

# An Attentive Sequence Model for Adverse Drug Event Extraction from Biomedical Text

Suriyadeepan Ramamoorthy    Selvakumar Murugan

SAAMA AI Research Lab

Chennai, India

{suriyadeepan.ramamoorthy, selvakumar.murugan}@saama.com

**Abstract**— Adverse reaction caused by drugs is a potentially dangerous problem which may lead to mortality and morbidity in patients. Adverse Drug Event (ADE) extraction is a significant problem in biomedical research. We model ADE extraction as a Question-Answering problem and take inspiration from Machine Reading Comprehension (MRC) literature, to design our model. Our objective in designing such a model, is to exploit the local linguistic context in clinical text and enable intra-sequence interaction, in order to jointly learn to classify drug and disease entities, and to extract adverse reactions caused by a given drug. Our model makes use of a self-attention mechanism to facilitate intra-sequence interaction in a text sequence. This enables us to visualize and understand how the network makes use of the local and wider context for classification.

## I. INTRODUCTION

Adverse reactions caused by drugs is a potentially life-threatening problem. Extracting such adverse effects is a non-trivial problem. It has attracted the interest of researchers and Health-care Providers. ADE (Adverse Drug Event) is an umbrella term that includes adverse reaction to drugs, unintended side effects and effects from discontinuation or overdose of prescribed drugs. 94% of ADE cases are under-reported by the official systems used in the pharmaceutical industry [4]. It is reported that approximately 7000 deaths [5] were caused by ADE annually, in a study conducted in 2000. ADE are hard to discover because it manifests on certain group of people in certain circumstances, and it may take a long time to expose. Recent research in post-market drug surveillance [4], [6], [7], [8] make use of informal text from social media platforms like twitter. Our research is entirely focused on biomedical text extracted from PubMed reports [1].

The objective of ADE task is to identify drug and disease mentions in a sequence and possible ADE relations between them. Most of the prior work model this task as a entity and relation extraction problem. There are two broad categories of models used to solve this task. The first category uses traditional pipeline method which consists of two steps - Named Entity Recognition (NER) followed by relation classification. Most early works [9], [10] in this area, use pipeline models for ADE task. Pipeline models heavily rely on manual feature engineering.

The second category of models are largely based on End-to-End Deep Neural Networks (DNN). Recent advances in deep learning have spawned diverse groups of neural

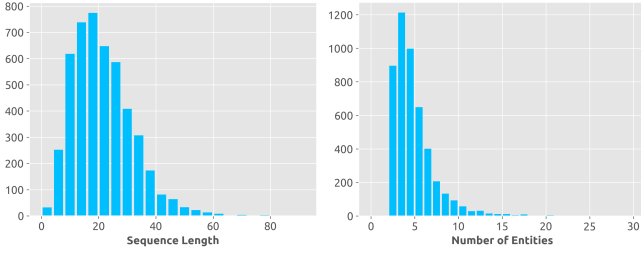
network architectures tailored for Computer Vision, Natural Language Processing (NLP), Speech Recognition, etc.. The promise of deep learning is automatic learning of significant features from large data. The NLP community has adopted Recurrent Neural Networks (RNN) for processing large, unstructured, variable length text sequences. Miwa et al., [11] used Long Short Term Memory (LSTM) [20] based RNNs for relation classification. Li et al., [13] used a feed forward neural network to jointly extracting drug-disease entity mentions and their ADE relations.

Li et al.'s work [3], which appears to be the state of the art in ADE task, employs a neural joint model to extract entities and their relations. They use a bidirectional LSTM based RNN for learning entity representations from text sequences. The output of this network is fed as input to another bidirectional LSTM RNN, which learns ADE representations. The LSTM parameters are shared by both the networks. The ADE network processes Shortest Dependency Paths (SDPs) between possible entities, based on the dependency graph of the text sequence. Our approach makes use of several techniques employed in [3].

In this work, we present a simple sequence model for the combined task of Entity Recognition (ER) and Adverse Drug Event (ADE) extraction. The design of our model is inspired by modern end-to-end sequence processing Recurrent Neural Network architectures. It is considerably simple and easier to train and extensible to similar tasks in NLP. The model makes use of a self-attention mechanism to facilitate intra-sequence interaction i.e., interaction between constituents of a sequence. We postulate that the attention mechanism, widely employed by the deep learning community, could be an apt replacement for SDP-based methods used in [3], [11] and [14]. Additionally, we make use of the heat maps generated by the attention mechanism, to visualize and understand how the network makes use of the local linguistic context and global semantic context, for entity recognition and ADE extraction.

