

Evidence-aware Fake News Detection with Graph Neural Networks

Weizhi Xu^{1,2,*}, Junfei Wu^{3,*†‡}, Qiang Liu^{1,2}, Shu Wu^{1,2,†}, Liang Wang^{1,2*†‡}

¹Center for Research on Intelligent Perception and Computing (CRIPAC)

National Laboratory of Pattern Recognition (NLPR)

Institute of Automation, Chinese Academy of Sciences

²School of Artificial Intelligence, University of Chinese Academy of Sciences

³School of Computer Science and Technology, Beijing Institute of Technology

weizhi.xu@cripac.ia.ac.cn, junfei.wu@bit.edu.cn, {qiang.liu, shu.wu, wangliang}@nlpr.ia.ac.cn

ABSTRACT

The prevalence and perniciousness of fake news has been a critical issue on the Internet, which stimulates the development of automatic fake news detection in turn. In this paper, we focus on the evidence-based fake news detection, where several evidences are utilized to probe the veracity of news (i.e., a claim). Most previous methods first employ sequential models to embed the semantic information and then capture the claim-evidence interaction based on different attention mechanisms. Despite their effectiveness, they still suffer from two main weaknesses. Firstly, due to the inherent drawbacks of sequential models, they fail to integrate the relevant information that is scattered far apart in evidences for veracity checking. Secondly, they neglect much redundant information contained in evidences that may be useless or even harmful. To solve these problems, we propose a unified **G**raph-based **s**emantic **s**tructure mining framework, namely GET in short. Specifically, different from the existing work that treats claims and evidences as sequences, we model them as graph-structured data and capture the long-distance semantic dependency among dispersed relevant snippets via neighborhood propagation. After obtaining contextual semantic information, our model reduces information redundancy by performing graph structure learning. Finally, the fine-grained semantic representations are fed into the downstream claim-evidence interaction module for predictions. Comprehensive experiments have demonstrated the superiority of GET over the state-of-the-arts.

CCS CONCEPTS

- **Computing methodologies** → **Natural language processing**;
- **Information systems** → **Data mining**.

*The first two authors contributed equally to this work.

†To whom correspondence should be addressed.

‡Work done while Junfei Wu interned at CRIPAC, CASIA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WWW '22, April 25–29, 2022, Virtual Event, Lyon, France

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9096-5/22/04...\$15.00

<https://doi.org/10.1145/3485447.3512122>

KEYWORDS

evidence-based fake news detection, graph neural networks

ACM Reference Format:

Weizhi Xu, Junfei Wu, Qiang Liu, Shu Wu, Liang Wang. 2022. Evidence-aware Fake News Detection with Graph Neural Networks. In *Proceedings of the ACM Web Conference 2022 (WWW '22)*, April 25–29, 2022, Virtual Event, Lyon, France. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3485447.3512122>

1 INTRODUCTION

Fake news, which is always fabricated by making some minor changes to the correct statement, is highly deceptive and indistinguishable. The widespread of fake news in diverse domains, such as politics [2] and public health [27], has posed a huge threat to web security and human society. Therefore, the research on automatic fake news detection is challenging but in demand.

Generally, previous methods could be roughly categorized into two groups, i.e., pattern-based approaches and evidence-based approaches [32]. The former methods regard the fake news detection as a feature recognition task, where language models are employed to verify the veracity of news solely according to the text pattern, e.g., writing styles. However, pattern-based methods usually suffer from the poor generalization and interpretability. The latter approaches model the task as a reasoning process, where external evidences are provided to probe the veracity of a claim. Models are required to discover and integrate useful information in given evidences for claim verification.

In this paper, we concentrate on the evidence-based pipeline. Existing methods usually follow a two-step paradigm: 1) they first capture the semantics of claims and evidences separately. 2) Next, they model the claim-evidence interaction to explore the semantic coherence or conflict for more accurate and interpretable verdict. To name a few representative models, the pioneering work De-ClarE [30] utilizes bidirectional LSTMs to model textual features, followed by a word-level attention mechanism to capture the claim-evidence interaction. HAN [26] further considers the sentence-level interaction to explore more general semantic coherence. To obtain multi-level semantic interaction, some recent works [37, 41] employ hierarchical attention networks.

Nevertheless, existing work focuses on the specific design of different interaction models (the second step) while neglecting exploring fine-grained semantics of claims and evidences (the first step). To be specific, we argue that there are two main weaknesses in previous methods.

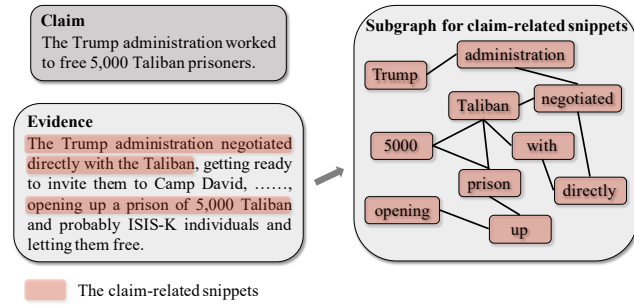


Figure 1: A toy example where a claim and its relevant evidence are given. Two significant snippets for verifying the claim are highlighted (“.....” represents that we omit several sentences for conciseness). The right graph is constructed according to the highlighted snippets. Such two snippets have a long distance in the plain text while they are pulled close on the constructed semantic graph via the shared keyword “Taliban”. Besides, there is much redundant information (texts except the highlighted parts), which is useless for claim verification.

Firstly, the complex, long-distance semantic dependency is less explored. Taking Figure 1 as an example, two highlighted snippets are separated by plenty of words, which induces a long distance between them. Such snippets both contain important information for verifying the claim, i.e., the subject “The Trump administration” and the action “opening up a prison of 5,000 Taliban”. Therefore, fusing the information is indispensable and beneficial for claim veracity prediction. However, the long-distance semantic dependency between such information is hard to be captured due to the inherent drawbacks of sequential models utilized in previous methods.

Secondly, existing methods neglect the redundant information involved in semantics. Such redundancy is useless or even harmful for fake news detection, e.g., as depicted in Figure 1, a large number of text segments, such as “getting ready to invite them to Camp David”, have no substantial contribution to the news veracity checking. Though previous models employ attention mechanisms to reduce the effect of unrelated words, these irrelevant texts are still preserved, which may introduce noises to the downstream claim-evidence interaction, deteriorating the final performance of veracity checking. An intuitive solution is to discard words with low attentive scores based on previous methods. However, they compute the score for each word independently, ignoring the complex semantic structure among words. We argue that it is significant to modeling the redundancy with rich semantic structural information, as the redundancy is not only related to the self-information, but also induced by its contexts, e.g., if a claim can be verified by a snippet in an evidence, the snippet’s context will be redundant.

To tackle the aforementioned problems, we propose a unified **Graph-based sEmantic sTructure mining framework**, namely GET for exploring fine-grained semantics. Specifically, modeling sequential data as graphs has benefited many tasks, such as text classification [47, 54] and sequential recommendation [44], owing to its capability of capturing long-distance structural dependency. To this

end, we utilize graph structure to model both claims and evidences, where nodes indicate words and edges represent the co-occurrence between two words. Thereafter, the dispersed claim-related snippets are pulled close on graphs, thus the useful information could be better fused via neighborhood propagation. For example, in Figure 1, after constructing the graph for two highlighted snippets distant from each other in plain texts, they are pulled close via the shared keyword “Taliban” so that the long-distance semantic dependency can be captured. Moreover, to alleviate the negative impact of redundant information, within our graph-based framework, we treat the redundancy mitigation as a graph structure learning process, where unimportant nodes are discarded according to complex semantic structures, i.e., both self-features (node attributes) and their contexts (graph topology). To date, our graph-based framework has captured the fine-grained semantics via long-distance dependency modeling and redundancy mitigation. Based on such semantics, we can apply the widely used attention mechanism in previous work to readout node features and form the claim- and evidence-level representations, followed by claim-evidence interactions to integrate information for the final veracity prediction.

