

# Document-Level Relation Extraction with Reconstruction

Wang Xu<sup>1</sup>, Kehai Chen<sup>2</sup>, and Tiejun Zhao<sup>1</sup>

<sup>1</sup> Harbin Institute of Technology, Harbin, China

<sup>2</sup> National Institute of Information and Communications Technology (NICT), Kyoto, Japan  
xuwang@hit-mtlab.net, khchen@nict.go.jp, tjzhao@hit.edu.cn

## Abstract

In document-level relation extraction (DocRE), graph structure is generally used to encode relation information in the input document to classify the relation category between each entity pair, and has greatly advanced the DocRE task over the past several years. However, the learned graph representation universally models relation information between all entity pairs regardless of whether there are relationships between these entity pairs. Thus, those entity pairs without relationships disperse the attention of the encoder-classifier DocRE for ones with relationships, which may further hinder the improvement of DocRE. To alleviate this issue, we propose a novel encoder-classifier-reconstructor model for DocRE. The reconstructor manages to reconstruct the ground-truth path dependencies from the graph representation, to ensure that the proposed DocRE model pays more attention to encode entity pairs with relationships in the training. Furthermore, the reconstructor is regarded as a relationship indicator to assist relation classification in the inference, which can further improve the performance of DocRE model. Experimental results on a large-scale DocRE dataset show that the proposed model can significantly improve the accuracy of relation extraction on a strong heterogeneous graph-based baseline. The code is publicly available at <https://github.com/xwjim/DocRE-Rec>.

## 1 Introduction

Graph structure plays an important role in the document relation extraction (DocRE) (Christopoulou, Miwa, and Ananiadou 2019; Sahu et al. 2019; Nan et al. 2020; Tang et al. 2020). Typically, one unstructured input document is first organized as a structure input graph (i.e., homogeneous or heterogeneous graphs) based on syntactic trees, co-reference, or heuristics rules, thereby building relationships between entity pairs within and across multiple sentences of the input document. Neural networks (i.e., graph network) are used to iteratively encode the structure input graph as a graph representation to model relation information in the input document. The graph representation is fed into one classifier to classify the relation category between each entity pair, which has achieved the state-of-the-art performance in DocRE (Christopoulou, Miwa, and Ananiadou 2019; Nan et al. 2020).

However, during the training of DocRE model, the graph representation universally encodes relation information between all entity pairs regardless of whether there are relationships between these entity pairs. For example, Figure 1 shows three entities in an input document: *X-Files*, *Chris Carter*, and *Fox Mulder*. Intuitively, they are three entity pairs:  $\{X\text{-Files}, \text{Chris Carter}\}$ ,  $\{X\text{-Files}, \text{Fox Mulder}\}$ , and  $\{\text{Chris Carter}, \text{Fox Mulder}\}$ . The DocRE model learns the node representations of each entity pair to classify their relation. As seen, there exists relationship between  $\{\text{Chris Carter}, \text{Fox Mulder}\}$  in the reference, indicating that there is naturally a reliable reasoning path from *Chris Carter* to *Fox Mulder*. In comparison, there do not exist relationships between  $\{X\text{-Files}, \text{Chris Carter}\}$  and between  $\{X\text{-Files}, \text{Fox Mulder}\}$ , indicating that there are not reasoning paths between  $\{X\text{-Files}-\text{Chris Carter}\}$  or  $\{X\text{-Files}, \text{Fox Mulder}\}$ . However, the learned graph representation models the three path dependencies universally and does not consider whether there is a path dependency between one target entity pair. As a result,  $\{X\text{-Files}, \text{Chris Carter}\}$  and  $\{X\text{-Files}, \text{Fox Mulder}\}$  without relationships disperse the attention of the DocRE model for the learning of  $\{\text{Fox Mulder}, \text{Chris Carter}\}$  with relationship, which may further hinder the improvement of the DocRE model.

To alleviate this issue, we propose a novel reconstructor method to enable the DocRE model to model path dependency between one entity pair with the ground-truth relationship. To this end, the reconstructor generates a sequence of node representations on the path from one entity node to another entity node and thereby maximizes the probability of its path if there is a ground-truth relationship between one entity pair and minimizes the probability otherwise. This allows the proposed DocRE model to pay more attention to the learning of entity pairs with relationships in the training, thereby learning an effective graph representation for the subsequent relation classification. Furthermore, the reconstructor is regarded as a relationship indicator to assist relation classification in the inference, which can further improve the performance of DocRE model. Experimental results on a large-scale DocRE dataset show that the proposed method gained improvement of 1.7 F1 points over a strong heterogeneous graph-based DocRE model, especially outperformed the recent state-of-

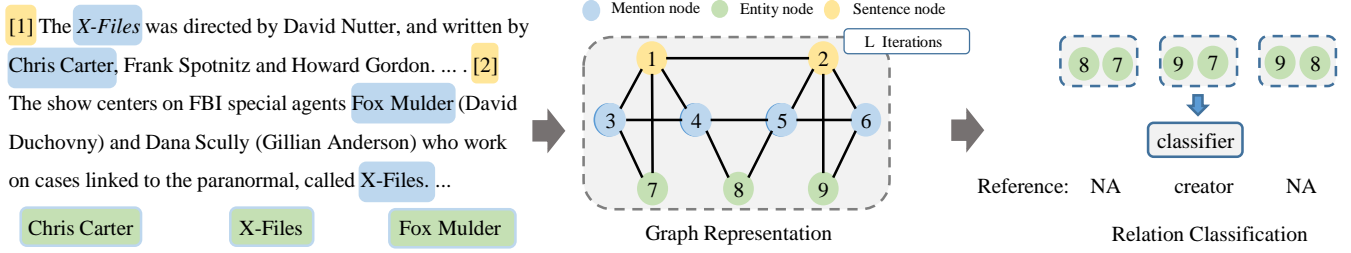


Figure 1: Heuristic rules are used to convert the input document into a heterogeneous graph. Then graph attention network is applied to learn the graph representation. Finally the node representations of entity pairs are used to classify their relationships.

the-art LSR model for DocRE (Nan et al. 2020).

