# Dialogue Relation Extraction with Document-Level Heterogeneous Graph Attention Networks

**Hui Chen, Pengfei Hong, Wei Han, Navonil Majumder, Soujanya Poria**
DeCLaRe Lab, Singapore University of Technology and Design, Singapore
`hui_chen@mymail.sutd.edu.sg,`
`{hongpengfei.emrys,henryhan88888}@gmail.com,`
`{navonil_majumder,sporia}@sutd.edu.sg`

## Abstract

Dialogue relation extraction (DRE) aims to detect the relation between pairs of entities mentioned in a multi-party dialogue. It plays an essential role in constructing knowledge graphs from conversational data increasingly abundant on the internet and facilitating intelligent dialogue system development. The prior methods of DRE do not meaningfully leverage speaker information—they just prepend the utterances with the respective speaker names. Thus, they fail to model the crucial inter-speaker relations that may provide additional context to relevant argument entities through pronouns and triggers. We present a graph attention network-based method for DRE where a graph that contains meaningfully connected speaker, entity, type, and utterance nodes is constructed. This graph is fed to a graph attention network for context propagation among relevant nodes, which effectively captures the dialogue context. We empirically show that this graph-based approach quite effectively captures the relations between different argument pairs in a dialogue as it outperforms the state-of-the-art approaches by a significant margin on the benchmark dataset DialogRE. Our code is released at: https://github.com/declare-lab/dialog-HGAT.

## 1 Introduction

The relation extraction (RE) task aims to identify relations between pairs of entities that exist in a document. It plays a pivotal role in understanding unstructured text and constructing knowledge bases (Peng et al., 2017; Quirk and Poon, 2017). Although the task of document-level relation extraction has been studied extensively in the past, the task of relation extraction from dialogues has yet to receive extensive study.

Most previous works in this field focus on the professional and formal literature like biomedical documents (Li et al., 2016; Wu et al., 2019) and Wikipedia articles (Elsahar et al., 2018; Yao
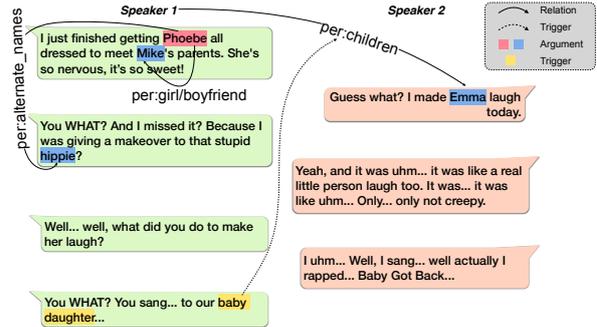


Figure 1: An example adapted from DialogRE dataset. Words with red and blue background represent subject and object entities. Words with yellow background represent triggers that facilitate the relation inference. Solid and dash lines stand for intra- and inter-utterance relations.

et al., 2019; Mesquita et al., 2019). These kinds of datasets are well-formatted and logically coherent with clear referential semantics. Hence, for most NLP tasks, analyzing a few continuous sentences is enough to grasp pivotal information. However, in dialogue relation extraction, conversational text is sampled from daily chat, which is more casual in nature. Hence its logic is simpler but more entangled, and the referential ambiguity always occurs to an external reader. Compared with formal literature, it has lower information density (Wang and Liu, 2011) and thus is more difficult for models to understand. Moreover, compared with other document-level RE datasets such as DocRED, dialogue text has much more cross-sentence relations (Yu et al., 2020).

Fig. 1 presents an example of the target dialogue, taken from DialogRE (Yu et al., 2020) dataset. In order to infer the relation between *Speaker1* and *Emma*, we may need to find some triggers to recognize the characteristics of *Emma*. Triggers are shreds of evidence that can support the inference. As we can see, the following utterances are talking about *Emma*, and the keyword

*baby daughter* mentioned by *Speaker1* is a trigger, which provides evidence that *Emma* is *Speaker1*'s daughter.

Prior works show that triggers of arguments facilitate the document-level relation inference. Thus, the DocRED dataset (Yao et al., 2019) provides several supporting evidence for each argument pair. Some efforts utilize the dependency paths of arguments to find possible triggers. For example, LSR model (Nan et al., 2020) constructs meta dependency paths of each argument pair and aggregates all the word representations located in these paths to their model to enhance the model's reasoning ability. Sahu et al. (2019) uses syntactic parsing and coreference resolution to find intra- and inter-related words of each argument. Christopoulou et al. (2019) proposes an edge-oriented graph to synthesize argument-related information. These models are graph-based and have proven powerful in encoding long-distance information. However, for dialogue relation extraction, interlocutors exist in every utterance of the dialogue, and they are often considered as an argument. Although these previous approaches have utilized entity features of arguments, most of them employ meta dependency paths to find the related words, which neglect necessary information related to speakers, since the speaker references have very little dependency features in each utterance. In this work, we formulate the dialogue relation extraction task as a classification problem, where we design a graph attention network to model semantic, syntactic, and speaker information. Compared with other graph-based models in the relation extraction task, our model is lightweight, without any costly matrix operation, and it can generalize to completely unseen graphs.