## II. DATA

The data consists of annotated sentences extracted from PubMed abstracts [1]. 6821 sentences contain at least one ADE relation. 16695 unlabeled sentences with no ADE



(a) Sequence Length Histogram (b) Entity Count Histogram

Fig. 1. ADE corpus statistics

relations, are ignored. Each sample consists of an index and a text sequence, followed by a drug entity from the sequence and an ADE entity caused by the drug. We have modified the data samples such that each sample is annotated with one drug and a list of ADE’s caused by that drug. This allows us to condition the model on a given drug and focus on the segments of text corresponding to the effect of the drug. As suggested in prior works [2], [3], 120 sentences with overlapping entities (e.g., ”lithium intoxication”, where ”lithium” is a drug that causes ”lithium intoxication”) are removed from the dataset. We present the sequence length and average entity count in a sequence, as histograms in figure 1.

### III. PROBLEM DEFINITION

The problem we consider is two-fold. Given a sequence, task one is to recognize drug and disease entities in the sequence i.e., Entity Recognition. Task two is identifying adverse drug events in the text sequence i.e., ADE Extraction. Task one is a simple sequence labeling problem, where we predict the category of each token in the sequence.

$$p(y_i^{entity} | X_{seq}) = f_1(x_i, X_{seq}) \quad (1)$$

We have considered task two as a question answering problem. In question answering, we are given (context, query) tuples and tasked to select an answer from a vocabulary. Most architectures in reading comprehension literature construct a query-aware context representation for selecting answers to query. Based on this idea, we considered ADE extraction as a question answering problem, where the text sequence becomes the context and the drug whose adverse effects are to be predicted, becomes the query. Rather than selecting an answer (adverse effect) from a vocabulary, we consider each token in the sequence as a potential ADE. We pose the question, ”Is the  $t^{th}$  word in the sequence an adverse effect of the given drug?”. Thus we have loosened the constraints, to change a relation extraction problem, to a sequential binary classification problem.

$$p(y_i^{ADE} | X_{seq}, X_{drug}) = f_2(x_i, X_{seq}, X_{drug}) \quad (2)$$

### IV. MODEL DEFINITION

Our model can be partitioned into 4 stages - *embedding*, *encoding*, *interaction* and *prediction*. We use pretrained word

TABLE I  
ARCHITECTURE OF CHARCNN

Filter Attributes	Values
Embedding Dimensions	25
Filter widths ( $w$ )	[1,2,3,4,5,6]
Number of filter matrices	[25. $w$ ]

embeddings to obtain distributed, continuous representation of tokenized sequences. Character-level word representations are obtained by using a combination of different sized convolutional filters, similar to Kim et al.’s work [24]. The word-level and character-level representations are concatenated and fed as input to the encoder. The encoder is a bi-directional LSTM, which processes the embedded sequence in either directions. The forward and backward hidden states of the LSTM are concatenated together to get an encoded sequence  $\{\{\vec{h}_i, \overleftarrow{h}_i\}\}$ . We use the same encoder to obtain an encoded representation of the drug. In the interaction stage, we iterate through the timesteps of the sequence and obtain a weighted representation of the whole encoded sequence, conditioned on current state. We combine the weighted representation with the current state to predict the entity type of the token corresponding to current state. We apply an affine transformation to the combined representation and then combine it with the encoded drug, to predict if the token corresponding to current state is an ADE of the given drug. Each stage is discussed in detail in the following sections.

#### A. Embedding

Distributed representation of words introduced by Bengio et al., [15] has become competent a replacement for the traditional bag-of-words encoding technique. Similar to [3], we adopt multiple embedding matrices to represent different representations of text sequence. Word-level and Character-level representations of input sequence are obtained by techniques explained in sections IV-A.1 and IV-A.2 respectively. In addition to these, we make use of Parts of Speech (PoS) and entity label information, by creating PoS embeddings and label embeddings. Both of these embedding matrices are randomly initialized.

