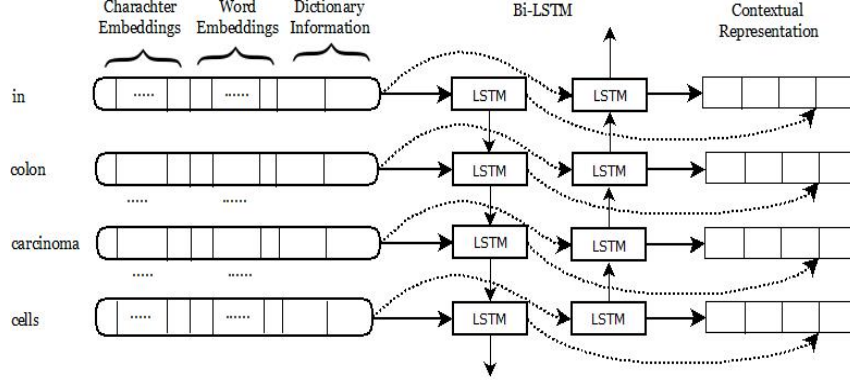


Figure 11: The contextual representation learning model

score matrix, where  $T$  is number of predefined tags and  $s_{i,j} \in S$  is the score of  $i^{\text{th}}$  tag for the  $j^{\text{th}}$  word. A global score,  $Score \in R$  of sequence of tags  $\mathbf{Y} = (y_1, y_2, \dots, y_m)$  is defined as:

$$Score(y_1, y_2, \dots, y_m) = \sum_{t=1}^m (s_{y_t, t} + Tr[y_t, y_{t+1}])$$

where,  $Tr[y_t, y_{t+1}]$  is the score of assigning the tag  $y_{t+1}$  given the tag  $y_t$ .

The sum of scores of all possible sequences of the tags for an input sentence  $\mathbf{X}$  is calculated as:

$$Z = \sum_{y'_i \in Y(x), 1 \leq i \leq m} e^{Score(y'_1, y'_2, \dots, y'_m)}$$

Then given the sentence  $\mathbf{X} = (x_1, x_2, \dots, x_m)$ , the conditional probability of a label sequence  $\mathbf{Y} = (y_1, y_2, \dots, y_m)$  is defined as:

$$P(\mathbf{Y}|\mathbf{X}) = \frac{e^{Score(y_1, y_2, \dots, y_m)}}{Z}$$

The predicted tag sequence  $\hat{y} \in Y(\mathbf{x})$  is calculated as:

$$\hat{y} = \mathbf{argmax}\{P(\mathbf{Y}|\mathbf{X})\}$$

Example of decoding a contextual representation for text fragment “In colon carcinoma cells” is given in Figure 12. In this example, numbers in columns are scores of assigning corresponding tags to the word. CRF is used to calculate the score of all paths and then the path with maximum score will be chosen. The global scores of two sequences are calculated as follows:

$$Score(O, B\text{-Disease}, I\text{-Disease}, O) = 4+3+2+5+4+7+2+8+1=36$$

$$Score(O, B\text{-Disease}, B\text{-Disease}, O) = 4+3+2+5+1+9+1+8+1=34$$

The first path having higher score will be selected by CRF.

In addition to this approach, a local score can be used to decode the output vectors to tags as follows:

$$local\_score(y_1, y_2, \dots, y_m) = \sum_{t=1}^m s_{y_t, t}$$

The predicted tag sequence  $\hat{y} \in Y(\mathbf{x})$  is calculated as:

$$\hat{y} = \mathbf{argmax}\{local\_score(y_1, y_2, \dots, y_m)\}$$

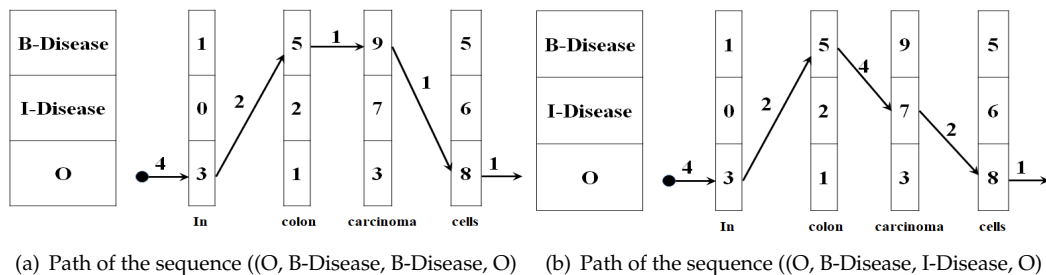


Figure 12: An example of decoding a text fragment

## 5 Datasets

### 5.1 NCBI Dataset

NCBI disease corpus was introduced for disease name recognition and normalization [32]. It is the most inclusive publicly available dataset annotated with disease mentions. The corpus was manually annotated by 14 medical practitioners. General statistics of the dataset is given Table 1.

	Training	Dev	Test	Total
<b>No. of Abstracts</b>	593	100	100	793
<b>No. of Sentences</b>	5661	939	961	7261
<b>Total Disease mentions</b>	5145	787	960	6892
<b>Unique Disease mentions</b>	1710	368	427	2136

Table 1: Statistics of NCBI dataset

### 5.2 BC5CDR dataset

BC5CDR dataset was created for BioCreative V Chemical Disease Relation (CDR) task and consists of 1500 PubMed articles with 5818 disease mentions [33]. The corpus was randomly split into three subsets: 500 each for training, testing and development sets. A BioNE class label named DISEASE and O (for non-BioNEs) are used to annotate the dataset.

## 6 Experiments

### 6.1 Pre-processing

Pre-trained word embeddings using Skip-gram model are used to represent tokens in the data set. This word embeddings<sup>1</sup> model combines domain-specific texts (PMC and PubMed texts) with generic ones (English Wikipedia dump) for better representation of tokens. We have used pre-trained word embeddings trained by Sampo Pyysalo et al. [34], as constructing word embeddings requires a huge corpus and machines with high specifications. Using skip-gram model for training word embeddings improves the semantic and syntactic representation. Character-level embeddings are created by initializing vector representations for every character in the corpus and then the character representation corresponding to every character in a word are given in direct and reverse order to BiLSTM as shown in Figure 13. All numbers in the input are replaced

<sup>1</sup><http://bio.nlplab.org/>

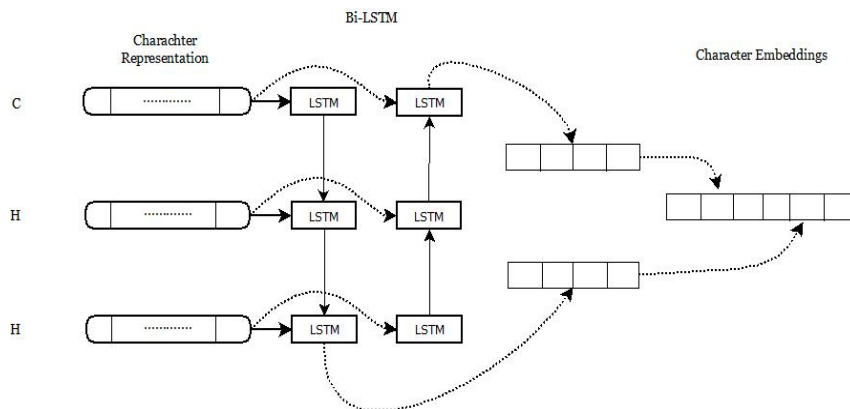


Figure 13: Character representation learning model

with the value `NUM` (number), all letters are converted to lowercase and words that are not represented in the pre-trained word embeddings are marked as `UNK` (unknown). As the pre-trained word embeddings are of large size, a look-up table containing the word embeddings of all words in the dataset are extracted from the pre-trained word embeddings which used as input. To evaluate the impact of the pre-trained word embeddings, randomly initialized word embeddings have been used separately.

