# Adversarial Contrastive Learning for Evidence-aware Fake News Detection with Graph Neural Networks

Junfei Wu, Weizhi Xu, Qiang Liu, *Member, IEEE*,
Shu Wu, *Senior Member, IEEE*, and Liang Wang, *Fellow, IEEE*

**Abstract**—The prevalence and perniciousness of fake news have 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 three 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 underestimate much redundant information contained in evidences that may be useless or even harmful. Thirdly, insufficient data utilization limits the separability and reliability of representations captured by the model, which are sensitive to local evidence. To solve these problems, we propose a unified **G**raph-based s**E**mantic structure mining framework with Con**TRA**stive **L**earning, namely GETRAL in short. Specifically, different from the existing work that treats claims and evidences as sequences, we first 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. Then the fine-grained semantic representations are fed into the downstream claim-evidence interaction module for predictions. Finally, the supervised contrastive learning accompanied with adversarial augmented instances is applied to make full use of data and strengthen the representation learning. Comprehensive experiments have demonstrated the superiority of GETRAL over the state-of-the-arts and validated the efficacy of semantic mining with graph structure and contrastive learning.

**Index Terms**—Evidence-based Fake News Detection, Graph Neural Networks, Contrastive Learning.

✦

## 1 INTRODUCTION

SOCIAL media has facilitated the dissemination and exchange of information, thus profoundly reshaping the convention of people to consume information. However, due to the inability to verify lots of real-time information, social media has also become a hotbed of fake news, which is always fabricated by making some minor changes to the correct statement. Fake news is not only highly deceptive but also inflammatory, potentially influencing real-world events. The widespread of fake news in diverse domains, such as politics [2] and public health [3], has posed a huge threat to web security and human society. Therefore, the research on automatic fake news detection is challenging and in demand.

Generally, previous methods could be roughly categorized into two groups, i.e., pattern-based approaches and evidence-based approaches [4]. 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 focus 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 DeClarE [5] utilizes bidirectional LSTMs to model textual features, followed by a word-level attention mechanism to capture the claim-evidence interaction. HAN [6] further considers the sentence-level interaction to explore more general semantic coherence. To obtain multi-level semantic interaction, some recent works [7], [8] 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). In addition, they ignore sufficient data utilization for capturing separable and reliable representations. To be specific, we argue that there are three main

**Claim**
The Trump administration worked to free 5,000 Taliban prisoners.

**Evidence**
The Trump administration negotiated directly with the Taliban, getting ready to invite them to Camp David, ……, opening up a prison of 5,000 Taliban and probably ISIS-K individuals and letting them free.

The claim-related snippets

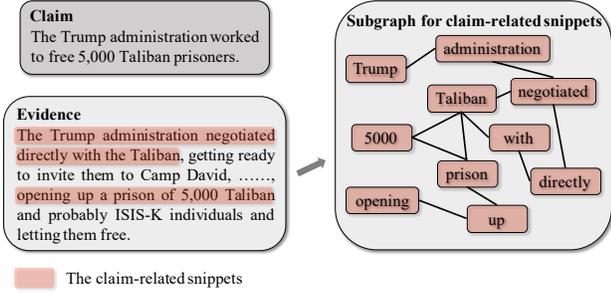**Subgraph for claim-related snippets**

Fig. 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.

weaknesses in previous methods,

(1) 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.

(2) Existing methods pay little attention to 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, we argue that it is significant to model the redundancy with rich semantic structural information, as the redundancy is not only related to the self-information, but also induced by its contexts.

(3) Previous works mainly focus on the learning of mapping from the interaction representations to veracity labels, without sufficient data utilization. This is likely to limit the separability of representations learned in the model. Moreover, insufficient data utilization also causes the sensitivity of the high-level representations to slight changes in the retrieved evidence. In real-world scenes, as the evidence retrieved online may contain much noise and change under different conditions, the captured representations are not reliable and may degrade the detection performance. In addition, the attention mechanism is always utilized for selecting the most useful evidence, which aggravates the

sensitivity of the captured representations to some significant evidence. Therefore, we argue that taking full use of data for efficient training can obtain more separable and reliable representations, which contain core clues instead of being sensitive to local evidence.

To tackle the aforementioned problems, we propose a unified **G**raph-based s**E**mantic structure mining framework with Con**TRA**stive **L**earning, namely GETRAL for exploring fine-grained semantics and capturing enhanced representations. Firstly, modeling sequential data as graphs has benefited many tasks, such as text classification [9], [10] and sequential recommendation [11], owing to its capability of capturing long-distance structural dependency. To this end, we firstly utilize graph structure to model both claims and evidences, where nodes indicate words and edges represent the co-occurence 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 including both self-features and their contexts. The former is related to its own information and its relevance to the claim, and the latter is related to its graph topology. To date, our graph-based framework has captured the fine-grained semantics via long-distance dependency modeling and redundancy mitigation.

Finally, inspired by the success of recent work [12], [13] which integrates contrastive learning, we introduce a supervised contrastive learning auxiliary task to strengthen the representation learning. In addition, to reduce the sensitivity to local evidence, we employ adversarial gradient perturbation [14], [15] to augment the contrastive instances at the feature level. In specific, with the veracity label utilized as the supervised signal, the representations of claim-evidence interactions of the same class are pulled close, while those of different classes are pulled apart. Subsequently, the reliable representations to discriminate different claim-evidence interactions can be better captured.

Our main contributions can be summarized as follows:

- We model claims and evidences as graph-structured data and introduce a simple and effective graph structure learning approach for redundancy mitigation. The captured long-distance and fine-grained semantics based on the structure can boost the performance of downstream interaction models.
- We introduce a supervised contrastive learning task and integrate the adversarial gradient perturbation for efficient training. Then the captured representations are more separable and reliable for detection.
- Comprehensive experiments are conducted to demonstrate the superiority of GETRAL and the effectiveness of each component. Our code is available at https://github.com/CRIPAC-DIG/GETRAL

## 2 RELATED WORK

### 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 [16], [17] and spatial approaches [18], [19]. 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 [11], [20], [21], [22], text classification [9], [10], and sentiment analysis [23], [24].