1) *Word Embedding*: We use two instances of pretrained word embeddings from PubMed word2vec [16]. We disable gradient updates on the first instance - *Fixed Embedding* ( $emb_f$ ) while we enable gradient updates on the second instance - *Variable Embedding* ( $emb_v$ ). Fixed embeddings provide good, stable representations for drug and disease entities in the sequence. Variable embeddings learn representations for common words and phrases in the corpus that highlight their role as linguistic supporting structures. We concatenate the fixed and variable embeddings for word-level representation of sequence tokens.

$$emb(x) = [emb_f(x); emb_v(x)] \quad (3)$$

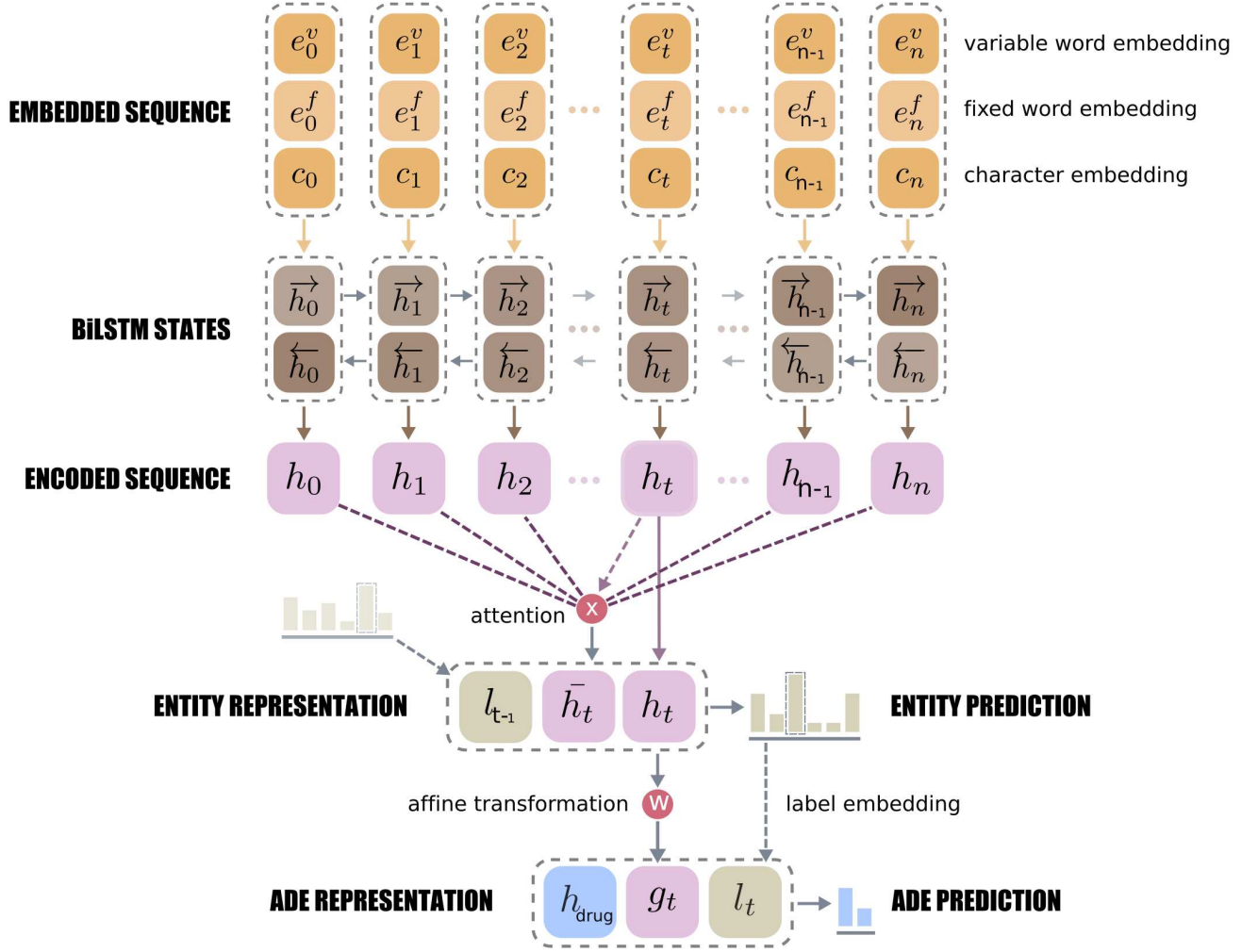


Fig. 2. Architecture of Character-aware Attentive ADENet

2) *Character-level Embedding:* Convolutional Neural Networks (CNNs) have been shown to perform well on NLP tasks. [24] and [25] have used CNNs for learning character-level features in text. Kim et al., [24] propose a character-aware neural language model, that learns character-level word representations using CNNs - *CharCNN*. They employ multiple convolutional filters of varying widths to obtain feature vectors for words in the sequence. We follow the same mechanism for character embedding.