## 2 Background

In this section, based on (Christopoulou, Miwa, and Ananiadou 2019)’s work, we used heuristic rules to convert the input document into a heterogeneous graph without external syntactic knowledge. Moreover, a graph attention network is used to encode the heterogeneous graph instead of the edge-oriented graph network (Christopoulou, Miwa, and Ananiadou 2019), thereby implementing a strong and general baseline for DocRE.

### 2.1 Heterogeneous Graph Construction

Formally, given an input document that consists of  $L$  sentences  $\{S^1, S^2, \dots, S^L\}$ , each of which is a sequence of words  $\{x_1^l, x_2^l, \dots, x_j^l\}$  with the length  $J=|S^l|$ . A bidirectional long short-term memory (BiLSTM) reads word by word to generate a sequence of word vectors to represent each sentence in the input document. Also, we apply the heterogeneous graph (Christopoulou, Miwa, and Ananiadou 2019) to the input document to build relationships between all entity pairs. Specifically, the heterogeneous graph includes three defined distinct types of nodes: Mention Node, Entity Node, and Sentence Node. For example, Figure 1 shows an input document including two sentences (yellow color index) in which there are four mentions (blue color) and three entities (green color). The representation of each node is the average of the words in the concept, thereby forming a set of node representations  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}$ , where  $N$  is the number of nodes. For edge connections, there are five distinct types of edges between pairs of nodes following (Christopoulou, Miwa, and Ananiadou 2019)’s work, Mention-Mention(MM) edge, Mention-Sentence (MS) edge, Mention-Entity (ME) edge, Sentence-Sentence (SS) edge, Sentence-Sentence (SS) edge, Entity-Sentence (ES) edge respectively. In addition, we add a Mention-Coreference (CO) edges between the two mentions which are referred to the same entity. According to these above definitions, there is a  $N \times N$  adjacency matrix  $\mathbb{E}$  denoting edge connections. Finally, the heterogeneous graph can be denoted as  $G=\{\mathbf{V}, \mathbb{E}\}$ , to keep relation information between all entity pairs in the input document.

### 2.2 Encoder

To learn an effective graph representation, we used the graph attention network (Guo, Zhang, and Lu 2019) to encode the feature representation of each node in the heterogeneous graph. Formally, given the outputs of all previous hop reasoning operations  $\{\mathbf{s}_n^1, \mathbf{s}_n^2, \dots, \mathbf{s}_n^{l-1}\}$ , they are concatenated and then transformed to a fixed dimensional vector as the input of the  $l$  hop reasoning:

$$\mathbf{v}_n^l = \mathbf{W}_e^l \cdot [\mathbf{v}_n : \mathbf{s}_n^1 : \mathbf{s}_n^2 : \dots : \mathbf{s}_n^{l-1}], \quad (1)$$

where  $\mathbf{s}_n^{l-1} \in \mathbb{R}^{d_0}$  and  $\mathbf{W}_e^l \in \mathbb{R}^{d_0 \times (l \times d_0)}$ . Also, according to edge matrix  $\mathbb{E}[n][a_c]=k$  ( $0 \leq a_c < N, k > 0$ ),  $C$  direct adjacent nodes of  $\mathbf{v}_n$  are  $\{\mathbf{z}_{a_1}^l, \mathbf{z}_{a_2}^l, \dots, \mathbf{z}_{a_C}^l\}$ . We then use the self-attention mechanism (Vaswani et al. 2017) to capture the feature information of  $\mathbf{v}_n$  between  $\mathbf{z}_n^l$  and  $\{\mathbf{z}_{a_1}^l, \mathbf{z}_{a_2}^l, \dots, \mathbf{z}_{a_C}^l\}$ :

$$\mathbf{s}_n^l = \text{softmax}\left(\frac{\mathbf{z}_n^l \mathbf{K}^\top}{\sqrt{d_0}}\right) \mathbf{V}, \quad (2)$$

where  $\{\mathbf{K}, \mathbf{V}\}$  are key and value matrices that are transformed from the direct adjacent nodes representations  $\{\mathbf{z}_{a_1}^l, \mathbf{z}_{a_2}^l, \dots, \mathbf{z}_{a_C}^l\}$  according to the edge type.

After performing  $L$  hop reasonings, there is a sequence of annotations  $\{\mathbf{s}_n^1, \mathbf{s}_n^2, \dots, \mathbf{s}_n^L\}$  to encode relation information in the input document. Finally, another no-linear layer is applied to integrate the reason information  $\{\mathbf{s}_n^1, \mathbf{s}_n^2, \dots, \mathbf{s}_n^L\}$  and the node information  $\mathbf{v}_n$ :

$$\mathbf{q}_n = \text{Relu}(\mathbf{W}_o \cdot [\mathbf{v}_n : \mathbf{s}_n^1 : \dots : \mathbf{s}_n^L]), \quad (3)$$

where  $\mathbf{W}_o \in \mathbb{R}^{d_1 \times (d_0 \times (L+1))}$ ,  $\mathbf{q}_n \in \mathbb{R}^{d_1}$ . As a result, the heterogeneous graph  $G$  is represented as  $\{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N\}$ .