Our main contributions can be summarized as follows:

- We model claims and evidences as graph-structured data and design a graph-based framework to explore the complex semantic structure. To the best of our knowledge, this is the first work to propose a unified graph-based method for evidence-based fake news detection.
- We introduce a simple and effective graph structure learning approach for redundancy mitigation. By capturing long-distance semantic dependency and reducing redundancy, we obtain the fine-grained semantics, which can boost the performance of downstream interaction models.
- Comprehensive experiments are conducted to verify the effectiveness of GET, where the results demonstrate its superiority.

2 RELATED WORK

In this section, we briefly review previous work in two related domains: graph neural networks and fake news detection.

2.1 Graph Neural Networks

Graph neural networks (GNNs) learn the node representation by gathering information from the neighborhood, i.e., neighborhood propagation/aggregation. Current GNNs can be roughly divided into two groups, namely spectral approaches [8, 19] and spatial approaches [14, 33]. Owing to the capability of capturing long-distance structural relationship on graphs, GNNs have been widely utilized and achieved satisfactory performance in several tasks, such as recommender system [5, 44, 51], text classification [47, 54], and sentiment analysis [22, 39]. Recently, researchers have observed that graphs inevitably contain noises that may deteriorate the training of GNNs [17]. To handle this problem, graph structure learning (GSL) is proposed, aiming to jointly learn an optimized graph structure and node embeddings. Existing GSL methods mainly fall into three groups [58]: 1) *the metric-learning-based methods* where the

adjacency matrices are built as metrics coupled with node embeddings. Therefore, the graph topology is updated with node embeddings being optimized. The metrics are mainly defined as the attention-based function [6, 7, 16] or kernel function [23, 45]. 2) *the probabilistic methods* assume that the adjacency matrix is generated by sampling from a specific probabilistic distribution [10, 11, 53]. 3) *the direct-optimized methods* treat the graph topology as learnable parameters that are updated together with task-specific parameters simultaneously, without depending on preset priors (node embeddings and distributions in the first two groups, respectively). The topology is optimized with the guidance of task-specific objectives (and some normalization constraints) [17, 46]. It is worth noting that existing graph pooling methods [12, 20, 48] could also be regarded as GSL algorithms, since the pooling target is to keep the most valuable nodes that preserve the graph structural information well, where the graph structure is optimized via merging or dropping nodes. Besides, GNNs are widely employed in the domain of fact verification, which have achieved promising performance [25, 56, 57]. Though fact verification is similar to fake news detection on the task setting, the latter requires more fine-grained semantics since the texts consist of more redundancy.

2.2 Fake News Detection

Several fake news detection methods have been proposed in recent years, which can be roughly grouped into two categories.

The first is the pattern-based pipeline where models solely consider the text pattern involved in the news itself. Different work always focus on different kinds of patterns. Popat et al. [28] classify a claim as true or fake in accordance with stylistic features and the article stance. Besides, some researchers attempt to verify the truthiness via the feedback in social media, such as reposts, likes, and comments [3, 4, 18, 24, 36, 38, 49]. Recently, more attention has been paid to the emotional pattern mining, where they hold an assumption that there are probably obvious sentiment biases in fake news [1, 13, 52].

The second is the evidence-based pipeline where researchers propose to explore the semantic similarity (conflict) in claim-evidence pairs to check the news veracity. Evidences are usually retrieved from the knowledge graph [35] or fact-checking websites [34] by giving unverified claims as queries. DeClarE [30] is the first work to utilize evidences in fake news detection. They employ BiLSTMs to embed the semantics of evidences and obtain the claim's sentence-level representation via average pooling. Next, they introduce an attention-based interaction to compute the claim-aware score for each word in evidences. Similar to the pioneering work, the following methods utilize the sequential models to obtain the semantic embeddings, followed by attention mechanisms performed on different granularities. HAN [26] compute the sentence-level coherence and entailment scores between claims and evidences. EHIAN [43] employs the self-attention mechanism to obtain word-level interaction scores. Recent work [37, 41, 42] hierarchically integrates both word-level and sentence-level interactions into the final representation for verification. In summary, they all employ sequential models to embed the semantics and apply attention mechanisms to capture the claim-evidence semantic relationship.

Different from the existing work, we propose a unified graph-based model, where the long-distance semantic dependency is captured via constructed graph structures and the redundancy is reduced by performing graph structure learning.

3 METHOD

3.1 Task Formulation

Evidence-based fake news detection is a classification task, where the model is required to output the prediction of news veracity. Specifically, the inputs are a claim c , several related evidences $\mathcal{E} = \{e_1, e_2, \dots, e_n\}$, and their corresponding speakers $\mathbf{s} \in \mathbb{R}^{1 \times b}$ or publishers $\mathbf{p} \in \mathbb{R}^{n \times b}$, where n is the number of evidences and b is the dimension of speaker and publisher embeddings. The output is the predicted probability of veracity $\hat{y} = f(c, \mathcal{E}, \mathbf{s}, \mathbf{p}, \Theta)$, where f is the verification model and Θ is its trainable parameters.

3.2 The Proposed Model: GET

In this part, we elaborate our unified graph-based model GET, which can be mainly separated into four modules: 1) *Graph Construction*, 2) *Graph-based Semantics Encoder*, 3) *Semantic Structure Refinement*, and 4) *Attentive Graph Readout Layer*.

3.2.1 Graph Construction. In order to capture the long-distance dependency of relevant information, we first convert the original claims and evidences to graphs. Like previous graph-based methods in other NLP tasks [47, 50, 54, 55], we use a fix-sized sliding window to screen out the connectivity for each word on graphs. In detail, the center words in every window will be connected with the rest of words in it (if connected, the corresponding entry in adjacency matrix is 1, otherwise 0), which captures the local context in center word's neighborhood. Furthermore, to model the long-distance dependency, we merge all the same words into one node on graph, which explicitly gathers their local contexts (e.g., the word e_2 in evidence text 1 in Figure 2). Therefore, several relevant snippets that scatter far apart is close on graphs, which can be explored via the high-order message propagation. In addition, the initial node representations are the corresponding word embeddings. Note that we also try to construct a graph in a fully-connected or semantic-similarity-based manner, but these two ways are all inferior to the sliding-window-based method, which may due to the redundant noises induced by the dense connection.

To ensure the numerical stability, we perform Laplacian normalization on adjacency matrices, denoted as $\tilde{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})\mathbf{D}^{-\frac{1}{2}}$, where \mathbf{D} is the diagonal degree matrix (i.e., $\mathbf{D}_{ii} = \sum_j \mathbf{A}_{ij}$) and \mathbf{I} is the identical matrix. Finally, we denote the initial normalized adjacency matrices and node feature matrices of claim and evidence as $\tilde{\mathbf{A}}_c^{(0)} \in \mathbb{R}^{N_c \times N_c}$, $\tilde{\mathbf{A}}_e^{(0)} \in \mathbb{R}^{N_e \times N_e}$ and $\mathbf{H}_c^{(0)} \in \mathbb{R}^{N_c \times d}$, $\mathbf{H}_e^{(0)} \in \mathbb{R}^{N_e \times d}$, respectively. N_c and N_e is the number of nodes in initial claim and evidence graphs, d is the dimension of word embeddings.

