

Knowledge Graph Enhanced Event Extraction in Financial Documents

Kaihao Guo, Tianpei Jiang, Haipeng Zhang*

School of Information Science and Technology, ShanghaiTech University

Shanghai, China

{guokh, jiangtp, zhanghp}@shanghaitech.edu.cn

Abstract—Event extraction is a classic task in natural language processing with wide use in handling large amount of yet rapidly growing financial, legal, medical, and government documents which often contain multiple events with their elements scattered and mixed across the documents, making the problem much more difficult. Though the underlying relations between event elements to be extracted provide helpful contextual information, they are somehow overlooked in prior studies.

We showcase the enhancement to this task brought by utilizing the knowledge graph that captures entity relations and their attributes. We propose a first event extraction framework that embeds a knowledge graph through a Graph Neural Network and integrates the embedding with regular features, all at document-level. Specifically, for extracting events from Chinese financial announcements, our method outperforms the state-of-the-art method by 5.3% in F1-score.

Index Terms—Knowledge Graph, Event Extraction, Graph Neural Network, Financial Documents, Financial Events

I. INTRODUCTION

Event extraction is a classic yet challenging task, which obtains key elements with attributes and relations between elements from documents. Take the equity pledge event in financial announcements for example, it contains various elements that are scattered across the document, including pledged shares, percentage of shares, start and end time, pledgers, and pledgees. Besides, there are relations between entities that can be crucial and yet error-prone. For instance, the pledger-pledgee relation is easily mistaken because the two words pledger and pledgee have multiple aliases that may cause confusion in Chinese. Furthermore, it is common that a document contains multiple events which makes the task much more complicated – one event’s elements can be confused with the elements of other events.

Intuitively, external information can help event extraction tasks. Knowledge bases, such as the FrameNet lexical database, are utilized to label events [1]. In a more specific domain, tree-structured long short-term memory networks have been used to capture gene tree dependency information in order to extract biomedical events [2]. However, the rich relational and structural information from knowledge graphs that links elements has not been fully exploited for enhancement in this scenario.

On the basis of a state-of-the-art document-level event extraction method [3], we propose a framework that integrates a prior knowledge graph to extract events from long documents. It consists of three steps. First, we use a graph encoder to embed the knowledge graph information for entities in our graph. Second, we perform named entity recognition (NER) to identify possibly relevant entities. Finally, we use two different approaches, Transformer layer [4] and linear layer, to input the graph embedding of the extracted entity to the event extraction model to help the tasks.

We apply the framework to perform event extraction on the Chinese financial announcements of the listed companies in China, in which we identify and extract six types of events (equity freeze, equity repurchase, equity overweight, equity underweight, equity pledge and lawsuit) that may indicate risks and require immediate actions from market regulators and financial institutions that have corresponding companies in their portfolios. Our dataset consists of 21,871 financial announcements from over 3,000 listed companies. With regard to the knowledge graph, we collect directional relation data for these companies, covering 347,570 entities (companies and persons) and 441,127 relations (e.g. share-holding relation, managing relation) which are publicly available.

Zheng *et al.* [3] extract 5 types of the 6 aforementioned events and provide corresponding dataset labeled with event details. The lawsuit events, which they have not covered are often accompanied by risks. In order to better apply the model to the real-world scenario, we use the distant supervision (DS) method based on a knowledge graph to label the announcement data and construct a lawsuit type event dataset.

The main contributions in this paper are two-folded:

1. We propose a framework that uses prior knowledge from basic knowledge graph to assist automated information extraction that outperforms state-of-the-art baselines;
2. We construct and release a new lawsuit event dataset for event extraction¹.

The paper is organized as follows. Section II shows the related work in event extraction, knowledge graph and graph neural network. Section III shows our definition of the problem, including the definition of the knowledge graph embedding, document-level event extraction task, and the description

*Haipeng Zhang is the corresponding author.

¹<http://bitly.ws/auMb>

of the predefined event roles and types. Section IV shows the main algorithms in our framework, including the graph embedding method and the encoder, event extraction and how to use distant supervision to label the data. Section V describes our data, the experiment setup, and our results compared with a state-of-the-art method. Finally, we conclude our work in Section VI.

II. RELATED WORK

Our work is inspired by three threads of research: event extraction, knowledge graph and graph neural networks.

A. Event Extraction

There has been a lot of studies on event extraction (EE), but the previous research focuses on sentence-level extraction (SEE) with most experiments performed on the ACE2005 dataset², which contains labeled events in English, Arabic and Chinese. As a first systematic SEE framework, Ahn [5] divides EE tasks into multiple sub-tasks, and proved that a series simple machine learning approaches can achieve the performance of best systems. Liao and Grishman [6] use document-level information for EE tasks, but their method still cannot handle the scenario where the event arguments are scattered in different sentences. Hong *et al.* [7] propose a method using cross-entity inference, but the authors only consider the entity types information and ignore the other information of entities. With the development of deep learning, Chen *et al.* [8] propose a multi-pooled convolution neural network EE method, which become a standard method for SEE tasks. However, these aforementioned methods ignore the external information that may be useful in the tasks.