In this paper, we propose a simple yet effective attention-based heterogeneous graph neural network to tackle the dialogue relation extraction task in an inductive manner. We use multi-type features to create the graph and employ graph attention mechanism to propagate contextual information. Different from most of the previous works, our proposed model is customized for the relation extraction task in dialogue background, as we have specially modeled speaker information and designed a mechanism to propagate messages among different sentences for better inter-sentence representation learning.

The remainder of this paper is organized as follows: Section 2 briefly discusses relevant works of heterogeneous graph neural networks; Section 3 elaborates on our proposed framework; Section 4 introduces the used dataset and baseline models; Section 5 lays out the experiment results and analysis; Section 6 concludes the paper.

## 2 Related Work

Graph-based models have raised widespread attention from NLP researchers, as it is demonstrated as a powerful mathematical tool to represent complicated syntactic and semantic relations among structured language data. Early work applies classic graph processing algorithms onto language graphs. Pang and Lee (2004) constructed a text graph and adopt the minimum-cut method to cluster the nodes for sentiment analysis. Agirre and Soroa (2009) leveraged PageRank algorithm on personalized subgraphs of a wordnet to disambiguate polysemous words according to connected context words.

Recently, graph neural networks (GNN) (Kipf and Welling, 2017) becomes popular in relation extraction tasks. For example, Peng et al. (2017) tried to build a computation graph from syntactic parsing trees and employed graph LSTM to obtain better word embeddings for multi-ary relation extraction. Zhang et al. (2018) designed a pruning algorithm for syntactic graphs and add a graph convolution layer on top of the sequential LSTM encoder in the learning process. The combination with typical attention-based language models such as transformer (Vaswani et al., 2017) is also studied. Cai and Lam (2020) and Yao et al. (2020) used transformer-based graph convolutional networks to explicitly encode relations among distant syntactic nodes, to address the long-distance propagation issue.

Based on GNN, heterogeneous graph neural networks are proposed and have been applied in many NLP tasks, like text classification (Linmei et al., 2019), text summarization (Wang et al., 2020), user profiling (Chen et al., 2019), and event categorization (Peng et al., 2019). The prior work proves that heterogeneous graph neural network is a powerful tool in NLP. For the relation extraction task, Christopoulou et al. (2019) constructed an edge-oriented heterogeneous graph that contains sentence, mention, and entity information. However, syntactic information is neglected in their model. Differently, homogeneous nodes in our

graph are all independent, and we take syntactic features to initialize sentence information as well as edges features.

# 3 Method

## 3.1 Task Definition

Given a dialogue containing $N$ utterances $\mathcal{D} = \{u_1, u_2, ..., u_N\}$ and a couple of argument pairs $\mathcal{A} = \{(x_1, y_1), (x_2, y_2), ...\}$, where subject $x_i$ and object $y_i$ are entities mentioned in the dialogue, the goal is to identify the relation between argument pairs $(x_i, y_i)$. For document-level relation extraction task, there are many cross-sentence relations which are supported by various sentences.

## 3.2 Model Overview

In this work, we introduce an attention-based graph network to tackle the problem where each conversation is represented as a heterogeneous graph. We first utilize an utterance encoder, which is composed of two Bidirectional long short-term memory networks to encode conversational information. These utterance encodings, along with word embeddings, speaker embeddings, argument embeddings, and type embedding, are logically connected to form a heterogeneous graph, which will be discussed in detail later in this section. Further, this graph is fed through five graph attention layers (Veličković et al., 2018) that aggregate information from the neighboring nodes. Lastly, we concatenate the learned argument embeddings and feed them to a classifier. An overview of the proposed model is shown in Fig. 2.

## 3.3 Data Preprocessing

In the data preprocessing period, we use spaCy[1] to tokenize utterances, and at the same time, we obtain part-of-speech (POS) tags as well as named entity types of each token.

## 3.4 Utterance Encoder

Given a dialogue $\mathcal{D} = \{u_1, u_2, ..., u_N\}$, we use GloVe (Pennington et al., 2014) to initialize the word embeddings and then feed them to a contextual Bidirectional Long Short-Term Memory network (BiLSTM) to obtain contextualized representations. The operation of BiLSTM can be de-

---

[1] https://spacy.io

fined as:

$$\overleftarrow{h}_j^i = LSTM_l(\overleftarrow{h}_{j+1}^i, e_j^i) \tag{1}$$

$$\overrightarrow{h}_j^i = LSTM_r(\overrightarrow{h}_{j-1}^i, e_j^i) \tag{2}$$

$$h_j^i = [\overleftarrow{h}_j^i; \overrightarrow{h}_j^i] \tag{3}$$

where $\overleftarrow{h}_j^i$ and $\overrightarrow{h}_j^i$ denote the hidden representations in the $j$-th layer of utterance $u_i$ from two directions, $h_j^i$ is the contextual representation which is the concatenation of $\overleftarrow{h}_j^i$ and $\overrightarrow{h}_j^i$, and $e_j^i$ stands for the embedding of the $j$-th token in utterance $u_i$. Unlike the previous approaches (Christopoulou et al., 2019; Nan et al., 2020) that only adopt semantic contextual features in utterance encoding, we add syntactic features such as POS tags and named entity types to the contextual representations. The embedding of each token in the utterance can be described as:

$$e = [e_w; e_p; e_t] \tag{4}$$

where we concatenate word embedding $e_w$ initialized by GloVe (Pennington et al., 2014), syntactic POS embedding $e_p$, and type embedding $e_t$ to form the token embedding $e$.

