# Drug-drug Interaction Extraction via Recurrent Neural Network with Multiple Attention Layers

Zibo Yi, Shasha Li, Jie Yu, Qingbo Wu

*College of Computer, National University of Defense Technology*
*Changsha, Hunan, China*
{yizibo14, shashali, yj, qingbo.wu}@nudt.edu.cn

*Abstract*—Drug-drug interaction (DDI) is a vital information when physicians and pharmacists intend to co-administer two or more drugs. Thus, several DDI databases are constructed to avoid mistakenly combined use. In recent years, automatically extracting DDIs from biomedical text has drawn researchers' attention. However, the existing work utilize either complex feature engineering or NLP tools, both of which are insufficient for sentence comprehension. Inspired by the deep learning approaches in natural language processing, we propose a recurrent neural network model with multiple attention layers for DDI classification. We evaluate our model on 2013 SemEval DDIExtraction dataset. The experiments show that our model classifies most of the drug pairs into correct DDI categories, which outperforms the existing NLP or deep learning methods.

## I. Introduction

Drug-drug interaction (DDI) is a situation when one drug increases or decreases the effect of another drug [1]. Adverse drug reactions may cause severe side effect, if two or more medicines were taken and their DDI were not investigated in detail. DDI is a common cause of illness, even a cause of death [2]. Thus, DDI databases for clinical medication decisions are proposed by some researchers. These databases such as SFINX [3], KEGG [4], CredibleMeds [5] help physicians and pharmacists avoid most adverse drug reactions.

Traditional DDI databases are manually constructed according to clinical records, scientific research and drug specifications. For instance, The sentence "With combined use, clinicians should be aware, when **phenytoin** is added, of the potential for reexacerbation of pulmonary symptomatology due to lowered serum **theophylline** concentrations [6]", which is from a pharmacotherapy report, describe the side effect of phenytoin and theophylline's combined use. Then this information on specific medicines will be added to DDI databases. As drug-drug interactions have being increasingly found, manually constructing DDI database would consume a lot of manpower and resources.

There has been many efforts to automatically extract DDIs from natural language [1], [7]–[13], mainly medical literature and clinical records. These works can be divided into the following categories:

- Text analysis and statistics based approach [1], [7], [14]. This kind of work utilizes NLP tools to analysis biomedical text's semantics or statistics features (such as TF-IDF) before the DDI decision. However, the semantics and statistics features are insufficient for understanding the whole text. Even worse, NLP toolkits are imperfect and may propagate error to the classification.
- Feature based machine learning approach [8], [11], [12], [15], [16]. Such method always need complex feature engineering. In addition, the quality of feature engineering have a great effect on the precision of DDI classification, which becomes the shortcoming of such method.
- Deep learning based approach [9], [10], [13]. Deep learning neural networks, such as convolutional neural networks (CNN) and long short-term memory networks (LSTM), have been utilized for DDI extraction. Deep learning method avoids complicated feature engineering since CNN and LSTM can extract semantics features automatically through a well designed network.

To avoid complex feature engineering and NLP toolkits' usage, we employ deep learning approaches for sentence comprehension as a whole. Our model takes in a sentence from biomedical literature which contains a drug pair and outputs what kind of DDI this drug pair belongs. This assists physicians refrain from improper combined use of drugs. In addition, the word and sentence level attentions are introduced to our model for better DDI predictions.

We train our language comprehension model with labeled instances. Figure 1 shows partial records in DDI corpus [17]. We extract the sentence and drug pairs in the records. There are 3 drug pairs in this example thus we have 3 instances. The DDI corpus annotate each drug pair in the sentence with a DDI type. The DDI type, which is the most concerned information, is described in table I. The details about how we train our model and extract the DDI type from text are described in the remaining sections.

## II. Related Work

In DDI extraction task, NLP methods or machine learning approaches are proposed by most of the work. Chowdhury [15] and Thomas *et al.* [12] proposed methods that use linguistic phenomenons and two-stage SVM to classify DDIs. FBK-irst [11] is a follow-on work which applies kernel method to the existing model and outperforms it.