Document-level event extraction (DEE) is the method of extracting elements from the full text and combining the elements into events. Compared with SEE, DEE is designed to solve the problem of argument dispersal. At present, there is little research on DEE. One of the major difficulties is associated with the lack of enough document-level labeled datasets. The DCFEE model [9] uses a distant supervised method to label data, and utilizes the sequence labeling model combined with the key sentence extraction model to extract the events. The latest Doc2EDAG model [3] performs named entity recognition (NER) and encodes the entity and the sentence which contains the entity. An entity's path generation method, which is a Transformer, is then used to extract events. Though these DEE methods mitigate the argument dispersal problem, the external knowledge which indicates the relations between entities is somehow also overlooked.

B. Knowledge Graph Enhanced Information Extraction

Knowledge graph has been a powerful tool for representing the knowledge by entities and relations between them, which is usually stored in the form of triples [10]. A knowledge graph with rich information can enhance the performance of the task in event extraction. Ruan *et al.* [11] incorporate cross-genre knowledge to improve the information extractor performance

for social media by 15% in F-score. Li and Ji [12] propose a distance supervised relation extraction method, where a knowledge graph with various semantic correlations among entities is embedded to help extracting information from long-tailed and imbalanced data. Zhang *et al.* [13] develop a system of knowledge-driven information extraction in Russian with a set of linguistic resources introduced. These works have proved the effectiveness of external knowledge in helping information extraction, but their potential improvements in event extraction, which is a quite different task, are barely touched.

C. Graph Convolution Network

Graph Convolution Network (GCN) is a common method to embed the knowledge stored in graph structure. The recent advance in the GCN can be categorized as spectral and non-spectral approaches.

Spectral approaches use a spectral transformation to deal with graph structure data. Bruna *et al.* [14] propose a construction based on the spectrum of the graph Laplacian [15] to learn convolution layers on low-dimension graphs. Kipf *et al.* [16], design a model which only operates the filter on the target node and its neighbor nodes to learn a hidden layer representation that encodes the local graph structure and node features. The disadvantage of spectral approaches is that they depend on the structure of the graph. As a result, a model trained on a specific structure cannot be directly applied to graphs with different structures.

Non-spectral approaches of GCN perform convolution by only considering close neighbor nodes of the target node. For instance, Duvenaud *et al.* [17] use a series of weight matrices to allow end-to-end learning on graphs, while Hamilton *et al.* [18] propose an inductive framework that learns a function to generate embeddings by aggregating features. These non-spectral approaches lack flexibility when the numbers of neighbors vary and a general operator is desirable. In order to solve this problem, Veličković *et al.* [19], implement the graph attention networks (GATs) model, which introduces an attention mechanism to extract the features of all neighbor nodes. The GATs model only depends on the edge instead of the complete graph structure and can be flexibly applied on directed and undirected graph.

III. PROBLEM DEFINITION

Our task is to extract structured event information from the financial announcement text in Chinese by a knowledge enhancement framework. In event extraction (EE) tasks, events are defined as structured data event types, event arguments, and event roles. Event role is a field which is predefined for each event type. In addition, we define the key roles for each event type which are indispensable in the corresponding event type. Table II shows an equity pledge event example. Event argument is the entity which plays an event role in a specific event. We obtain a dataset consisting of six different event types. For each type, there is a list of predefined event roles as shown in Table I. When dealing with an announcement,

²<https://www ldc.upenn.edu/collaborations/past-projects/ace>

TABLE I
EVENT TYPES AND THEIR CORRESPONDING EVENT ROLES.

Event Type	Event Roles
Equity Freeze	Equity Holder, Froze Shares, Legal Institution, Total Holding Shares, Total Holding Ratio, Date, Unfroze Date
Equity Repurchase	Company Name, Highest Trading Price, Lowest Trading Price, Repurchased Shares, Closing Date, Repurchase Amount
Equity Overweight	Equity Holder, Trading Shares, Date, Later Holding Shares, Average Price
Equity Underweight	Equity Holder, Trading Shares, Date, Later Holding Shares, Average Price
Equity Pledge	Pledger, Pledged Shares, Pledgee, Total Holding Shares, Total Holding Ratio, Total Pledged Shares, Date
Lawsuit	Plaintiff, Defendant, Legal Institution, Date

TABLE II
AN EQUITY PLEDGE EVENT EXAMPLE

Event Roles	Event Argument
Pledger*	Zhu Mingqiu
PledgedShares*	4,000,000 shares
Pledgee*	Huatai Securities Asset Management Co., Ltd.
TotalHoldingShares	131,573,433 shares
TotalHoldingRatio	8.47%
TotalPledgedShares	55,265,780 shares
StartDate	2017-07-19
EndDate	Null
ReleasedDate	Null

we determine the event type of the document and then fill the extracted entities into the event role list of the corresponding event type. To be specific, we use the dataset from the Doc2EDAG paper [3] and we further complement it with the lawsuit event data which is obtained by distant supervision labelling to be described in Section IV.

