# Multi-Level Structured Self-Attentions for Distantly Supervised Relation Extraction

**Jinhua Du[†], Jingguang Han[§], Andy Way[†], Dadong Wan[§]**
[†]ADAPT Centre, School of Computing, Dublin City University, Ireland
[§]Accenture Labs Dublin, Ireland
{jinhua.du, andy.way}@adaptcentre.ie
{jingguang.han, dadong.wan}@accenture.com

## Abstract

Attention mechanisms are often used in deep neural networks for distantly supervised relation extraction (DS-RE) to distinguish valid from noisy instances. However, traditional 1-$D$ vector attention models are insufficient for the learning of different contexts in the selection of valid instances to predict the relationship for an entity pair. To alleviate this issue, we propose a novel multi-level structured (2-$D$ matrix) self-attention mechanism for DS-RE in a multi-instance learning (MIL) framework using bidirectional recurrent neural networks. In the proposed method, a structured word-level self-attention mechanism learns a 2-$D$ matrix where each row vector represents a weight distribution for different aspects of an instance regarding two entities. Targeting the MIL issue, the structured sentence-level attention learns a 2-$D$ matrix where each row vector represents a weight distribution on selection of different valid instances. Experiments conducted on two publicly available DS-RE datasets show that the proposed framework with a multi-level structured self-attention mechanism significantly outperform state-of-the-art baselines in terms of PR curves, P@N and F1 measures.

## 1 Introduction

Relation extraction is a fundamental task in information extraction (IE), which studies the issue of predicting semantic relations between pairs of entities in a sentence (Zelenko et al., 2003; Bunescu and Mooney, 2005; Zhou et al., 2005). One crucial problem in RE is the relative lack of large-scale, high-quality labeled data. In recent years, one commonly used and effective technique for dealing with this challenge is the distant supervision method via knowledge bases (KBs) (Mintz et al., 2009; Riedel et al., 2010; Hoffmann et al., 2011), which assumes that if one entity pair appearing in some sentences can be observed in a KB with a certain relationship, then these sentences will be labeled as the context of this entity pair and this relationship. The distant supervision strategy is an effective and efficient method for automatically labeling large-scale training data. However, it also introduces a severe mislabelling problem due to the fact that a sentence that mentions two entities does not necessarily express their relation in a KB (Surdeanu et al., 2012; Zeng et al., 2015).

Plenty of research work has been proposed to deal with distantly supervised data and has achieved significant progress, especially with the rapid development of deep neural networks (DNN) for relation extraction in recent years (Zeng et al., 2014, 2015; Lin et al., 2016, 2017a; Wang et al., 2016; Zhou et al., 2016; Ji et al., 2017; Yang et al., 2017; Zeng et al., 2017). DNN models under an MIL framework for DS-RE have become state-of-the-art, replacing statistical methods, such as feature-based and graphical models (Riedel et al., 2010; Hoffmann et al., 2011; Surdeanu et al., 2012). In the MIL framework for distantly supervised RE, each entity pair often has multiple instances where some are noisy and some are valid. The attention mechanism in DNNs, such as convolutional (CNN) and recurrent neural networks (RNN), is an effective way to select valid instances by learning a weight distribution over multiple instances. However, there are two important representation learning problems in DNN-based distantly supervised RE: (1) **Problem I**: entity pair-targeted context representation learning from an instance; and (2) **Problem II**: valid instance selection representation learning over multiple instances. The former can use a word-level attention mechanism to learn a weight distribution on words and then a weighted sentence representation regarding two entities; the latter can employ a sentence-level attention mechanism to

learn a weight distribution on multiple instances so that valid sentences with higher weights can be focused and selected, and noisy instances with lower weights are suppressed.

Both the word-level and sentence-level attention mechanisms in previous work on the RE task are simple 1-$D$ vectors which are learned using the hidden states of the RNN, or via pooling from either the RNNs' hidden states or convolved $n$-grams (Zeng et al., 2014, 2015; Zhou et al., 2016; Wang et al., 2016; Ji et al., 2017; Yang et al., 2017). The deficiency of the 1-$D$ attention vector is that it only focuses on one or a small number of aspects of the sentence, or one or a small number of instances (Lin et al., 2017b), with the result that different semantic aspects of the sentence, or different multiple valid sentences are ignored, and cannot be utilised.

Inspired by the structured self-attentive sentence embedding in Lin et al. (2017b), we propose a novel multi-level structured (2-$D$) self-attention mechanism (MLSSA) in a bidirectional LSTM-based (BiLSTM) (Hochreiter and Schmidhuber, 1997) MIL framework to alleviate two **problems** in the distantly supervised RE. Regarding **Problem I**, we propose a 2-$D$ matrix-based word-level attention mechanism, which contains multiple vectors, each focusing on different aspects of the sentence for better context representation learning. In terms of **Problem II**, we propose a 2-$D$ sentence-level attention mechanism for multiple instance learning, where it contains multiple vectors, each focusing on different valid instances for a better sentence selection. "**structured**" indicates that the weight vectors in the learned 2-$D$ matrix try to construct a structural dependency relationship by learning different weight distributions for different contexts or instances given the entity pair. We can see that our structured attention mechanism is different from that in Kim et al. (2017) which incorporates richer structural distributions and are simple extensions of the basic attention procedure. We verify the proposed framework on two distantly supervised RE datasets, namely the New York Times (NYT) dataset (Riedel et al., 2010) and the DBpedia Portuguese dataset (Batista et al., 2013). Experimental results show that our MLSSA framework significantly outperforms state-of-the-art baseline systems in terms of different evaluation metrics.