### 2.3 Classifier

Given the heterogeneous graph representation  $\{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N\}$ , two node representations of each entity pair are as the input to the classifier to classify their relationship. Specifically, the classifier is a multi-layer perceptron (MLP) layer with sigmoid function to calculate the relationship probability:

$$R(r) = P(r|\{e_i, e_j\}) = \text{sigmoid}(\text{MLP}([\mathbf{q}_i : \mathbf{q}_j])). \quad (4)$$

To train the DocRE model, the binary cross-entropy is used to optimize parameters of neural networks

over the triple examples (subject, object, relation) on the training data set (including  $T$  documents), that is,  $\{\{e1_n^t, e2_n^t, r_n^t\}_{n=1}^{N_t}\}_{t=1}^T$ :

$$Loss_c = -\frac{1}{\sum_{t=0}^T N_t} \sum_{t=1}^T \sum_{n=1}^{N_t} \{r_n^t \log(R(r_n^t)) + (1 - r_n^t) \log(1 - R(r_n^t))\}, \quad (5)$$

where  $r_n^t \in \{0, 1\}$  indicates whether the entity pair has relation label  $r$  and  $N_t$  is the number of relations in the  $t$ -th document.

### 3 Methodology

Intuitively, when a human understands a document with relationships, he or she often pays more attention to learn entity pairs with relationships rather than ones without relationships. Motivated by this observation, we proposed a novel DocRE model with reconstruction (See Figure 2) to pay more attention to entity pairs with relationships, thus enhancing the accuracy of relationship classification.

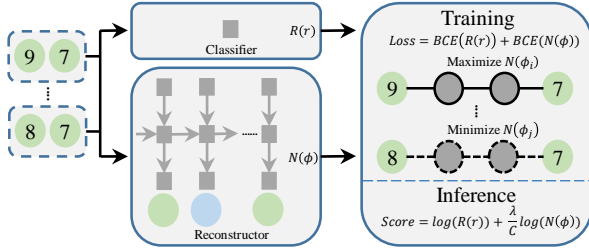


Figure 2: Model overview. The reconstructor manages to reconstruct the ground-truth path dependencies from the graph representation to ensure that the model to pay attention to model entity pairs with relationships. Furthermore, the reconstructor is regarded as a relationship indicator to assist relation classification in the inference.

#### 3.1 Meta Path of Entity Pair

Generally, when there is a relationship between two entities, they should have one strong path dependency in the graph structure (or representation). In comparison, when there is not a relationship between two entities, there is a weak path dependency.<sup>1</sup> Thus, we explore to reconstruct the path dependency between each entity pair from the learned graph representation. To this end, we first define three type paths between two entity nodes in the graph representation as reconstructed candidates according to the meta-path information (Sun and Han 2013).

- 1) Meta Path1 of Pattern Recognition: Two entities are connected through a sentence in this reasoning type. The relation schema is  $EM \circ MM \circ EM$ , for example node sequence  $\{7, 3, 4, 8\}$  in Figure 1.

<sup>1</sup>If there is no path dependency between two target entities without a relationship, this may weaken the understanding of relationship information in the document.

- 2) Meta Path2 of Logical Reasoning: the relation between two entities is indirectly established by a bridge entity. The bridge entity occurs in a sentence with the two entities separately. The relation schema is  $EM \circ MM \circ CO \circ MM \circ EM$ , for example node sequence  $\{7, 3, 4, 5, 6, 9\}$  in Figure 1.
- 3) Meta Path3 of Coreference Reasoning: Coreference resolution must be performed first to identify target entities. A reference word refers to an entity that appear in the previous sentence. The two entities occur in the same sentence implicitly. The relation schema is  $ES \circ SS \circ ES$ , for example node sequence  $\{7, 1, 2, 9\}$  in Figure 1.

Actually, all the entity pairs have at least one of the three meta-paths. We select one meta-path type according to the priority, meta-path1 > meta-path2 > meta-path3. Generally, several instance paths may exist corresponding to the meta path, we select the instance path that appears firstly in the document.

#### 3.2 Path Reconstruction

For each entity pair, one instance path is selected as the supervision of the reconstruction of the path dependency. In other words, there is only one supervision path  $\phi_n = \{\mathbf{v}_{b_1}, \mathbf{v}_{b_2}, \dots, \mathbf{v}_{b_C}\}$  between each target pair  $\{e1_n, e2_n\}$ , where  $b_C$  is the number of nodes.

To reconstruct the path dependency of each entity pair, we model the reconstructor as the sequence generation. Specifically, we use a LSTM to compute a path hidden state  $\mathbf{p}_{b_c}$  for each node  $\mathbf{q}_{b_{c-1}}$  on the path  $\phi_n$ :

$$\mathbf{p}_{b_c} = \text{LSTM}(\mathbf{p}_{b_{c-1}}, \mathbf{q}_{b_{c-1}}). \quad (6)$$

Note that  $\mathbf{p}_{b_0}$  is initialized as the transform of  $\mathbf{o}_{ij}$ , since it plays a key role in classification.  $\mathbf{p}_{b_c}$  is fed into a softmax layer to compute the probability of node  $\mathbf{v}_{b_c}$  on the path:

$$\mathcal{P}(\mathbf{v}_{b_c} | \mathbf{v}_{<b_c}) = \frac{\exp(\mathbf{p}_{b_c} \mathbf{W}_r \mathbf{q}_{b_c})}{\sum_n \exp(\mathbf{p}_{b_c} \mathbf{W}_r \mathbf{q}_n)}, \quad (7)$$

where  $\mathbf{W}_r \in \mathbb{R}^{d_1 \times d_1}$ . Also, there is a set of node probabilities  $\{\mathcal{P}(\mathbf{v}_{b_1} | \mathbf{v}_{<b_1}), \mathcal{P}(\mathbf{v}_{b_2} | \mathbf{v}_{<b_2}), \dots, \mathcal{P}(\mathbf{v}_{b_C} | \mathbf{v}_{<b_C})\}$  for the path  $\phi_n$ . Finally, the probability of this path  $\phi_n$  is computed:

$$\mathcal{N}(\phi_n) = \prod_{c=1}^C (\mathcal{P}(\mathbf{v}_{b_c} | \mathbf{v}_{<b_c})). \quad (8)$$