To evaluate the performance of the event extraction algorithms, we use the micro precision, recall and F1-score as evaluation metrics. To be specific, we compare the extracted event with the ground-truth event which shares the largest overlap in terms of number of event arguments. For each event role, we calculate the true positive rate, false positive rate and false negative rate. Besides, we further evaluate the model on multi-event documents and single-event documents, respectively.

IV. METHODOLOGY

The workflow of the framework is described in this section. We first show the process of employing the distant supervision (DS) method to extract lawsuit events from the announcements. We then describe two different methods to learn the embeddings of entities from the main knowledge graph: one-hot encoder and graph attention network. Finally, we demonstrate the event extraction model and the way of using the previously learned graph embeddings in the model.

Recall that one aim of the framework is to incorporate external knowledge to enhance the performance of the event extractor. Our main knowledge graph stores the external knowledge including key information of China’s listed A-share companies.

As illustrated in Figure 1, the framework consists of four major steps:

(1) We embed the entities in the main knowledge graph by their relation and then get a vector that focuses on describing the relations between entities.

(2) We perform named entity recognition (NER) to get the list of candidate entities.

(3) We correspond the extracted entities to the entities on the knowledge graph to obtain the embeddings of these entities. Then we use an encoder to process the embedding of these entities.

(4) We embed the entities and sentences in the article, and use a series of path-expanding tasks to extract events. In path expanding tasks, we use the embeddings obtained from the knowledge graph.

A. Unstructured text labeling based on knowledge graph

The task of document-level event extraction requires labelled data. To increase the various of the dataset, we use a knowledge graph based Distance Supervision (DS) method to label new type of data [20].

It is worth mentioning that in this part, we use a new DS knowledge graph constructed from data matched by the template. The relations in the DS knowledge graph and the main knowledge graph used in the event extraction are exclusive.

Here we describe the three steps to labeling dataset based on the DS knowledge graph. (1) The first step is template matching, which is to ensure high accuracy on the event extraction tasks. We create some templates for events of the selected type. Through these templates, we can extract events from the structured text in the corresponding announcement. (2) The second step is to construct a DS knowledge graph. The extracted information in the first step is transformed into triples to describe the entities, relations, and attributes of events. Disambiguation is a main problem in constructing the DS knowledge graph. That is, a same entity can have different names in different announcements or even in a same one. We employ a series of rules, such as removing location names (e.g. Shanghai) and suffixes of companies or courts (e.g. Co., Ltd.) to create a new expression of the name. We also crawl a list of company abbreviations to resolve the entity disambiguation problem. (3) In the third step, we use a DS method based on the DS knowledge graph to label new type of data. For each event type, the parameters of the event is defined as:

$$\text{Event}_p = \{a_{p1} \dots a_{pi}, b_{p1} \dots, b_{pj}\} \quad (1)$$

where a and b represent key and non-key roles of the event type respectively, p is the event type, and i and j ’s are indices