The main contributions of this paper include:

(1) we propose a novel multi-level structured (2-$D$) self-attention mechanism for DS-RE which can make full use of input sequences to learn different contexts, without integrating extra resources; (2) we propose a 2-$D$ matrix-based word-level attention for better context representation learning targeting two entities; (3) we propose a 2-$D$ sentence-level attention mechanism over multiple instances to select different valid instances; and (4) we verify the proposed framework on two publicly available distantly supervised datasets.

## 2 Related Work

Most existing work on distant supervision data mainly focuses on denoising the data under the MIL strategy by learning a valid sentence representation or features, and then selecting one or more valid instances for relation classification (Riedel et al., 2010; Hoffmann et al., 2011; Surdeanu et al., 2012; Zeng et al., 2015; Lin et al., 2016, 2017a; Zhou et al., 2016; Ji et al., 2017; Zeng et al., 2017; Yang et al., 2017).

Riedel et al. (2010) and Surdeanu et al. (2012) use a graphical model and MIL to select the valid sentences and classify the relations. However, these models are based on statistical methods and feature engineering, i.e. extracting sentence features using other NLP tools. Zeng et al. (2015) proposed a piece-wise CNN (PCNN) method to automatically learn sentence-level features and select one valid instance for the relation classification. The one-sentence-selection strategy does not make full use of the supervision information among multiple instances.

Lin et al. (2016) and Ji et al. (2017) introduce an attention mechanism to the PCNN-based MIL framework to select informative sentences, which outperforms all baseline systems on the NYT data set. However, their attention mechanism is only a sentence-level model without incorporating word-level attention. Zhou et al. (2016) introduce a word-level attention model to the BiLSTM-based MIL framework and obtain significant improvements on the SemEval2010 (Hendrickx et al., 2010) data set. Wang et al. (2016) extend the single word-level attention model to multiple word levels in CNNs to discern patterns in heterogeneous contexts of the input sentence, and achieve best performance on the SemEval2010 data set. However, these two works were not targeting the distantly supervised RE problem.

Yang et al. (2017) experiment with word-level and sentence-level attention models in the bidirectional RNN on the NYT dataset on the basis of the open source DS-RE system,[1] and verify that a two-level attention mechanism achieves best performance compared to PCNN/CNN models. Both the word-level and sentence-level attention models are 1-$D$ vectors.

From previous work, we can see that the attention mechanism in DNNs has made significant progress on the RE task. However, both word-level and sentence-level attention models are still based on 1-$D$ vectors which have the following insufficiencies: (1) although the 1-$D$ attention model can learn weights for different contexts, it only focuses on one or very few aspects of a single sentence (Lin et al., 2017b), or one or very few instances; (2) in order to allow the attention mechanism to learn more aspects of the sentence, or different instances, extra knowledge needs to be integrated, such as the work in Ji et al. (2017) and Lin et al. (2017a). The former integrates entity descriptions generated from Freebase and Wikipedia as supplementary background knowledge to disambiguate the entity. The latter introduces a multilingual framework which employs a monolingual attention mechanism to utilize the information within monolingual texts, and further uses a cross-lingual attention mechanism to consider the information consistency and complementarity among cross-lingual texts. However, extra resources are difficult to obtain in many practical scenarios.

In order to alleviate the burden of integrating extra knowledge, and make full use of the input sentence (i.e. learning different aspects of context and focusing on different valid instances), we propose a multi-level structured self-attention mechanism in a BiLSTM-based MIL framework without integrating extra resources.

## 3 Approach

The distantly supervised RE can be formalised as follows: given an entity pair $(e_1, e_2)$, a bag $\mathcal{G}$ containing $\mathcal{J}$ instances, and the relation label $r$ for $\mathcal{G}$, the goal of the training process is to denoise these instances by selecting valid candidates based on $r$, and the goal of the testing process is to denoise multiple instances by selecting valid candidates to

predict the relation $r$ for $\mathcal{G}$.

To alleviate the aforementioned two problems, improving the following two representation learning issues is clearly important for a DNN-based RE classifier:

- *Entity pair-targeted context representation*: The model should have the capability to learn a better context representation from the input sentence targeting the entity pair;

- *Instance selection representation*: The model should have the capability to learn a better weight distribution over multiple instances to select valid instances regarding an entity pair.

Motivated by these two issues, we propose a multi-level structured self-attention framework.

### 3.1 Architecture

The proposed framework consists of three parts as shown in Figure 1. The first part includes the input layer, embedding layer and BiLSTM layer which transform the input sequence at different time steps to LSTM hidden states.

The second part implements the entity pair-targeted context representation learning, including:

- *a structured word-level self-attention layer*: this generates a set of summation weight vectors (or a 2-$D$ matrix) taking the LSTM hidden states as input. Each vector in the 2-$D$ matrix represents the weights for different aspects of the input sentence.