Recently, researchers have observed that graphs inevitably contain noises that may deteriorate the training of GNNs [25]. 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 [26]: 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 [27], [28], [29] or kernel function [30], [31]. 2) *the probabilistic methods* assume that the adjacency matrix is generated by sampling from a specific probabilistic distribution [32], [33], [34]. 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 (namely 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) [25], [35]. It is worth noting that existing graph pooling methods [36], [37], [38] could also be viewed 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 [39], [40], [41]. 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 recently, 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 works always focus on different kinds of patterns. Popat et al. [42] 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 [43], [44], [45], [46], [47], [48], [49]. Recently, more attention has been paid to the emotional pattern mining, where it holds an assumption that there are probably obvious sentiment biases in fake news [50], [51], [52], [53].

The second is the evidence-based pipeline where researchers propose to explore the semantic similarity (con-flict) in claim-evidence pairs to check the news veracity. Evidences are usually retrieved from the knowledge graph [54] or fact-checking websites [55] by giving unverified claims as queries. DeClarE [5] is the first work to utilize evidences in fake news detection. It employs BiLSTMs to embed the semantics of evidences and obtains the claim's sentence-level representation via average pooling. Next, it introduces 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 [6] computes the sentence-level coherence and entailment scores between claims and evidences. EHIAN [56] employs the self-attention mechanism to obtain word-level interaction scores. Recent works [7], [8], [57] hierarchically integrate both word-level and sentence-level interactions into the final representation for verification. In summary, they all employ sequential models to embed semantics and apply attention mechanisms to capture the claim-evidence interactions.

Different from existing works, 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 graph structure learning.

### 2.3 Contrastive Learning

Contrastive learning is an effective training paradigm that captures separable and distinguishable representations which can bring significant improvement for downstream tasks. Specifically, it utilizes InfoNCE loss [58] to pull the representations of positive samples closer while pulling the negative samples apart, forming a representation space with alignment and uniformity. Nowadays, contrastive learning has been applied to several tasks [58], [59], [60], [61].

Contrastive learning was first introduced for training visual representations [59], [60], [62], by conducting visual augmentations including cropping, resizing and other operations at the visual level to obtain positive pairs, while different instances naturally form the negative pairs with each other. Subsequently, supervised contrastive learning [61] has also gained much attention. It integrates the class relationship to construct contrastive instances to calibrate representations. For graph data, perturbations are imposed on the graph topology and node features to generate corrupted views as the contrastive instances. Then, the graph representations which better contain structure and semantic information are obtained by maximizing the agreement between either global graph embeddings or local node embeddings [63], [64], [65], [66]. In addition, as it is critical for natural language processing to compress dense semantics, contrastive learning has been introduced to leverage abundant textual resources to learn better representations. SimCSE [67] uses only dropout as minimal data augmentation to sentence embeddings and boost the performance of pretrained model significantly. In distantly supervised relation extraction, CIL [12] and HiCLRE [15] exploit the abundant instance relations and propose contrastive instance learning to obtain accurate representations under noise efficiently. Contrastive learning has also been adopted by pattern-based fake news

detection for domain adaptation [68], [69]. Different from these works, we combine a supervised contrastive learning task with the classification task to capture more separable representations for evidence-aware fake news detection.

# 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 $\mathcal{E} = \{e_1, e_2, \ldots, e_n\}$, 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 $\Theta$ is its trainable parameters.

## 3.2 The Proposed Model: GETRAL

In this part, we elaborate our unified graph-based model GETRAL, which can be mainly separated into five modules: 1) *Graph Construction*, 2) *Graph-based Semantics Encoder*, 3) *Semantic Structure Refinement*, 4) *Attentive Graph Readout Layer*, and 5) *Adversarial Contrastive Learning Module*.

### 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 [9], [10], [70], [71], 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 the adjacency matrix is 1, otherwise 0), which captures the local context in the 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 inferior to the sliding-window-based method, which may be due to the redundant noises induced by the dense connection.

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.

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.

### 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 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:

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

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

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

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

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

where $\mathcal{C}$ 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. $\sigma$ is the non-linear activation unit and we utilize the Sigmoid function in our model. For brevity, we denote Eq. (1) - (5) as $\mathbf{GGNN}(\tilde{\mathbf{A}}, \mathbf{H})$[1].

### 3.2.3 Semantic Structure Refinement

As evidences always contain redundant information that may mislead the 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 [28], [34], [38].

In particular, we propose to compute a redundancy score for each node, based on which we obtain a ranking list and the nodes with the top-$k$ redundancy scores will be discarded, then we adjust the aggregation weight for the rest nodes. For each node, its redundancy can be determined by its own information and its relevance to the claim. Hence, we evaluate the independent information redundancy of each node from the view of the node itself and claim relevance, respectively. Specifically, the node-self redundancy is directly measured by a linear projection. To obtain the claim-related redundancy, we utilize the Gaussian kernel to measure the fine-grained relevance of each token in the evidence to all tokens in claims. It has been proved in [40], [48] that Gaussian kernel can summarize matching features effectively. We construct a fine-grained translation matrix $\mathbf{M}$ based on the cosine similarity, where $\mathbf{M}_{ij} = cos(\mathbf{H}_{ei}, \mathbf{H}_{cj})$. Subsequently, the Gaussian kernel is applied to transform

---

1. When generally describing the module that will be repeatedly utilized in the model, we omit the superscripts indicating layer number for brevity.
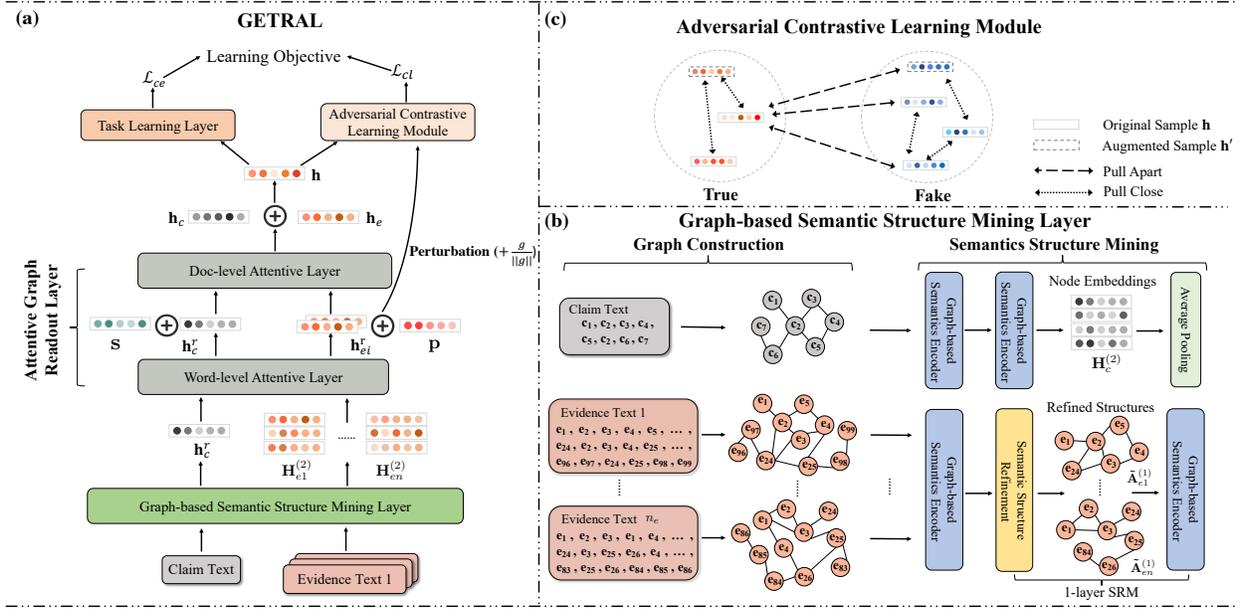
Fig. 2. (a) The overall architecture of GETRAL. It consists of Graph-based Semantic Structure Mining Layer, Attentive Graph Readout Layer, Task Learning Layer and Adversarial Contrastive Learning Module. (b) The Graph-based Semantic Structure Mining Layer first transform plain texts into graphs, then perform neighbourhood aggregation and structure learning to obtain fine-grained semantics. (c) The Adversarial Contrastive Learning Module pulls the representations of positive pairs close while pushing representations of negative pairs apart.