Taking the established graph structures and node embeddings as inputs, we design a graph-based model to better capture complex semantics and obtain refined semantic structures.

3.2.2 Graph-based Semantics Encoder. To mine the long-distance semantic dependency, we propose to utilize GNNs as the semantics encoder. In particular, as we expect to adaptively keep a balance

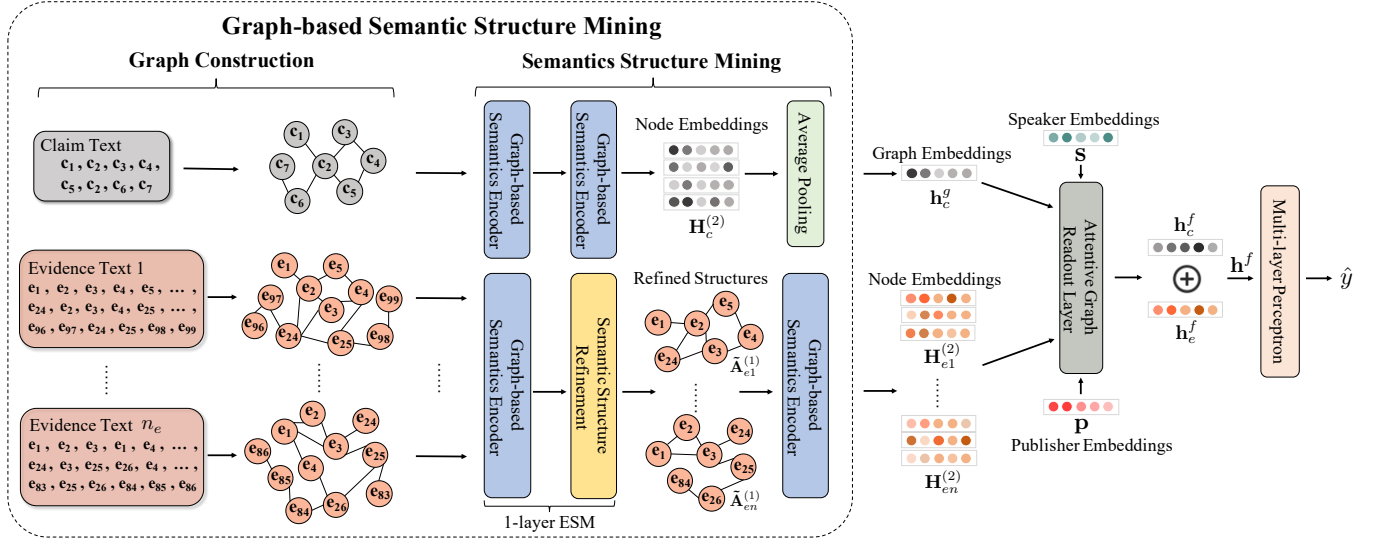


Figure 2: The architecture of GET. The plain texts are first transformed into graphs using a sliding window (the window size is 2 in the figure). The same words repeatedly appear in texts are merged into one node. Next, we introduce graph-based semantics encoder to capture long-distance structural dependencies and generate high-order representations via neighborhood aggregation. Furthermore, the semantic structure refinement layer is proposed to generate optimized structures $\{\tilde{A}_{e1}^{(1)}, \dots, \tilde{A}_{en}^{(1)}\}$ for n evidences, where redundant nodes are discarded (The 1-layer ESM consists of a graph-based semantics encoder and a semantic structure refinement layer). Thereafter, the fine-grained semantics is obtained by performing neighborhood propagation on refined graphs. Finally, claim and evidence embeddings along with their speaker and publisher information are fed into the attentive graph readout layer to output the final prediction \hat{y} .

between self-features and the information of neighboring nodes, we employ graph gated neural networks (GGNN) to perform neighborhood propagation on both claim and evidence graphs, enabling nodes to capture their contextual information, which is significant for learning high-level semantics. Formally, it can be written as follows:

$$\mathbf{a}_i = \sum_{(w_i, w_j) \in C} \tilde{A}_{ij} \mathbf{W}_a \mathbf{H}_j \quad (1)$$

$$\mathbf{z}_i = \sigma(\mathbf{W}_z \mathbf{a}_i + \mathbf{U}_z \mathbf{H}_i + \mathbf{b}_z) \quad (2)$$

$$\mathbf{r}_i = \sigma(\mathbf{W}_r \mathbf{a}_i + \mathbf{U}_r \mathbf{H}_i + \mathbf{b}_r) \quad (3)$$

$$\tilde{\mathbf{H}}_i = \tanh(\mathbf{W}_h \mathbf{a}_i + \mathbf{U}_h (\mathbf{r}_i \odot \mathbf{H}_i) + \mathbf{b}_h) \quad (4)$$

$$\hat{\mathbf{H}}_i = \tilde{\mathbf{H}}_i \odot \mathbf{z}_i + \mathbf{H}_i \odot (1 - \mathbf{z}_i) \quad (5)$$

where C denotes the edge set, \mathbf{W}_* , \mathbf{U}_* , and \mathbf{b}_* are trainable parameters, which control the proportion of the neighborhood information and self-information. σ is the non-linear activation unit and we utilize the Sigmoid function in our model. For brevity, we denote Eq. (1) - (5) as $\text{GGNN}(\tilde{\mathbf{A}}, \mathbf{H})^1$.

3.2.3 Semantic Structure Refinement. As evidences always contain redundant information that may mislead model to focus on unimportant features, it is beneficial to discover and filter out the redundancy, thus obtaining refined semantic structures. To this end,

in our graph-based framework, we treat the redundancy mitigation as a graph structure learning process, whose aim is to learn the optimized graph topology along with better node representations. Previous GSL methods generally optimize the topology in three ways, i.e., dropping nodes, dropping edges, and adjusting edge weights. Since the redundancy information is mainly involved in words denoted as nodes in evidence graphs, we attempt to refine evidence graph structures via discarding redundant nodes, inspired by previous GSL methods [6, 20, 53].

In particular, we propose to compute a redundancy score for each node, based on which we obtain a ranking list and nodes with the top- k redundancy scores will be discarded. The redundancy is not only related to the self-information contained in each node, but also induced by the contextual information, which is involved in the neighborhood on graphs. For example, if a claim can be verified by a snippet in an evidence, the rest of segments (including the snippet's context) will be redundant. Therefore, we utilize a 1-layer GGNN to compute the redundancy scores, which takes into account both self- and context-information in score computation. Mathematically, it can be formulated as:

$$\mathbf{s}_r = \text{GGNN}(\tilde{\mathbf{A}}, \hat{\mathbf{H}}_e \mathbf{W}_s) \quad (6)$$

$$\text{idx} = \text{topk_index}(\mathbf{s}_r) \quad (7)$$

$$\tilde{\mathbf{A}}_{\text{idx},:} = \tilde{\mathbf{A}}_{:, \text{idx}} = 0 \quad (8)$$

where $\mathbf{W}_s \in \mathbb{R}^{d \times 1}$ is the trainable weights that project node representations into the 1-dimension score space. idx denotes the indices

¹When generally describing the module that will be repeatedly utilized in the model, we omit the superscripts indicating layer number for brevity.

of node with top- k redundancy scores which are discarded by masking their degrees as 0 (c.f., Eq. (8)). Note that $\text{GGNN}(\cdot)$ in Eq. (6) does not share parameters with the semantics encoder due to their different targets. Besides, we only perform semantic structure refinement on evidences since claims are usually short (less than 10 words) so that the semantic structures are simple and unnecessary to be refined.