Moreover, we believe conversation-level contextual features play an important role in understanding a conversation. To encode non-local contextual information between each utterance, we apply max pool operation to the hidden states of each utterance-level BiLSTM (local LSTM), and then feed the sequence $c = \{c_1, c_2, ..., c_N\}$ to a conversation-level BiLSTM (global LSTM). The operation of global LSTM is the same as Eqs. (1) to (3).

## 3.5 Graph Construction

### 3.5.1 Node Construction

In our model, we design a heterogeneous graph network containing five types of nodes: utterance nodes, type nodes, word nodes, speaker nodes, and argument nodes. Each type of node is used to encode a type of information in the dialogue. In the task, only word nodes, speaker nodes and argument nodes are probable to attend the final classification process. In other words, only these types of nodes are possible arguments. For simplicity, we name them as *basic nodes* in our illustration.

**Utterance and Type Nodes** Utterance nodes are initialized by the utterance embeddings which we
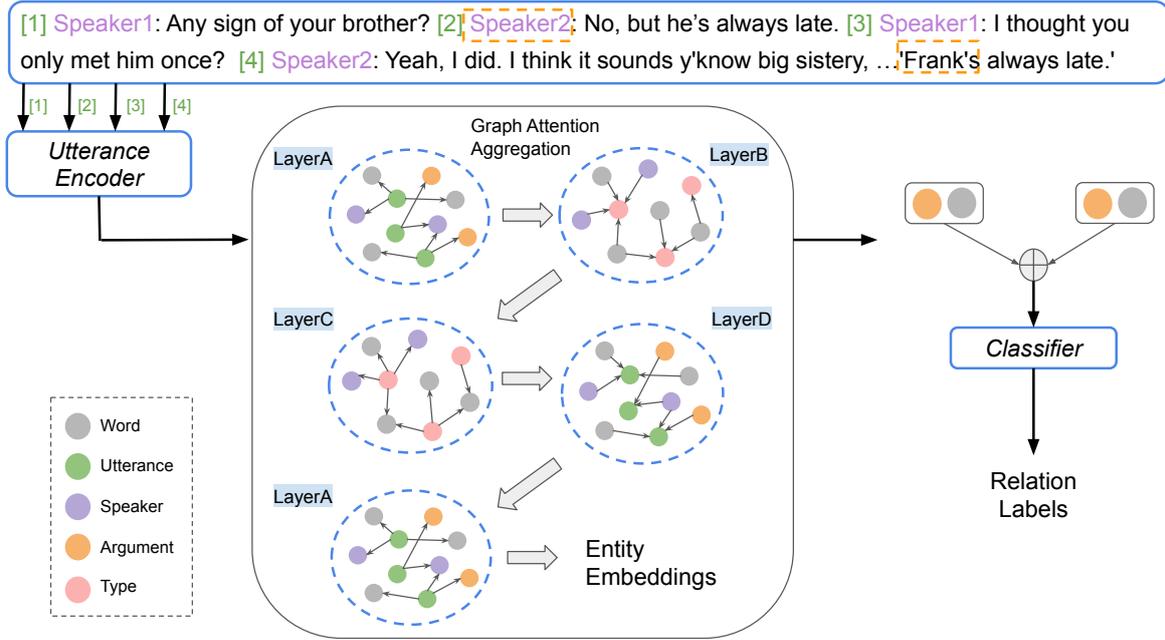
Figure 2: An overview of the proposed model.

obtain from the utterance encoder. They are connected with the basic nodes which constitute the utterance. Type nodes represent the entity types of words in an utterance, where a variety of named and numeric entities, such as PERSON or LOCATION, are included. Since one mention may have different types in one conversation, type nodes can facilitate information integration. For example, 'Frank' can be a string if it represents an alternative name, and at the same time, it can be a person if it refers to a speaker in the conversation. Type nodes are connected with the basic nodes having the type attribute in the conversation. Each type node is initialized with a specific type of embedding. We believe that type of information has a positive influence on the relation inference process.

**Basic Nodes** Word nodes represent the vocabulary of a conversation. Each word node is connected with the utterance, which contains the word and it is also connected with all the possible types that the word may have in the conversation. We initialize the states of word nodes with GloVe (Pennington et al., 2014).

Speaker nodes represent each unique speaker in the conversation. Each speaker node is connected with the utterances uttered by the speaker himself/herself. This type of node is initialized with some specific embeddings and can gather informa-

tion from different speakers.

Argument nodes are two special nodes that are used to encode relative positional information of argument pairs. There are two argument nodes in each graph in total. One stands for the subject argument and the other represents the object argument. Similarly, argument nodes are also encoded by specific embeddings.

### 3.5.2 Edge Construction

The proposed graph is undirected but the propagation has directions. There are five types of edges: utterance-word, utterance-argument, utterance-speaker, type-word, and type-argument edges. Each edge has its own type. These edges are randomly initialized except the utterance-word edge.

For the edge between utterance and word nodes, we adopt POS tags to initialize the edge features. This type of edge aggregates not only global semantic features of the conversation but also local syntactic features to the word nodes.