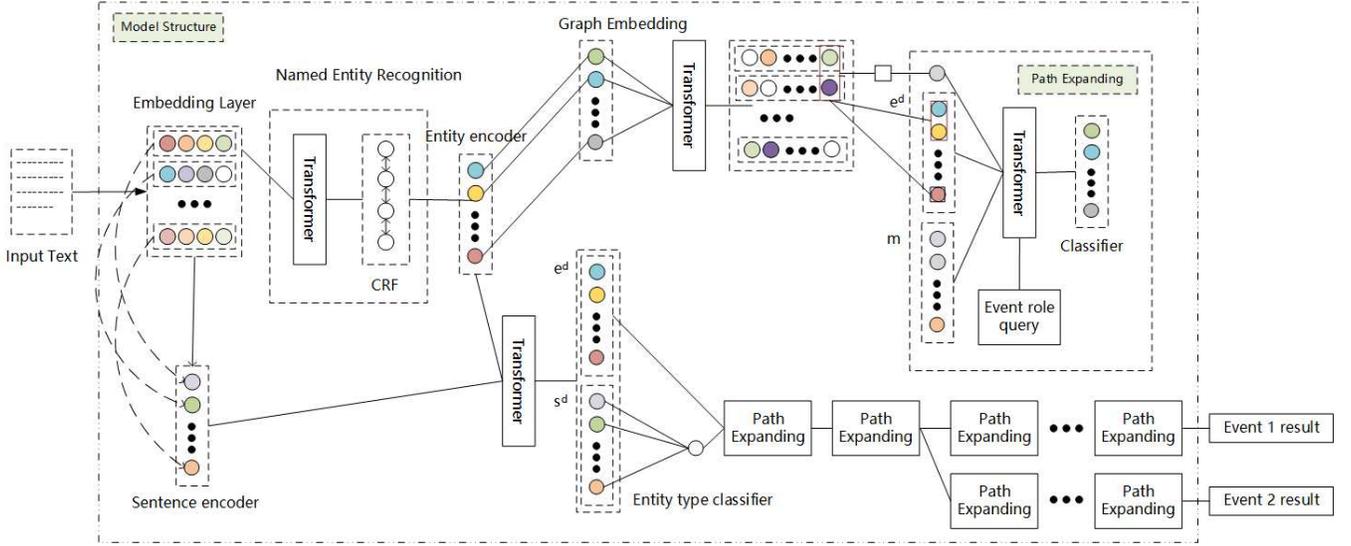


Fig. 1. The Transformer encoder approach of the event extraction model. The input of our model is a piece of text, and the output is the event results. In this example, two branches are generated on the second path-expanding task, and finally two different events are extracted.

of event arguments. We set two constraints in this step. First, we check whether all the predefined key roles of the event appear in the text. Second, we only keep the results whose ratios between extracted arguments and the amount of all possible event roles are above a certain threshold.

B. Knowledge graph embedding

We use two different approaches to embed the information of the main knowledge graph. The first is one-hot encoder based on the connecting paths between two entity nodes and the second is a graph neural network.

1) *One-hot Encoder*: For each pair of two target nodes, we select the first k shortest paths $\{P_1, P_2, \dots, P_k\}$. $P_i = \{t = (n_a) - [r_{ab}] - (n_b) | t \in Graph\}$ is a set of triples of nodes and relations, where n_a represents the node a and r_{ab} is the relation between node a and node b . If the number of paths between the two target nodes is less than k , we add empty paths to the set until the number meets the requirement.

For each path P_i , we extract all the relations r_{iab} and use one-hot encoder to encode them as vectors \vec{r}_{iab} . We then define the graph embedding between the two target nodes as:

$$\vec{h}_{ij} = \left[\sum \vec{r}_1, \sum \vec{r}_2, \dots, \sum \vec{r}_k \right] \quad (2)$$

For each empty path, we use an all -1 vector to represent the summation result of its relation vector.

2) *Graph Neural Network*: In this approach, we use the Graph Neural Network based on self-attention to calculate the embedding of each entity node in the main knowledge graph. We use the masked attention mechanism which designed by Veličković *et al.* in GAT [19] to solve the lack of flexibility for the non-spectral graph convolution networks approaches. Figure 2 shows the model structure of the graph self-attention model.

We randomly initialize $\mathbf{h} = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_N\}$, $\vec{h}_i \in R^F$ as a set of node feature vectors, where N is the number of nodes,

and F is a constant. We use one-hot encoder for relations. The relation between node i and node j are defined as r_{ij} , $r_{ij} \in R^S$, and S is the number of relation categories.

We use the graph self-attention model based on link prediction to train the embedding of nodes. Our input is the embeddings \vec{h}_i of the target node i and a set of tuples where each tuple represents one of i 's neighbor nodes and the its relation to i denoted as $\{(\vec{h}_j, \vec{r}_{ij}) | j \in N(i)\}$, where $N(i)$ represents i 's neighbor nodes. The relation embedding layer is used to encode the embedding of the neighbor node itself and the one-hot vector of the relation into one feature vector. It works as:

$$\vec{h}'_j = W_a \left[\vec{h}_j, \vec{r}_{ij} \right] \quad (3)$$

where $W_a \in R^{F \times (F+S)}$. Another linear transformation layer, which is parameterized by a weight matrix $W_b \in R^{F \times 2F}$, is then used to combine the embedding of i with the previous layer outputs.

$$\vec{h}''_j = \sigma \left(W_b \left[\vec{h}_i, \vec{h}'_j \right] \right) \quad (4)$$

We then use a shared attention mechanism to compute the attention coefficients e_{ij} for each neighbor node j to measure the importance of node j f node i :

$$e_{ij} = \vec{a} \vec{h}''_j \quad (5)$$

where $\vec{a} \in R^{1 \times F}$ is a single-layer feed-forward neural network. We use the softmax function to normalize the coefficients as:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in N(i)} \exp(e_{ik})} \quad (6)$$

We compute the linear combination of the features of each neighbor node by the attention coefficients as the result of