the translation matrix into kernel features, attending to different levels of relevance. It can be denoted as:

$$\mathbf{S}_{se} = \hat{\mathbf{H}}_e \mathbf{W}_{se} \tag{6}$$

$$\mathbf{K}_i = [\mathbf{K}_{i1}; \mathbf{K}_{i2}; \ldots; \mathbf{K}_{ik}] \tag{7}$$

$$\mathbf{K}_{it} = \log \sum_j exp(-\frac{(\mathbf{M}_{ij} - \mu_t)^2}{2\sigma_t^2}) \tag{8}$$

$$\mathbf{S}_{sc} = \mathbf{K}\mathbf{W}_{sc} \tag{9}$$

where $\mathbf{W}_{se} \in \mathbb{R}^{d \times 1}, \mathbf{W}_{sc} \in \mathbb{R}^{k \times 1}$ are trainable weights that project representations into the shared 1-dimension score space. $k$ denotes the number of kernel with corresponding mean $\mu_t$ and width $\sigma_t$ that captures different specific similarity regions [72].

However, the redundancy is not only related to the information contained for claim verification 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 context-aware redundancy scores, which takes into account both node-self and context information. Finally the two scores are fused by a simple linear combination. Mathematically, it can be formulated as:

$$\mathbf{s}_{se} = \mathbf{GGNN}(\tilde{\mathbf{A}}, \mathbf{S}_{sc}) \tag{10}$$

$$\mathbf{s}_{sc} = \mathbf{GGNN}(\tilde{\mathbf{A}}, \mathbf{S}_{se}) \tag{11}$$

$$\mathbf{s}_r = (1 - \beta)\mathbf{s}_{se} + \beta\mathbf{s}_{sc} \tag{12}$$

$$idx = topk\_index(\mathbf{s}_r) \tag{13}$$

$$\tilde{\mathbf{A}}_{idx,:} = \tilde{\mathbf{A}}_{:,idx} = 0 \tag{14}$$

where $\beta$ is an introduced coefficient that controls the fusion proportion of node-self and claim-related score. $idx$ denotes

the indices of node with top-$k$ redundancy scores which are discarded by masking their degrees as 0 (c.f., Eq. (14)). Note that $\mathbf{GGNN}(\cdot)$ in Eq. (10) 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.

Notably, as the trainable parameters related to redundancy scoring need to be updated by back-propagation, we modify the Eq.(1) in $\mathbf{GGNN}(\cdot)$ when it is just after semantic structure refinement by multiplying the normalized score to scale features:

$$\mathbf{a}_i = \sum_{(w_i, w_j) \in \mathcal{C}} \tilde{\mathbf{A}}_{ij} \mathbf{W}_a \mathbf{H}_j (1 - \sigma(\mathbf{s}_{rj})) \tag{15}$$

where $\sigma$ is the Sigmoida function.

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

### 3.2.4 Attentive Graph Readout Layer

So far, we have obtained refined structures $\tilde{\mathbf{A}}_e^{(T_R)}$ for each evidence and fine-grained node embeddings $\mathbf{H}_c^{(T_E)}$, $\mathbf{H}_e^{(T_R+1)}$ for claims and evidences separately[2], where $T_R$

---

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

and $T_E$ are the numbers of the SRM layer and semantics encoder layer of the 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 [7], we propose to obtain claim-aware evidence representations via a word-level attentive layer. In detail, we compute the attention score of the $j$-th word $\mathbf{H}_{ej}$ in the refined evidence graph with the claim representation $\mathbf{h}_c^r$. Thereafter, the evidence representation $\mathbf{h}_e^r$ is obtained via weighted summation:

$$\mathbf{h}_c^r = \frac{1}{l_c} \sum_{i=1}^{l_c} \mathbf{H}_{ci} \tag{16}$$

$$\mathbf{p}_j = \tanh\left([\mathbf{H}_{ej}; \mathbf{h}_c^r] \mathbf{W}_c\right) \tag{17}$$

$$\alpha_j = \frac{\exp\left(\mathbf{p}_j \mathbf{W}_p\right)}{\sum_{i=1}^{l_e} \exp\left(\mathbf{p}_i \mathbf{W}_p\right)} \tag{18}$$

$$\mathbf{h}_e^r = \sum_{j=1}^{l_e} \alpha_j \mathbf{H}_{ej} \tag{19}$$

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. (17) - (19) as $\mathbf{ATTN}(\mathbf{H}_e, \mathbf{h}_c^r)$ and the attention modules can be easily extended to multi-head ones by concatenating outputs of each head.