Finally, we stack the semantic structure refinement layer over one semantics encoder to form a unified module, namely *evidence semantics miner* (ESM in short), where the long-distance semantic dependency is captured and the redundant information is reduced. In general, we can stack T_R layers of ESM to refine the semantic structures T_R times, eventually followed by a semantics encoder to perform neighborhood propagation on refined semantic graphs, obtaining the fine-grained representations.

3.2.4 Attentive Graph Readout Layer. So far, we have obtained refined structures $\hat{\mathbf{A}}_e^{(T_R)}$ for each evidence and fine-grained node embeddings $\mathbf{H}_c^{(T_E)}$, $\mathbf{H}_e^{(T_R+1)}$ for claims and evidences separately², where T_R and T_E are the numbers of ESM layer and semantics encoder layer of claim, respectively ($T_R = 1$ and $T_E = 2$ in Figure 2). Next, to perform the claim-evidence interaction, we first need to integrate all node embeddings (word embeddings) into general graph embeddings (claim and evidence embeddings). Following previous work [37], we propose to obtain claim-aware evidence representations via the attention mechanism. In detail, we compute the attention score of the j -th word \mathbf{H}_{ej} in the refined evidence graph with the claim representation \mathbf{h}_c^g . Thereafter, the evidence representation \mathbf{h}_e^g is obtained via weighted summation:

$$\mathbf{h}_c^g = \frac{1}{l_c} \sum_{i=1}^{l_c} \mathbf{H}_{ci} \quad (9)$$

$$\mathbf{p}_j = \tanh\left([\mathbf{H}_{ej}; \mathbf{h}_c^g] \mathbf{W}_c\right) \quad (10)$$

$$\alpha_j = \frac{\exp(\mathbf{p}_j \mathbf{W}_p)}{\sum_{i=1}^{l_e} \exp(\mathbf{p}_i \mathbf{W}_p)} \quad (11)$$

$$\mathbf{h}_e^g = \sum_{j=1}^{l_e} \alpha_j \mathbf{H}_{ej} \quad (12)$$

where $[\cdot; \cdot]$ denotes the concatenation of two vectors and $\mathbf{W}_c \in \mathbb{R}^{2d \times d}$ and $\mathbf{W}_p \in \mathbb{R}^{d \times 1}$ are the trainable parameters. l_c and l_e are the length of claim and evidence, respectively. We denote Eq. (10) - (12) as $\text{ATTN}(\mathbf{H}_e, \mathbf{h}_c^g)$ and the attention modules can be easily extended to multi-head ones by concatenating outputs of each head. It is worth noting that based on the fine-grained representations our graph-based model outputs, the above attention mechanism can be replaced by any interaction method in previous work, which we further discuss in Section 4.4.

As Vo and Lee [37] empirically demonstrate that claim speaker and evidence publisher information is important for verification, we extend claim and evidence representations by concatenating

Dataset	# True	# False	# Evi.	# Spe.	# Pub.
Snopes	1164	3177	29242	N/A	12236
PolitiFact	1867	1701	29556	664	4542

Table 1: The statistics of two datasets. The symbol “#” denotes “the number of”. “True” and “False” stand for true claims and false claims, respectively. “Evi.”, “Spe.”, and “Pub.” denote evidences, speakers and publishers.

them with corresponding information vectors, i.e., $\mathbf{h}_c^f = [\mathbf{h}_c^g; \mathbf{s}]$ and $\mathbf{h}_e^f = [\mathbf{h}_e^g; \mathbf{p}]$.

After obtaining the claim and evidence representations, we further employ another attentive network, which is of the same structure as the above, to capture the document-level interaction between a claim and several evidences:

$$\mathbf{H}_e^r = [\mathbf{h}_{e1}^r; \mathbf{h}_{e2}^r; \dots; \mathbf{h}_{en}^r] \quad (13)$$

$$\mathbf{h}_e^f = \text{ATTN}(\mathbf{H}_e^r, \mathbf{h}_c^f) \quad (14)$$

where \mathbf{H}_e^r denotes the concatenation of embeddings of n evidences. Eventually, we integrate claim and evidence embeddings into one unified representation via concatenation, followed by a multi-layer perceptron to output the veracity prediction \hat{y} .

$$\mathbf{h}^f = [\mathbf{h}_c^f; \mathbf{h}_e^f] \quad (15)$$

$$\hat{y} = \text{Softmax}(\mathbf{W}_f \mathbf{h}^f + \mathbf{b}_f) \quad (16)$$

3.2.5 Training Objective. As it is fundamentally a classification task, we utilize the standard cross entropy loss as the objective function, which can be written as:

$$\mathcal{L}_\Theta(y, \hat{y}) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y})) \quad (17)$$

where $y \in \{0, 1\}$ denotes the label of each unverified news.

4 EXPERIMENTS

In this section, we conduct comprehensive experiments to answer the following research questions:

- RQ1: How does GET perform compared to previous fake news detection baselines?
- RQ2: How does the redundant information involved in evidences affect the fake news detection?
- RQ3: How is the performance of different semantic encoders?
- RQ4: How does GET perform with different interaction modules?
- RQ5: How does GET perform under different hyperparameter settings?

4.1 Experimental Setup

4.1.1 Datasets. We utilize two widely used datasets to verify our proposed model. The detailed statistics is summarized in Table 1.

- Snopes [29]. Claims and their corresponding labels (*true* or *false*) are collected from the fact-checking website³. Taking

²We omit the index subscript of evidences for brevity, as they are all fed into the same networks.

³<https://www.snopes.com/>

Method	Snopes								PolitiFact							
	F1-Ma	F1-Mi	F1-T	P-T	R-T	F1-F	P-F	R-F	F1-Ma	F1-Mi	F1-T	P-T	R-T	F1-F	P-F	R-F
LSTM	0.621	0.719	0.430	0.484	0.397	0.812	0.791	0.837	0.606	0.609	0.618	0.632	0.613	0.593	0.590	0.604
TextCNN	0.631	0.720	0.450	0.482	0.430	0.812	0.799	0.826	0.604	0.607	0.615	0.630	0.610	0.592	0.591	0.604
BERT	0.621	0.716	0.431	0.477	0.407	0.810	0.793	0.830	0.597	0.598	0.608	0.619	0.599	0.586	0.577	0.597
DeClarE	0.725	0.786	0.594	0.610	0.579	0.857	0.852	0.863	0.653	0.652	0.675	0.667	0.683	0.631	0.637	0.625
HAN	0.752	0.802	0.636	0.625	0.647	0.868	0.876	0.861	0.661	0.660	0.679	0.676	0.682	0.643	0.650	0.637
EHIAN	0.784	0.828	0.684	0.617	0.768	0.885	0.882	0.890	0.676	0.679	0.689	0.686	0.693	0.655	0.675	0.636
MAC	0.786	0.833	0.687	0.700	0.686	0.886	0.886	0.887	0.672	0.673	0.718	0.675	0.735	0.643	0.676	0.617
CICD	0.789	0.837	0.691	0.632	0.775	0.893	0.890	0.895	0.682	0.685	0.702	0.689	0.714	0.657	0.691	0.629
GET	0.800 [‡]	0.846 [‡]	0.705 [‡]	0.721 [‡]	0.694	0.895 [‡]	0.890	0.902 [‡]	0.691 [‡]	0.694 [‡]	0.723 [‡]	0.687	0.764 [‡]	0.660 [‡]	0.708 [‡]	0.629

Table 2: The model comparison on two datasets Snopes and PolitiFact. “F1-Ma” and “Fi-Mi” denote the metrics F1-Macro and F1-Micro, respectively. “-T” represents “True News as Positive” and “-F” denotes “Fake news as Positive” in computing the precision and recall values. The best performance is highlighted in boldface. ‡ indicates that the performance improvement is significant with p-value ≤ 0.05 .

each claim as a query, the evidences and their publishers are retrieved via the search engine.