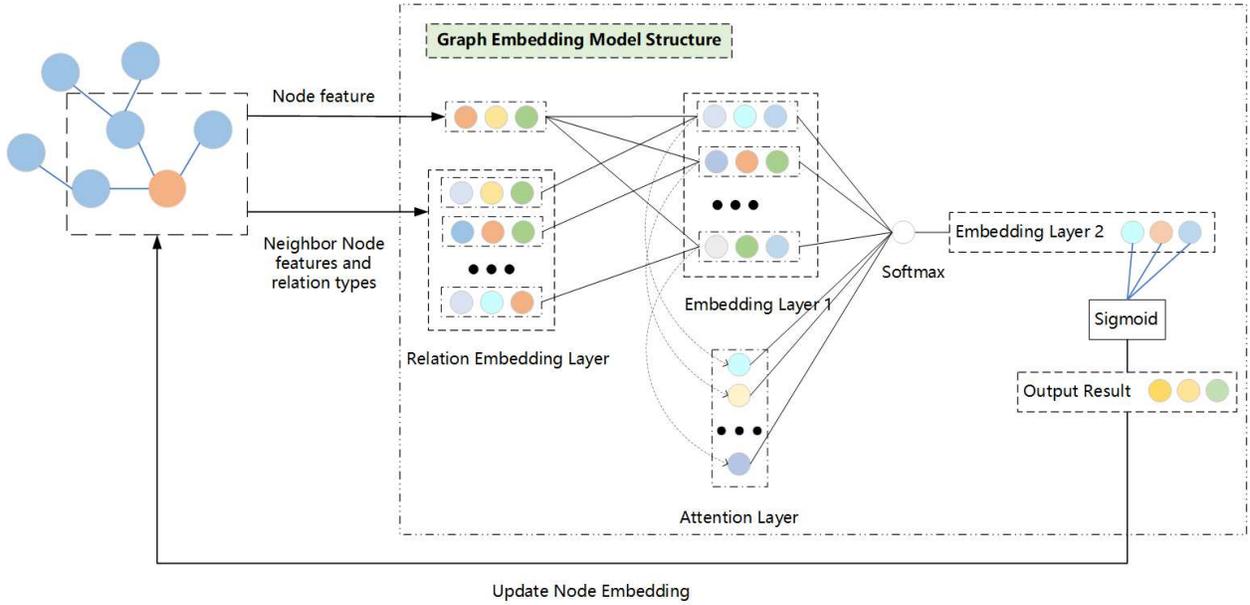


Fig. 2. The Graph embedding model structure. The input of the model is the embedding of a node, the embeddings of its neighbor nodes and the one-hot vectors representing the relation types. After completing the workflow, the model will update the new node embedding to the original graph.

embedding. To avoid the problem of covariate shift, we apply a layer normalization function.

$$\vec{h}'_i = \sigma \left[W_c \text{LayerNorm} \left(\sum_{j \in N(i)} \alpha_{ij} [\vec{h}_i, \vec{h}'_j] \right) \right] \quad (7)$$

where $W_c \in R^{F \times 2F}$.

To train the model, we use a simple linear classifier to predict the relation between two nodes. To mitigate the problem of imbalanced data, we use the focal loss [21]:

$$FL(p_t) = -\alpha_t (1 - p_t)^\lambda \log(p_t) \quad (8)$$

where p_t is the predicted probability of the model for category t , λ is the modulating factor and α_t is the weighting factor of category t . We calculate the factor α_t by the equation:

$$\alpha_t = \frac{\exp(n_{sum}/n_t)}{\sum_{s \in c} \exp(n_{sum}/n_s)} \quad (9)$$

where n_t is the number of occurrences for relation t in the graph, and $n_{sum} = \sum n_t$ represents the total number of relations in the graph.

C. DEE model based on knowledge graph

In this section, we describe the event extraction model. The whole model is divided into three parts: named entity recognition, entity and sentence information embedding, and event path construction.

1) *Named Entity Recognition*: We use a Transformer+CRF [4] model in the NER part. We input the document into the model in the form of a sequence of sentences $[s_1; s_2; \dots; s_{N_s}]$. The sentence s_i is represented by a sequence of token embeddings $[w_{i,1}, w_{i,2}, \dots, w_{i,N_w}]$, where $w_{i,j} \in R^{d_w}$ corresponds to the j th token in i th sentence. The

hyperparameters N_s and N_w are the maximum lengths of the sentence sequence and the token sequence respectively, and d_w is the embedding size of tokens.

After that, the Transformer layer embeds the input into a vector $h_i \in R^{d_w \times N_w}$. Finally, we use a conditional random field (CRF) layer to get the output, and apply Viterbi decoding on the output to get the NER result.

2) *Entity Embedding*: Because the information of the document is sparse, we encode entities at document-level to obtain as much information as possible. We encode different entities from three dimensions: the label of entities, sentence, and document.

We input entities and sentences into a max-pooling layer for dimensionality reduction. In this stage, a sentence or an entity is regarded as a list of tokens. We apply the max-pooling method to each entity e_l and sentence s_l to get the new embeddings $e'_l, s'_l \in R^{d_w}$.

We also encode the label type for each entity. The output of our NER contains different label types such as date, company name, and company id. The encoding process is

$$e'_i = \text{LayerNorm}(e_i + W_l [l_i]) \quad (10)$$

where $l_i \in R^{N_l}$ represents the one-hot vector of the entity's label type, N_l is the total number of label type categories and $W_l \in R^{d_w \times N_l}$ is a weighted matrix.

To obtain document-level embeddings, we use a Transformer layer to encode entities and sentences. After Transformer encoding, we use the max-pooling operation to merge the entities with the same entity name into a single embedding. Finally, we get a list of entities $e^d = [e_1^d, \dots, e_{N_e}^d]$ and a list of sentences $s^d = [s_1^d, \dots, s_{N_s}^d]$, where N_e is the number of distinct entities.