Neural network based approaches have been proposed by several works. Liu *et al.* [10] employ CNN for DDI extraction for the first time which outperforms the traditional machine learning based methods. Limited by the convolutional kernel size, the CNN can only extracted features of continuous 3

## TABLE I
### THE DDI TYPES AND CORRESPONDING EXAMPLES

| DDI types | Definition | Example sentence | Drug pair |
|---|---|---|---|
| False | An interaction between the two drugs is not shown in the sentence. | Concomitantly given thiazide diuretics did not interfere with the absorption of a tablet of digoxin. | thiazide diuretics, digoxin |
| Mechanism | An pharmacokinetic mechanism is shown in the sentence. | Additional iron significantly inhibited the absorption of cobalt. | iron, cobalt |
| Effect | The effect of two drugs' combination use is shown in the sentence. | Methotrexate: An increased risk of hepatitis has been reported to result from combined use of methotrexate and etretinate. | methotrexate, etretinate |
| Advise | An advise about two drugs is given in the sentence. | UROXATRAL should NOT be used in combination with other alpha-blockers. | UROXATRAL, alpha-blockers |
| Int | A drug interaction without any further information is mentioned in the sentence. | Clinical implications of warfarin interactions with five sedatives. | warfarin, sedatives |

to 5 words rather than distant words. Liu *et al.* [9] proposed dependency-based CNN to handle distant but relevant words. Sahu *et al.* [13] proposed LSTM based DDI extraction approach and outperforms CNN based approach, since LSTM handles sentence as a sequence instead of slide windows. To conclude, Neural network based approaches have advantages of 1) less reliance on extra NLP toolkits, 2) simpler preprocessing procedure, 3) better performance than text analysis and machine learning methods.

Drug-drug interaction extraction is a relation extraction task of natural language processing. Relation extraction aims to determine the relation between two given entities in a sentence. In recent years, attention mechanism and various neural networks are applied to relation extraction [18]–[22]. Convolutional deep neural network are utilized for extracting sentence level features in [20]. Then the sentence level features are concatenated with lexical level features, which are obtained by NLP toolkit WordNet [23], followed by a multilayer perceptron (MLP) to classify the entities' relation. A fixed work is proposed by Nguyen *et al.* [22]. The convolutional kernel is set various size to capture more *n*-gram features. In addition, the word and position embedding are trained automatically instead of keeping constant as in [20]. Wang *et al.* [21] introduce multi-level attention mechanism to CNN in order to emphasize the keywords and ignore the non-critical words during relation detection. The attention CNN model outperforms previous state-of-the-art methods.

Besides CNN, Recurrent neural network (RNN) has been applied to relation extraction as well. Zhang *et al.* [19] utilize long short-term memory network (LSTM), a typical RNN model, to represent sentence. The bidirectional LSTM chronologically captures the previous and future information, after which a pooling layer and MLP have been set to extract feature and classify the relation. Attention mechanism is added to bidirectional LSTM in [18] for relation extraction. An attention layer gives each memory cell a weight so that classifier can catch the principal feature for the relation detection. The Attention based bidirectional LSTM has been proven better than previous work.

```
<sentence id="DDI-DrugBank.d353.s8" text="Methotrexate: An increased risk
of hepatitis has been reported to result from combined use of methotrexate
and etretinate.">
        <entity id="DDI-DrugBank.d353.s8.e0" charOffset="0-11"
            type="drug" text="Methotrexate"/>
        <entity id="DDI-DrugBank.d353.s8.e1" charOffset="94-105"
            type="drug" text="methotrexate"/>
        <entity id="DDI-DrugBank.d353.s8.e2" charOffset="111-120"
            type="drug" text="etretinate"/>
        <pair id="DDI-DrugBank.d353.s8.p0" e1="DDI-DrugBank.d353.s8.e0"
            e2="DDI-DrugBank.d353.s8.e1" ddi="false"/>
        <pair id="DDI-DrugBank.d353.s8.p1" e1="DDI-DrugBank.d353.s8.e0"
            e2="DDI-DrugBank.d353.s8.e2" ddi="false"/>
        <pair id="DDI-DrugBank.d353.s8.p2" e1="DDI-DrugBank.d353.s8.e1"
            e2="DDI-DrugBank.d353.s8.e2" ddi="true" type="effect"/>
</sentence>
```

Fig. 1. Partial records in DDI corpus

## III. PROPOSED MODEL

In this section, we present our bidirectional recurrent neural network with multiple attention layer model. The overview of our architecture is shown in figure 2. For a given instance, which describes the details about two or more drugs, the model represents each word as a vector in embedding layer. Then the bidirectional RNN layer generates a sentence matrix, each column vector in which is the semantic representation of the corresponding word. The word level attention layer transforms the sentence matrix to vector representation. Then sentence level attention layer generates final representation for the instance by combining several relevant sentences in view of the fact that these sentences have the same drug pair. Followed by a softmax classifier, the model classifies the drug pair in the given instance as specific DDI.