#### 3.3 Training with Reconstruction Loss

We use the reconstructed path probability to compute an additional reconstruction loss over the triple examples of the training data set  $\{\{e1_n^t, e2_n^t, r_n^t\}_{n=1}^{N_t}\}_{t=1}^T$ :

$$Loss_r = -\frac{1}{\sum_{t=0}^T N_t} \sum_{t=1}^T \sum_{n=1}^{N_t} \{r_n^t \log \mathcal{N}(\phi_n) + (1 - r_n^t) \log(1 - \mathcal{N}(\phi_n))\}, \quad (9)$$

where  $r_n^t$  is one of  $\{0, 1\}$ , that is, we maximize the probability of the path  $\mathcal{N}(\phi_n)$  if the entity pair has relation, and minimize the probability otherwise. To simplify the

Eq.(9), we use  $\prod_{c=1}^C(1 - P_{b_c})$  to replace with the  $(1 - \mathcal{N}(\phi_n))$ , where  $P_{b_c} = \mathcal{P}(\mathbf{v}_{b_c} | \mathbf{v}_{<b_c})$ . The reconstruction loss is modified as Eq.(10):

$$Loss_r = -\frac{1}{\sum_{t=0}^T N_t} \sum_{t=1}^T \sum_{n=1}^{N_t} \left\{ \sum_{b_c=1}^{b_C} \{ (r_n^t \log P_{b_c}) + (1 - r_n^t) \log(1 - P_{b_c}) \} \right\}. \quad (10)$$

Finally, the reconstructor loss and the existing classification loss in Eq.(5) is added as the training objective of the proposed DocRE model:

$$Loss = Loss_c + Loss_r. \quad (11)$$

### 3.4 Inference with Path Reconstruction

Intuitively, the proposed reconstructor encourages the DocRE model to pay more attention to model entity pairs with ground-truth relationships. Furthermore, we maximized the path probability between one entity pair if there is indeed a relation and we minimized it otherwise when computing the reconstruction loss in Eq.(10). In other words, the higher the probability of this path is, the greater the likelihood of a relationship between the entity pair is. Naturally, we treat this path probability as a relational indicator to assist relation classification in the inference:

$$S(r) = \log(R(r)) + \lambda \cdot \frac{1}{C} \sum_{b_c=1}^{b_C} \log(P_{b_c}), \quad (12)$$

where  $\lambda$  is a hyper-parameter to control the importance of reconstruction probability in the inference.

## 4 Experiments

### 4.1 Setup

The proposed methods were evaluated on a large-scale human-annotated dataset for document-level relation extraction (Yao et al. 2019). DocRED contains 3,053 documents for the training set, 1,000 documents for the development set, and 1,000 documents for the test set, totally with 132,375 entities, 56,354 relational facts, and 96 relation types. More than 40% of the relational facts require the reading and reasoning over multiple sentences. Following settings of (Nan et al. 2020)’s work, we used the GloVe embedding (100d) and BiLSTM (128d) as word embedding and encoder. The hop number  $L$  of the encoder was set to 2. The learning rate was set to  $1e-4$  and we trained the model using Adam as the optimizer. For the BERT representations, we used uncased BERT-Based model (768d) as the encoder and the learning rate was set to  $1e^{-5}$ . For evaluation, we used  $F_1$  and Ign  $F_1$  as the evaluation metrics. Ign  $F_1$  denotes  $F_1$  score excluding relational facts shared by the training and development/test sets. In particular, the predicted results were ranked by their confidence and traverse this list from top to bottom by  $F_1$  score on development set, and the score value corresponding to the maximum  $F_1$  is picked as threshold  $\theta$ . All hyper-parameters were tuned based on the development set. In addition, the results on the test set were evaluated through CodaLab<sup>2</sup>.

<sup>2</sup><https://competitions.codalab.org/competitions/20717>

### 4.2 Baseline Systems

According to Section 2, there is a baseline heterogeneous-based graph self-attention network model (HeterGSAN). Also, there are some recent DocRE methods as our comparison systems:

- **Sequence-based Models:** These models used different neural architectures to encode sentences in the document, including including convolution neural networks (CNN) (Yao et al. 2019), bidirectional LSTM (BiLSTM) (Yao et al. 2019) and Context-Aware LSTM (Yao et al. 2019).

- **Graph-based Models.** GCNN (Sahu et al. 2019), GAT (Veličković et al. 2018), AGGCN (Guo, Zhang, and Lu 2019) constructed the graph from syntactic parsing and sequential information, or non-local dependencies from coreference resolution and other semantic dependencies, and then uses the GCN based method to calculate the node embedding. EoG (Christopoulou, Miwa, and Ananiadou 2019) defined several node types and edges to construct a heterogeneous graph of the input document without external syntactic knowledge. EoG uses an iterative algorithm to learn new edge representations between different nodes in the heterogeneous graph and classify relationships between entity pairs. Instead of constructing a static graph representation, LSR (Nan et al. 2020) empowered the relational reasoning across sentences by automatically inducing the latent document-level graph.

- **BERT.** It applied a pre-trained language model to learn the representations of the input document (Wang et al. 2019; Devlin et al. 2019). Furthermore, it used a two-phase training process to enhance the performance of DocRE model. Specifically, it first predicts whether a pair of entities has a relation or not and classifies the relation for each entity pair.