## 6.2 Parameters

In this work, the dimension of character-level embeddings is set to 100 and that of the pre-trained word embeddings to 200. Different values for these parameters have been experimented and we set the values which reported the best performance for our experiments. Table 2 shows the parameters and their values used in our experiments.

Parameter	Value
No. of epochs	15
Dropout	0.5
Batch size	20
Learning method	Adam
Learning decay	0.9
No. of LSTM units for character embeddings	100
No. of LSTM units for word embeddings	300
No. of distinct words in NCBI dataset	8147
No. of distinct words in BC5CDR dataset	13427
No. of distinct characters in NCBI dataset	84
No. of distinct characters in BC5CDR dataset	80

Table 2: Parameters and their values used for training the proposed model

## 7 Results and Discussion

NCBI and BC5CDR datasets have been used to evaluate the proposed model. SR schemes namely IOB2 and IOBES have been used to train the model as these two schemes reported high performance among other schemes for BioNER as discussed in [35].

Results shown in Table 3 illustrate that the character embeddings significantly increase the performance for both datasets. Decoding the output vectors using CRF with global scores improves the performance rather than using local scores. The pre-trained word embedding results in better performance than randomly initialized word embeddings. The reason is that these pre-trained embeddings are trained over a collection of huge texts including texts from biomedical domain as well as generic domain. In both datasets, the best results are reported using character embeddings, pre-trained word embeddings, dictionary information and CRF with global scores. IOBES and IOB2 schemes do not show significant difference with NCBI dataset. However, using the same schemes with BC5CDR dataset shows a significant difference of 1.13 in the f1-measure. Ta-

Dataset	SR scheme	V1	V2	V3	V4	f1-measure
NCBI	IOB2	x	x	x	x	71.16%
		x	x	x	✓	80.46%
		x	x	✓	✓	83.13%
		x	✓	✓	✓	84.24%
		✓	✓	✓	✓	85.19%
	IOBES	x	x	x	x	73.23%
		x	x	x	✓	81.60%
		x	x	✓	✓	83.44%
		x	✓	✓	✓	84.40%
		✓	✓	✓	✓	<b>85.40%</b>
BC5CDR	IOB2	x	x	x	x	73.26%
		x	x	x	✓	75.81%
		x	x	✓	✓	78.30%
		x	✓	✓	✓	78.36%
		✓	✓	✓	✓	78.49%
	IOBES	x	x	x	x	73.64%
		x	x	x	✓	76.98%
		x	x	✓	✓	77.95%
		x	✓	✓	✓	79.33%
		✓	✓	✓	✓	<b>79.62%</b>

Table 3: Results of the proposed model with following variations:

**V1:** (x) without dictionary information (✓) with dictionary information

**V2:** (x) randomly initialized word embeddings (✓) pre-trained word embeddings

**V3:** (x) using local score for decoding (✓) using global score for decoding

**V4:** (x) without character embeddings (✓) with character embeddings

Table 4 and Table 5 give the comparisons between our model and the state-of-the-art models for

NCBI and BC5CDR datasets respectively. For NCBI dataset, the performance of our model outperforms the state-of-the-art works.

Model	f1-measure
ANN (BiLSTM + CNN) [23]	79.13%
Pairwise learning [20]	80.90%
ANN (ReLU + Softmax activation) [21]	80.74%
TaggerOne (Semi-Markov model) [22]	82.90%
BANNER (CRF) [25]	81.80%
CNN (Dictionary + Postprocessing) [27]	85.17%
<b>Our model</b> (Best result)	<b>85.40%</b>

Table 4: Comparison of our model with related work on NCBI dataset

Model	f1-measure
LSTM [36]	76.50%
Ensemble (RNN + CRF) [24]	78.04%
CD-REST (CRF + External resource) [26]	84.43%
CRF (dictionary Information)[29]	86.46%
Ensemble(SVM+CRF) [28]	86.93%
CNN (Dictionary + Postprocessing) [27]	87.83%
<b>Our model</b> (Best result)	<b>79.62%</b>