### 3.5.3 Graph Attention Mechanism

We use graph attention mechanism (Veličković et al., 2018) to aggregate neighboring information to the target node. Suppose we have a node $i$ and some neighborhood nodes $j$, the graph attention mechanism can be described as:

$$\mathcal{F}(h_i, h_j) = \text{LeakyReLU}(\mathbf{a}^T(\mathbf{W}_i h_i; \mathbf{W}_j h_j; \mathbf{E}_{ij}))$$

$$\tag{5}$$

$$\alpha_{ij} = \text{softmax}(\mathcal{F}(h_i, h_j)) \tag{6}$$

$$= \frac{\exp(\mathcal{F}(h_i, h_j))}{\sum_k \exp(\mathcal{F}(h_i, h_k))} \tag{7}$$

$$h'_i = ||_{k=1}^{K} \sigma(\sum_j \alpha_{ij}^k \mathbf{W}_q^k h_j) \tag{8}$$

where $h_i$ and $h_j$ are representations of node $i$ and nodes $j$, $\mathbf{W}_i$, $\mathbf{W}_j$, $\mathbf{W}_q$ and $\mathbf{a}^T$ are trainable weight matrices, $\mathbf{E}_{ij}$ is the edge weight matrix that is mapped to the multi-dimensional embedding space, $\alpha_{ij}$ is the attention weight between $i$ and $j$, $\sigma$ is an activation function, and $||$ is concatenation operation.

### 3.5.4 Message Propagation

As shown in Fig. 2, there are five layers in our proposed graph module, where each layer represents an aggregation. There are four types of layers that we mark in the figure. LayerA and LayerD contain the message propagation between utterance nodes and basic nodes, and similarly, LayerB and LayerC are the message propagation between basic nodes and type nodes. We would call the whole message propagation path meta path. Different meta path strategies may lead to different performance.

Our meta path in this work can be described as follows: First, we use utterance nodes to update word nodes, speaker nodes, and argument nodes; secondly, the updated word nodes and argument nodes pass messages to type nodes; then type nodes conversely update the word nodes and argument nodes; next we use word nodes, speaker nodes, and argument nodes to update utterance nodes; and lastly the updated utterance nodes update word nodes, speaker nodes and argument nodes. The path can be denoted as $V_u - V_b - V_t - V_b - V_u - V_b$, where $V_u$, $V_b$, and $V_t$ refer to utterance nodes, basic nodes, and type nodes.

Following Wang et al. (2020), we add a residual connection (He et al., 2016) to avoid gradient vanishing during updating:

$$\hat{h}_i = \bar{h}_i + h'_i \tag{9}$$

where $\bar{h}_i$ is the output learned in the graph attention layer, and $h'_i$ is the original input of the graph attention layer.

In message passing, except for graph attention operation, there is also a two-layer feed-forward network which can be denoted as:

$$h_i^{new} = \text{FFN}(\hat{h}_i) \tag{10}$$

Suppose we have the initial embeddings of utterance nodes, basic nodes and type nodes, denoted as embedding matrices $\mathbf{H}_u = \{\mathbf{H}_u, \mathbf{H}_b, \mathbf{H}_t\}$, the message propagating process can be written as:

$$\mathbf{H}_b^1 = \text{GAT}(\mathbf{H}_b^0, \mathbf{H}_u^0) \tag{11}$$

$$\mathbf{H}_t^1 = \text{GAT}(\mathbf{H}_t^0, \mathbf{H}_b^1) \tag{12}$$

$$\mathbf{H}_b^2 = \text{GAT}(\mathbf{H}_b^1, \mathbf{H}_t^1) \tag{13}$$

$$\mathbf{H}_u^1 = \text{GAT}(\mathbf{H}_u^0, \mathbf{H}_b^2) \tag{14}$$

$$\mathbf{H}_b^3 = \text{GAT}(\mathbf{H}_b^2, \mathbf{H}_u^1) \tag{15}$$

where the GAT operation is the same as Eqs. (5) to (10). The superscripts represent the $n^{th}$ update of the matrix and 0 marks the initial state.

### 3.6 Relation Classifier

After the message propagation in the heterogeneous graph, we obtain new representations of all entities. We select the argument nodes $\tau_x$ and $\tau_y$, as well as the corresponding word nodes $e_x$ and $e_y$ from basic nodes, and concatenate them. Finally, they are fed to a linear transformation and a sigmoid function to get the predictions:

$$e'_x = [\text{maxpool}(\tau_x); \text{maxpool}(e_x)] \tag{16}$$

$$e'_y = [\text{maxpool}(\tau_y); \text{maxpool}(e_y)] \tag{17}$$

$$e' = [e'_x; e'_y] \tag{18}$$

$$P(r|e_x, e_y) = \sigma(\mathbf{W}_e e' + b_e)_r \tag{19}$$

where $P(r|e_x, e_y)$ is the probability of relation type $r$ given argument pair $(e_x, e_y)$, $\mathbf{W}_e$ and $b_e$ are linear transformation weight and bias vector, $maxpool$ is max pooling operation, and $\sigma$ is sigmoid function.