- PolitiFact [34]. Claim-label pairs are collected from another fact-checking website⁴ about US politics and evidences are obtained in a similar way to that in Snopes. Aside from publisher information, claim promulgators are added into the dataset. Following previous work [30, 31, 37], we merge *true, mostly true, half true* into the unified class *true* and *false, mostly false, pants on fire* into *false*.

4.1.2 Baselines. To demonstrate the effectiveness of our proposed model GET, we compare it with several existing methods, including both pattern- and evidence-based models, the specific description is listed as follows:

Pattern-based methods.

- LSTM [15]. They utilize LSTM to encode the semantics with the news as input and obtain the final representation of claim via the average pooling.
- TextCNN [40]. They apply a 1D-convolutional network to embed the semantics of claim.
- BERT [9]. They employ BERT to learn the representation of claim. A linear layer is stacked over the special token [CLS] to output the final prediction.

Evidence-based methods.

- DeClarE [30]. They employ BiLSTMs to embed the semantics of evidences and obtain the claim’s representation via average pooling, followed by an attention mechanism performing among claim and each word in evidences to generate the final claim-aware representation.
- HAN [26]. They use GRUs to embed semantics and design two modules named topic coherence and semantic entailment to model the claim-evidence interaction, which are based on sentence-level attention mechanism.
- EHIAN [43]. They utilize self-attention mechanism to learn semantics and concentrate on the important part of evidences for interaction.
- MAC [37]. They introduce a hierarchical attentive framework to model both word- and evidence-level interaction.

- CICD [41]. They introduce individual and collective cognition view-based interaction to explore both local and global opinions towards a claim.

4.1.3 Experimental Environment. We conduct all experiments using PyTorch 1.5.1 on a Linux server equipped with GeForce RTX 3090 GPUs (with 24GB memory each) and AMD EPYC 7742 (256) CPUs.

4.2 Model Comparison (RQ1)

We compare our model GET with eight baselines⁵, including three pattern-based methods and five evidence-based methods. The overall results are shown in Table 2, from which we have the following observations:

Firstly, our proposed model GET outperforms all existing methods on most of metrics on both two datasets by a significant margin, demonstrating the effectiveness of GET. It is worth noting that GET stands out from the recent three sequential-based baselines (EHIAN, MAC, and CICD) whose performance is close, indicating the positive impact of introducing graph-based models to evidence-based fake news detection. In detail, compared to the strongest baselines CICD on two datasets, GET advances the performance about 1 percent on F1-Macro and F1-Micro, which are the evaluation metrics better reflect the overall detection capability of models. With regard to the more fine-grained evaluation, i.e., ‘True news as Positive’ and ‘Fake news as Positive’, GET also achieve the best results on the F1 score on two datasets, where the F1 score is more representative than Precision and Recall since it takes into account both of them synthetically.

Secondly, compared to the pattern-based methods (i.e., the first three methods in Table 2), evidence-based approaches have a substantial performance improvement. This is probably due to the better generalization of evidence-based methods, where the external information is utilized to probe the claim veracity, avoiding the over-reliance on text patterns. In addition, the performance of BERT is similar to that of other pattern-based approaches. We suspect the reason is probably that claims are short and contain lots of noises (e.g., spelling errors and domain-specific abbreviations),

⁴<https://www.politifact.com/>

⁵As some evidence-based methods do not release codes, we reproduce results carefully following settings reported in their original paper.

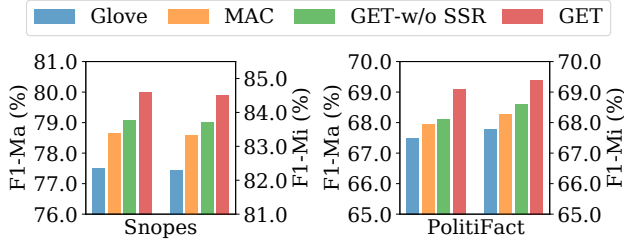


Figure 3: The performance comparison between GET and model variants with different semantic encoders (Glove and MAC) and without structure refinement (GET-w/o SSR).

which are rarely appeared in the pretraining corpus, thus it is hard for BERT to transfer the contextual information learned from the pretrained stage.

Thirdly, among five evidence-based baselines, the performance of DeClarE and HAN is inferior to other three models, which is mainly because they lack exploring the different grain-sized semantics. Specifically, DeClarE only considers word-level semantic interaction and HAN solely relies on document-level representations to model claim-evidence interaction. However, the rest of evidence-based methods all consider multi-level semantics, thus achieving better performance.

4.3 Ablation Study (RQ2, RQ3)

To verify the positive effect of structure refinement for reducing the useless redundancy in evidences, we conduct the ablation study where the structure learning layer is removed and other parts are kept unchanged. We name this model variant as GET-w/o SSR. As shown in Figure 3, we can observe an obvious decline on both datasets regarding the F1-Micro and F1-Macro. This demonstrates the necessity of performing structure refinement on semantic graphs and confirm the effectiveness of our structure learning method. Furthermore, it also indicates that reducing the effect of unimportant information via attention mechanisms will lead to suboptimal results, since they still maintain the noisy semantic structure unchanged [6] (i.e., specifically, all words will participate in the claim-evidence interaction). Therefore, the effect of structure refinement is not overlapped with the attention mechanism, but further goes beyond.

To demonstrate the superiority of the proposed graph-based semantics encoder, we further conduct experiments on two model variants. One is named Glove, where the pretrained word embeddings are directly fed into the attentive readout layer; the other is named MAC, where the semantics encoder is a BiLSTM the same as the baseline [37]. As shown in Figure 3, Glove has the poorest performance since the contextual information is not captured. Moreover, the performance of GET-w/o SSR is superior to that of MAC, indicating that the long-distance structural dependency involved in semantic structure, which is less explored in sequential models, is significant for veracity checking. Note that we choose GET-w/o SSR instead of GET to be compared with MAC fairly, since the only difference between GET-w/o SSR and MAC is the semantics encoder.

Dataset	Metric	DeC	GET-DeC	EHl	wGET-EHl
Snopes	F1-Ma	0.725	0.761	0.784	0.795
	F1-Mi	0.786	0.813	0.828	0.841
	F1-T	0.594	0.649	0.684	0.693
	F1-F	0.857	0.873	0.885	0.897
PolitiFact	F1-Ma	0.653	0.681	0.676	0.688
	F1-Mi	0.652	0.685	0.679	0.690
	F1-T	0.675	0.714	0.689	0.713
	F1-F	0.631	0.647	0.655	0.663

Table 3: The performance of GET with different claim-evidence interaction modules, compared to their corresponding baselines DeClarE (DeC) and EHIAN (EHl). The superior results are highlighted in boldface.