TABLE III
MAIN KNOWLEDGE GRAPH DATA DISTRIBUTION

relation	Quantity	Ratio
Branch	98,428	22.68%
Creditor	1,014	0.23%
Share Holder	138,853	32.00%
Invest	4,873	1.12%
Legal Person	95,745	22.07%
Pledge	61,242	14.11%
Managing member	33,735	7.78%

3) *Event Path construction*: We transform the event extraction problem into a problem of constructing an entity-based directed acyclic graph (EDAG) which is similar to [3]. Compared with the simple table filling method [9], the method of constructing a graph can be divided into some path-expanding task, which is easier to handle. At the same time, the method of constructing a graph is more suitable for incorporating graph embedding information. For each event type, we define a field list to describe the event roles, and connect the fields in order to form a directed acyclic graph. This method allows us to extract the target event by path-expanding on each node in the graph.

In this part, we first use max-pooling to encode the sentences of the document to get a vector $t \in R^{d_w}$ and use a linear classifier to get the event type of this document.

Next, we construct an initialized EDAG node according to the event type. Then we expand the path on this node based on the predefined sequence of event role. In order to remember the entities in the previous path, we construct a memory vector m . After expanding the current path, we append the current entity to the memory vector. If no entity is selected, we use a zero-padding vector to update m .

As shown in Figure 1, we add the knowledge graph embedding information in the path-expanding task. We use a list e^p to record the graph embedding of all entities expanded in the previous path-expanding task. If there is no such entity in our main knowledge graph, a zero-padding vector is used to represent it. For each candidate entity, we input its graph embedding and the list e^p to an encoder to obtain a vector $g \in R^F$. We use the Transformer model with max-pooling as the encoder. We also implement a single linear layer with max-pooling approach for comparisons.

For each path-expanding task, we concatenate all the entities e^d , memory vector m and the graph embedding vector v , and feed them to a Transformer model. Finally, we use a linear classifier to determine which entity will be added as a node under the current path. If multiple entities are added to the graph in the same path-expanding task, the graph generates multiple branches at the current node.

V. EXPERIMENTS

In this section, we first describe the data used in our framework, including main knowledge graph data and the labeled financial announcement dataset. We then compare the

TABLE IV
LABELED DATA QUALITY

Method	Precision	Recall	F1-score
Template Matching	76.36%	64.70%	70.04%
Distant Supervision	92.86%	85.71%	89.14%

TABLE V
DEE DATASET STATISTICS AND SPLIT.

Event Type	Train	Val.	Test	Total
Equity Freeze	407	183	106	696
Equity Repurchase	2,331	547	534	3,412
Equity Overweight	3,292	240	350	3,882
Equity Underweight	3,278	297	217	3,774
Equity Pledge	7,657	908	913	9,478
Lawsuit	498	64	59	621
Total	17,463	2,179	2,179	21,871

experiment results from our models with these from a state-of-the-art model and further divide the test data into single event data and multi-event data to examine the differences in performance. We also carefully discuss the advantages and disadvantages of our framework.

A. Knowledge Graph

We use business information of 3,000 listed companies from China’s National Enterprise Credit Information Publicity System³, and construct a directed knowledge graph with 347,570 entities (companies or persons) and 441,127 relations (7 in total, such as Share Holder and Branch, mainly formed by commercial activities). In our framework, we focus on the relations among entities. The details about the relations are listed in Table III.

B. Data Labelling

In addition to the dataset in [3], we use the distant supervision approach described in Section IV-A to obtain 2,528 labelled lawsuit type financial announcements. To be specific, we start with a total of 14,052 lawsuit type financial announcements from Shanghai Stock Exchange and Shenzhen Stock Exchange in over 10 years (2010-2020). After applying the template matching method, 2,937 announcements contain enough event information to be used in our experiments and they include 3,197 events in total. From these announcements, we construct a basic lawsuit knowledge graph with 9,839 nodes and 12,887 relations that serves as the knowledge base for the DS method mentioned in Section IV-A, which yields the aforementioned 2,528 announcements.

Table IV shows the results of template matching and the DS labelling experiment. We randomly select 100 documents and manually check the documents to test the quality of our dataset. It suggests that the DS method improves template matching and achieves a satisfying F1-score of almost 90%. It is worth mentioning that our released dataset based on this has been further manually checked and corrected to facilitate relevant research.

³<http://www.gsxt.gov.cn/>

TABLE VI
EVENT EXTRACTION RESULTS (%).