## 4 Experiments

### 4.1 Dataset Used

We evaluate the proposed framework on the DialogRE dataset (Yu et al., 2020), which contains 1,788 dialogues and 10,168 relational triples. The data statistics are shown in Table 1. DialogRE is

adapted from the complete transcripts of *Friends*, a widely used corpus in dialogue research these years (Chen et al., 2017; Zhou and Choi, 2018; Yang and Choi, 2019; Poria et al., 2019), and there are 36 possible relation types, most of which focus on biographical attributes of person entities. Each dialogue contains several relational triples $(x, y, r)$, and the task is to predict the relation $r$ between each argument pair $(x, y)$. In the experiments, the dataset is partitioned into train, dev, and test set with a roughly 60/20/20 ratio. Following the evaluation metrics of DialogRE, we report macro $F1$ scores of the proposed model and all the baselines in both the standard and conversational settings. In the following sections, we use $F1_c$ to represent $F1$ scores in the conversational setting.

| DialogRE | Train | Dev | Test |
|---|---|---|---|
| #Conversations | 1073 | 358 | 357 |
| #Argument Pairs | 5963 | 1928 | 1858 |
| Average dialogue length | 229.5 | 224.1 | 214.2 |
| Average # of turns | 13.1 | 13.1 | 12.4 |
| Average # of speakers | 3.3 | 3.2 | 3.3 |

Table 1: DialogRE dataset statistics.

## 4.2 Baseline models

### 4.2.1 Sequence-based Models

We select convolutional neural networks (CNN) (Zeng et al., 2014), LSTM, and BiL-STM (Cai et al., 2016) as the sequence-based baselines. These models take word embeddings, mention embeddings, and type embeddings as features. Concretely, they use GloVe and spaCy to get word embeddings and label named-entity types, and then take an average of all the embeddings of mention names for each entity to get mention embeddings.

### 4.2.2 Graph-based Models

As our proposed model is graph-based, we also select two graph-based models AGGCN (Guo et al., 2019) and LSR (Nan et al., 2020) as the baselines. AGGCN directly feeds the full dependency tree of each sentence to a graph convolutional network, which takes self-attention weights as soft edges. It achieves state-of-the-art results in various relation extraction tasks. LSR adopts an adaptation of Kirchhoff's Matrix-Tree Theorem (Tutte, 1984; Koo et al., 2007) to induce the latent dependency

structure of each document and then feeds the latent structure to a densely connected graph convolutional network to inference the relations. These graph-based models both utilize dependency information to construct the inference graph.

## 5 Result and Analysis

### 5.1 Comparison with Baselines

We present our main results on DialogRE dataset in Table 2. As shown in Table 2, our model surpasses the state-of-the-art method by 9.6%/7.5% $F1$ scores, and 8.4%/5.7% $F1_c$ scores in both validation and test sets, which demonstrates the effectiveness of the information propagation along task-specific functional meta-paths in the heterogeneous graph. As a result, inter-sentence communication usually passes through a long distance, which causes information loss or degradation. However, this kind of information transmission is critically important for dialog-style text, because logical connections are not locally compact within adjacent sentences, instead, they are spread over the whole conversations. Our proposed model constructs a heterogeneous graph with shorter distances between logically closed but syntactically faraway word pairs. Hence the long-distance issue is mitigated.

We also compare the model sizes as an efficiency indicator. Although creating numerous nodes and edges inevitably brings overhead, the total number of parameters is still moderate.

### 5.2 Ablation Study

To understand the impact of our model's components, we perform ablation studies using our proposed model on the DialogRE dataset. The ablation results are shown in Table 3. First, we remove local LSTM and global LSTM. The dropping accuracy proves that the contextual encoder plays an important role in semantic feature extraction. Second, we remove the specific argument nodes and have observed that $F1$ and $F1_c$ scores decrease to 55.0% and 50.2% on test set. This proves that our design on argument nodes effectively synthesizes argument features to the model. Further, we test the performance of the syntactic features we inject by removing POS embedding, NER embedding, and POS edge features. The scores record a decrease under all these experiment settings. Notably, removing POS embedding leads to even about 2% drops in all the evaluation

| Model | #params | Dev (%) | | Test (%) | |
|---|---|---|---|---|---|
| | | $F1$ | $F1_c$ | $F1$ | $F1_c$ |
| Majority (Yu et al., 2020) | - | 38.9 | 38.7 | 35.8 | 35.8 |
| CNN (Yu et al., 2020) | - | 46.1 | 43.7 | 48.0 | 45.0 |
| LSTM (Yu et al., 2020) | - | 46.7 | 44.2 | 47.4 | 44.9 |
| BiLSTM (Yu et al., 2020) | 4.1M | 48.1 | 44.3 | 48.6 | 45.0 |
| AGGCN (Guo et al., 2019) | 3.7M | 46.6 | 40.5 | 46.2 | 39.5 |
| LSR (Nan et al., 2020) | 20.5M | 44.5 | - | 44.4 | - |
| This work | 4.0M | **57.7** | **52.7** | **56.1** | **50.7** |

Table 2: Main results on DialogRE dataset. Values in the #params column refer to parameter sizes of the models. $F1$ and $F1_c$ are macro $F1$ scores under standard setting and conversational setting, respectively. Word embeddings of the models are captured by GloVe (Pennington et al., 2014).

metrics.

| Model | Dev (%) | | Test (%) | |
|---|---|---|---|---|
| | $F1$ | $F1_c$ | $F1$ | $F1_c$ |
| Full model | 57.7 | 52.7 | 56.1 | 50.7 |
| w/o Local BiLSTM | 54.9 | 50.0 | 55.3 | 50.3 |
| w/o Global BiLSTM | 54.7 | 50.2 | 53.5 | 48.7 |
| w/o Argument nodes | 56.0 | 51.3 | 55.0 | 50.2 |
| w/o POS embedding | 54.6 | 50.9 | 53.0 | 48.5 |
| w/o NER embedding | 56.8 | 51.5 | 54.2 | 49.2 |
| w/o POS edge weights | 56.9 | 52.4 | 54.7 | 50.4 |

Table 3: Ablation results on DialogRE dataset.