### A. Preprocessing

The DDI corpus contains thousands of XML files, each of which are constructed by several records. For a sentence containing $n$ drugs, there are $C_n^2$ drug pairs. We replace the interested two drugs with "drug1" and "drug2" while the other drugs are replaced by "durg0", as in [10] did. This step is called drug blinding. For example, the sentence in figure 1 generates 3 instances after drug blinding: "drug1: an increased risk of hepatitis has been reported to result from combined use of drug2 and drug0", "drug1: an increased risk of hepatitis has been reported to result from combined use of drug0 and

drug2", "drug0: an increased risk of hepatitis has been reported to result from combined use of drug1 and drug2". The drug blinded sentences are the instances that are fed to our model.

We put the sentences with the same drug pairs together as a set, since the sentence level attention layer (will be described in Section III-E) will use the sentences which contain the same drugs.

### B. Embedding Layer

Given an instance $S = (w_1, w_2, ..., w_t)$ which contains specified two drugs $w_u = $ "drug1", $w_v = $ "drug2", each word is embedded in a $d = d_{WE} + 2d_{PE}$ dimensional space ($d_{WE}$, $d_{PE}$ are the dimension of word embedding and position embedding). The look up table function $LT.(\cdot)$ maps a word or a relative position to a column vector. After embedding layer the sentence is represented by $S = (\boldsymbol{x}_1, \boldsymbol{x}_2, ..., \boldsymbol{x}_t)$, where

$$\boldsymbol{x}_i = (LT_W(w_i)^{\mathrm{T}}, (LT_P(i-u)^{\mathrm{T}}, (LT_P(i-v)^{\mathrm{T}})^{\mathrm{T}} \quad (1)$$

The $LT.(\cdot)$ function is usually implemented with matrix-vector product. Let $\overline{w_i}$, $\overline{k}$ denote the one-hot representation (column vector) of word and relative distance. $E_w$, $E_p$ are word and position embedding query matrix. The look up functions are implemented by

$$LT_W(w_i) = E_w\overline{w_i}, LT_P(k) = E_p\overline{k} \quad (2)$$

Then the word sequence $S = (\boldsymbol{x}_1, \boldsymbol{x}_2, ..., \boldsymbol{x}_t)$ is fed to the RNN layer. Note that the sentence will be filled with $\boldsymbol{0}$ if its length is less than $t$.

### C. Bidirectional RNN Encoding Layer

The words in the sequence are read by RNN's gated recurrent unit (GRU) one by one. The GRU takes the current word $\boldsymbol{x}_i$ and the previous GRU's hidden state $h_{i-1}$ as input. The current GRU encodes $h_{i-1}$ and $\boldsymbol{x}_i$ into a new hidden state $h_i$ (its dimension is $d_h$, a hyperparameter), which can be regarded as informations the GRU remembered.

Figure 3 shows the details in GRU. The reset gate $r_i$ selectively forgets informations delivered by previous GRU. Then the hidden state becomes $\tilde{h}_i$. The update gate $z_i$ updates the informations according to $\tilde{h}_i$ and $h_{i-1}$. The equations below describe these procedures. Note that $\otimes$ stands for element wise multiplication.

$$r_i = \sigma(W_r\boldsymbol{x}_i + U_r h_{i-1}) \quad (3)$$

$$\tilde{h}_i = \Phi(W\boldsymbol{x}_i + U(r_i \otimes h_{i-1})) \quad (4)$$

$$z_i = \sigma(W_z\boldsymbol{x}_i + U_z h_{i-1}) \quad (5)$$

$$h_i = z_i \otimes h_{i-1} + ((1, 1, ..., 1)^{\mathrm{T}} - z_i) \otimes \tilde{h}_i \quad (6)$$

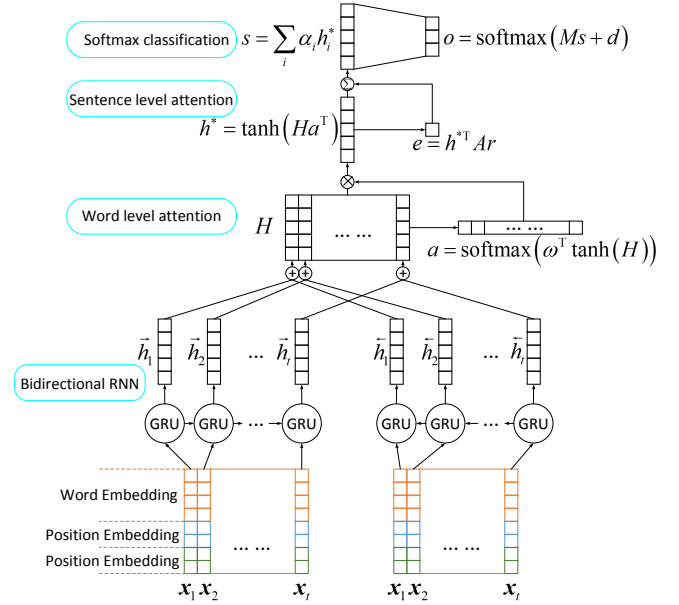The bidirectional RNN contains forward RNN and backward RNN. Forward RNN reads sentence from $\boldsymbol{x}_1$ to $\boldsymbol{x}_t$,