Initially the text sequence is split into a nested sequence of characters. A lookup of character embeddings is performed and the vectors corresponding to individual characters in a word, are stacked together to form a matrix. This is followed by a series of convolutions with multiple filters of different widths. Max-over-time pooling is applied to the output of convolutional layer, to obtain character-level word embeddings for each word in the sequence. The details of filter size, number of filters and embedding size are tabulated in Table I.

### B. Encoder

The embedded text sequence which is comprised of word embedding (fixed and variable), character embedding and PoS feature embedding, is processed by the encoder, a bidirectional LSTM based RNN. The forward  $\{\vec{h}_i\}$  and backward  $\{\overleftarrow{h}_i\}$  hidden states of the RNN are concatenated together. Now we have a primary representation of the text sequence - the encoded sequence  $\{h_i\}$ . The drug, in its embedded form, is also processed by the same encoder network. The final state of the RNN is considered the encoded drug representation  $h_{drug}$ .

$$h_t = [\vec{h}_t; \overleftarrow{h}_t] \quad (4)$$

$$h_{drug} = [\vec{h}_{drug}; \overleftarrow{h}_{drug}] \quad (5)$$

### C. Interaction Layer

During each step  $t$  of encoder, information from the local neighborhood flows from either directions. The hidden state  $\bar{h}_t$  represents the local context corresponding to  $t^{th}$  time

step. In [3], ADE relation extraction is done by building SDPs between drug entities and potential ADE candidates in the sequence. We have designed the interaction layer as an alternative to this approach. The objective is to facilitate interaction between different parts of the sequence. Our model considers each token in the sequence as a potential entity and an ADE candidate.

In the *interaction* stage, we iterate through the encoded sequence and learn a entity representation and an ADE representation for each sequence. During each time step  $t$ , the model creates a weighted representation over all the encoded states  $\{\bar{h}_i\}$ , conditioned on the current state  $h_t$ . We choose multiplicative attention mechanism proposed in [22], over additive attention [21], since the former is faster and more space-efficient [23].

$$\bar{h}_i = \sum_j a_{ij} h_j \quad (6)$$

$$a_i = \text{softmax}(f_{att}(h_i, h_j)) \quad (7)$$

$$f_{att}(h_i, h_j) = h_i^T W_a h_j \quad (8)$$

#### D. Prediction

At each step  $t$  of the interaction layer, we construct entity and ADE representations corresponding to token  $t$ . For entity representation, we consider the embedding of predicted label of the previous token,  $l_{t-1}$ . We concatenate previous label embedding  $l_{t-1}$  with weighted representation of encoded sequence  $\bar{h}_t$  and the current state  $h_t$ . This gives us the entity representation  $[l_{t-1}; \bar{h}_t; h_t]$ . A *tanh* non-linearity is applied, followed by an affine transformation. The output of this operation is normalized with softmax, which provides a probability distribution over the entity labels.

$$p(y_t^{ent} | X_{seq}) = \text{softmax}(\text{tanh}(W^{ent} r^{ent} + b^{ent})) \quad (9)$$

$$r^{ent} = [l_{t-1}; \bar{h}_t; h_t] \quad (10)$$

We apply an affine transformation over  $r^{ent}$  to get  $\bar{r}^{ent}$  to enable variation between entity and ADE representation. We concatenate encoded drug state  $h_{drug}$  with  $\bar{r}^{ent}$  and the embedding of label predicted at current step  $l_t$ , to obtain the ADE representation,  $r^{ade}$ . By applying another affine transformation over  $r^{ade}$ , followed by softmax, we get the probability distribution over ADE labels, 0 or 1 corresponding to *not an ADE* and *ADE*.

$$p(y_t^{ade} | X_{seq}, X_{drug}) = \text{softmax}(W^{ade} r^{ade} + b^{ade}) \quad (11)$$

$$r^{ade} = [h_{drug}; \bar{r}^{ent}; l_t] \quad (12)$$

$$\bar{r}^{ent} = W^c r^{ent} + b^c \quad (13)$$

$$l_t = \text{emb}_l(\hat{y}_t) \quad (14)$$

$$\hat{y}_t = \text{argmax}(p(y_t^{ent} | X_{seq})) \quad (15)$$

TABLE II  
COMPARISON WITH THE STATE OF THE ART

Model	Entity Recognition			ADE Extraction		
	P	R	F1	P	R	F1
Li [3]	82.70	86.70	84.60	67.50	75.80	71.40
Our Model	88.41	82.41	85.30	86.28	87.29	86.78