## 5.3 Effect of the Meta Path

We test the performance of our message propagation strategy via changing meta-path strategies. In our proposed model, there are five layers in the heterogeneous graph. Those basic nodes, corresponding to different types of words, speakers, and arguments, are updated totally three times, i.e., they are first updated by utterance nodes, second updated by type nodes, and ultimately updated by utterance again. To investigate the meta path's effect, we compare our proposed five-layer graph module with different strategies where the numbers of layers are one, seven, and nine in Table 4. In Strategy1, we only set up one LayerA, where the basic nodes are updated by the initialized utterance nodes once. We observe that all the macro $F1$ scores drop dramatically, showing the one-layer structure is not deep enough to grasp complex dependencies. To make node features more informative, we would add more layers. At this time, we may be curious about how many layers the module should have to induct an optimal structure in this task. In Strategy2 and Strategy3,

we design a seven-layer module and a nine-layer module, respectively. For Strategy2, the order of layers is A-B-C-D-A-D-A, where A,B,C, and D are layer labels introduced in Fig. 2. Compared with our proposed module, scores on validation set decrease about 1% and scores on test set decrease 1.7% and 0.6 % with the standard-setting and the conversational setting, respectively. However, the module with nine layers in Strategy3 shows a larger gap between itself and the best performance, where the order of layers is A-B-C-D-A-B-C-D-A. We think this is probably because the structure is so complicated, which causes an over-smooth problem and prevents itself from learning meaningful hidden representations.

| Strategy | Dev (%) | | Test (%) | |
|---|---|---|---|---|
| | $F1$ | $F1_c$ | $F1$ | $F1_c$ |
| This work (L=5) | 57.7 | 52.7 | 56.1 | 50.7 |
| Strategy1 (L=1) | 48.3 | 46.5 | 48.4 | 45.9 |
| Strategy2 (L=7) | 56.8 | 51.6 | 54.4 | 50.1 |
| Strategy3 (L=9) | 53.8 | 49.1 | 52.2 | 47.2 |

Table 4: Comparison with different meta-path strategies on DialogRE dataset. 'L' means the number of layers in the graph module.

## 5.4 Case Studies

In the dataset, 95% of argument pairs span in at least two consecutive sentences instead of being restricted to the same sentence. Therefore, it is crucial that the model can tackle long-distance learning issues. Compared with the LSTM model, direct connections among different types of nodes in HGNN reduce the length of information propagation paths between pairs of argument nodes. Considering the following example in Fig. 3, sub-
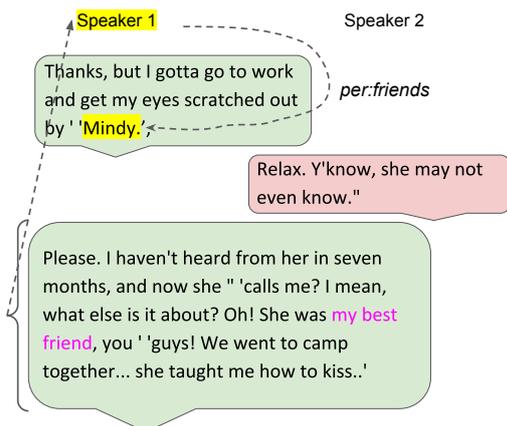
Figure 3: An example to show the effective message propagation between argument pairs

ject a - 'Mindy' and object b - 'Speaker 1' share the relationship 'per:friends', which is indicated by the trigger 'my best friend' in the first utterance. The entity information is relayed from 'Mindy' to 'Speaker 1' in the update process: 'speaker 1' node aggregates utterance level information from its neighbor nodes containing a. the relation trigger 'best friend'. b. in BiLSTM model, the key information has to travel a long journey from the subject entity word to the object one as there are too many words between them in the context.

### 5.5 Error Analysis

Type information involves in the information propagation process and thus affects the contents of output embeddings. The model is prone to make incorrectly and biased predictions. If it fails to receive enough certainty from other information sources and then can only rely on the entity types of the two arguments. For example, given an argument pair of two human names, both are named entity type 'PERSON'. Sometimes the model inclines to deem the relationship between the two arguments to be 'per:alternate_name' instead of the correct answer 'per:alumni' or 'per:roommate'. This is because among all of these classes, 'PERSON-PERSON' is a preferable type pair. However, the class 'per:alternate_name' (22.01%) presents more frequently than 'per:alumni' (1.83%) and 'per:roommate' (1.29%) in the dataset. When information aggregated from all sources other than the argument pair is not evident for judgment, en-

tity bias misguides the model to the wrong classification results.