Fig. 2. The bidirectional recurrent neural network with multiple attentions

generating $\overrightarrow{h}_1$, $\overrightarrow{h}_2$, ..., $\overrightarrow{h}_t$. Backward RNN reads sentence from $\boldsymbol{x}_t$ to $\boldsymbol{x}_1$, generating $\overleftarrow{h}_t$, $\overleftarrow{h}_{t-1}$, ..., $\overleftarrow{h}_1$. Then the encode result of this layer is

$$H = (\overrightarrow{h}_1 + \overleftarrow{h}_1, \overrightarrow{h}_2 + \overleftarrow{h}_2, ..., \overrightarrow{h}_t + \overleftarrow{h}_t) \quad (7)$$

We apply dropout technique in RNN layer to avoid overfitting. Each GRU have a probability (denoted by $Pr_{dp}$, also a hyperparameter) of being dropped. The dropped GRU has no output and will not affect the subsequent GRUs. With bidirectional RNN and dropout technique, the input $S = (\boldsymbol{x}_1, \boldsymbol{x}_2, ..., \boldsymbol{x}_t)$ is encoded into sentence matrix $H$.

### D. Word Level Attention

The purpose of word level attention layer is to extract sentence representation (also known as feature vector) from encoded matrix. We use word level attention instead of max pooling, since attention mechanism can determine the importance of individual encoded word in each row of $H$. Let $\omega$ denotes the attention vector (column vector), $a$ denotes the filter that gives each element in the row of $H$ a weight. The following equations shows the attention operation, which is also illustrated in figure 2.

$$a = \mathrm{softmax}(\omega^{\mathrm{T}}\mathrm{tanh}(H)) \quad (8)$$

$$h^* = \mathrm{tanh}(Ha^{\mathrm{T}}) \quad (9)$$

The softmax function takes a vector $v = [v_1, v_2, ..., v_n]$ as input and outputs a vector,

$$\mathrm{softmax}(v) = [\frac{e^{v_1}}{\sum_{i=1}^n e^{v_i}}, \frac{e^{v_2}}{\sum_{i=1}^n e^{v_i}}, ..., \frac{e^{v_n}}{\sum_{i=1}^n e^{v_i}}] \quad (10)$$
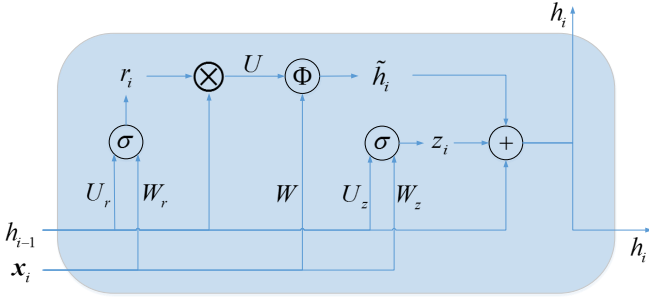
Fig. 3. The Gated Recurrent Unit

$h^*$ denotes the feature vector captured by this layer. Several approaches [13], [18] use this vector and softmax classifier for classification. Inspired by [24] we propose the sentence level attention to combine the information of other sentences for a improved DDI classification.

### E. Sentence Level Attention

The previous layers captures the features only from the given sentence. However, other sentences may contains informations that contribute to the understanding of this sentence. It is reasonable to look over other relevant instances when determine two drugs' interaction from the given sentence. In our implementation, the instances that have the same drug pair are believed to be relevant. The relevant instances set is denoted by $\mathcal{L} = \{h_1^*, h_2^*, ..., h_N^*\}$, where $h_i^*$ is the sentence feature vector. $e_i$ stands for how well the instance $h_i^*$ matches its DDI $r$ (Vector representation of a specific DDI). $A$ is a diagonal attention matrix, multiplied by which the feature vector $h_i^*$ can concentrate on those most representative features.

$$e_i = h_i^{*\mathrm{T}} A r \tag{11}$$

$$\alpha_i = \frac{\exp(e_i)}{\sum_{k=1}^{N} \exp(e_k)} \tag{12}$$

$\alpha_i$ is the softmax result of $e_i$. The final sentence representation is decided by all of the relevant sentences' feature vector, as Equation 13 shows.

$$s = \sum_{i=1}^{N} \alpha_i h_i^* \tag{13}$$

Note that the set $\mathcal{L}$ is gradually growing as new sentence with the same drugs pairs is found when training. An instance $S = (\boldsymbol{x}_1, \boldsymbol{x}_2, ..., \boldsymbol{x}_t)$ is represented by $h^*$ before sentence level attention. The sentence level attention layer finds the set $\mathcal{L}$, instances in which have the same drug pair as in $S$, and put $S$ in $\mathcal{L}$. Then the final sentence representation $s$ is calculated in this layer.