#### V. MODEL VARIANTS AND FEATURES

We have experimented with different variants of the model architecture presented in figure 2. The performance of these models are tabulated in Table III. We started with a baseline model - **ADENet**, which consists of an LSTM-based bidirectional encoder. The final state of encoder which processes the drug, is taken as encoded drug representation. At each step of encoding the sequence, the hidden state  $\bar{h}_t$  is combined with the embedding of entity label  $\hat{y}_{t-1}^{entity}$  predicted in the previous step. This becomes the entity representation corresponding to token  $X_t^{seq}$ . An affine transformation, followed by a *tanh* non-linearity gives the logits for entity classification. By combining the current hidden state  $\bar{h}_t$  with label embedding of current entity prediction  $\hat{y}_t^{entity}$  and the encoded drug  $\bar{h}^{drug}$ , we obtain the ADE representation of  $X_t^{seq}$ . Likelihood of ADE is obtained by applying softmax over affine transformed ADE representation.

Parts of speech tags of the text sequence are obtained using *nlk*'s [17] PoS tagger. PoS embeddings of the sequence are concatenated with the word embeddings to get a richer primary sequence representation. This model - *Baseline + PoS Features* outperforms the baseline by a small margin.

The Interaction layer, defined in section IV-C, is added to the baseline - *Baseline + Attention*. We can observe a sharp increase in the performance of this model, in both the tasks, compared to the baseline. This is due to the fact that during each step of prediction, the model has access to the encoded representation of the whole sequence. By conditioning the weighted representation on the current state, the model performs a search on the sequence, looking for information relevant to current prediction. The hidden state bottleneck [21] is overcome by allowing the model to peek at the whole sequence.

When PoS features are added to this model - *Baseline + Attention + PoS Features*, there is an increase in performance on both the tasks. We theorize that the attention mechanism learns to make better use of new information provided by the additional features (PoS). This theory holds true when we add character embeddings to the model.

We then introduce character embedding as an additional feature to the model - *Baseline + Attention + Character Embedding*. This results in a much better performance. The character embeddings contain rich morphological information, which are crucial in differentiating a drug or disease

TABLE III  
PERFORMANCE COMPARISON OF MODEL AUGMENTED WITH VARIOUS FEATURES

Model	Entity Recognition			ADE Extraction		
	P	R	F1	P	R	F1
Baseline	79.91	67.62	72.25	75.30	80.27	77.70
Baseline + PoS Features	79.93	68.73	73.90	76.76	81.14	78.89
Baseline + Attention	83.57	74.29	78.65	81.14	81.49	81.31
Baseline + Attention + PoS Features	84.23	75.35	79.53	81.93	83.14	82.53
Baseline + Attention + Character Embedding	82.69	77.74	80.13	79.35	86.20	82.62
Baseline + Attention + Character Embedding + PoS Features	88.41	82.41	85.30	86.28	87.29	86.78

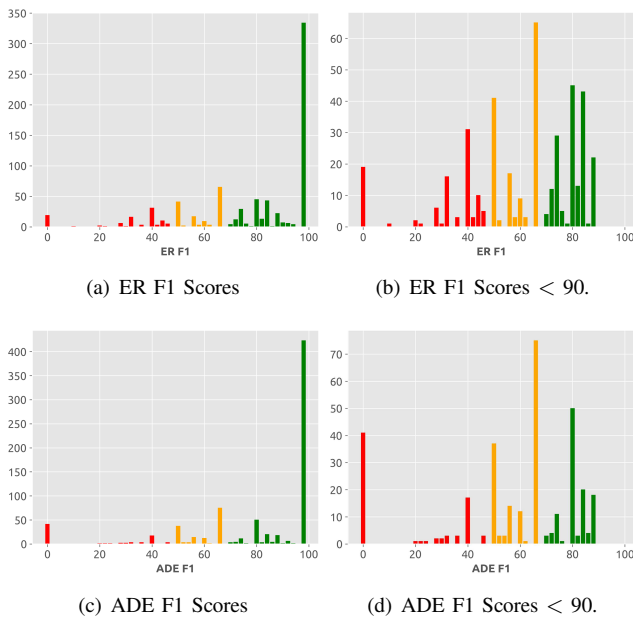


Fig. 3. Histogram of F1 scores in Entity Recognition and ADE Extraction tasks are presented in figures (a) and (c) respectively. (b) and (d) ignore F1 scores greater than 90.

entity from a regular English word.

And finally when we add PoS features to this model - *Baseline + Attention + Character Embedding + PoS Features*, we obtain our best performing model by far (illustrated in figure 2). Based on these observations, we add that by including domain-specific or context-specific features and facilitating interaction between the features corresponding to individual tokens in the sequence, the performance on NLP tasks could be significantly improved.