4.4 GET with Different Claim-evidence Interaction Modules (RQ4)

The GET mainly consists of two parts, i.e., graph-based semantic structure mining and attentive graph readout, where the refined semantic structure is obtained in the former stage and the claim-evidence interactions are captured in the latter. As we have mentioned in Section 3.2.4, the semantic structure mining framework can be adaptively connected with any interaction module. Therefore, to further verify the positive impact of graph-based structure mining, we replace the concatenation attention mechanism in our base model with different interaction modules used in existing work. In detail, we choose two modules in representative work: one is the word-level attention mechanism employed in DeClarE [30], the other is the self-attention mechanism utilized to obtain global claim-evidence interactions in EHIAN [43]. We name such two model variants as GET-DeC and GET-EHl, respectively. Thereafter, we can compare them with DeClarE (DeC) and EHIAN (EHl) to see whether the optimized semantic structure can boost model performance with different downstream interaction modules.

The experimental results are shown in Table 3. It is obvious that GET-DeC and GET-EHl both surpass their corresponding competitors, which indicates the effectiveness of our unified graph-based semantic structure mining framework, with being agnostic to the downstream interaction modules. In other words, we can employ such graph-based framework in any evidence-based fake news detection model in a plug-in-play manner, obtaining the fine-grained representations on optimized semantic structures and advancing the model performance.

4.5 Sensitivity Analysis (RQ5)

In this section, we conduct experiments to analyse the performance fluctuation of GET with respect to different values of key hyperparameters.

4.5.1 The number of semantics encoder layer for claims T_E . This hyperparameter decides propagation field on graphs, since stacking T_E -layer encoder (GGNN) makes each node aggregate information within T_E -hop neighborhood. We report the model performance when $T_E = 0, 1, 2, 3$ (See Figure 4) and summarize the observations as follows:

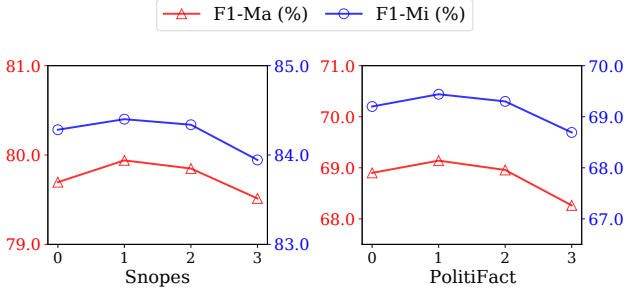


Figure 4: The influence of different semantics encoder layers T_E for claims on model performance.

There is no drastic rise and fall when T_E is changed from 0 to 3. Specifically, the model with $T_E = 1$ slightly outperforms its counterparts. We suspect that the close results are due to the short length of claims (the average lengths of claim are about 6 and 8 in Snopes and PolitiFact, respectively), where the semantic structure can be well-explored merely via 1-hop propagation.

Only one obvious decline is observed between $T_E = 2$ and $T_E = 3$, which is probably caused by the inappropriate propagation field. When the layer number is 3, each node on graphs aggregate information from 3-hop neighborhood, which may cover all nodes since the claims are short, thus failing to model the local semantic structure and leading to the poor performance.

4.5.2 The discarding rate r . This rate is also an important hyperparameter in our proposed model GET. It decides the proportion of redundant information in evidences we filter out. We test the model with r ranging from 0 to 0.6 (See Figure 5) and have the following observations:

When $r = 0$, the model is the same as GET-w/o SSR in the ablation study, where structure refinement layer is removed and no words are dropped. We can see that the performance is not satisfactory since redundant information is preserved that may mislead the model.

The performance grows with r increasing and peaks at the best when $r = 0.4$, which indicates that reducing redundant information plays a positive role in improving the model performance. When r is larger than 0.6, a obvious performance decline can be seen. The probable reason is that some useful information for veracity prediction is mistakenly discarded, so that the model fails to capture the rich semantics in evidences, as the r is too large.

4.5.3 The number of ESM layer T_R . It is a key hyperparameter that controls the information propagation field on graphs and the extent of structure refinement. We observe some phenomena when T_R increases from 0 to 3 (See Figure 6):

The performance is first improved from $T_R = 0$ to $T_R = 1$. Note that when $T_R = 0$, the model downgrades into the one with only a semantics encoder layer. The inferior performance is mainly due to two aspects: 1) It is unable to capture the high-order semantics of long evidences since only features from 1-hop neighborhood are aggregated. 2) Moreover, no redundancy reduction may affect other claim-relevant useful information, since they are fused via

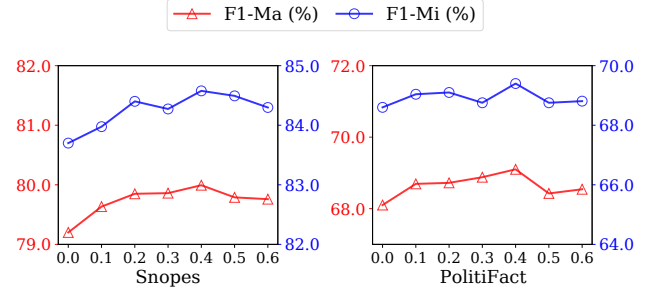


Figure 5: The influence of different discarding rates r on model performance.

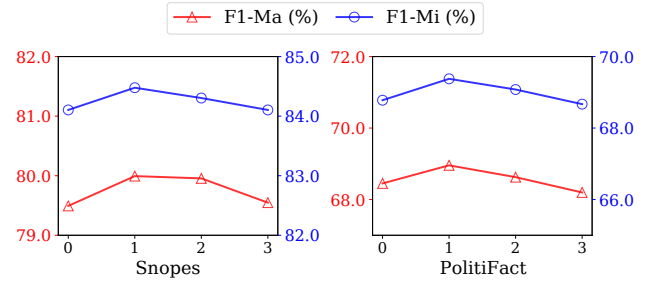


Figure 6: The influence of different evidence semantics miner layers T_R on model performance.

neighborhood propagation. Therefore, these drawbacks, in turn, demonstrate the significance of high-order semantics and structure refinement.

A moderate fall of performance can be seen when T_R ranges from 1 to 3. This is probably because the networks suffer from the over-smoothing problem, which is common in GNNs [21]. Besides, the information is overly discarded so that the evidence semantics is not well modeled, which is also a main reason.

5 CONCLUSION

In this paper, we have proposed a unified graph-based fake news detection model named GET to explore the complex semantic structure. Based on constructed claim and evidence graphs, the long-distance semantic dependency is captured via the information propagation. Moreover, a simple and effective structure learning module is introduced to reduce the redundant information, obtaining fine-grained semantics that are more beneficial for the downstream claim-evidence interaction. We have also validated the performance of GET with different interaction methods, where results demonstrate its ability of acting as a plug-in-play module to boost the performance of other fake news detection models.

ACKNOWLEDGMENTS

This work is supported by National Natural Science Foundation of China (U19B2038, 61772528, 62141608).