### 4.3 Main Results

Table 1 presents the detailed results on the development set and the test set of DocRED. As seen, our baseline HeterGSAN model achieved 53.52  $F_1$  score on the test set and outperformed the EoG model which is also a heterogeneous-based graph DocRE model by 1.7 points in terms of  $F_1$ . Meanwhile, HeterGSAN is consistently superior to the most of comparison methods, including CNN, BiLSTM, ContextAware, GCNN, GAT, and AGGCN. This indicates that the graph self-attention network can give a strong baseline in the heterogeneous-based methods of DocRE. HeterGSAN+reconstruction achieved 55.23  $F_1$ , which outperformed the baseline HeterGSAN by 1.71  $F_1$  score. In particular, HeterGSAN+reconstruction outperformed the existing state-of-the-art LSR model by 1.05  $F_1$  score, which is a new state-of-the-art result on the DocRED dataset without the pre-trained model (BERT). This means that the proposed reconstructor is beneficial to encode relation information in the input document, thereby enhancing the relation extraction.

In addition, we evaluated the proposed HeterGSAN model with a pre-trained language model as shown in Table 1. First, HeterGSAN+BERT model consistently outperformed the comparison BERT model, Two-Phase BERT

Groups	Methods	Dev		Test	
		Ign F1	F1	Ign F1	F1
w/o BERT	CNN* (Yao et al. 2019)	41.58	43.45	40.33	42.26
	BiLSTM* (Yao et al. 2019)	48.87	50.94	48.78	51.06
	ContexAware* (Yao et al. 2019)	48.94	51.09	48.40	50.07
	GCNN <sup>†</sup> (Sahu et al. 2019)	46.22	51.52	49.59	51.62
	EoG <sup>†</sup> (Christopoulou, Miwa, and Ananiadou 2019)	45.94	52.15	49.48	51.82
	GAT <sup>†</sup> (Veličković et al. 2018)	45.17	51.44	47.36	49.51
	AGGCN <sup>†</sup> (Guo, Zhang, and Lu 2019)	46.29	52.47	48.89	51.45
	LSR* (Nan et al. 2020)	48.82	55.17	52.15	54.18
	HeterGSAN +Reconstruction	52.17	54.40	52.07	53.52
w/ BERT	BERT* (Wang et al. 2019)	-	54.16	-	53.20
	Two-Phase BERT* (Wang et al. 2019)	-	54.42	-	53.92
	BERT+LSR* (Nan et al. 2020)	52.43	59.00	56.97	59.05
	HeterGSAN +BERT	52.17	54.40	52.07	53.52
	HeterGSAN +BERT +Reconstruction	<b>57.00</b>	<b>59.13</b>	<b>56.21</b>	<b>58.54</b>
		<b>58.13</b>	<b>60.18</b>	<b>57.12</b>	<b>59.45</b>

Table 1: Results on the development set and the test set. Results with \* are reported in their original papers. Results with † are reported in (Nan et al. 2020). Bold results indicate the best performance of the current method.

model, BERT+LSR model. This confirms the effectiveness of the BERT method, which we believe makes the evaluation convincing. Moreover, HeterGSAN+BERT+Reconstruction model outperformed HeterGSAN+BERT model by 0.91  $F_1$  score, indicating that our approach is complementary to BERT, and combining them is able to further improve the accuracy of relation extraction. Meanwhile, HeterGSAN+BERT+Reconstruction model ( $F_1$  59.45) outperformed BERT+LSR model ( $F_1$  59.05) by 0.40  $F_1$  score on the test set, which is a new state-of-the-art result.

#### 4.4 Effect of Reconstruction

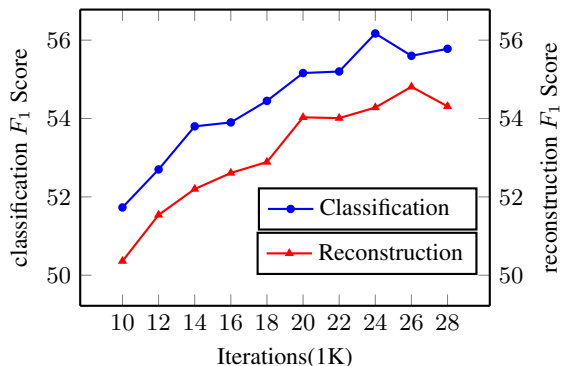


Figure 3: Learning curves of classification (left y-axis) and reconstruction (right y-axis) performances (in  $F_1$  scores) on the development set during the training.

To valid the effect of reconstruction, Figure 3 showed learning curves of classification and reconstruction performances (in  $F_1$  scores) on the development set during the training. For reconstruction, we used the reconstructor to generate the source path for each entity pair and calculated the probability of the reconstructed path to indicate how

much there is a relationship. As seen, the reconstruction  $F_1$  scores went up with the improvement of reconstruction over time. When the classification performance reached a peak at iteration 24K, the proposed model achieved a balance between classification and reconstruction scores. Therefore, we use the trained model at iteration 24K in Table 1.

#### 4.5 Ablation in Training and Inference

	Reconstructor used in		Metric	
	Training	Inference	Ign F1	F1
#1	×	×	52.17	54.40
#2	✓	×	53.69	55.66
#3	✓	✓	54.27	56.22

Table 2: Ablation of Reconstructor in training and inference.

To further explore the effect of Reconstructor, we incrementally introduced it into the training and inference phases in turn. Table 2 shows the results of the ablation experiment on the development set. As seen, when Reconstructor was only introduced into the training phase (#2), there was 1.26  $F_1$  improvement over the baseline HeteGASN model (#1) in which there are not Reconstructor in the training and inference phases. Moreover, Reconstructor was introduced into the inference as a relation indicator to assist relation classification, that is, there are Reconstructor in both training and inference contain the Reconstructor (#3), As a result, there gained 0.56  $F_1$  further improvement. This shows that the proposed Reconstructor can not only encode relation information of the input document efficient but also indicate how much there is a relationship, to enhance relation classification between entity pair.