Table 5: Comparison of our model with related work on BC5CDR dataset

## 8 Conclusion

In this paper, an ANN-based model for Disease-NER is presented. Instead of hand engineering the features for tokens, character-level embeddings are used to represent orthographical features of tokens in addition to using word embeddings and dictionary information. Two different SR schemes namely IOB2 and IOBES are used for annotating the corpus. Results show that using character embeddings, pre-trained word embeddings, dictionary information and CRF with global scores improves the performance of BioNER. In addition, IOBES scheme outperforms IOB2 scheme.

## References

- [1] Hamada A. Nayel. *Biomedical Named Entity Recognition*. PhD thesis, Mangalore Univesity, Karnataka, India, 2018.
- [2] Hamada Nayel and H L Shashirekha. Improving NER for Clinical Texts by Ensemble Approach using Segment Representations. In *Proceedings of the 14th International Conference*

- on *Natural Language Processing (ICON-2017)*, pages 197–204, Kolkata, India, December 2017. NLP Association of India.
- [3] Haodi Li, Qingcai Chen, Buzhou Tang, Xiaolong Wang, Hua Xu, Baohua Wang, and Dong Huang. Cnn-based ranking for biomedical entity normalization. *BMC Bioinformatics*, 18(11):385, Oct 2017.
  - [4] Kai Xu, Zhanfan Zhou, Tianyong Hao, and Wenyin Liu. *A Bidirectional LSTM and Conditional Random Fields Approach to Medical Named Entity Recognition*, pages 355–365. Springer International Publishing, Cham, 2018.
  - [5] Hamada A. Nayel, Hiroyuki Shindo, H. L. Shashirekha, and Yuji Matsumoto. Improving multi-word entity recognition for biomedical texts. *International Journal of Pure and Applied Mathematics*, 118(16):301–319, 2018.
  - [6] Fei Li, Meishan Zhang, Guohong Fu, and Donghong Ji. A neural joint model for entity and relation extraction from biomedical text. *BMC Bioinformatics*, 18(1):198, Mar 2017.
  - [7] Sudha Rao, Daniel Marcu, Kevin Knight, and Hal Daumé III. Biomedical event extraction using abstract meaning representation. In *BioNLP 2017*, pages 126–135, Vancouver, Canada,, August 2017. Association for Computational Linguistics.
  - [8] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *Transactions on Neural Networks*, 5(2):157–166, March 1994.
  - [9] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, ICML’13*, pages III–1310–III–1318. JMLR.org, 2013.
  - [10] Sepp Hochreiter and Jrgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
  - [11] Alex Graves and Jrgen Schmidhuber. Frameworkise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks*, 18(5):602 – 610, 2005. IJCNN 2005.
  - [12] Peter D. Turney and Patrick Pantel. From frequency to meaning: Vector space models of semantics. *J. Artif. Int. Res.*, 37(1):141–188, January 2010.
  - [13] Debasis Ganguly, Dwaipayan Roy, Mandar Mitra, and Gareth J.F. Jones. Word embedding based generalized language model for information retrieval. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’15*, pages 795–798, New York, NY, USA, 2015. ACM.
  - [14] Cicero dos Santos and Maira Gatti. Deep convolutional neural networks for sentiment analysis of short texts. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 69–78, 2014.
  - [15] Mahmood Yousefi-Azar and Len Hamey. Text summarization using unsupervised deep learning. *Expert Systems with Applications*, 68:93 – 105, 2017.
  - [16] Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, et al. End-to-end memory networks. In *Advances in neural information processing systems*, pages 2440–2448, 2015.
  - [17] 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 3111–3119, 2013.
  - [18] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