### F. Classification and Training

A given sentence $S = (w_1, w_2, ..., w_t)$ is finally represented by the feature vector $s$. Then we feed it to a softmax classifier. Let $C$ denotes the set of all kinds of DDI. The output $o \in R^{|C|}$ is the probabilities of each class $S$ belongs.

$$o = \mathrm{softmax}(Ms + d) \tag{14}$$

We use cross entropy cost function and $L^2$ regularization as the optimization objective. For $i$-th instance, $Y_i$ denotes the one-hot representation of it's label, where the model outputs $o_i$. The cross entropy cost is:

$$l_i = -\ln Y_i^{\mathrm{T}} o_i \tag{15}$$

For a mini-batch $\mathcal{M} = \{S_1, S_2, ..., S_M\}$, the optimization objective is:

$$J(\theta) = -\frac{1}{|\mathcal{M}|} \sum_{i=1}^{|\mathcal{M}|} \ln Y_i^{\mathrm{T}} o_i + \lambda ||\theta||_2^2 \tag{16}$$

All parameters in this model is:

$$\theta = \{E_w, E_p, W_r, U_r, W, U, W_z, U_z, \omega, A, r, M, d\} \tag{17}$$

We optimize the parameters of objective function $J(\theta)$ with Adam [25], which is a variant of mini-batch stochastic gradient descent. During each train step, the gradient of $J(\theta)$ is calculated. Then $\theta$ is adjusted according to the gradient. After the end of training, we have a model that is able to predict two drugs' interactions when a sentence about these drugs is given.

### G. DDI Prediction

The model is trained for DDI classification. The parameters in list $\theta$ are tuned during the training process. Given a new sentence with two drugs, we can use this model to classify the DDI type.

The DDI prediction follows the procedure described in Section III-A - III-F. The given sentence is eventually represented by feature vector $s$. Then $s$ is classified to a specific DDI type with a softmax classifier. In next section, we will evaluate our model's DDI prediction performance and see the advantages and shortcomings of our model.

## IV. EXPERIMENTS

### A. Datasets and Evaluation Metrics

We use the DDI corpus of the 2013 DDIExtraction challenge [17] to train and test our model. The DDIs in this corpus are classified as five types. We give the definitions of these types and their example sentences, as shown in table I. This standard dataset is made up of training set and testing set. We use the same metrics as in other drug-drug interaction extraction literature [9]–[13], [26]: the overall precision, recall, and F1 score on testing set. $C$ denotes the set of {False, Mechanism, Effect, Advise, Int}. The precision and recall of each $c \in C$ are calculated by
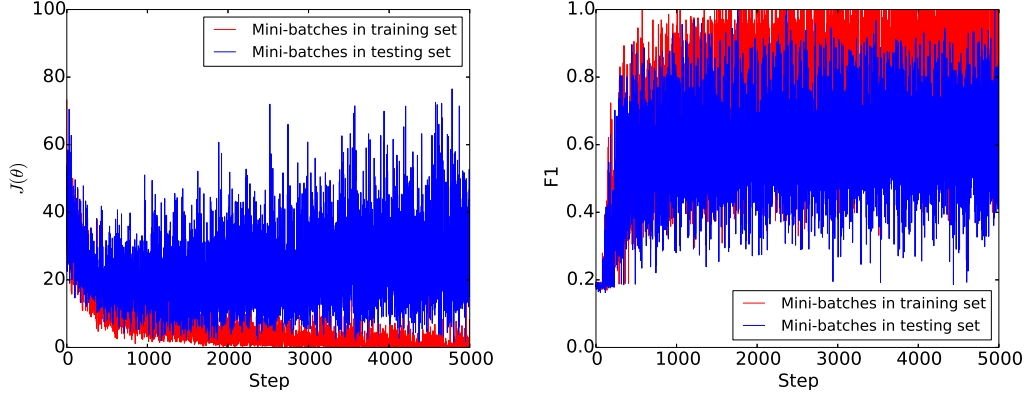
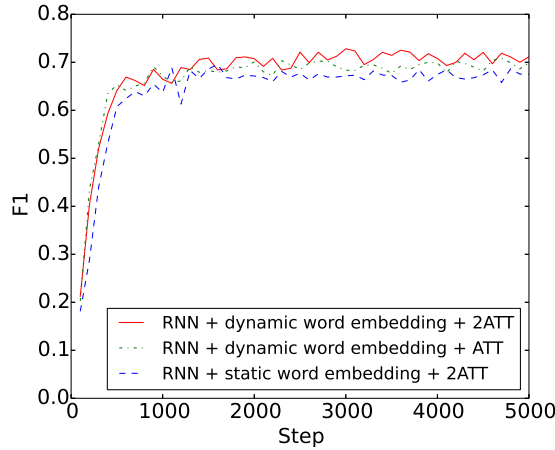Fig. 4. The objective function and F1 in the train process