Model	EF			ER			EO			EU			EP			LA			Total		
	P.	R.	F1	P.	R.	F1															
Doc2EDAG	78.4	70.2	74.1	81.3	80.9	81.6	84.0	73.2	78.2	77.3	64.8	70.5	77.5	73.0	75.2	65.9	59.3	62.4	79.0	73.9	76.4
One-hot+Transformer	82.2	76.2	79.1	86.4	87.4	86.9	84.5	70.5	76.9	81.7	76.7	79.1	83.9	76.3	79.9	59.0	66.3	62.4	79.6	75.6	77.4
GNN+Linear	82.5	70.2	75.8	83.4	79.8	81.6	84.0	71.8	77.4	71.6	64.1	67.7	74.4	76.7	75.5	55.4	56.6	56.0	77.2	75.2	76.2
GNN+Transformer	85.9	82.3	84.1	86.1	87.4	86.7	83.9	69.5	76.0	85.0	76.5	80.5	83.9	77.9	80.8	67.4	65.6	66.5	84.4	79.1	81.7

TABLE VII
SINGLE EVENT (S.) VS MULTI-EVENT (M.) RESULTS IN F1-SCORES (%)

Model	EF		ER		EO		EU		EP		LA		Total	
	S.	M.												
Doc2EDAG	74.9	64.0	87.7	62.7	80.1	74.4	76.5	60.9	83.5	69.0	67.6	58.3	78.4	64.9
GNN+Transformer	88.1	81.9	92.1	74.7	85.3	76.3	79.8	70.0	87.8	76.4	69.6	55.0	87.8	76.2

C. Dataset Construction

We keep the announcements from the Doc2EDAG dataset [3] and the lawsuit dataset (Section V-B) which include at least one entity that appears in the main knowledge graph such that the knowledge graph can potentially help. As a result, we get a new dataset with 21,871 announcements, covering 6 event types: equity freeze (EF), equity repurchase (ER), equity overweight (EO), equity underweight (EU), equity pledge (EP) and lawsuit (LA). We split the data into training, validation, and test sets following a ratio of 8:1:1 and use this dataset for all of our DEE experiments and evaluations. Details of this dataset can be found in Table V.

D. DEE Experiment

In this part, we compare our models with the state-of-the-art Doc2EDAG model. Our models include the main GNN+Transformer model as well as the One-hot+Transformer and GNN+Linear models, all elaborated in Section IV. GNN embedding and one-hot encoder are compared as two different methods of incorporating priori information from external knowledge graphs. With regard to embedding merging, Single-Linear layer and Transformer are evaluated.

1) *Main result:* As we can see in Table VI, the overall performance of our models that incorporate knowledge graph information beats Doc2EDAG, with the GNN+Linear model as the only exception that has a slightly lower F1-score and a higher recall. This indicates that proper integration of relational knowledge can improve the performance of the event extraction task.

It is worth noting that GNN+Transformer, as the best-performing model with an overall F1-score of 81.7%, surpasses Doc2EDAG by a relatively large margin (5.4% for precision, 5.2% for recall, and 5.3% for F1-score). In particular, for the EF, ER, EU, and EP categories, the GNN+Transformer model demonstrates significant improvements by 10.0%, 6.1%, 10.0%, and 5.6%, respectively. This may be related with the characteristics of the events in these categories – entities within a document are often related intrinsically. For instance, in the case of equity pledge, the pledgee is more likely to be a bank that already has pledge relations with many companies,

which is recorded in our relational knowledge graph. Another example is equity overweight which happens when a listed company’s major stake-holders plan to purchase more of its shares and these stake-holders are usually captured in the knowledge graph.

When compared with the other two graph based models, the GNN+Transformer model also shows improvements in F1-scores. Its 4.3% improvement over One-hot+Transformer implies that GNN is more effective in capturing graph information, as apposed to the intuitive feature engineering over the graph, while the 5.5% improvement over GNN+Linear suggests that the Transformer model with its self-attention mechanism better encodes the information, compared with the simple linear layer.

2) *Single Event vs Multi-Event:* We examine the performance of GNN+Transformer and Doc2EDAG on single event documents (each containing only one event) and multi-event documents (each containing more than one event), by further splitting the test data and calculating the corresponding F1-scores at document level (shown in Table VII). Overall, our model improves Doc2EDAG by 9.4% and 11.3% in the single event dataset and the multi-event dataset, respectively. For the multi-event documents, the improvement is higher. A possible reason is that these documents with multiple events usually contain multiple entities and entities across events may cause confusions for the Doc2EDAG model, while this might be mitigated by our model based on a relational knowledge base.

3) *Discussion:* All of the above experiments suggest that the combination of GNN and Transformer can more effectively embed knowledge graph information, comparing with the other two combinations. Different from the plain one-hot encoding of the path between two nodes, GNN can implicitly learn the relation between a node and all its neighbor nodes, as well as the position of the node in the graph. In addition, one-hot encoder cannot describe sequence information, but our GNN model over come this problem through self-attention. Transformer, as a seq2seq model, can encode the relations between two nodes well when there are multiple nodes in between and the paths are sequential. Besides, compared with other Recurrent Neural Networks, the Transformer model can

be easily paralleled and offers higher interpretability, making it more applicable in solving real-world problems [4].