#### 4.6 Ablation of Reconstruction Loss

In the reconstruction phase, we maximized (max) the path probability if the entity pair has the ground-truth relationship

and minimized (min) the path probability otherwise. Therefore, we performed the ablation of the above two reconstruction paths. Specifically, we gradually introduced them into the proposed HeterGSAN with Reconstruction to verify the effect of two reconstruction paths, as shown in Table 3. Here, “relation” denotes entity pairs with ground-truth relationships while “no-relation” denotes entity pairs without ground-truth relationships. As seen, when one of “no-relation” (#2) and “relation” (#3) entity pairs were used to compute the reconstruction loss, their  $F_1$  scores were better than the baseline HeterGSAN (#1). This means that reconstructing one of two paths is beneficial to improve the performance of DocRE model. Meanwhile, “relation” (#3) was superior to “no-relation” (#2). In particular, both of them can complement each other to further improve  $F_1$  score (#4). This indicates that two path reconstruction methods help the DocRE model capture more diverse useful information from the input document.

	entity pair		Metric	
	relation	no-relation	Ign F1	F1
#1	–	–	52.17	54.40
#2	–	min	53.29	55.04
#3	max	–	53.53	55.55
#4	max	min	54.27	56.22

Table 3: Ablation experiments of reconstruction loss for the proposed HeterGSAN+Reconstruction model.

#### 4.7 Effect of Path Probability in Inference

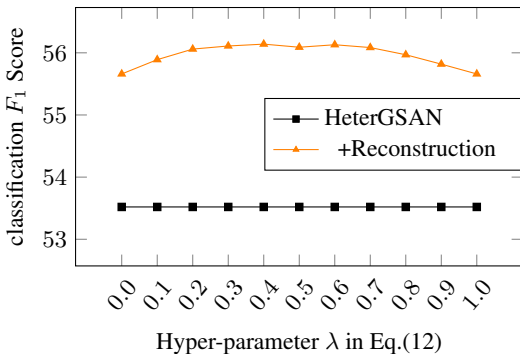


Figure 4: Classification  $F_1$  scores of different hyper-parameter  $\lambda$  for the reconstructed path probability of DocRE models (HeterGSAN and +Reconstruction) in inference.

In inference, the reconstructor is regarded as a relationship indicator to assist relation classification. The hyper-parameter  $\lambda$  in Eq.(12) keeps a trade-off between the classification scores and the construction scores when classifying the relation of each entity pair. Figure 4 shows classification  $F_1$  scores of different hyper-parameter  $\lambda$  for the reconstructed path probability of HeterGSAN and +Reconstruction models in inference. As seen,  $F_1$  scores of +Reconstruction model increased with the increasing of

$\lambda$  until 0.4, indicating that the probability of reconstructed path is useful for improving the relation classification. Subsequently, larger values of  $\lambda$  reduced the  $F_1$  scores, suggesting that excessive biased path information may be weak at keeping the gained improvement. Therefore, we set the hyper-parameter  $\lambda$  to 0.4 to control the effect of reconstructed path information in our experiments (Table 1).

#### 4.8 Evaluating Different Meta Path

To evaluate our defined three candidate meta paths, we divide entity pairs of the same type of meta path in the development set to three groups, for example, “MP1” indicates that the path representation of entity pairs are from the defined Meta Path1 (See section 3.1) during the reconstruction. Table 4 showed  $F_1$  scores of three groups (MP1, MP2, and MP3) for HeterGSAN and +Reconstruction models. As seen,  $F_1$  scores of +Reconstruction outperformed that of HeterGSAN in all three groups. This means that our defined meta paths can efficient capture path dependency between entity pairs in the reconstruction processing.

	MP1	MP2	MP3
HeterGSAN (%)	60.67	50.29	46.30
+Reconstruction (%)	61.73	52.19	47.57

Table 4:  $F_1$  scores of three groups (MP1, MP2, and MP3) with different meta paths.

#### 4.9 Path Attention Scores

To study how the reconstructor (Rec) affect the distribution of attention scores along the path in the HeterGSAN, we divided attention scores into five intervals (i.e, 0-0.2, 0.2-0.4, etc) and showed the percent of attention distribution on HeterGSAN and +Reconstruction on the development set as shown in the Table 5. The attention scores of HeterGSAN are mainly concentrated in interval 0-.2, which may indicate the hypothesis of universally learning relationship information. Thus, +Reconstruction significantly reduced the percent of attention scores in interval 0-.2 and increased the percent of remaining intervals with higher attention scores. This means that the reconstructor guides the DocRE model to pay more attention to model meta-path dependencies for the ground-truth relationships.

	0-.2	.2-.4	.4-.6	.6-.8	.8-1.0
HeterGSAN(%)	84.43	5.28	2.28	2.53	5.48
+Rec (%)	69.04	9.21	10.78	1.34	9.63

Table 5: Changes of the distribution of path attention scores

#### 4.10 Ablation of Different Meta-Paths

we reconstruct one of three meta-paths (MP1, MP2 and MP3) in each DocRE model and not consider the reconstructor in inference. The results are as follows in

[0] *Lark Force* was an Australian Army formation established in March 1941 during World War II for service in New Britain and New Ireland.

[1] Under the command of Lieutenant Colonel John Scanlan, it was raised in *Australia* and deployed to *Rabaul* and Kavieng, aboard *SS Katoomba*, *MV Neptuna* and *HMAT Zealania*, to defend their strategically important harbours and airfields.