- *a structured context representation layer*: the weight vectors learned by the 2-$D$ word-level self-attention are dotted with the BiLSTM hidden states. Accordingly, a 2-$D$ matrix or a set of weighted LSTM hidden state vectors, denoted as "$M_{L1}$" in Figure 1, is obtained. Each weighted vector represents a sentence embedding reflecting a different aspect of the sentence targeting the entity pair. By this means, a dependency parsing-like structure of the input sentence can be constructed, obtaining different semantic representations of the sentence for the two entities in question.

- *a flattened representation layer*: this concatenates each vector in the 2-$D$ matrix of the sentence embedding to one vector. Then, the

**MLSSA-ATT-2**
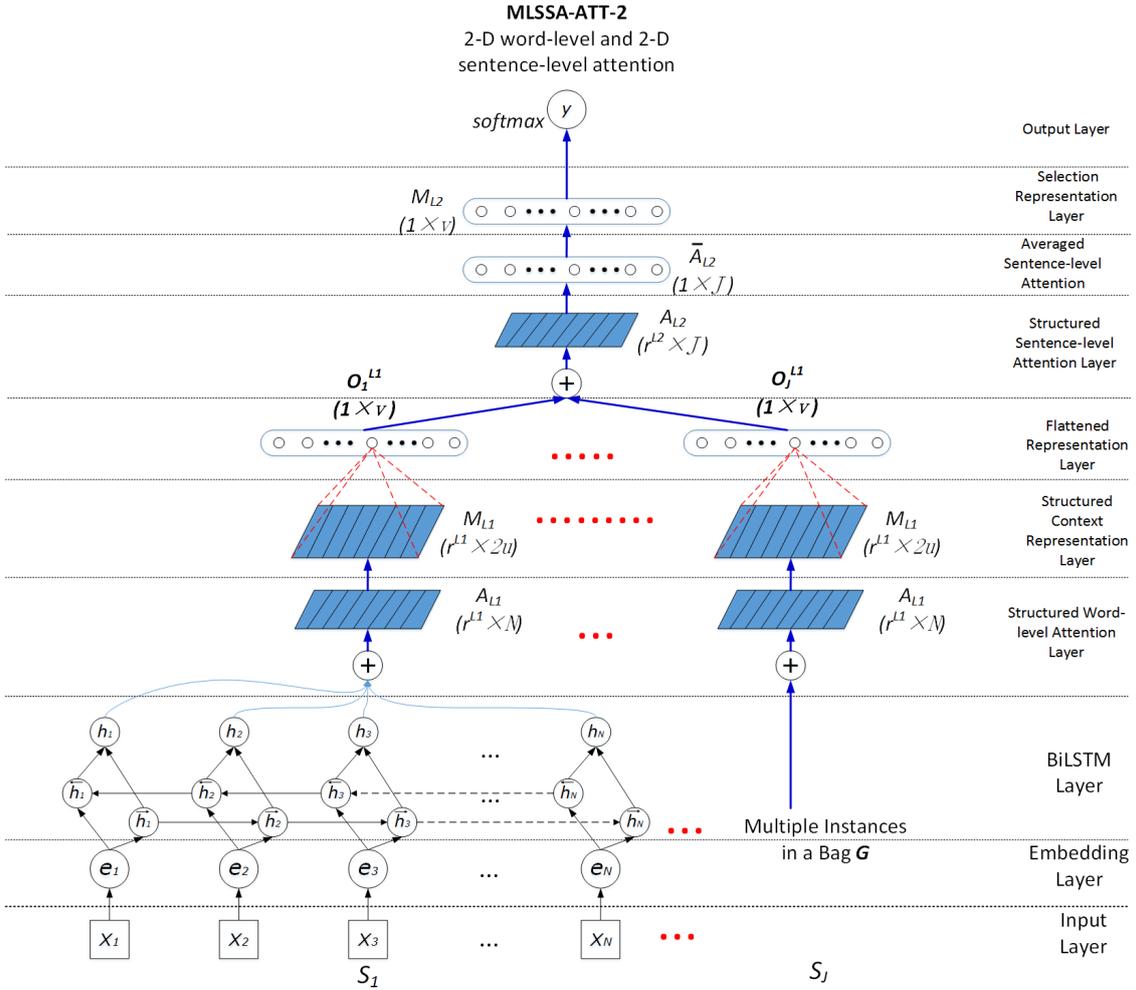2-D word-level and 2-D
sentence-level attention

Figure 1: Multi-level structured self-attention framework for distantly supervised RE

flattened vector connects to a 1-layer multi-layer perceptron (MLP) with ReLU activation function, generating an aggregated sentence representation.

The first and second parts operate on the single instance level, i.e. given a bag $\mathcal{G}$ and feeding each instance into the framework, the structured word-level self-attention mechanism will construct $J$ individual structured sentence representations corresponding to $J$ input instances.

The third part targets the instance selection representation learning issue, and operates on the bag level, i.e. considering weighted context representations of all instances in the bag $\mathcal{G}$ and learning probability distributions to distinguish informative from noisy sentences. This part includes:

- *a structured sentence-level attention model*: this has a similar structure to the structured word-level attention mechanism, except that

it generates a set of summation weight vectors for all input instances in the same bag $\mathcal{G}$. Each vector is a weight distribution over all instances. Accordingly, the 2-$D$ sentence-level matrix is expected to learn a set of different weight distributions focusing on different informative instances. As a result, informative sentences are expected to contribute more with higher weights, and noisy sentences are expected to contribute less with smaller weights, to the relation classification.