As MAC [7] empirically demonstrates that claim speaker and evidence publisher information is important for verification, we extend claim and evidence representations by concatenating them with corresponding information vectors:

$$\mathbf{h}_c = [\mathbf{h}_c^r; \mathbf{s}] \tag{20}$$

$$\mathbf{h}_e^g = [\mathbf{h}_e^r; \mathbf{p}] \tag{21}$$

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

$$\mathbf{H}_e^g = [\mathbf{h}_{e1}^g; \mathbf{h}_{e2}^g; \ldots; \mathbf{h}_{en}^g] \tag{22}$$

$$\mathbf{h}_e = \mathbf{ATTN}(\mathbf{H}_e^g, \mathbf{h}_c) \tag{23}$$

where $\mathbf{H}_e^g$ 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 $\hat{y}$.

$$\mathbf{h} = [\mathbf{h}_c; \mathbf{h}_e] \tag{24}$$

$$\hat{y} = \text{Softmax}(\mathbf{W}_f \mathbf{h} + \mathbf{b}_f) \tag{25}$$

As it is fundamentally a classification task, we utilize the standard cross entropy loss $\mathcal{L}_{ce}$ as the task loss, which can be written as:

$$\mathcal{L}_{ce} = -(y \log \hat{y} + (1-y) \log(1-\hat{y})) \tag{26}$$

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

### 3.2.5 Adversarial Contrastive Learning Module

To make full use of data, we propose an auxiliary supervised contrastive learning task to help calibrate representations. The supervised contrastive learning aims to pull samples of the same class close and samples of different classes apart, exploiting the common properties within the class. Corresponding contrastive auxiliary loss $\mathcal{L}_{sup}$ takes the following form:

$$\mathcal{L}_{cl} = \frac{-1}{|\mathcal{P}(\mathbf{h})|} \sum_{\mathbf{h}_p \in \mathcal{P}(\mathbf{h})} \log \frac{\exp(cos(\mathbf{h}, \mathbf{h}_p)/\tau)}{\sum_{\mathbf{h}_n \in \mathcal{N}(\mathbf{h})} \exp(cos(\mathbf{h}, \mathbf{h}_n)/\tau)} \tag{27}$$

where $\mathcal{P}(\mathbf{h})$ is the set of positive samples and $\mathcal{N}(\mathbf{h})$ is the set of negative samples to $\mathbf{h}$ in the same batch. $\mathbf{h}_p$ is the positive sample with the same label and $\mathbf{h}_n$ is the negative sample with a different label. $\tau \in \mathbb{R}$ is a temperature parameter and $cos(\cdot)$ denotes the cosine similarity function.

Moreover, to further reduce the sensitivity of models to local evidence, we employ adversarial gradient perturbation [14], [15] to construct augmented samples for more efficient contrastive learning. It can be viewed as an intentional injected noise. Compared to discrete augmentations on textual content like word deletion, insertion, and substitution which may hurt the semantics of complex sentences, the gradient-based method can maintain semantics to the greatest extent and simulate a worst case. These augmented instances are integrated to enrich the representations in the representation space.

Specific to the augmentation process, we first select the piece of evidence $\mathbf{h}_{ek}^g$ with the highest attention score in $\mathbf{H}_e^g$, which contributes most to the synthetic evidence representation $\mathbf{h}_e$. Then we utilize the gradient adversarial perturbation to $\mathbf{h}_{ek}^g$ in the representation level to obtain the perturbed evidence $\mathbf{h}_{ek}^{g'}$. Finally, the perturbed evidence along with other evidence representations is input into the same document-level attentive network $\mathbf{ATTN}()$ to obtain the final augmented view $\mathbf{h}'$. So far, the positive set $\mathcal{P}(\mathbf{h})$ contains the original samples and adversarial augmented samples of the same class, and the negative set $\mathcal{N}(\mathbf{h})$ contains those of different classes. The perturbation process above can be mathematically expressed as:

$$\mathbf{g}_{ek} = \nabla_{\mathbf{h}_{ek}^g} \mathcal{L}_{ce} \tag{28}$$

$$\mathbf{h}_{ek}^{g'} = \mathbf{h}_{ek}^g + \epsilon \frac{\mathbf{g}_{ek}}{||\mathbf{g}_{ek}||} \tag{29}$$

$$\mathbf{H}_e^{g'} = [\mathbf{h}_{e1}^g; \ldots; \mathbf{h}_{ek}^{g'}; \ldots] \tag{30}$$

$$\mathbf{h}_e' = \mathbf{ATTN}(\mathbf{H}_e^{g'}, \mathbf{h}_c) \tag{31}$$

$$\mathbf{h}' = [\mathbf{h}_c; \mathbf{h}_e'] \tag{32}$$

where $\mathbf{g}_{ek}$ is the first-order derivation of target loss at variable $\mathbf{h}_{ek}^g$. $\epsilon$ is the norm parameter to control the normalized gradient as a valid perturbation. Then $\mathbf{h}'$ is the unified representation of an augmented view that shares the same label with the original sample. These instances are used to simulate the instances with noise, thus improving the difficulty of contrastive learning.

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.

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

### 3.2.6  Training Objective

Finally, we combine the cross entropy loss and supervised contrastive loss as a joint optimization target loss:

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda \mathcal{L}_{cl} \tag{33}$$

where $\lambda$ is a hyper-parameter to control the extent of contrastive learning.

## 4  EXPERIMENTS

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

- RQ1: How does GETRAL perform compared to previous fake news detection baselines?
- RQ2: How effective are the structure modeling component and contrastive learning component to GETRAL?
- RQ3: How does GETRAL 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 [73]. Claims and their corresponding labels ($true$ or $false$) are collected from the fack-checking website[3]. Taking each claim as a query, the evidences and their publishers are retrieved via the search engine.
- PolitiFact [55]. Claim-label pairs are collected from another fact-checking website[4] 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 [5], [7], [74], 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 GETRAL, we compare it with several existing methods, including both pattern- and evidence-based models, and the specific description is listed as follows:

**Pattern-based methods.**

- **LSTM** [75] utilizes LSTM to encode the semantics with the news as input and obtains the final representation of claim via average pooling.

3. https://www.snopes.com/
4. https://www.politifact.com/

- **TextCNN** [76] applies a 1D-convolutional network to embed the semantics of claim.
- **BERT** [77] employs 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** [5] employs BiLSTMs to embed the semantics of evidences and obtains 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** [6] uses GRUs to embed semantics and designs two modules named topic coherence and semantic entailment to model the claim-evidence interaction, which are based on sentence-level attention mechanism.
- **EHIAN** [56] utilizes self-attention mechanism to learn semantics and concentrates on the important part of evidences for interaction.
- **MAC** [7] introduces a hierarchical attentive framework to model both word- and evidence-level interaction.
- **CICD** [8] introduces individual and collective cognition view-based interaction to explore both local and global opinions towards a claim.