<i>(Lark Force, Australia, P17)</i>		Reference: P17
HeterGSAN		Prediction: NA
$\log(R(r))$ : -1.1271		$\theta_1$ : -0.9828
HeterGSAN +Reconstruction		Prediction: P17
$\log(R(r))$ : -0.9760		$\theta_2$ : -1.0270
<hr/>		
<i>(Rabaul, Australia, P137)</i>		Reference: NA
HeterGSAN		Prediction: P137
$\log(R(r))$ : -0.7095		$\theta_1$ : -0.9828
HeterGSAN +Reconstruction		Prediction: NA
$\log(R(r))$ : -1.4012		$\lambda \log(R(\phi))$ : -0.6017
		$S(r)$ : -2.0029
		$\theta_2$ : -1.7340

Figure 5: Case Study

Table 6. First, reconstruction of each meta-path is beneficial to enhance the DocRE model, confirming our motivation. Thus, the improved range of each meta-path is in descending order: MP1, MP2, MP3, confirming the priority for the reconstruction meta-path in Sec 3.1. It is a statistic that the percentage of MP1, MP2, and MP3 are 22.39%, 23.15%, and 54.46%. Then, when two different meta-paths are considered, their F1 values are higher than the single path which is reconstructed, indicating that more ground-truth path relationships are reconstructed to enhance the training of the DocRE model. Similarly, considering three meta-paths gain the highest F1 on development/test sets.

type of meta-path	Dev F1	Test F1	type of meta-path	Dev F1	Test F1
None	54.40	53.52	MP1&MP2	55.26	54.40
MP1	54.79	54.22	MP1&MP3	55.12	54.37
MP2	54.78	54.20	MP2&MP3	54.96	54.28
MP3	54.54	53.88	All	55.66	54.91

Table 6: Ablation experiments of different Meta-Paths.

#### 4.11 Case Study

Figure 5 shows a case study of HeterGSAN and +Reconstruction models. For the entity pair  $\{Lark Force, Australia\}$ , HeterGSAN classified its relation to “NA” which is inconsistent with the Reference “P17” because of its classifier score -1.1271 is less than the threshold  $\theta_1$  -0.9828. In comparison, the classifier score of +Reconstruction classified its relation to “P17” which is consistent with the Reference “P17” because of its classifier score -0.9760 was greater than the threshold  $\theta_2$  -1.0270. This means that the proposed Reconstructor can better guide the training of DocRE model. For another entity pair  $\{Rabaul, Australia\}$ , the classifier scores of HeterGSAN and +Reconstruction models were greater than  $\theta_1$  and  $\theta_2$ , respectively. However, they gained a relation category P137 which is inconsistent with the Reference “NA”. When the path score -0.6017 was considered in the inference, +Reconstruction classified its relation to “NA” which is consistent with the Reference “NA”. This indicates that the inference with Reconstructor can further improve the accuracy of relation classification.

## 5 Related Work

**DocRE** Early efforts focus on classifying relationships between entity pair within a single sentence or extract entity and relations jointly in a sentence (Zeng et al. 2014; Wang et al. 2016; Wei et al. 2020; Song et al. 2019). These approaches do not consider interactions across mentions and ignore relations expressed across sentence boundaries. Recently, the extraction scope has been expanded to the entire document in the biomedical domain by only considering a few relations among chemicals (Peng et al. 2017; Quirk and Poon 2017; Gupta et al. 2019; Zhang, Qi, and Manning 2018; Christopoulou, Miwa, and Ananiadou 2019). In particular, Yao et al. (2019) proposed a large-scale human-annotated DocRED dataset. The dataset requires understanding a document and performing multi-hop reasoning and several works (Wang et al. 2019; Nan et al. 2020) have been done on the dataset.

**Reconstruction** Reconstructor was used to solve the problem that translations generated by neural network translation (NMT) often lack adequacy (Tu et al. 2016; Cheng et al. 2016). (Cheng et al. 2016) reconstructs the monolingual corpora with two separate source-to-target and target-to-source NMT models. (Tu et al. 2016) aims at enhancing adequacy of unidirectional (i.e., source-to-target) NMT via a target-to-source objective on parallel corpora. Besides, (Hu et al. 2020) uses reconstructor to pre-train a graph neural network on the unlabeled data with self-supervision to reduce the cost of labeled data.

## 6 Conclusion

This paper proposed a novel reconstruction method to guide the DocRE model to pay more attention to the learning of entity pairs with the ground-truth relationships, thereby learning an effective graph representation to classify relation category. In inference, the reconstructor is further regarded as a relation indicator to assist relation classification between entity pair. Experimental results on a large-scale DocRED dataset show that our method can greatly advance the DocRE task. In the future, we will explore more information related to relationship classification in the input document, for example, syntax constraint (Chen et al. 2018), diverse information (Chen et al. 2020), and knowledge reasoning (Cohen et al. 2020).

## Acknowledgments

We are grateful to the anonymous reviewers, senior program Committee and area chair for their insightful comments and suggestions. The corresponding authors are Kehai Chen and Tiejun Zhao. This work is supported by the National Key R&D Program of China (No. 2018YFC0830700) and Huawei Technologies CO., Ltd (No. YBN2019115122).