REFERENCES

- [1] Oluwaseun Ajao, Deepayan Bhowmik, and Shahrzad Zargari. 2019. Sentiment Aware Fake News Detection on Online Social Networks. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2507–2511.
- [2] Hunt Allcott and Matthew Gentzkow. 2017. Social Media and Fake News in the 2016 Election. *CSN: Politics (Topic)* (2017).
- [3] Adrien Benamira, Benjamin Devillers, Etienne Lesot, Ayush Ray, Manal Saadi, and Fragkiskos D. Malliaros. 2019. Semi-Supervised Learning and Graph Neural Networks for Fake News Detection. *2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)* (2019), 568–569.
- [4] Shantanu Chandra, Pushkar Mishra, Helen Yannakoudakis, and Ekaterina Shutova. 2020. Graph-based Modeling of Online Communities for Fake News Detection. *ArXiv abs/2008.06274* (2020).
- [5] Tianwen Chen and Raymond Chi-Wing Wong. 2020. Handling Information Loss of Graph Neural Networks for Session-based Recommendation. *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (2020).
- [6] Yu Chen, Lingfei Wu, and Mohammed Zaki. 2020. Iterative Deep Graph Learning for Graph Neural Networks: Better and Robust Node Embeddings. In *NIPS*. 19314–19326.
- [7] Luca Cosmo, Anees Kazi, Seyed-Ahmad Ahmadi, Nassir Navab, and Michael M. Bronstein. 2020. Latent Patient Network Learning for Automatic Diagnosis. (2020).
- [8] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In *NIPS*.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL*.
- [10] Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. 2018. Bilevel Programming for Hyperparameter Optimization and Meta-Learning. In *Proceedings of the 35th International Conference on Machine Learning*. 1568–1577.
- [11] Luca Franceschi, Mathias Niepert, Massimiliano Pontil, and Xiao He. 2019. Learning Discrete Structures for Graph Neural Networks. In *Proceedings of the 36th International Conference on Machine Learning*. 1972–1982.
- [12] Hongyang Gao and Shuiwang Ji. 2019. Graph U-Nets. In *Proceedings of the 36th International Conference on Machine Learning*.
- [13] Anastasia Giachanou, Paolo Rosso, and Fabio Crestani. 2019. Leveraging Emotional Signals for Credibility Detection. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '19)*. 877–880.
- [14] William L. Hamilton, Zitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *NIPS*.
- [15] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation* 9 (1997), 1735–1780.
- [16] Bo Jiang, Ziyang Zhang, Doudou Lin, Jin Tang, and Bin Luo. 2019. Semi-Supervised Learning With Graph Learning-Convolutional Networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 11305–11312.
- [17] Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. 2020. Graph Structure Learning for Robust Graph Neural Networks. *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (2020).
- [18] Yiqiao Jin, Xiting Wang, Ruichao Yang, Yizhou Sun, Wei Wang, Hao Liao, and Xing Xie. 2021. Towards Fine-Grained Reasoning for Fake News Detection. *ArXiv abs/2110.15064* (2021).
- [19] Thomas Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. *ArXiv abs/1609.02907* (2017).
- [20] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. 2019. Self-attention graph pooling. In *36th International Conference on Machine Learning, ICML 2019*. 6661–6670.
- [21] Qimai Li, Zhichao Han, and Xiao-Ming Wu. 2018. Deeper Insights into Graph Convolutional Networks for Semi-Supervised Learning. *ArXiv abs/1801.07606* (2018).
- [22] Ruifan Li, Hao Chen, Fangxiang Feng, Zhanyu Ma, Xiaojie Wang, and Eduard H. Hovy. 2021. Dual Graph Convolutional Networks for Aspect-based Sentiment Analysis. In *ACL/IJCNLP*.
- [23] Ruoyu Li, Sheng Wang, Feiyun Zhu, and Junzhou Huang. 2018. Adaptive Graph Convolutional Neural Networks. *ArXiv abs/1801.03226* (2018).
- [24] Qiang Liu, Feng Yu, Shu Wu, and Liang Wang. 2018. Mining significant microblogs for misinformation identification: an attention-based approach. *ACM Transactions on Intelligent Systems and Technology (TIST)* 9, 5 (2018), 1–20.
- [25] Zhenghao Liu, Chenyan Xiong, Maosong Sun, and Zhiyuan Liu. 2020. Fine-grained Fact Verification with Kernel Graph Attention Network. In *ACL*.
- [26] Jing Ma, Wei Gao, Shafiq Joty, and Kam-Fai Wong. 2019. Sentence-Level Evidence Embedding for Claim Verification with Hierarchical Attention Networks. In *ACL*. 2561–2571.
- [27] Salman Bin Naeem and Rubina Bhatti. 2020. The Covid-19 'infodemic': a new front for information professionals. *Health Information and Libraries Journal* (2020).
- [28] Kashyap Popat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. 2016. Credibility Assessment of Textual Claims on the Web. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management (CIKM '16)*. 2173–2178.
- [29] Kashyap Popat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. 2017. Where the Truth Lies: Explaining the Credibility of Emerging Claims on the Web and Social Media. *Proceedings of the 26th International Conference on World Wide Web Companion* (2017).
- [30] Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. 2018. DeClarE: Debunking Fake News and False Claims using Evidence-Aware Deep Learning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. 22–32.
- [31] Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. 2017. Truth of Varying Shades: Analyzing Language in Fake News and Political Fact-Checking. In *EMNLP*.
- [32] Qiang Sheng, Xuexiao Zhang, Juan Cao, and Lei Zhong. 2021. Integrating Pattern- and Fact-based Fake News Detection via Model Preference Learning. *ArXiv abs/2109.11333* (2021).
- [33] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph Attention Networks. *ArXiv abs/1710.10903* (2018).
- [34] Andreas Vlachos and Sebastian Riedel. 2014. Fact Checking: Task definition and dataset construction. In *Proceedings of the ACL 2014 Workshop on Language Technologies and Computational Social Science*. 18–22.
- [35] Andreas Vlachos and Sebastian Riedel. 2015. Identification and Verification of Simple Claims about Statistical Properties. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. 2596–2601.
- [36] Nguyen Vo and Kyumin Lee. 2018. The Rise of Guardians: Fact-Checking URL Recommendation to Combat Fake News. In *The 41st International ACM SIGIR Conference on Research Development in Information Retrieval (SIGIR '18)*. 275–284.
- [37] Nguyen Vo and Kyumin Lee. 2021. Hierarchical Multi-head Attentive Network for Evidence-aware Fake News Detection. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. 965–975.
- [38] Svitlana Volkova, Kyle Shaffer, Jin Yea Jang, and Nathan Hodas. 2017. Separating Facts from Fiction: Linguistic Models to Classify Suspicious and Trusted News Posts on Twitter. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. 647–653.
- [39] Kai Wang, Weizhou Shen, Yunyi Yang, Xiaojun Quan, and Rui Wang. 2020. Relational Graph Attention Network for Aspect-based Sentiment Analysis. In *ACL*.
- [40] William Yang Wang. 2017. "Liar, Liar Pants on Fire": A New Benchmark Dataset for Fake News Detection. In *ACL*.
- [41] Lianwei Wu, Yuan Rao, Yuqian Lan, Ling Sun, and Zhaoyin Qi. 2021. Unified Dual-view Cognitive Model for Interpretable Claim Verification. *arXiv preprint arXiv:2105.09567* (2021).
- [42] Lianwei Wu, Yuan Rao, Ling Sun, and Wangbo He. 2021. Evidence Inference Networks for Interpretable Claim Verification. *Proceedings of the AAAI Conference on Artificial Intelligence* (2021), 14058–14066.
- [43] Lianwei Wu, Yuan Rao, Xiong Yang, Wanzhen Wang, and Ambreen Nazir. 2020. Evidence-Aware Hierarchical Interactive Attention Networks for Explainable Claim Verification. In *IJCAI*. 1388–1394.
- [44] Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. 2019. Session-based Recommendation with Graph Neural Networks. In *AAAI*.
- [45] Xuan-Wei Wu, Lingxiao Zhao, and Leman Akoglu. 2018. A Quest for Structure: Jointly Learning the Graph Structure and Semi-Supervised Classification. *Proceedings of the 27th ACM International Conference on Information and Knowledge Management* (2018).
- [46] Liang Yang, Zesheng Kang, Xiaochun Cao, Di Jin, Bo Yang, and Yuanfang Guo. 2019. Topology Optimization based Graph Convolutional Network. In *IJCAI*.
- [47] Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Graph Convolutional Networks for Text Classification. *ArXiv abs/1809.05679* (2019).
- [48] Zitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. 2018. Hierarchical graph representation learning with differentiable pooling. In *NIPS*. 4800–4810.
- [49] Feng Yu, Q. Liu, Shu Wu, Liang Wang, and Tieniu Tan. 2017. A Convolutional Approach for Misinformation Identification. In *IJCAI*.
- [50] Xueli Yu, Weizhi Xu, Zeyu Cui, Shu Wu, and Liang Wang. 2021. Graph-based Hierarchical Relevance Matching Signals for Ad-hoc Retrieval. *Proceedings of the Web Conference 2021* (2021).
- [51] Jinghao Zhang, Yanqiao Zhu, Qiang Liu, Shu Wu, Shuhui Wang, and Liang Wang. 2021. Mining Latent Structures for Multimedia Recommendation. *Proceedings of the 29th ACM International Conference on Multimedia* (2021).
- [52] Xuexiao Zhang, Juan Cao, Xirong Li, Qiang Sheng, Lei Zhong, and Kai Shu. 2021. Mining Dual Emotion for Fake News Detection (WWW '21). 3465–3476.