#### 4.1.3  Implementation Details

Following previous work [5], [7], we utilize the same data split[5] 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. The fusion rate $\beta$ is 0.5. The kernel size is set to 11 for Snopes and 21 for PolitiFact. Specifically, one kernel with $\mu = 1$ and $\sigma = 10^{-3}$ can capture exact matches [48], while $\mu$ of the other kernels is spaced evenly between $[-1, 1]$ and $\sigma$ is set to 0.1. We set the redundancy discarding rate $r = 0.3$ for Snopes and $r = 0.2$ for PolitiFact, 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 $T_E$ and evidence semantics miner layer $T_R$ is both 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 [7]. We use the Glove pretrained embedding with the dimension $d = 300$ for all baselines for a fair comparison. We conduct all experiments using PyTorch 1.5.1 on a Linux server equipped with GeForce RTX 3090 GPUs and AMD EPYC 7742 CPUs.

5. https://github.com/nguyenvo09/EACL2021/tree/main/formatted_data/declare

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 $\leq 0.05$.

| 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 | 62.10 | 71.87 | 42.95 | 48.42 | 39.69 | 81.25 | 79.14 | 83.67 | 60.56 | 60.87 | 61.82 | 63.19 | 61.27 | 59.31 | 59.05 | 60.43 |
| TextCNN | 63.08 | 72.01 | 45.00 | 48.16 | 43.04 | 81.16 | 79.88 | 82.62 | 60.38 | 60.74 | 61.52 | 63.01 | 61.03 | 59.24 | 59.05 | 60.42 |
| BERT | 62.05 | 71.62 | 43.07 | 47.73 | 40.65 | 81.04 | 79.31 | 82.97 | 59.71 | 59.81 | 60.81 | 61.95 | 59.90 | 58.60 | 57.73 | 59.70 |
| DeClarE | 72.54 | 78.61 | 59.43 | 61.03 | 57.93 | 85.67 | 85.25 | 86.39 | 65.31 | 65.25 | 67.49 | 66.71 | 68.32 | 63.11 | 63.70 | 62.46 |
| HAN | 75.21 | 80.23 | 63.58 | 62.50 | 64.69 | 86.83 | 87.64 | 86.11 | 66.12 | 66.01 | 67.92 | 67.58 | 68.20 | 64.33 | 64.97 | 63.73 |
| EHIAN | 78.43 | 82.83 | 68.41 | 61.69 | 76.79 | 88.47 | 88.18 | 89.04 | 67.22 | 67.95 | 68.92 | 68.64 | 69.34 | 65.52 | 67.49 | 63.60 |
| MAC | 78.66 | 83.32 | 68.74 | 70.00 | 68.60 | 88.58 | 88.62 | 88.71 | 68.03 | 68.25 | 71.78 | 67.54 | 73.49 | 64.28 | 67.61 | 61.68 |
| CICD | 78.92 | 83.73 | 69.07 | 63.20 | **77.48** | 89.30 | **88.99** | 89.54 | 68.18 | 68.48 | 70.24 | 68.92 | 71.44 | 65.72 | 69.12 | 62.93 |
| GETRAL | **80.61**‡ | **85.12**‡ | **71.26**‡ | **74.18**‡ | 68.79 | **89.96**‡ | 88.90 | **91.04**‡ | **69.53**‡ | **69.81**‡ | **72.21**‡ | **69.73**‡ | **75.10**‡ | **66.84**‡ | **70.26**‡ | **64.01**‡ |

## 4.2 Model Comparison (RQ1)

We compare our model GETRAL with eight baselines[6], 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 model GETRAL outperforms all existing methods on most of metrics on both two datasets by a significant margin, demonstrating the effectiveness of GE-TRAL. It is worth noting that GETRAL stands out from the recent three sequential-based baselines (EHIAN, MAC, and CICD) whose performances are close, indicating the positive impact of introducing graph-based models and contrastive learning to evidence-based fake news detection. In detail, compared to the strongest baselines CICD on two datasets, GETRAL advances the performance by about 1.5 percent on F1-Macro and F1-Micro, which can 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', GETRAL 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.

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), 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

---

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

---

TABLE 3
The performance comparison between GETRAL and model variants.

| Method | Snopes | | PolitiFact | |
|---|---|---|---|---|
| | F1-Ma | F1-Mi | F1-Ma | F1-Mi |
| GETRAL-SE-CL | 77.51 | 82.31 | 67.47 | 67.77 |
| GETRAL-GSE-CL | 78.66 | 83.32 | 68.03 | 68.25 |
| GETRAL-SSR-CL | 79.49 | 84.10 | 68.45 | 68.78 |
| GETRAL-CL | 80.12 | 84.52 | 69.25 | 69.60 |
| GETRAL-AD | 80.32 | 84.81 | 69.40 | 69.69 |
| GETRAL | 80.61 | 85.12 | 69.53 | 69.81 |

methods all consider multi-level semantics, thus achieving better performance.

## 4.3 Ablation Study (RQ2)

To verify the effect of components related to fine-grained semantics mining or contrastive learning of GETRAL, we conduct the ablation study for several variants by removing the specific component: **-SE** removes any semantics encoder and feeds the pretrained word embeddings e.g. Glove, into the attentive readout layer directly; **-GSE** removes the graph semantic encoder and utilizes the BiLSTM as the semantics encoder like the baseline [7]; **-SSR** removes the structure learning layer which is proposed to reduce the useless redundancy in evidences; **-CL** only uses the task specific classification loss $\mathcal{L}_{ce}$ for training; **-AD** performs a simplified supervised contrastive learning without constructing the augmented instances by adversarial gradient perturbation.

The experimental results are shown in Table 3, from which we can observe that each variant suffers form an obvious decline on both datasets regarding the F1-Micro and F1-Macro. Specific analyses are as follows:

- In terms of different semantics encoders, GETRAL-SE-CL has the poorest performance since the contextual information is not captured. Moreover, the performance of GETRAL-SSR-CL is superior to that of GETRAL-GSE-CL, 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 GETRAL-SSR-CL instead of GETRAL-CL to be compared with GETRAL-GSE-CL fairly, since the only difference between GETRAL-SSR-CL and GETRAL-GSE-CL is the semantics encoder.
- The performance degradation of GETRAL-SSR-CL demonstrates the necessity of performing structure

refinement on semantic graphs and confirms the effectiveness of our structure learning method. Furthermore, it indicates that reducing the effect of unimportant information via attention mechanisms will lead to suboptimal results, since it still maintains the noisy semantic structure unchanged [28] (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.