Fig. 5. The F1 scores on the whole testing set

$$P_c = \frac{\# \; DDI \; is \; c \; and \; is \; classified \; as \; c}{\# \; Classified \; as \; c} \quad (18)$$

$$R_c = \frac{\# \; DDI \; is \; c \; and \; is \; classified \; as \; c}{\# \; DDI \; is \; c} \quad (19)$$

Then the overall precision, recall, and F1 score are calculated by

$$P = \frac{1}{|C|}\sum_{c \in C} P_c, \;\; R = \frac{1}{|C|}\sum_{c \in C} R_c, \;\; F1 = \frac{2PR}{P+R} \quad (20)$$

Besides, we evaluate the captured feature vectors with t-SNE [27], a visualizing and intuitive way to map a high dimensional vector into a 2 or 3-dimensional space. If the points in a low dimensional space are easy to be split, the feature vectors are believed to be more distinguishable.

### B. Hyperparameter Settings and Training

We use TensorFlow [28] r0.11 to implement the proposed model. The input of each word is an ordered triple (*word,*

*relative distance from drug1, relative distance from drug2*). The sentence, which is represented as a matrix, is fed to the model. The output of the model is a $|C|$-dimensional vector representing the probabilities of being corresponding DDI. It is the network, parameters, and hyperparameters which decides the output vector. The network's parameters are adjusted during training, where the hyperparameters are tuned by hand. The hyperparameters after tuning are as follows. The word embedding's dimension $d_{WE} = 100$, the position embedding's dimension $d_{PE} = 10$, the hidden state's dimension $d_h = 230$, the probability of dropout $Pr_d = 0.5$, other hyperparameters which are not shown here are set to TensorFlow's default values.

The word embedding is initialized by pre-trained word vectors using GloVe [29], while other parameters are initialized randomly. During each training step, a mini-batch (the mini-batch size $|\mathcal{M}| = 60$ in our implementation) of sentences is selected from training set. The gradient of objective function is calculated for parameters updating (See Section III-F).

Figure 4 shows the training process. The objective function $J(\theta)$ is declining as the training mini-batches continuously sent to the model. As the testing mini-batches, the $J(\theta)$ function is fluctuating while its overall trend is descending. The instances in testing set are not participated in training so that $J(\theta)$ function is not descending so fast. However, training and testing instances have similar distribution in sample space, causing that testing instances' $J(\theta)$ tends to be smaller along with the training process. $J(\theta)$ has inverse relationship with the performance measurement. The F1 score is getting fluctuating around a specific value after enough training steps. The reason why fluctuating range is considerable is that only a tiny part of the whole training or testing set has been calculated the F1 score. Testing the whole set during every step is time consuming and not necessary. We will evaluate the model on the whole testing set in Section IV-C.

### C. Experimental Results

We save our model every 100 step and predict all the DDIs of the instances in the testing set. These predictions'

| Systems | Methods | Performance | | |
|---|---|---|---|---|
| | | P | R | F1 |
| WBI [12] | Two stage SVM classification | 0.6420 | 0.5790 | 0.6090 |
| FBK-ist [11] | Hand crafted features + SVM | 0.6460 | 0.6560 | 0.6510 |
| SCNN [26] | Two stage syntax CNN | 0.725 | 0.651 | 0.686 |
| Liu *et al.* [10] | CNN + Pre-trained WE | 0.7572 | 0.6466 | 0.6975 |
| DCNN [9] | Dependency-based CNN + Pretrained WE | **0.7721** | 0.6435 | 0.7019 |
| Sahu *et al.* [13] | bidirectional LSTM + ATT | 0.7341 | 0.6966 | 0.7148 |
| This paper | RNN + dynamic WE + 2ATT | 0.7367 | **0.7079** | **0.7220** |



(a) Static word embedding + 2ATT    (b) Dynamic word embedding + ATT    (c) Dynamic word embedding + 2ATT