## VI. TRAINING

We train our models with an initial learning rate of 0.001 and a batch size of 16, with Adagrad Optimizer [18]. We have split the dataset in 8:1:1 ratio into training set, test set and validation set. We use a dropout [19] of 0.5 for the LSTM encoder. All the trainable parameters of the model are initialized with a random uniform initializer within range (-0.01, 0.01). All our models converge within 10 epochs. The dimensions of pretrained PubMed word embeddings is 200. Dimensions of character embedding, PoS embedding and label embedding, along with the rest of the hyperparameter settings, are listed in Table IV.

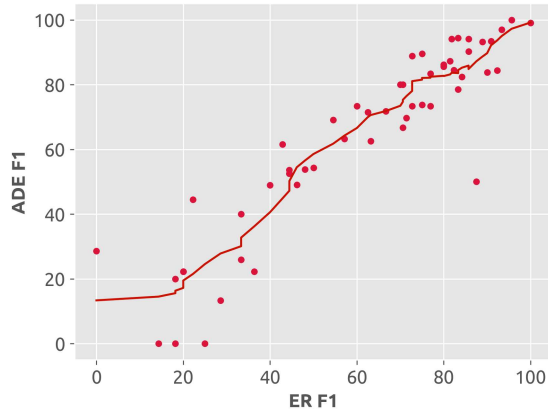


Fig. 4. Correlation between F1 scores of ER task and ADE task

TABLE IV  
HYPERPARAMETER SETTINGS

Hyperparameter	Value
Word Embedding	400
Character Embedding	525
PoS Embedding	25
Label Embedding	25
BiLSTM hidden dimensions	[150, 150]
Learning rate	0.01

## VII. EXPERIMENT RESULTS

Table II compares the performance of our best performing model with that of Li et al.'s state of the art model [3]. Our model's performance in Entity Recognition is slightly higher compared to the state of the art. We also present the results of our ADE extraction system. It is unfair to compare its performance with that of [3], as ADE extraction is treated as a relation extraction problem in the latter. Although, in a practical setting, we could augment an end-to-end neural network with a drug lexicon, to condition the model to extract ADE's caused by a specific drug. By constraining our model to focus on ADE's only relevant to a particular drug, we have achieved an F1 score of 86.78, which is an improvement of approximately 15% over the state of the art.

## VIII. ANALYSIS

Our analysis is focused on predictions made by our best performing model which includes Character Embedding, PoS features and Attention mechanism. We have clustered the