- Additionally, GETRAL gains significant improvements on GETRAL-CL and GETRAL-AD. This indicates that contrastive learning and adversarial gradient perturbation are both beneficial. Contrastive learning can mine the underlying relations to help capture the core difference between classes, while the adversarial augmented instances can further boost this process. It is worth noting that gradient adversarial augmentation has a more significant effect on the Snopes dataset compared to PolitiFact. We attribute this improvement to the supplement of obviously unbalanced samples in Snopes, which promotes efficient contrastive learning.

## 4.4 Sensitivity Analysis (RQ3)

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

### 4.4.1 The number of semantics encoder layer for claims $T_E$

This hyperparameter decides the 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 3) and summarize the observations as follows:

There is a significant improvement when $T_E$ is changed from 0 to 1. Specifically, the model with $T_E = 1$ outperforms its counterparts. We suspect that the gains 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.

An obvious decline are observed between $T_E = 1$ and $T_E = 3$, which is probably caused by the inappropriate propagation field. When the layer number is increased, each node on graphs aggregates information from the multi-hop neighborhood, which may cover all nodes since the claims are short, thus failing to model the local semantic structure and leading to poor performance.

### 4.4.2 The fusion coefficient $\beta$

The fusion coefficient $\beta$ controls the fusion rate of node-self and claim related score in the graph structure learning. $\beta = 0$ denotes GETRAL only considers the information within evidence to determine redundancy, and $\beta = 1$ denotes GETRAL only considers the fine-grained claim-related information to determine redundancy. From the results in Figure 4, we can observe that:
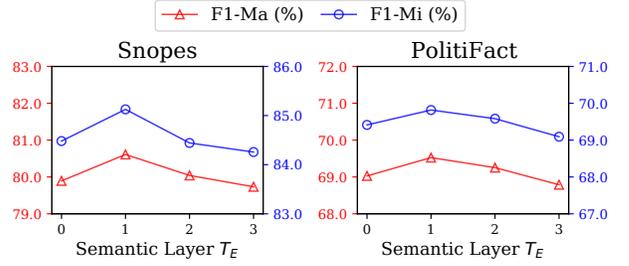


Fig. 3. The influence of different semantics encoder layers $T_E$ for claims on model performance.
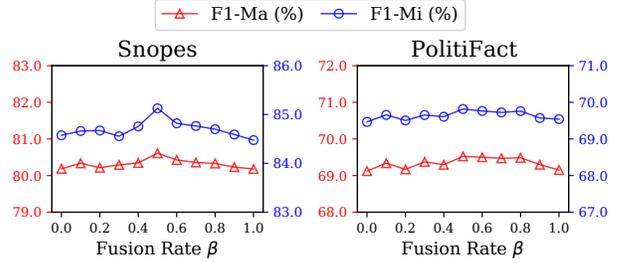


Fig. 4. The influence of different fusion rate $\beta$ on model performance.

When $\beta$ is set to 0 or 1, the performance is relatively poor. Because it only determines redundancy from a single perspective, failing to exploit the rich information contained in each node and its relation to the claim. And when $\beta$ changing from 0 to 1, the performance increases first and then decreases, indicating that a moderate fusion score is useful. It is worth noting that the best performance is achieved when $\beta = 0.5$. It denotes that the claim-related redundancy score is as important as the node-self score and they both should be taken into consideration.

With varied fusion coefficient $\beta$, GETRAL always obtains a competitive performance on both datasets. It indicates that our method is relatively insensitive to the change of $\beta$ and proves the validity of the fused redundancy score.

### 4.4.3 The discarding rate $r$

This discarding rate 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 can be viewed as integrating the redundancy signals when encoding semantics, without dropping nodes. It is worth noting that this degraded version of GETRAL differs from the one without semantic structure refinement and can still maintain competitive performance. It is mainly because the former can adaptively reduce the weight of redundant nodes while the latter can't distinguish the difference and treat each node equally.

The performance grows with $r$ increasing and peaks at the best when $r = 0.3$ in Snopes and $r = 0.2$ in PolitiFact, which indicates that reducing redundant information plays a positive role in improving the model performance. When $r$ continues to increase, an obvious performance decline can be seen. The probable reason is that some useful information
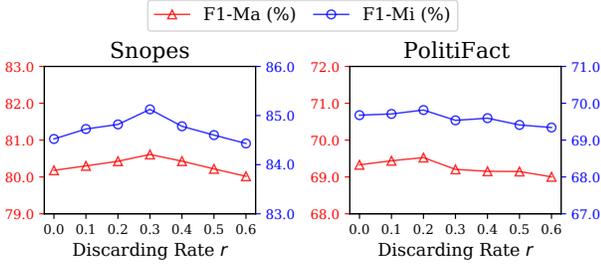
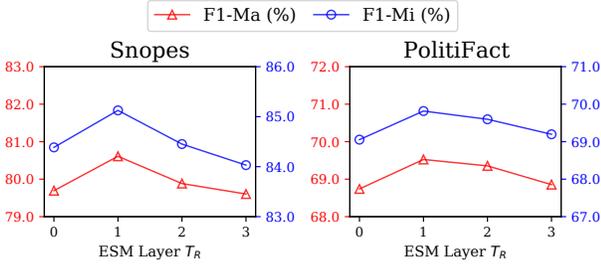Fig. 5. The influence of different discarding rates $r$ on model performance.



Fig. 6. The influence of different semantic refinement and miner layers $T_R$ on model performance.

for veracity prediction is mistakenly discarded, so that the model fails to capture the rich semantics in evidences, as the $r$ is too large.

### 4.4.4 The number of SRM 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 this redundant information are fused via neighborhood propagation. Therefore, these drawbacks, in turn, demonstrate the significance of high-order semantics and structure refinement.

A significant 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 [78]. Besides, the information is overly discarded so that the evidence semantics is not well modeled.

### 4.4.5 The contrastive coefficient $\lambda$

We also conduct experiments to study the impact of contrastive coefficient $\lambda$, with different values ranging from 0.0 to 0.3 (See Figure 7). This hyperparameter decides the extent of the auxiliary contrastive learning task besides classification task. We have the following observations:
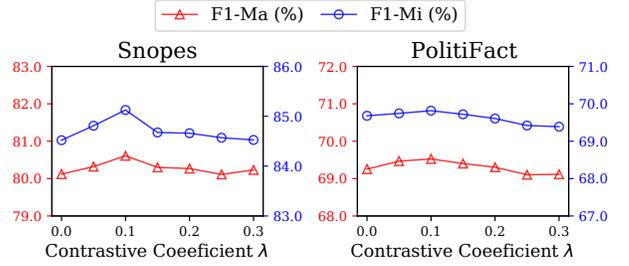


Fig. 7. The influence of different contrastive coefficient $\lambda$ on model performance.