- *an averaged sentence-level attention layer*: the 2-$D$ sentence-level attention matrix is averaged and converted to a 1-$D$ vector.

- *a selection representation layer*: the 1-$D$ averaged attention vector is dotted with the output of the flattened representation layer. Accordingly, a 1-$D$ vector, denoted as "$M_{L2}$" in Figure 1, is obtained which represents an av-

eraged weighted selection representation of multiple sentences.

- *an output layer*: this connects to a *softmax* layer and produces a probability distribution corresponding to relation classes.

## 3.2 Structured Word-Level Self-Attention and its Penalisation Function

Given a bag $\mathcal{G} = (S_1, S_2, \ldots, S_J)$ containing $J$ instances, and a sentence $S_j$ in $\mathcal{G}$ consisting of $N$ tokens, $S_j$ can be represented using a sequence of word embeddings, as in (1):

$$S_j = (e_1, e_2, \ldots, e_N) \qquad (1)$$

where $e_i$ is a $d$-dimension vector for the $i$-th word, and $S_j$ is the $j$-th instance in $\mathcal{G}$.

We denote the hidden state of the BiLSTM as in (2):

$$H = (\mathbf{h_1}, \mathbf{h_2}, \ldots, \mathbf{h_N})^T \qquad (2)$$

where $\mathbf{h_t}$ is a concatenation of the forward hidden state $\overrightarrow{h}_t$ and the backward hidden state $\overleftarrow{h}_t$ at time step $t$. $T$ is the transpose operation. If the size of each unidirectional LSTM is $u$, then $H$ has the size $2u$-by-$N$.

Then, the structured word-level self-attention mechanism is defined as in (3):

$$A_{L1} = softmax(W_{s2}^{L1} tanh(W_{s1}^{L1} H)) \qquad (3)$$

where $L1$ stands for the first-level attention mechanism, i.e. the word-level; $W_{s1}^{L1}$ is a weight matrix of size $d_a^{L1} \times 2u$, where $d_a^{L1}$ is a hyper-parameter for the number of neurons in the attention network; $W_{s2}^{L1}$ is a weight matrix with the shape $r^{L1} \times d_a^{L1}$, where $r^{L1}$ ($r^{L1} > 1$) is the hyper-parameter representing the size of multiple vectors in the 2-$D$ attention matrix. The size of $r^{L1}$ is defined based on how many different aspects of the sentence need to be focused on; $A_{L1}$ is the annotation matrix of size $r^{L1} \times N$. We can see that in $A_{L1}$, there are $r^{L1}$ attention vectors for the $N$-token input sentence.

Finally, we compute the $r^{L1}$ weighted sums by multiplying the annotation matrix $A_{L1}$ and BiL-STM hidden states $H$. The resulting structured sentence representation $M_{L1}$ is (4):

$$M_{L1} = A_{L1} H^T \qquad (4)$$

where $M_{L1}$ has the shape $r^{L1} \times 2u$. It can be seen that the traditional 1-$D$ sentence representation is extended to a 2-$D$ representation ($r^{L1} > 1$).

Subsequently, the output of the flattened representation layer for the instance $S_j$ in $\mathcal{G}$ is (5):

$$O_j^{L1} = ReLU(W_o^{L1} M_{L1}^{FT} + b_o^{L1}) \qquad (5)$$

where $W_o^{L1}$ is the weight matrix that has the shape $v$-by-$r^{L1} * 2u$, where $v$ is the amount of neurons in the $ReLU$-based MLP layer; $M_{L1}^{FT}$ is the flattened structured sentence representation which is a concatenated vector of each row in $M_{L1}$ and has the dimension $r^{L1} * 2u$; $b_o^{L1}$ is the bias vector of size $v$; $O_j^{L1}$ is the aggregated sentence representation of the $j$-th instance in the bag $\mathcal{G}$ with size $v$.

Then, the output of all instances in $\mathcal{G}$ from the flattened representation layer is denoted as in (6):

$$O^{L1} = (O_1^{L1}, O_2^{L1}, \ldots, O_J^{L1})^T \qquad (6)$$

where $O^{L1}$ has the shape of $v \times J$.

As in Lin et al. (2017b), the penalisation term for the structured word-level attention is as in (7):

$$P_{L1} = ||(A_{L1} A_{L1}^T - I)||_F^2 \qquad (7)$$

where $|| \cdot ||_F$ is the Frobenius norm of a matrix. $I$ is an identity matrix. Minimising this penalisation term means that we learn an orthogonal matrix for $A_{L1}$ so that each row in $A_{L1}$ only focuses on a single aspect of semantics.

## 3.3 Structured Sentence-Level Self-Attention and Averaged Selection Representation

Taking $O^{L1}$ as the input to the structured 2-$D$ sentence-level attention model, the annotation matrix $A_{L2}$ is calculated as in (8):

$$A_{L2} = softmax(W_{s2}^{L2} tanh(W_{s1}^{L2} O^{L1})) \qquad (8)$$

where $W_{s1}^{L2}$ is the weight matrix of size $d_a^{L2} \times v$, and $d_a^{L2}$ is the number of neurons in the attention network; $W_{s2}^{L2}$ is the weight matrix of shape $r^{L2} \times d_a^{L2}$, where $r^{L2}$ ($r^{L2} > 1$) is the hyper-parameter representing the size of multiple vectors in the 2-$D$ sentence-level attention matrix. The $r^{L2}$ multiple vectors are expected to focus on different informative instances for the relation classification; $A_{L2}$ is the sentence-level annotation matrix of size $r^{L2} \times J$. We can see that the traditional 1-$D$ sentence-level attention model is expanded to a multi-vector attention ($r^{L2} > 1$).