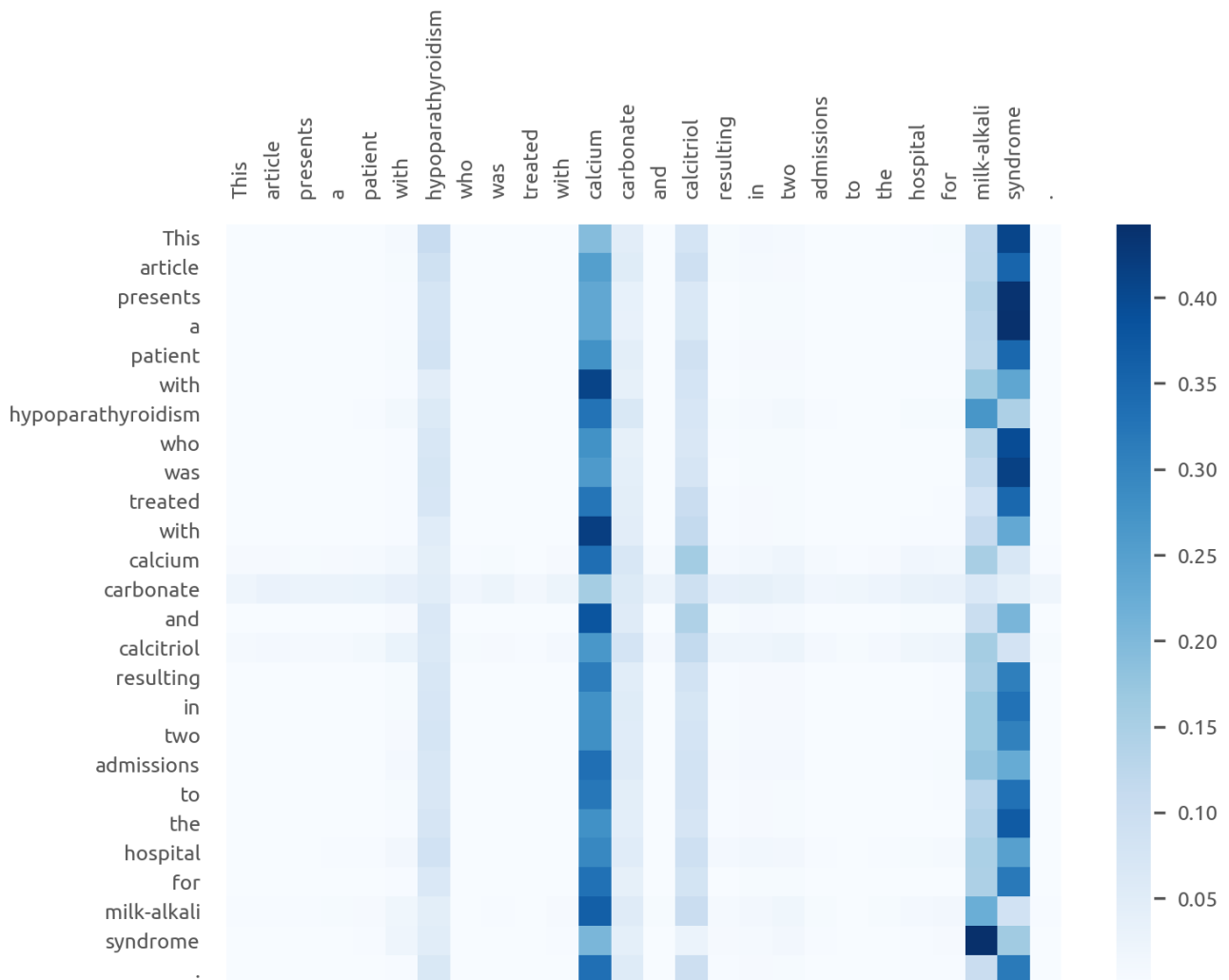


Fig. 5. Visualization of Attention scores. To understand the figure, the sequence on the left should be read from top to bottom. Every row is a probability distribution (normalized attention score) across the whole sequence conditioned on the token on the left. Notice the peak in attention in positions corresponding to drug or disease entities.

data samples into three categories based on our model’s performance on the sample. We studied the properties of each category of samples and present our observations in this section. The histograms of F1 scores of Entity Recognition and ADE extraction tasks are presented in figure 3.

We can clearly observe that more than 60% of samples achieve an F1 score greater than 95 in ADE extraction task and almost 50% of samples achieve an F1 score greater than 95 in Entity Recognition task. Figures 3(b) and 3(d) provide a closer look at the rest of the samples. The performance of the rest of the samples seem to be almost evenly distributed between scores of 20 and 90, with a few outliers. We can also notice a sharp peak around the score of zero.

We further investigate these samples by studying the effect of text sequence length, number of entities in the sequence and number of ADEs in the sequence, on performance.

Based on which we conclude that there is no observable correlation between sequence length and performance. A similar comparison based on number of entities and number of drugs in the sequence, does not reveal any significant distinguishing features that separate good samples from bad samples. We then, compare the performance of our model in individual tasks. It is clear from figure 4 that there is an increase in ADE performance with increase in ER performance. This kind of correlation is to be expected in a joint model, where sub-systems designed to solve different problems share parameters.

As discussed in section IV-C, we employ an attention mechanism across all the encoded states, corresponding to word tokens in the text sequences. This mechanism allows the system at each prediction step, to focus on information from the sequence relevant to current prediction. We have

mapped the attention scores  $a_i$  obtained from equation 7, to the tokens in the text sequence and present a sample visualization of the interaction layer in figure 5. We can observe clear peaks in attention scores, corresponding to drug and disease entities. This pattern is observed in every good sample (high performance in both the tasks) from the test set. To add to that, in most bad samples, the attention seems to be scattered across the tokens in the sequence. This could be interpreted as the model searching the encoder states, for information that could be useful for prediction and failing to find it.

## IX. CONCLUSION

In this paper, we propose a joint model for Entity Recognition and Adverse Drug Event extraction in biomedical text. By loosening the constraints on ADE extraction, our model outperforms the state of the art by a large margin. The performance of our Entity Recognition sub-system is slightly better than that of state of the art. We observe that by using an attention mechanism to facilitate intra-sequence interaction, we could replace SDP based methods used in prior work. Additionally, the use of attention mechanism allows us to investigate the dependence of classifier on different segments of the text sequence.