There is a significant improvement when $\lambda$ ranging from 0.0 to 0.10, and the model achieves the best performance when $\lambda = 0.1$. In specific, $\lambda = 0.0$ denotes the simplified version GETRAL-CL which only focuses on the classification task. It indicates that the auxiliary contrastive learning task can improve the performance of our method.

However, the performance begins to decline obviously when $\lambda$ continues to increase. This is due to a moderate $\lambda$ can help the model to capture separable and reliable representations which are beneficial to detection, while a larger value may distract it from the main task. The problem of sensitivity to the proportion of the auxiliary contrastive learning task has recently been noticed by existing works in [79], [80].

## 5 CONCLUSION

In this paper, we have proposed a unified graph-based fake news detection model with adversarial contrastive learning named GETRAL to explore the complex semantic structure and enhance the representation learning. 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. Finally, we integrate an adversarial contrastive learning task to capture separable representations to help fake news detection. We have conducted empirical experiments to demonstrate the superiority of our proposed method.

## REFERENCES

[1] W. Xu, J. Wu, Q. Liu, S. Wu, and L. Wang, "Evidence-aware fake news detection with graph neural networks," in *WWW*, 2022, p. 2501–2510.
[2] H. Allcott and M. Gentzkow, "Social media and fake news in the 2016 election," *CSN: Politics (Topic)*, 2017.
[3] S. B. Naeem and R. Bhatti, "The covid-19 'infodemic': a new front for information professionals," *Health Information and Libraries Journal*, 2020.

[4] Q. Sheng, X. Zhang, J. Cao, and L. Zhong, "Integrating pattern- and fact-based fake news detection via model preference learning," *ArXiv*, vol. abs/2109.11333, 2021.

[5] K. Popat, S. Mukherjee, A. Yates, and G. Weikum, "Declare: Debunking fake news and false claims using evidence-aware deep learning," in *EMNLP*, 2018, pp. 22–32.

[6] J. Ma, W. Gao, S. Joty, and K.-F. Wong, "Sentence-level evidence embedding for claim verification with hierarchical attention networks," in *ACL*, 2019, pp. 2561–2571.

[7] N. Vo and K. Lee, "Hierarchical multi-head attentive network for evidence-aware fake news detection," in *EACL*, 2021, pp. 965–975.

[8] L. Wu, Y. Rao, Y. Lan, L. Sun, and Z. Qi, "Unified dual-view cognitive model for interpretable claim verification," *ArXiv*, 2021.

[9] L. Yao, C. Mao, and Y. Luo, "Graph convolutional networks for text classification," *ArXiv*, vol. abs/1809.05679, 2019.

[10] Y. Zhang, X. Yu, Z. Cui, S. Wu, Z. Wen, and L. Wang, "Every document owns its structure: Inductive text classification via graph neural networks," *ArXiv*, vol. abs/2004.13826, 2020.

[11] S. Wu, Y. Tang, Y. Zhu, L. Wang, X. Xie, and T. Tan, "Session-based recommendation with graph neural networks," in *AAAI*, 2019.

[12] T. Chen, H. Shi, S. Tang, Z. Chen, F. Wu, and Y. Zhuang, "Cil: Contrastive instance learning framework for distantly supervised relation extraction," in *ACL*, 2021, pp. 6191–6200.

[13] D. Zhang, F. Nan, X. Wei, S.-W. Li, H. Zhu, K. Mckeown, R. Nallapati, A. O. Arnold, and B. Xiang, "Supporting clustering with contrastive learning," in *NAACL*, 2021, pp. 5419–5430.

[14] T. Miyato, A. M. Dai, and I. Goodfellow, "Adversarial training methods for semi-supervised text classification," *ArXiv*, 2016.

[15] D. Li, T. Zhang, N. Hu, C. Wang, and X. He, "Hiclre: A hierarchical contrastive learning framework for distantly supervised relation extraction," in *ACL findings*, 2022, pp. 2567–2578.

[16] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *NIPS*, 2016.

[17] T. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *ArXiv*, vol. abs/1609.02907, 2017.

[18] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Lio', and Y. Bengio, "Graph attention networks," *ArXiv*, vol. abs/1710.10903, 2018.

[19] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *NIPS*, 2017.

[20] T. Chen and R. C.-W. Wong, "Handling information loss of graph neural networks for session-based recommendation," in *KDD*, 2020.

[21] M. Zhang, S. Wu, M. Gao, X. Jiang, K. Xu, and L. Wang, "Personalized graph neural networks with attention mechanism for session-aware recommendation," *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 2020.

[22] J. Zhang, Y. Zhu, Q. Liu, S. Wu, S. Wang, and L. Wang, "Mining latent structures for multimedia recommendation," in *ACM Multimedia*, 2021.

[23] K. Wang, W. Shen, Y. Yang, X. Quan, and R. Wang, "Relational graph attention network for aspect-based sentiment analysis," in *ACL*, 2020.

[24] R. Li, H. Chen, F. Feng, Z. Ma, X. Wang, and E. H. Hovy, "Dual graph convolutional networks for aspect-based sentiment analysis," in *ACL/IJCNLP*, 2021.

[25] W. Jin, Y. Ma, X. Liu, X. Tang, S. Wang, and J. Tang, "Graph structure learning for robust graph neural networks," *KDD*, 2020.

[26] Y. Zhu, W. Xu, J. Zhang, Q. Liu, S. Wu, and L. Wang, "Deep graph structure learning for robust representations: A survey," *CoRR*, 2021.

[27] B. Jiang, Z. Zhang, D. Lin, J. Tang, and B. Luo, "Semi-supervised learning with graph learning-convolutional networks," in *CVPR*, 2019, pp. 11 305–11 312.

[28] Y. Chen, L. Wu, and M. Zaki, "Iterative deep graph learning for graph neural networks: Better and robust node embeddings," in *NIPS*, 2020, pp. 19 314–19 326.

[29] L. Cosmo, A. Kazi, S. Ahmadi, N. Navab, and M. M. Bronstein, "Latent patient network learning for automatic diagnosis," 2020.