Then, we average the 2-$D$ $A_{L2}$ to a 1-$D$ vector $\bar{A}_{L2}$ which has the dimension of $J$.

Accordingly, we calculate the averaged weighted sum by multiplying $\bar{A}_{L2}$ and the aggregated sentence representation $O^{L1}$, with the

resulting instance selection representation $M_{L2}$ being (9):

$$M_{L2} = \bar{A}_{L2} \cdot (O^{L1})^T \qquad (9)$$

where $M_{L2}$ has the size of $v$.

The probability distribution of the predicted relation type, i.e. the final output for relation prediction, can be calculated as in (10):

$$p(\hat{y}|G) = softmax(W_o^{L2} tanh(M_{L2}) + b_o^{L2}) \,(10)$$

### 3.4 Loss Function and Optimisation

The total loss of the network is the summation of the penalisation term $P_{L1}$, *softmax* loss in Eq. (10) and the $L2$ regularisation loss.

We use the ADAM optimiser (Kingma and Ba, 2014) to minimize the loss function on the mini-batch basis which is randomly selected from the training set.

## 4 Experiments

### 4.1 Datasets

We use two distantly supervised datasets, namely the NYT corpus (NYT) and the DBpedia Portuguese dataset (PT),[2] to verify our method.

In the NYT dataset, there are 53 relationships including a special relation *NA* which indicates a *None Relation* between two entities. The training set contains 580,888 sentences, 292,484 entity pairs and 19,429 relational facts (Non-NA). The test set contains 172,448 sentences, 96,678 entity pairs and 1,950 relational facts (Non-NA). There are 19.24% and 22.57% entity pairs corresponding to multiple instances in the training set and test set, respectively.

The DBpedia Portuguese dataset is smaller, containing just 10 relationships including a special relation *Other*. After preprocessing the original dataset, we obtain 96,847 sentences, 85,528 entity pairs and 77,321 relational facts (Non-Other). There are 8.61% entity pairs corresponding to multiple instances in the whole dataset. As in Batista et al. (2013), we use two different settings for the training and test sets: (1) a manually

reviewed subset that contains 602 sentences (PT-MANUAL) as the test set; and (2) 70%–30% out of the whole data as the training set and test set, respectively (PT-SPLIT).

### 4.2 Word Embeddings and Relative Position Features

For the NYT dataset, we use the 200-dimensional word vectors pre-trained using the NYT corpus;[3] for the PT dataset, we use a pre-trained 300-dimensional word vector model.[4] For the two-word entities in the data set, we use *underscore* to connect them as one word. The word embeddings of unknown words are intialised using the normal distribution with the standard deviation 0.05. Similar to previous work, we also use position embeddings specified by entity pairs. It is defined as the combination of the relative distances from the current word to head or tail entities (Zeng et al., 2014, 2015; Lin et al., 2016).

### 4.3 Baselines and Our MLSSA Systems

Neural RE systems have become the state-of-the-art, such as CNN-based (Zeng et al., 2014; Lin et al., 2017a), Piecewise CNN-based (Zeng et al., 2015; Lin et al., 2016; Ji et al., 2017), and BiLSTM-based (Zhou et al., 2016) models with or without an attention mechanism. In order to carry out a fair comparison, we select CNN+ATT, PCNN+ATT, BiGRU+ATT (bidirectional gated recurrent unit) and BiGRU+2ATT models as baselines on the NYT data, PCNN+ATT and BiGRU+2ATT as baselines on the PT data, where ATT indicates that the model has a sentence-level attention mechanism, and 2ATT indicates that the model has a 1-$D$ word-level and a 1-$D$ sentence-level attention.[5]

To show the incremental effectiveness of structured 2-$D$ word-level and 2-$D$ sentence-level self-attention mechanisms, we use two different settings for our MLSSA system: (1) **MLSSA-1**: this has a 2-$D$ word-level self-attention and a 1-$D$ sentence-level attention, i.e. $A_{L2}$ in Figure 1 is a 1-$D$ vector. This system is used to verify the context representation learning targeting **Problem I**;

---

[2]There are several reasons to use the Portuguese dataset: (i) the data sets reported in previous work, such as the KBP data, are not publicly available, or (ii) SemEval data sets which are not distantly supervised data. Google has also released a dataset (https://github.com/google-research-datasets/relation-extraction-corpus), but it is smaller and only has 4 relation types. For all these reasons, the Portuguese data is a better option to verify our method.

[3]https://catalog.ldc.upenn.edu/ldc2008t19
[4]https://s3-us-west-1.amazonaws.com/fasttext-vectors/wiki.pt.vec
[5]All the baseline systems are obtained from https://github.com/thunlp/NRE and https://github.com/thunlp/TensorFlow-NRE.

(2) **MLSSA-2**: both the word-level and sentence-level attentions are structured 2-$D$ matrices. This system verifies the instance selection representation learning targeting **Problem II**.