## REFERENCES

- [1] Gurulingappa, H., Rajput, A.M., Roberts, A., Fluck, J., Hofmann-Apitius, M. and Toldo, L., 2012. Development of a benchmark corpus to support the automatic extraction of drug-related adverse effects from medical case reports. *Journal of biomedical informatics*, 45(5), pp.885-892.
- [2] Kang, N., Singh, B., Bui, C., Afzal, Z., van Mulligen, E.M. and Kors, J.A., 2014. Knowledge-based extraction of adverse drug events from biomedical text. *BMC bioinformatics*, 15(1), p.64.
- [3] Li, F., Zhang, M., Fu, G. and Ji, D., 2017. A neural joint model for entity and relation extraction from biomedical text. *BMC bioinformatics*, 18(1), p.198.
- [4] Gupta, S., Pawar, S., Ramrakhiyani, N., Palshikar, G. and Varma, V., 2017. Semi-Supervised Recurrent Neural Network for Adverse Drug Reaction Mention Extraction. *arXiv preprint arXiv:1709.01687*.
- [5] Hazell, L. and Shakir, S.A., 2006. Under-reporting of adverse drug reactions. *Drug safety*, 29(5), pp.385-396.
- [6] Nikfarjam, A., Sarker, A., OConnor, K., Ginn, R. and Gonzalez, G., 2015. Pharmacovigilance from social media: mining adverse drug reaction mentions using sequence labeling with word embedding cluster features. *Journal of the American Medical Informatics Association*, 22(3), pp.671-681.
- [7] Lee, K., Qadir, A., Hasan, S.A., Datla, V., Prakash, A., Liu, J. and Farri, O., 2017, April. Adverse Drug Event Detection in Tweets with Semi-Supervised Convolutional Neural Networks. In *Proceedings of the 26th International Conference on World Wide Web* (pp. 705-714). International World Wide Web Conferences Steering Committee.
- [8] Huynh, T., He, Y., Willis, A. and Rger, S., 2016. Adverse drug reaction classification with deep neural networks. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers* (pp. 877-887).
- [9] Gurulingappa, H., MateenRajpu, A. and Toldo, L., 2012. Extraction of potential adverse drug events from medical case reports. *Journal of biomedical semantics*, 3(1), p.15.
- [10] Li, Q. and Ji, H., 2014, June. Incremental Joint Extraction of Entity Mentions and Relations. In *ACL (1)* (pp. 402-412).
- [11] Miwa, M. and Bansal, M., 2016. End-to-end relation extraction using lstms on sequences and tree structures. *arXiv preprint arXiv:1601.00770*.
- [12] Hochreiter, S. and Schmidhuber, J., 1997. Long short-term memory. *Neural computation*, 9(8), pp.1735-1780.
- [13] Li, F., Zhang, Y., Zhang, M. and Ji, D., 2016, July. Joint Models for Extracting Adverse Drug Events from Biomedical Text. In *IJCAI* (pp. 2838-2844).
- [14] Xu, Y., Mou, L., Li, G., Chen, Y., Peng, H. and Jin, Z., 2015, September. Classifying Relations via Long Short Term Memory Networks along Shortest Dependency Paths. In *EMNLP* (pp. 1785-1794).
- [15] Bengio, Y., Ducharme, R., Vincent, P. and Jauvin, C., 2003. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb), pp.1137-1155.
- [16] 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* (pp. 3111-3119).
- [17] Bird, S., 2006, July. NLTK: the natural language toolkit. In *Proceedings of the COLING/ACL on Interactive presentation sessions* (pp. 69-72). Association for Computational Linguistics.
- [18] Duchi, J., Hazan, E. and Singer, Y., 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul), pp.2121-2159.
- [19] Hinton, G.E., Srivastava, N., Krizhevsky, A., Sutskever, I. and Salakhutdinov, R.R., 2012. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*.
- [20] Hochreiter, S. and Schmidhuber, J., 1997. Long short-term memory. *Neural computation*, 9(8), pp.1735-1780.
- [21] Bahdanau, D., Cho, K. and Bengio, Y., 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- [22] Luong, M.T., Pham, H. and Manning, C.D., 2015. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- [23] Ruder, S., 2017. Deep Learning for NLP Best Practices. In: Sebastian Ruder's Blog. 25-07-2017. Available from: <http://ruder.io/deep-learning-nlp-best-practices/index.html>
- [24] Kim, Y., Jernite, Y., Sontag, D. and Rush, A.M., 2016, February. Character-Aware Neural Language Models. In *AAAI* (pp. 2741-2749).
- [25] Zhang, X. and LeCun, Y., 2015. Text understanding from scratch. *arXiv preprint arXiv:1502.01710*.