## References

- Chen, K.; Wang, R.; Utiyama, M.; Sumita, E.; and Zhao, T. 2018. Syntax-Directed Attention for Neural Machine Translation. In *AAAI Conference on Artificial Intelligence*, 4792–4798. New Orleans, Louisiana, USA. URL <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16060/16008>.
- Chen, K.; Wang, R.; Utiyama, M.; Sumita, E.; Zhao, T.; Yang, M.; and Zhao, H. 2020. Towards More Diverse Input Representation for Neural Machine Translation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 28: 1586–1597. doi:10.1109/TASLP.2020.2996077.
- Cheng, Y.; Xu, W.; He, Z.; He, W.; Wu, H.; Sun, M.; and Liu, Y. 2016. Semi-Supervised Learning for Neural Machine Translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1965–1974. Berlin, Germany: Association for Computational Linguistics. doi:10.18653/v1/P16-1185. URL <https://www.aclweb.org/anthology/P16-1185>.
- Christopoulou, F.; Miwa, M.; and Ananiadou, S. 2019. Connecting the Dots: Document-level Neural Relation Extraction with Edge-oriented Graphs. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 4925–4936. Hong Kong, China: Association for Computational Linguistics. doi:10.18653/v1/D19-1498. URL <https://www.aclweb.org/anthology/D19-1498>.
- Cohen, W. W.; Sun, H.; Hofer, R. A.; and Siegler, M. 2020. Scalable Neural Methods for Reasoning With a Symbolic Knowledge Base. In *International Conference on Learning Representations*. URL <https://openreview.net/forum?id=BJlguT4YPr>.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, 4171–4186.
- Guo, Z.; Zhang, Y.; and Lu, W. 2019. Attention Guided Graph Convolutional Networks for Relation Extraction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 241–251. Florence, Italy: Association for Computational Linguistics. doi:10.18653/v1/P19-1024. URL <https://www.aclweb.org/anthology/P19-1024>.
- Gupta, P.; Rajaram, S.; Schütze, H.; and Runkler, T. A. 2019. Neural Relation Extraction within and across Sentence Boundaries. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, 6513–6520. AAAI Press. URL <https://doi.org/10.1609/aaai.v33i01.33016513>.
- Hu, Z.; Dong, Y.; Wang, K.; Chang, K.-W.; and Sun, Y. 2020. GPT-GNN: Generative Pre-Training of Graph Neural Networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '20*, 1857–1867. New York, NY, USA: Association for Computing Machinery. ISBN 9781450379984. doi:10.1145/3394486.3403237. URL <https://doi.org/10.1145/3394486.3403237>.
- Nan, G.; Guo, Z.; Sekulic, I.; and Lu, W. 2020. Reasoning with Latent Structure Refinement for Document-Level Relation Extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 1546–1557. Online: Association for Computational Linguistics. doi:10.18653/v1/2020.acl-main.141. URL <https://www.aclweb.org/anthology/2020.acl-main.141>.
- Peng, N.; Poon, H.; Quirk, C.; Toutanova, K.; and Yih, W. 2017. Cross-Sentence N-ary Relation Extraction with Graph LSTMs. *Trans. Assoc. Comput. Linguistics* 5: 101–115. URL <https://transacl.org/ojs/index.php/tacl/article/view/1028>.
- Quirk, C.; and Poon, H. 2017. Distant Supervision for Relation Extraction beyond the Sentence Boundary. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, 1171–1182. Valencia, Spain: Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/E17-1110>.
- Sahu, S. K.; Christopoulou, F.; Miwa, M.; and Ananiadou, S. 2019. Inter-sentence Relation Extraction with Document-level Graph Convolutional Neural Network 4309–4316. doi:10.18653/v1/P19-1423. URL <https://www.aclweb.org/anthology/P19-1423>.
- Song, L.; Zhang, Y.; Gildea, D.; Yu, M.; Wang, Z.; and Su, J. 2019. Leveraging Dependency Forest for Neural Medical Relation Extraction. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* doi:10.18653/v1/d19-1020. URL <http://dx.doi.org/10.18653/v1/D19-1020>.
- Sun, Y.; and Han, J. 2013. Mining heterogeneous information networks: a structural analysis approach. *SIGKDD Explorations* 14: 20–28.
- Tang, H.; Cao, Y.; Zhang, Z.; Cao, J.; Fang, F.; Wang, S.; and Yin, P. 2020. HIN: Hierarchical Inference Network for Document-Level Relation Extraction. *Advances in Knowledge Discovery and Data Mining* 12084: 197 – 209.
- Tu, Z.; Liu, Y.; Shang, L.; Liu, X.; and Li, H. 2016. Neural Machine Translation with Reconstruction. *CoRR* abs/1611.01874. URL <http://arxiv.org/abs/1611.01874>.



- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L. u.; and Polosukhin, I. 2017. Attention is All you Need. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30*, 5998–6008. Curran Associates, Inc.
- Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Liò, P.; and Bengio, Y. 2018. Graph Attention Networks. In *International Conference on Learning Representations*. URL <https://openreview.net/forum?id=rJXMpikCZ>.
- Wang, H.; Focke, C.; Sylvester, R.; Mishra, N.; and Wang, W. W. J. 2019. Fine-tune Bert for DocRED with Two-step Process. *ArXiv* abs/1909.11898.
- Wang, L.; Cao, Z.; de Melo, G.; and Liu, Z. 2016. Relation Classification via Multi-Level Attention CNNs. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1298–1307. Berlin, Germany: Association for Computational Linguistics. doi:10.18653/v1/P16-1123. URL <https://www.aclweb.org/anthology/P16-1123>.
- Wei, Z.; Su, J.; Wang, Y.; Tian, Y.; and Chang, Y. 2020. A Novel Cascade Binary Tagging Framework for Relational Triple Extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 1476–1488. Online: Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.acl-main.136>.
- Yao, Y.; Ye, D.; Li, P.; Han, X.; Lin, Y.; Liu, Z.; Liu, Z.; Huang, L.; Zhou, J.; and Sun, M. 2019. DocRED: A Large-Scale Document-Level Relation Extraction Dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 764–777. Florence, Italy: Association for Computational Linguistics. doi:10.18653/v1/P19-1074. URL <https://www.aclweb.org/anthology/P19-1074>.
- Zeng, D.; Liu, K.; Lai, S.; Zhou, G.; and Zhao, J. 2014. Relation Classification via Convolutional Deep Neural Network. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, 2335–2344. Dublin, Ireland: Dublin City University and Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/C14-1220>.
- Zhang, Y.; Qi, P.; and Manning, C. D. 2018. Graph Convolution over Pruned Dependency Trees Improves Relation Extraction. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2205–2215. Brussels, Belgium: Association for Computational Linguistics. doi:10.18653/v1/D18-1244. URL <https://www.aclweb.org/anthology/D18-1244>.