## 4.4 Experiment Setup and Evaluation Metrics

Following previous work, we use different evaluation metrics on these two datasets. For the NYT dataset:

- Overall evaluation: all training data is used for the model training, and all test data is used for the evaluation in terms of Precision-Recall (PR) curves;

- P@N evaluation: we select those entity pairs that have more than one instance to carry out the comparison in terms of the *precision at* $n$ (P@N) measure.[6] As in Lin et al. (2016), there are three settings: (1) One: for each testing entity pair corresponding to multiple instances, we randomly select one sentence to predict the relation; (2) Two: for each testing entity pair with multiple instances, we randomly select two sentences for the relation extraction; and (3) All: for each entity pair having multiple instances, we use all of them to predict the relation. Note that these three selections are only applied to the test set, and we keep all sentences in the training data for model building.

For the PT dataset, we use Macro F1 to evaluate system performance.[7]

## 4.5 Hyper-parameter Settings

We use cross-validation to determine the hyper-parameters of our system regarding two different settings and datasets. The in-common and different parameters for our two systems and two datasets are shown in Table 1.

## 4.6 PR Curves on NYT Dataset

The comparison results for the NYT test set are shown in Figure 2. We have the following observations: (1) BiGRU+ATT outperforms CNN+ATT

---

[6]P@N considers only the topmost results returned by the model.

[7]Regarding the metric, we keep the evaluation consistent with the work in Batista et al. (2013) where they used F1 to measure their RE systems on the Portuguese dataset, in order to maintain a fair comparison with their work using the same metric.

| Parameters for MLSSA-1/2 | NYT | PT |
|---|---|---|
| Word embedding dimension $d$ | 200 | 300 |
| Position embedding dimension | 50 | 50 |
| Batch size $B$ | 64 | 50 |
| Time steps $T$ | 70 | 70 |
| Learning rate $\lambda$ | 0.001 | 0.001 |
| Hidden size in BiLSTM $u$ | 300 | 300 |
| $d_a^{L1}$ at word-level attention | 300 | 300 |
| $r^{L1}$ at word-level attention | 9 | 5 |
| MLP size $v$ | 1000 | 1000 |
| Coefficient of the penalisation term | 1.0 | 1.0 |
| **Parameters for MLSSA-2 only** | **NYT** | **PT** |
| $d_a^{L2}$ at sentence-level attention | 300 | 300 |
| $r^{L2}$ at sentence-level attention | 9 | 3 |

Table 1: Hyper-parameter settings

and PCNN+ATT in terms of the PR curve, showing that it can learn a better semantic representation from the sequential input; (2) BiGRU+2ATT has better overall performance compared to Bi-GRU+ATT, showing that word-level attention is beneficial to sentence-level attention compared to single-attention models, i.e. the sentence-level attention model can select more informative sentences based on a more reasonable sentence embedding learned by the word-level attention model; (3) MLSSA-1 outperforms all baseline systems in terms of the PR curve, which demonstrates that the structured 2-$D$ word-level attention model can learn a better sentence representation by focusing on different aspects of the sentence, so that the sentence-level attention has a better chance of selecting the most informative sentences; and (4) the PR curve of MLSSA-2 is higher than that of MLSSA-1, demonstrating that the 2-$D$ sentence-level attention model can better select the most informative sentences compared to the 1-$D$ sentence-level attention model targeting those entity pairs with multiple instances.

## 4.7 P@N Evaluation on NYT Dataset

The results on the NYT dataset regarding P@100, P@200, P@300 and the mean of three settings for each model are shown in Table 2. From the table, we have similar observations to the PR Curves: (1) BiGRU+2ATT outperforms CNN+ATT, PCNN+ATT and BiGRU+ATT in most cases in terms of all P@N scores; and (2) MLSSA-1 and MLSSA-2 significantly outperform all baselines for all measures. We observe that MLSSA-1 performs better than MLSSA-2 on tasks **One** and **Two**, but worse on **All**. We infer that in our 2-$D$ sentence-level attention model, we

| Test Settings | One | | | | Two | | | | All | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| P@N(%) | 100 | 200 | 300 | Mean | 100 | 200 | 300 | Mean | 100 | 200 | 300 | Mean |
| CNN+ATT | 72.0 | 67.0 | 59.5 | 66.2 | 75.5 | 69.0 | 63.3 | 69.3 | 74.3 | 71.5 | 64.5 | 70.1 |
| PCNN+ATT | 73.3 | 69.2 | 60.8 | 67.8 | 77.2 | 71.6 | 66.1 | 71.6 | 76.2 | 73.1 | 67.4 | 72.2 |
| BiGRU+ATT | 75.0 | 69.5 | 64.7 | 69.7 | 80.0 | 72.5 | 69.3 | 73.9 | 82.0 | 76.5 | 71.3 | 76.6 |
| BiGRU+2ATT | 81.0 | 74.0 | 67.3 | 74.1 | 81.0 | 75.5 | 70.7 | 75.7 | 81.0 | 76.0 | 72.7 | 76.6 |
| MLSSA-1 | **88.0** | **77.0** | **70.0** | **78.3** | 88.0 | **79.0** | **73.3** | **80.1** | 87.0 | **81.5** | 76.0 | 81.5 |
| MLSSA-2 | 87.0 | 76.0 | **70.0** | 77.7 | **89.0** | 78.5 | 72.3 | 79.9 | **90.0** | **81.5** | **77.0** | **82.8** |