Fig. 6.   The features which mapped to 2D

TABLE III

PREDICTION RESULTS

| | Classified as | | | | | Sum |
|---|---|---|---|---|---|---|
| | False | Mechanism | Effect | Advise | Int | |
| False | 4490 | 138 | 49 | 45 | 15 | 4737 |
| Mechanism | 68 | 229 | 2 | 3 | 0 | 302 |
| Effect | 101 | 12 | 230 | 15 | 2 | 360 |
| Advise | 49 | 5 | 0 | 165 | 2 | 221 |
| Int | 13 | 3 | 37 | 0 | 43 | 96 |
| Sum | 4721 | 387 | 318 | 228 | 62 | 5716 |

F1 score is shown in figure 5. To demonstrate the sentence level attention layer is effective, we drop this layer and then directly use $h^*$ for softmax classification (See figure 2). The result is shown with "RNN + dynamic word embedding + ATT" curve, which illustrates that the sentence level attention layer contributes to a more accurate model.

Whether a dynamic or static word embedding is better for a DDI extraction task is under consideration. Nguyen *et al.* [22] shows that updating word embedding at the time of other parameters being trained makes a better performance in relation extraction task. We let the embedding be static when training, while other conditions are all the same. The "RNN + static word embedding + 2ATT" curve shows this case. We can draw a conclusion that updating the initialized word embedding trains more suitable word vectors for the task, which promotes the performance.

We compare our best F1 score with other state-of-the-art approaches in table II, which shows our model has competitive advantage in dealing with drug-drug interaction extraction.

The predictions confusion matrix is shown in table III. The DDIs other than false being classified as false makes most of the classification error. It may perform better if a classifier which can tells true and false DDI apart is trained. We leave this two-stage classifier to our future work. Another phenomenon is that the "Int" type is often classified as "Effect". The "Int" sentence describes there exists interaction between two drugs and this information implies the two drugs' combination will have good or bed effect. That's the reason why "Int" and "Effect" are often obfuscated.

To evaluate the features our model captured, we employ scikit-learn [30]'s t-SNE class [1] to map high dimensional feature vectors to 2-dimensional vectors, which can be depicted on a plane. We depict all the features of the instances in testing set, as shown in figure 6. The RNN model using dynamic word embedding and 2 layers of attention is the most distinguishable one. Unfortunately, the classifier can not classify all the instances into correct classes. Comparing table III with figure 6(c), both of which are from the best performed model, we can observe some conclusions. The "Int" DDIs are often misclassified as "Effect", for the reason that some of the "Int" points are in the "Effect" cluster. The "Effect" points are too scattered so that plenty of "Effect" DDIs are classified to other types. The "Mechanism" points are gathered around two clusters, causing that most of the "mechanism" DDIs are classified to two types: "False" and "Mechanism". In short, the visualizability of feature mapping gives better explanations for the prediction results and the quality of captured features.

---

[1] http://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html

## V. Conclusion and Future Work

To conclude, we propose a recurrent neural network with multiple attention layers to extract DDIs from biomedical text. The sentence level attention layer, which combines other sentences containing the same drugs, has been added to our model. The experiments shows that our model outperforms the state-of-the-art DDI extraction systems. Task relevant word embedding and two attention layers improved the performance to some extent.

The imbalance of the classes and the ambiguity of semantics cause most of the misclassifications. We consider that instance generation using generative adversarial networks would cover the instance shortage in specific category. It is also reasonable to use distant supervision learning (which utilize other relevant material) for knowledge supplement and obtain a better performed DDI extraction system.

## References

[1] L. Tari, S. Anwar, S. Liang, J. Cai, and C. Baral, "Discovering drug-drug interactions: a text-mining and reasoning approach based on properties of drug metabolism," *Bioinformatics*, vol. 26, no. 18, pp. 547–53, 2010.

[2] J. Lazarou, B. H. Pomeranz, and P. N. Corey, "Incidence of adverse drug reactions in hospitalized patients: a meta-analysis of prospective studies," *Jama*, vol. 279, no. 15, pp. 1200–1205, 1998.

[3] Y. Böttiger, K. Laine, M. L. Andersson, T. Korhonen, B. Molin, M.-L. Ovesjö, T. Tirkkonen, A. Rane, L. L. Gustafsson, and B. Eiermann, "Sfinx: a drug-drug interaction database designed for clinical decision support systems," *European journal of clinical pharmacology*, vol. 65, no. 6, pp. 627–633, 2009.

[4] M. Takarabe, D. Shigemizu, M. Kotera, S. Goto, and M. Kanehisa, "Network-based analysis and characterization of adverse drug-drug interactions." *Journal of Chemical Information and Modeling*, vol. 51, no. 11, pp. 2977–2985, 2011.

[5] P. R. Shankar, "Crediblemeds: Independent information on medicines," *Australasian Medical Journal*, vol. 7, no. 1, p. 149, 2014.