- [53] Yingxue Zhang, Soumyasundar Pal, Mark J. Coates, and Deniz Üstebay. 2019. Bayesian graph convolutional neural networks for semi-supervised classification. In *AAAI*.
- [54] Yufeng Zhang, Xueli Yu, Zeyu Cui, Shu Wu, Zhongzhen Wen, and Liang Wang. 2020. Every Document Owns Its Structure: Inductive Text Classification via Graph Neural Networks. *ArXiv abs/2004.13826* (2020).
- [55] Yufeng Zhang, Jinghao Zhang, Zeyu Cui, Shu Wu, and Liang Wang. 2021. A Graph-based Relevance Matching Model for Ad-hoc Retrieval. In *AAAI*.
- [56] Wanjun Zhong, Jingjing Xu, Duyu Tang, Zenan Xu, Nan Duan, M. Zhou, Jiahai Wang, and Jian Yin. 2020. Reasoning Over Semantic-Level Graph for Fact Checking. In *ACL*.
- [57] Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2019. GEAR: Graph-based Evidence Aggregating and Reasoning for Fact Verification. In *ACL*.
- [58] Yanqiao Zhu, Weizhi Xu, Jinghao Zhang, Qiang Liu, Shu Wu, and Liang Wang. 2021. Deep Graph Structure Learning for Robust Representations: A Survey. *CoRR* (2021).

A IMPLEMENTATION DETAILS

We introduce the specific settings in our experiments including hyperparameters, training settings, and the experimental environment.

Following previous work [30, 37], we utilize the same data split⁶ to train and test our model. We also report 5-fold cross validation results, where 4 folds are used for training and the rest one fold is for testing. We utilize Adam optimizer with a learning rate $lr = 0.0001$ and weight decay $decay = 0.001$. The model early stops when F1-macro does not increase in 10 epochs and the maximum number of epoch is 100. We set the maximum length of claims and evidences in both datasets as 30 and 100, respectively. The number of evidences $n = 30$ and the batch size is 32. We set the redundancy discarding rate $r = 0.4$, i.e., $k = rl_e$ will be filtered out in a semantic refinement layer, where l_e is the length of evidence. The number of semantics encoder layer $T_E = 1$ and evidence semantics miner layer $T_R = 1$. The number of word-level and document-level attentive readout head as 5 and 2 for Snopes (3 and 1 for PolitiFact), the dimension of publisher and speaker embedding is both 128, following the work [37]. We use the Glove pretrained embedding with the dimension $d = 300$ for all baselines for a fair comparison.

B VISUALIZATION OF REFINEMENT

In order to better understand what redundant information is discarded by semantic structure refinement, we visualize examples

in both datasets depicted respectively by Figure 7 and Figure 8, where the discarded words are highlighted in grey. It is indicated that most dropped words are adverbs, conjunctions, and pronouns which have less valuable information or are relatively unrelated to the news. For instance, words like ‘it’ and ‘by’ occur frequently but contribute little to the semantic of text. And in the first example in PolitiFact, several nouns like ‘illinois state senate’ have weak connection with its topic and may interfere models’ judgement.

Therefore, the refinement layer can effectively distill important information and get rid of redundant noises.

⁶https://github.com/nguyen09/EACL2021/tree/main/formatted_data/declare

Claim	<i>A new executive order being circulated will severely limit firearms ownership by the elderly [False]</i>
Doc	... skip to content news local and beyond elderly elderly gun ban elderly gun ban claim a new executive order being circulated will severely limit firearms ownership by the elderly image image false image guns to be banned for elderly staff reports united press international washington obama deputy attorney general designate david ogden is circulating a draft of an executive order in which among other things firearms possession would be severely limited to people over 60 an assistant to ogden told us it appears that in these changing times it is no longer necessary to allow ...
Claim	<i>A chinese coal miner recently found alive in an abandoned mine 17 years trapped inside by earthquake [False]</i>
Doc	... images 4 a chinese coal miner was recently found alive in an abandoned mine 17 years after he had been trapped inside it by an earthquake okay i feel pretty silly because i genuinely thought this was real the first time i saw it it doesn't have the wackiness factor of the others that said it false origin ...

Figure 7: Visualization of discarded words in the examples in Snopes Dataset. [True/False] indicates veracity of claims.

Claim	<i>Sen. Obama has always had a 100 percent prochoice rating [True]</i>
Doc	... illinois state senate choice advocates asked strong prochoice legislators like senator obama to vote present on bills like a ban on partial birth abortion to protect a womans right to choose she wrote senator obama has always had a 100 percent prochoice rating and he is the only candidate running for president who stood up and spoke out when south dakota passed an incredibly restrictive ban on abortion the email was one of a handful of gambits ...
Claim	<i>We reduced abortion. We increased adoptions by 135 percent [False]</i>
Doc	...an theory to be sure possibly even more noteworthy than recent research indicating that liberalizing abortion increased premarital sex increased births reduced adoptions and ended so called shotgun marriages but a thorough analysis of abortion and crime statistics leads to the opposite conclusion that abortion increases crime the question about abortion and crime was greatly influenced by a swedish study published in 1966 by hans and inga ...

Figure 8: Visualization of discarded words in the examples in PolitiFact Dataset. [True/False] indicates veracity of claims.