Table 2: Precision values for the top-100, top-200, and top-300 relation instances that are randomly selected in terms of one, two and all sentences.
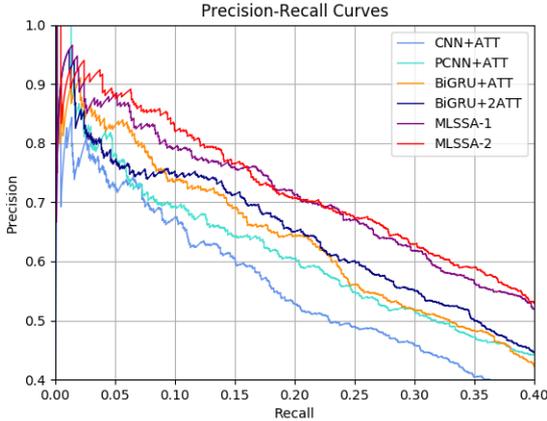


Figure 2: Comparison results of a variety of methods in terms of precision/recall curves.

set $r^{L2}$ to 9, but there are only *one* and *two* instances for selection in tasks **One** and **Two**, so the 2-$D$ matrix cannot demonstrate its full potential. However, in **All**, many entity pairs contain multiple or more than 9 instances, so it can learn a better 2-$D$ matrix to focus on different instances.

### 4.8 Results on PT Dataset

Based on results from the NYT dataset, we choose PCNN+ATT and BiGRU+2ATT as representative baselines to compare against our MLSSA-1/2 systems on the PT test sets. The results in terms of Macro F1 are shown in Table 3.

It can be seen that on both test sets, our MLSSA-2 model achieved the best performance which shows that the structured 2-$D$ word-level and sentence-level self-attention models can be well applied to datasets of a smaller scale and with a smaller ratio of multiple instances.

### 4.9 Examples and Analysis

In order to show the effectiveness of structured self-attention mechanisms, we show some examples by visualising the attentions on different as-

| SYS | PT-MANUAL (%) | PT-SPLIT (%) |
|---|---|---|
| PCNN+ATT | 62.3 | 74.1 |
| BiGRU+2ATT | 63.5 | 75.3 |
| MLSSA-1 | 66.0 | 77.2 |
| MLSSA-2 | **69.6** | **78.1** |

Table 3: Results on the PT test sets

pects of a sentence, and on different sentences comparing with BiLSTM+2ATT model.

Figure 3 shows the comparison of word-level attention mechanism between BiGRU+2ATT and MLSSA-1 reflecting their capability of context representation learning (*Problem I*). MLSSA-2 has a similar probability distribution to MLSSA-1 in terms of this example.

The *pink* fonts indicate lower probability and *red* indicates higher probability. We observe that: (1) BiGRU+2ATT mainly focuses on one word *baltimore*. We can see that it has little attention on the entity word *maryland*. In this example, the comma implies a semantic relationship *location/location/contains* for the entity pair (*Maryland, Baltimore*). However, BiGRU+2ATT allocates quite a small probability to it; and (2) we can see that our model focuses on different words via different attention vectors (9 in total). Words with a *red* background have a high probability of 0.98 or so. For rows 5, 6, 8 and 9, the focus is on the *BLANK* tokens. In both systems, the maximum time step is set to 70, which indicates that shorter sentences are padded with *BLANK* tokens and longer sentences are cut off. The last row shows the summation of **9** annotation vectors, and it constructs a dependency-like context of the relation for the entity pair. Attentions on different words are attributed to the penalisation $P_{L1}$ which is optimised to learn orthogonal eigenvectors.

Figure 4 shows the comparison of sentence-level attentions between BiGRU+2ATT, MLSSA-1 and MLSSA-2. The first, second and third columns are probability distributions over multiple instances. The entity pair is (*vinod khosla,*

| | | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BiGRU+2ATT | | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| MLSSA-1 | 1 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | 2 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | 3 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | 4 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | 5 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | 6 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | 7 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | 8 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | 9 | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |
| | sum | born | in | baltimore | , | maryland | on | january | 6th | , | 1941 | . |

Figure 3: Comparison of word-level attentions.

| BiGRU+2ATT | MLSSA-1 | MLSSA-2 | Sentences (Relation: Business/Person/Company) |
|---|---|---|---|
| 0.26 | 0.19 | 0.37 | the most visible and one of the most outspoken is vinod khosla , a founder of sun microsystems and now a partner at khosla ventures . |
| 0.73 | 0.47 | 0.41 | in early 2006 , he joined khosla ventures , a silicon valley venture firm started by vinod khosla , a founder of sun microsystems . |
| 0.01 | 0.33 | 0.21 | vinod khosla , a co-founder of sun microsystems who formed khosla ventures in 2004 , has invested in more than a dozen start-ups involved in '' clean fuel '' technologies . |

Figure 4: Comparison of sentence-level attentions.

*sun microsystems*), and their relation is *Business/Person/Company*. From this figure, we observe that: (1) BiGRU+2ATT allocates high probabilities to **Sentences 1** and **2** by learning the context of "*a founder of*", but does not recognise that "*co-founder*" is semantically the same as "*founder*"; and (2) our two models almost evenly focus on all sentences because they express the same semantic concept of "*a person is a founder of a company*" in terms of the given entity pair. Therefore, the structured self-attention mechanism is helpful to learn a better representation and select informative sentences.