- [19] Cícero Nogueira Dos Santos and Bianca Zadrozny. Learning character-level representations for part-of-speech tagging. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32, ICML'14*, pages 1818–1826. JMLR.org, 2014.
- [20] Robert Leaman, Rezarta Islamaj Doan, and Zhiyong Lu. Dnorm: disease name normalization with pairwise learning to rank. *Bioinformatics*, 29(22):2909–2917, 2013.
- [21] Gamal Crichton, Sampo Pyysalo, Billy Chiu, and Anna Korhonen. A neural network multi-task learning approach to biomedical named entity recognition. *BMC Bioinformatics*, 18(1):368, Aug 2017.
- [22] Robert Leaman and Zhiyong Lu. Taggerone: Joint named entity recognition and normalization with semi-markov models. *Bioinformatics*, 2016.
- [23] Sunil Sahu and Ashish Anand. Recurrent neural network models for disease name recognition using domain invariant features. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2216–2225, Berlin, Germany, August 2016. Association for Computational Linguistics.
- [24] Qikang Wei, Tao Chen, Ruifeng Xu, Yulan He, and Lin Gui. Disease named entity recognition by combining conditional random fields and bidirectional recurrent neural networks. *Database*, 2016.
- [25] Robert Leaman, Graciela Gonzalez, et al. Banner: an executable survey of advances in biomedical named entity recognition. In *Pacific Symposium on Biocomputing*, volume 13, pages 652–663. Kohala Coast, Hawaii, USA, 2008.
- [26] Jun Xu, Yonghui Wu, Yaoyun Zhang, Jingqi Wang, Hee-Jin Lee, and Hua Xu. Cd-rest: a system for extracting chemical-induced disease relation in literature. *Database*, 2016:baw036, 2016.
- [27] Zhehuan Zhao, Zhihao Yang, Ling Luo, Lei Wang, Yin Zhang, Hongfei Lin, and Jian Wang. Disease named entity recognition from biomedical literature using a novel convolutional neural network. *BMC medical genomics*, 10(5):73, 2017.
- [28] Haodi Li, Buzhou Tang, Qingcai Chen, Kai Chen, Xiaolong Wang, Baohua Wang, and Zhe Wang. Hitsz\_cdr: an end-to-end chemical and disease relation extraction system for biocreative v. *Database*, 2016:baw077, 2016.
- [29] Hsin-Chun Lee, Yi-Yu Hsu, and Hung-Yu Kao. An enhanced crf-based system for disease name entity recognition and normalization on biocreative v dner task. In *Proceedings of the Fifth BioCreative Challenge Evaluation Workshop*, pages 226–233, 2015.
- [30] Allan Peter Davis, Thomas C. Wieggers, Michael C. Rosenstein, and Carolyn J. Mattingly. Medic: a practical disease vocabulary used at the comparative toxicogenomics database. *Database*, 2012:bar065, 2012.
- [31] John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning, ICML '01*, pages 282–289, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc.
- [32] Rezarta Islamaj Doan, Robert Leaman, and Zhiyong Lu. Ncbi disease corpus: A resource for disease name recognition and concept normalization. *Journal of Biomedical Informatics*, 47(Supplement C):1 – 10, 2014.
- [33] 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*, 2016:baw068, 2016.

- [34] SPFGH Moen and Tapio Salakoski<sup>2</sup> Sophia Ananiadou. Distributional semantics resources for biomedical text processing. In *Proceedings of the 5th International Symposium on Languages in Biology and Medicine, Tokyo, Japan*, pages 39–43, 2013.
- [35] H. L. Shashirekha and H. A. Nayel. A comparative study of segment representation for biomedical named entity recognition. In *2016 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, pages 1046–1052, Sept 2016.
- [36] Hongwei Liu and Yun Xu. A deep learning way for disease name representation and normalization. In Xuanjing Huang, Jing Jiang, Dongyan Zhao, Yansong Feng, and Yu Hong, editors, *Natural Language Processing and Chinese Computing*, pages 151–157, Cham, 2018. Springer International Publishing.