### 6 Conclusion

In this work, we present an attention-based heterogeneous graph network to deal with the dialogue relation extraction task in an inductive manner. This heterogeneous graph attention network has modeled multi-type features of the conversation, such as utterance, word, speaker, argument, and entity type information. On the benchmark DialogRE dataset, our proposed framework outperforms the strongest baselines and the state-of-the-art approaches by a significant margin, which proves the proposed framework can effectively capture relations between different entities in the conversation. Future work will focus on making use of latent relations between entities that exist in dialogue history to develop intelligent conversational agents.

### References

Eneko Agirre and Aitor Soroa. 2009. Personalizing pagerank for word sense disambiguation. In *Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009)*, pages 33–41.

Deng Cai and Wai Lam. 2020. Graph transformer for graph-to-sequence learning. In *AAAI*, pages 7464–7471.

Rui Cai, Xiaodong Zhang, and Houfeng Wang. 2016. Bidirectional recurrent convolutional neural network for relation classification. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 756–765.

Henry Y. Chen, Ethan Zhou, and Jinho D. Choi. 2017. Robust coreference resolution and entity linking on dialogues: Character identification on TV show transcripts. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 216–225, Vancouver, Canada. Association for Computational Linguistics.

Weijian Chen, Yulong Gu, Zhaochun Ren, Xiangnan He, Hongtao Xie, Tong Guo, Dawei Yin, and Yongdong Zhang. 2019. Semi-supervised user profiling with heterogeneous graph attention networks. In *IJCAI*, volume 19, pages 2116–2122.

Fenia Christopoulou, Makoto Miwa, and Sophia Ananiadou. 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)*, pages 4927–4938.

Hady Elsahar, Pavlos Vougiouklis, Arslen Remaci, Christophe Gravier, Jonathon Hare, Frederique Laforest, and Elena Simperl. 2018. T-rex: A large scale alignment of natural language with knowledge base triples. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.

Zhijiang Guo, Yan Zhang, and Wei Lu. 2019. Attention guided graph convolutional networks for relation extraction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 241–251.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.

Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations*.

Terry Koo, Amir Globerson, Xavier Carreras, and Michael Collins. 2007. Structured prediction models via the matrix-tree theorem. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 141–150.

Jiao Li, Yueping Sun, Robin J Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J Mattingly, Thomas C Wiegers, and Zhiyong Lu. 2016. Biocreative v cdr task corpus: a resource for chemical disease relation extraction. *Database*, 2016.

Hu Linmei, Tianchi Yang, Chuan Shi, Houye Ji, and Xiaoli Li. 2019. Heterogeneous graph attention networks for semi-supervised short text classification. 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)*, pages 4823–4832.

Filipe Mesquita, Matteo Cannaviccio, Jordan Schmidek, Paramita Mirza, and Denilson Barbosa. 2019. Knowledgenet: A benchmark dataset for knowledge base population. 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)*, pages 749–758.

Guoshun Nan, Zhijiang Guo, Ivan Sekulic, and Wei Lu. 2020. Reasoning with latent structure refinement for document-level relation extraction. In *Proceedings*

*of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1546–1557.

Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*, pages 271–278.

Hao Peng, Jianxin Li, Qiran Gong, Yangqiu Song, Yuanxing Ning, Kunfeng Lai, and Philip S Yu. 2019. Fine-grained event categorization with heterogeneous graph convolutional networks. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 3238–3245. AAAI Press.

Nanyun Peng, Hoifung Poon, Chris Quirk, Kristina Toutanova, and Wen-tau Yih. 2017. Cross-sentence n-ary relation extraction with graph lstms. *Transactions of the Association for Computational Linguistics*, 5:101–115.

Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.

Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2019. Meld: A multimodal multi-party dataset for emotion recognition in conversations. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 527–536.

Chris Quirk and Hoifung Poon. 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*, pages 1171–1182.

Sunil Kumar Sahu, Fenia Christopoulou, Makoto Miwa, and Sophia Ananiadou. 2019. Inter-sentence relation extraction with document-level graph convolutional neural network. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4309–4316.

William Thomas Tutte. 1984. Graph theory. In *Clarendon Press*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. In *6th International Conference on Learning Representations*.

Danqing Wang, Pengfei Liu, Yining Zheng, Xipeng Qiu, and Xuanjing Huang. 2020. Heterogeneous graph neural networks for extractive document summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6209–6219.

Dong Wang and Yang Liu. 2011. A pilot study of opinion summarization in conversations. In *Proceedings of the 49th annual meeting of the Association for Computational Linguistics: Human language technologies*, pages 331–339.

Ye Wu, Ruibang Luo, Henry CM Leung, Hing-Fung Ting, and Tak-Wah Lam. 2019. Renet: A deep learning approach for extracting gene-disease associations from literature. In *International Conference on Research in Computational Molecular Biology*, pages 272–284. Springer.

Zhengzhe Yang and Jinho D Choi. 2019. Friendsqa: Open-domain question answering on tv show transcripts. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 188–197.

Shaowei Yao, Tianming Wang, and Xiaojun Wan. 2020. Heterogeneous graph transformer for graph-to-sequence learning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7145–7154.

Yuan Yao, Deming Ye, Peng Li, Xu Han, Yankai Lin, Zhenghao Liu, Zhiyuan Liu, Lixin Huang, Jie Zhou, and Maosong Sun. 2019. Docred: A large-scale document-level relation extraction dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 764–777.

Dian Yu, Kai Sun, Claire Cardie, and Dong Yu. 2020. Dialogue-based relation extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4927–4940.