VI. CONCLUSION AND FUTURE WORK

In this paper, we propose a knowledge graph based extraction framework. While a straightforward way of engineering the graph into features can effectively improve the state-of-the-art baseline in financial announcement extraction, our GNN+Transformer model that better integrates the graph information with limited human intervention archives a considerable 5.3% improvement. In order to make the model more applicable in the real scenario, we construct a lawsuit type event dataset and merge it with a public dataset. While all our models that incorporate the knowledge graph outperform the state-of-the-art baseline, the GNN+Transformer model that learns more precise relational knowledge achieves the best performance.

As a future direction, we plan to apply the framework on tasks including relation extraction and sentence-level event extraction where relational prior knowledge can also be helpful. Besides, we will process documents in other domains with our framework augmented with general knowledge graphs including DBpedia and Freebase.

VII. ACKNOWLEDGMENTS

We thank the anonymous reviewers for their constructive suggestions. This project is sponsored by Shanghai Sail Program (No.19YF1433800) and the Key Projects of Shanghai Soft Science Research Program (No.20692194300). This article has been accepted to IEEE Big Data 2020 Conference.

REFERENCES

- [1] S. Liu, Y. Chen, K. Liu, and J. Zhao, "Exploiting argument information to improve event detection via supervised attention mechanisms," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2017, pp. 1789–1798.
- [2] D. Li, L. Huang, H. Ji, and J. Han, "Biomedical event extraction based on knowledge-driven tree-1stm," in *NAACL-HLT 2019: Annual Conference of the North American Chapter of the Association for Computational Linguistics*, 2019, pp. 1421–1430.
- [3] S. Zheng, W. Cao, W. Xu, and J. Bian, "Doc2edag: An end-to-end document-level framework for chinese financial event extraction," in *2019 Conference on Empirical Methods in Natural Language Processing*, 2019, pp. 337–346.
- [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [5] D. Ahn, "The stages of event extraction," in *Proceedings of the Workshop on Annotating and Reasoning about Time and Events*, 2006, pp. 1–8.
- [6] S. Liao and R. Grishman, "Using document level cross-event inference to improve event extraction," in *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, 2010, pp. 789–797.
- [7] Y. Hong, J. Zhang, B. Ma, J. Yao, G. Zhou, and Q. Zhu, "Using cross-entity inference to improve event extraction," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 2011, pp. 1127–1136.
- [8] Y. Chen, L. Xu, K. Liu, D. Zeng, and J. Zhao, "Event extraction via dynamic multi-pooling convolutional neural networks," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, vol. 1, 2015, pp. 167–176.
- [9] H. Yang, Y. Chen, K. Liu, and J. Zhao, "Dcfee: A document-level chinese financial event extraction system based on automatically labeled training data," in *ACL 2018: 56th Annual Meeting of the Association for Computational Linguistics*, 2018, pp. 50–55.
- [10] C. Bizer, J. Lehmann, G. Kobilarov, S. Auer, C. Becker, R. Cyganiak, and S. Hellmann, "Dbpedia - a crystallization point for the web of data," *Journal of Web Semantics*, vol. 7, no. 3, pp. 154–165, 2009.
- [11] T. Ruan, L. Xue, H. Wang, F. Hu, L. Zhao, and J. Ding, "Building and exploring an enterprise knowledge graph for investment analysis," in *International Semantic Web Conference (2)*, 2016, pp. 418–436.
- [12] H. Li and H. Ji, "Cross-genre event extraction with knowledge enrichment," in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2016, pp. 1158–1162.
- [13] N. Zhang, S. Deng, Z. Sun, G. Wang, X. Chen, W. Zhang, and H. Chen, "Long-tail relation extraction via knowledge graph embeddings and graph convolution networks," in *NAACL-HLT 2019: Annual Conference of the North American Chapter of the Association for Computational Linguistics*, 2019, pp. 3016–3025.
- [14] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," in *ICLR 2014 : International Conference on Learning Representations (ICLR) 2014*, 2014.
- [15] M. Belkin and P. Niyogi, "Laplacian eigenmaps and spectral techniques for embedding and clustering," in *Advances in Neural Information Processing Systems 14*, 2001, pp. 585–591.
- [16] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *ICLR 2017 : International Conference on Learning Representations 2017*, 2017.
- [17] D. Duvenaud, D. Maclaurin, J. Aguilera-Iparraguirre, R. Gómez-Bombarelli, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams, "Convolutional networks on graphs for learning molecular fingerprints," in *NIPS'15 Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, 2015, pp. 2224–2232.
- [18] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in Neural Information Processing Systems*, 2017, pp. 1024–1034.
- [19] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *ICLR 2018 : International Conference on Learning Representations 2018*, 2018.
- [20] A. Smirnova, J. Audiffren, and P. Cudré-Mauroux, "Distant supervision from knowledge graphs," *Encyclopedia of Big Data Technologies*, pp. 1–7, 2018.
- [21] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 2, pp. 318–327, 2020.