[6] S. J. Sklar and J. C. Wagner, "Enhanced theophylline clearance secondary to phenytoin therapy," *Annals of Pharmacotherapy*, vol. 19, no. 1, p. 34, 1985.

[7] Y. Lu, D. Shen, M. Pietsch, C. Nagar, Z. Fadli, H. Huang, Y. C. Tu, and F. Cheng, "A novel algorithm for analyzing drug-drug interactions from medline literature," *Scientific Reports*, vol. 5, p. 17357, 2015.

[8] Q. C. Bui, P. M. A. Sloot, E. M. V. Mulligen, and J. A. Kors, "A novel feature-based approach to extract drugdrug interactions from biomedical text," *Bioinformatics*, vol. 30, no. 23, pp. 3365–71, 2014.

[9] S. Liu, K. Chen, Q. Chen, and B. Tang, "Dependency-based convolutional neural network for drug-drug interaction extraction," in *IEEE International Conference on Bioinformatics and Biomedicine*, 2016, pp. 1074–1080.

[10] S. Liu, B. Tang, Q. Chen, and X. Wang, "Drug-drug interaction extraction via convolutional neural networks," *Computational and Mathematical Methods in Medicine*, vol. 2016, pp. 1–8, 2016.

[11] M. F. M. Chowdhury and A. Lavelli, "Fbk-irst: A multi-phase kernel based approach for drug-drug interaction detection and classification that exploits linguistic information," *Atlanta, Georgia, USA*, vol. 351, p. 53, 2013.

[12] P. Thomas, M. Neves, T. Rocktschel, and U. Leser, "Wbi-ddi: Drug-drug interaction extraction using majority voting," in *DDI Challenge at Semeval*, 2013.

[13] S. K. Sahu and A. Anand, "Drug-drug interaction extraction from biomedical text using long short term memory network," *arXiv preprint arXiv:1701.08303*, 2017.

[14] M. P. Melnikov and P. N. Vorobkalov, *Retrieval of Drug-Drug Interactions Information from Biomedical Texts: Use of TF-IDF for Classification*. Springer International Publishing, 2014.

[15] M. F. M. Chowdhury and A. Lavelli, "Exploiting the scope of negations and heterogeneous features for relation extraction: A case study for drug-drug interaction extraction." in *HLT-NAACL*, 2013, pp. 765–771.

[16] M. Rastegar-Mojarad, "Extraction and classification of drug-drug interaction from biomedical text using a two-stage classifier," 2013.

[17] M. Herrero-Zazo, I. Segura-Bedmar, P. Martinez, and T. Declerck, "The ddi corpus: an annotated corpus with pharmacological substances and drug-drug interactions," *Journal of Biomedical Informatics*, vol. 46, no. 5, p. 914, 2013.

[18] P. Zhou, W. Shi, J. Tian, Z. Qi, B. Li, H. Hao, and B. Xu, "Attention-based bidirectional long short-term memory networks for relation classification," in *Meeting of the Association for Computational Linguistics*, 2016, pp. 207–212.

[19] S. Zhang, D. Zheng, X. Hu, and M. Yang, "Bidirectional long short-term memory networks for relation classification." in *PACLIC*, 2015.

[20] D. Zeng, K. Liu, S. Lai, G. Zhou, J. Zhao, *et al.*, "Relation classification via convolutional deep neural network." in *COLING*, 2014, pp. 2335–2344.

[21] L. Wang, Z. Cao, G. D. Melo, and Z. Liu, "Relation classification via multi-level attention cnns," in *Meeting of the Association for Computational Linguistics*, 2016, pp. 1298–1307.

[22] T. H. Nguyen and R. Grishman, "Relation extraction: Perspective from convolutional neural networks," in *The Workshop on Vector Space Modeling for Natural Language Processing*, 2015, pp. 39–48.

[23] D. Lin, "Wordnet: An electronic lexical database," *Computational Linguistics*, vol. 25, no. 2, pp. 292–296, 1999.

[24] Y. Lin, S. Shen, Z. Liu, H. Luan, and M. Sun, "Neural relation extraction with selective attention over instances," in *Meeting of the Association for Computational Linguistics*, 2016, pp. 2124–2133.

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

[26] Z. Zhao, Z. Yang, L. Luo, H. Lin, and J. Wang, "Drug drug interaction extraction from biomedical literature using syntax convolutional neural network." *Bioinformatics*, vol. 32, no. 22, pp. 3444–3453, 2016.

[27] L. V. Der Maaten and G. E. Hinton, "Visualizing data using t-sne," *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008.

[28] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.

[29] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Conference on Empirical Methods in Natural Language Processing*, 2014, pp. 1532–1543.

[30] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.