## 5 Conclusion and Future Work

This paper has proposed a multi-level structured self-attention mechanism for distantly supervised RE. In this framework, the traditional 1-$D$ word-level and sentence-level attentions are extended to 2-$D$ structured matrices which can learn different aspects of a sentence, and different informative instances. Experimental results on two distant supervision data sets show that (1) the structured 2-$D$ word-level attention can learn a better sentence representation; (2) the structured 2-$D$ sentence-level attention and averaged selection can perform better selection from multiple instances for relation classification; (3) the proposed framework significantly outperforms state-of-the-art baseline systems for a range of different measures, which verifies its effectiveness on two representation learning issues. A subsequent manual investigation via examples also show its effectiveness on two representation learning issues.

In future work, we will build a domain-specific distant supervision dataset with a higher ratio of multiple instances and compare our system with others. Furthermore, we will consider not using RNNs or CNNs, but a deeper neural networks with only attentions for distantly supervised RE, similar to the work in Vaswani et al. (2017).

# References

David Soares Batista, David Forte, Rui Silva, Bruno Martins, and Mario J. Silva. 2013. Exploring dbpedia and wikipedia for portuguese semantic relationship extraction. *Linguamatica*, 5(1).

Razvan Bunescu and Raymond J Mooney. 2005. Subsequence kernels for relation extraction. In *In Proceedings of NIPS*, pages 171–178.

Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2010. Semeval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals. In *Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions*, pages 94–99.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Comput.*, 9(8):1735–1780.

Raphael Hoffmann, Congle Zhang, Xiao Ling, Luke Zettlemoyer, and Daniel S. Weld. 2011. Knowledge-based weak supervision for information extraction of overlapping relations. In *Proceedings of the 49th ACL-HLT*, pages 541–550.

Guoliang Ji, Kang Liu, Shizhu He, and Jun Zhao. 2017. Distant supervision for relation extraction with sentence-level attention and entity descriptions. In *Proceedings of the Thirty-First (AAAI) Conference on Artificial Intelligence*, pages 3060–3066.

Yoon Kim, Carl Denton, Luong Hoang, and Alexander M. Rush. 2017. Structured attention networks. *CoRR*, abs/1702.00887.

Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980.

Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2017a. Neural relation extraction with multi-lingual attention. In *Proceedings of the 55th ACL*, pages 34–43.

Yankai Lin, Shiqi Shen, Zhiyuan Liu, Huanbo Luan, and Maosong Sun. 2016. Neural relation extraction with selective attention over instances. In *Proceedings of the 54th ACL*.

Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017b. A structured self-attentive sentence embedding. In *International Conference on Learning Representations 2017*.

Mike Mintz, Steven Bills, Rion Snow, and Dan Jurafsky. 2009. Distant supervision for relation extraction without labeled data. In *Proceedings of the Joint Conference of the 47th ACL and the 4th AFNLP*, pages 1003–1011.

Sebastian Riedel, Limin Yao, and Andrew McCallum. 2010. Modeling relations and their mentions without labeled text. In *Proceedings of the 2010 European Conference on Machine Learning and Knowledge Discovery in Databases: Part III*, pages 148–163.

Mihai Surdeanu, Julie Tibshirani, Ramesh Nallapati, and Christopher D. Manning. 2012. Multi-instance multi-label learning for relation extraction. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 455–465, Jeju Island, Korea.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, pages 6000–6010.

Linlin Wang, Zhu Cao, Gerard de Melo, and Zhiyuan Liu. 2016. Relation classification via multi-level attention cnns. In *Proceedings of the 54th ACL*, pages 1298–1307.

Linyi Yang, Tin Lok, James Ng, Catherine Mooney, and Ruihai Dong. 2017. Multi-level attention-based neural networks for distant supervised relation extraction. In *Proceedings of Irish Conference on Artificial Intelligence and Cognitive Science (AICS)*.

Dmitry Zelenko, Chinatsu Aone, and Anthony Richardella. 2003. Kernel methods for relation extraction. *Journal of machine learning research*, 3:1083–1106.

Daojian Zeng, Kang Liu, Yubo Chen, and Jun Zhao. 2015. Distant supervision for relation extraction via piecewise convolutional neural networks. In *Proceedings of EMNLP*, pages 1753–1762.

Daojian Zeng, Kang Liu, Siwei Lai, Guangyou Zhou, and Jun Zhao. 2014. Relation classification via convolutional deep neural network. In *Proceedings of COLING 2014*, pages 2335–2344.

Wenyuan Zeng, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2017. Incorporating relation paths in neural relation extraction. In *Proceedings of EMNLP*, pages 1769–1778.

Guodong Zhou, Jian Su, Jie Zhang, and Min Zhang. 2005. Exploring various knowledge in relation extraction. In *In Proceedings of the 43rd ACL*, pages 427–434.

Peng Zhou, Wei Shi, Jun Tian, Zhenyu Qi, Bingchen Li, Hongwei Hao, and Bo Xu. 2016. Attention-based bidirectional long short-term memory networks for relation classification. In *Proceedings of the 54th ACL*, pages 207–212.