Daojian Zeng, Kang Liu, Siwei Lai, Guangyou Zhou, and Jun Zhao. 2014. Relation classification via convolutional deep neural network. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 2335–2344.

Yuhao Zhang, Peng Qi, and Christopher D Manning. 2018. Graph convolution over pruned dependency trees improves relation extraction. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2205–2215.

Ethan Zhou and Jinho D Choi. 2018. They exist! introducing plural mentions to coreference resolution and entity linking. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 24–34.

# A    Settings and Hyperparameters

In our experiments, we tune the parameters of batch size, learning rate, and BiLSTM hidden size by testing the performance on the validation set. Table 5 lists the major parameters used in our experiments.

| Parameter | Value |
|---|---|
| Word embedding dimension | 300 |
| NER embedding dimension | 30 |
| POS embedding dimension | 30 |
| Local BiLSTM hidden Size | 200 |
| Local BiLSTM layers | 2 |
| Global BiLSTM hidden Size | 128 |
| Global BiLSTM layers | 2 |
| # Multihead attention | 10 |
| Learning rate | 0.0005 |
| Batch size | 16 |
| Edge embedding dimension | 50 |

Table 5: Parameter settings.

# B    Statistics of Relation Labels

Table 6 shows statistics of relation labels in DialogRE dataset. In the train set and test set, there are 35 types of relations, while in the dev set, there are 37 types. 'gpe:birth_in_place' and 'per:place_of_birth' only exist in the dev set.

| Relation Type | Quantity | | | Percentage (%) | | |
|---|---|---|---|---|---|---|
| | train | dev | test | train | dev | test |
| per:alternate_names | 1319 | 410 | 409 | 22.12 | 21.26 | 22.01 |
| unanswerable | 1308 | 404 | 388 | 21.94 | 20.95 | 20.88 |
| per:girl/boyfriend | 502 | 170 | 136 | 8.42 | 8.82 | 7.32 |
| per:positive_impression | 476 | 149 | 138 | 7.98 | 7.73 | 7.43 |
| per:friends | 444 | 156 | 122 | 7.45 | 8.09 | 6.57 |
| per:title | 250 | 86 | 78 | 4.19 | 4.46 | 4.20 |
| per:spouse | 204 | 72 | 54 | 3.42 | 3.73 | 2.91 |
| per:siblings | 196 | 64 | 58 | 3.29 | 3.32 | 3.12 |
| per:children | 171 | 55 | 48 | 2.87 | 2.85 | 2.58 |
| per:parents | 171 | 55 | 48 | 2.87 | 2.85 | 2.58 |
| per:negative_impression | 156 | 46 | 56 | 2.62 | 2.39 | 3.01 |
| per:roommate | 140 | 44 | 24 | 2.35 | 2.28 | 1.29 |
| per:alumni | 110 | 38 | 34 | 1.84 | 1.97 | 1.83 |
| per:other_family | 66 | 29 | 30 | 1.11 | 1.50 | 1.61 |
| per:works | 58 | 12 | 19 | 0.97 | 0.62 | 1.02 |
| per:age | 53 | 15 | 10 | 0.89 | 0.78 | 0.54 |
| per:client | 52 | 18 | 18 | 0.87 | 0.93 | 0.97 |
| per:place_of_residence | 49 | 12 | 23 | 0.82 | 0.62 | 1.24 |
| gpe:residents_of_place | 49 | 12 | 23 | 0.82 | 0.62 | 1.24 |
| per:boss | 49 | 13 | 12 | 0.82 | 0.67 | 0.65 |
| per:subordinate | 49 | 13 | 12 | 0.82 | 0.67 | 0.65 |
| per:visited_place | 48 | 20 | 25 | 0.80 | 1.04 | 1.35 |
| gpe:visitors_of_place | 48 | 20 | 25 | 0.80 | 1.04 | 1.35 |
| per:employee_or_member_of | 46 | 11 | 15 | 0.77 | 0.57 | 0.81 |
| org:employees_or_members | 46 | 11 | 15 | 0.77 | 0.57 | 0.81 |
| per:neighbor | 40 | 14 | 12 | 0.67 | 0.73 | 0.65 |
| per:place_of_work | 37 | 9 | 25 | 0.62 | 0.47 | 1.35 |
| per:pet | 30 | 10 | 8 | 0.50 | 0.52 | 0.43 |
| per:acquaintance | 26 | 12 | 34 | 0.44 | 0.62 | 1.83 |
| per:origin | 21 | 4 | 1 | 0.35 | 0.21 | 0.05 |
| per:dates | 20 | 14 | 6 | 0.34 | 0.73 | 0.33 |
| per:schools_attended | 5 | 2 | 1 | 0.08 | 0.10 | 0.05 |
| org:students | 5 | 2 | 1 | 0.08 | 0.10 | 0.05 |
| per:major | 2 | 1 | 3 | 0.03 | 0.05 | 0.16 |
| per:date_of_birth | 1 | 2 | 3 | 0.02 | 0.10 | 0.16 |
| gpe:birth_in_place | 0 | 1 | 0 | 0 | 0.05 | 0 |
| per:place_of_birth | 0 | 1 | 0 | 0 | 0.05 | 0 |

Table 6: Statistics of relation labels in DialogRE dataset.