[30] R. Li, S. Wang, F. Zhu, and J. Huang, "Adaptive graph convolutional neural networks," *ArXiv*, vol. abs/1801.03226, 2018.

[31] X.-W. Wu, L. Zhao, and L. Akoglu, "A quest for structure: Jointly learning the graph structure and semi-supervised classification," *CIKM*, 2018.

[32] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil, "Bilevel programming for hyperparameter optimization and meta-learning," in *ICML*, 2018, pp. 1568–1577.

[33] L. Franceschi, M. Niepert, M. Pontil, and X. He, "Learning discrete structures for graph neural networks," in *ICML*, 2019, pp. 1972–1982.

[34] Y. Zhang, S. Pal, M. J. Coates, and D. Üstebay, "Bayesian graph convolutional neural networks for semi-supervised classification," in *AAAI*, 2019.

[35] L. Yang, Z. Kang, X. Cao, D. Jin, B. Yang, and Y. Guo, "Topology optimization based graph convolutional network," in *IJCAI*, 2019.

[36] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec, "Hierarchical graph representation learning with differentiable pooling," in *NIPS*, 2018, pp. 4800–4810.

[37] H. Gao and S. Ji, "Graph u-nets," in *ICML*, 2019.

[38] J. Lee, I. Lee, and J. Kang, "Self-attention graph pooling," in *ICML*, 2019, pp. 6661–6670.

[39] J. Zhou, X. Han, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun, "Gear: Graph-based evidence aggregating and reasoning for fact verification," in *ACL*, 2019.

[40] Z. Liu, C. Xiong, M. Sun, and Z. Liu, "Fine-grained fact verification with kernel graph attention network," in *ACL*, 2020.

[41] W. Zhong, J. Xu, D. Tang, Z. Xu, N. Duan, M. Zhou, J. Wang, and J. Yin, "Reasoning over semantic-level graph for fact checking," in *ACL*, 2020.

[42] K. Popat, S. Mukherjee, J. Strötgen, and G. Weikum, "Credibility assessment of textual claims on the web," in *CIKM*, 2016, p. 2173–2178.

[43] F. Yu, Q. Liu, S. Wu, L. Wang, and T. Tan, "A convolutional approach for misinformation identification," in *IJCAI*, 2017.

[44] S. Volkova, K. Shaffer, J. Y. Jang, and N. Hodas, "Separating facts from fiction: Linguistic models to classify suspicious and trusted news posts on Twitter," in *ACL*, 2017, pp. 647–653.

[45] N. Vo and K. Lee, "The rise of guardians: Fact-checking url recommendation to combat fake news," in *SIGIR*, 2018, p. 275–284.

[46] A. Benamira, B. Devillers, E. Lesot, A. Ray, M. Saadi, and F. D. Malliaros, "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)*, pp. 568–569, 2019.

[47] S. Chandra, P. Mishra, H. Yannakoudakis, and E. Shutova, "Graph-based modeling of online communities for fake news detection," *ArXiv*, vol. abs/2008.06274, 2020.

[48] Y. Jin, X. Wang, R. Yang, Y. Sun, W. Wang, H. Liao, and X. Xie, "Towards fine-grained reasoning for fake news detection," *ArXiv*, vol. abs/2110.15064, 2021.

[49] Q. Liu, F. Yu, S. Wu, and L. Wang, "Mining significant microblogs for misinformation identification: an attention-based approach," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 9, no. 5, pp. 1–20, 2018.

[50] O. Ajao, D. Bhowmik, and S. Zargari, "Sentiment aware fake news detection on online social networks," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 2507–2511.

[51] A. Giachanou, P. Rosso, and F. Crestani, "Leveraging emotional signals for credibility detection," in *SIGIR*, 2019, p. 877–880.

[52] X. Zhang, J. Cao, X. Li, Q. Sheng, L. Zhong, and K. Shu, "Mining dual emotion for fake news detection," ser. WWW, 2021, p. 3465–3476.

[53] Y. Zhu, Q. Sheng, J. Cao, Q. Nan, K. Shu, M. Wu, J. Wang, and F. Zhuang, "Memory-guided multi-view multi-domain fake news detection," *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 2022.

[54] A. Vlachos and S. Riedel, "Identification and verification of simple claims about statistical properties," in *EMNLP*, 2015, pp. 2596–2601.

[55] ——, "Fact checking: Task definition and dataset construction," in *Proceedings of the ACL 2014 Workshop on Language Technologies and Computational Social Science*, 2014, pp. 18–22.

[56] L. Wu, Y. Rao, X. Yang, W. Wang, and A. Nazir, "Evidence-aware hierarchical interactive attention networks for explainable claim verification." in *IJCAI*, 2020, pp. 1388–1394.

[57] L. Wu, Y. Rao, L. Sun, and W. He, "Evidence inference networks for interpretable claim verification," *AAAI*, pp. 14 058–14 066, 2021.

[58] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *ArXiv*, vol. abs/1807.03748, 2018.

[59] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*. PMLR, 2020, pp. 1597–1607.

[60] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *CVPR*, 2020, pp. 9729–9738.

[61] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," in *NIPS*, vol. 33, 2020, pp. 18 661–18 673.

[62] P. Bachman, R. D. Hjelm, and W. Buchwalter, "Learning representations by maximizing mutual information across views," in *NIPS*, vol. 32, 2019.

[63] P. Velickovic, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm, "Deep graph infomax." in *ICLR*, vol. 2, no. 3, 2019, p. 4.

[64] Y. You, T. Chen, Y. Sui, T. Chen, Z. Wang, and Y. Shen, "Graph contrastive learning with augmentations," in *NIPS*, vol. 33, 2020, pp. 5812–5823.

[65] Y. Zhu, Y. Xu, F. Yu, Q. Liu, S. Wu, and L. Wang, "Deep graph contrastive representation learning," in *GRL+@ICML*, 2020.

[66] ——, "Graph contrastive learning with adaptive augmentation," in *WWW*, 2021.

[67] T. Gao, X. Yao, and D. Chen, "Simcse: Simple contrastive learning of sentence embeddings," in *EMNLP*, 2021, pp. 6894–6910.

[68] H. Lin, J. Ma, L. Chen, Z. Yang, M. Cheng, and G. Chen, "Detect rumors in microblog posts for low-resource domains via adversarial contrastive learning," *ArXiv*, 2022.

[69] Z. Yue, H. Zeng, Z. Kou, L. Shang, and D. Wang, "Contrastive domain adaptation for early misinformation detection: A case study on covid-19," *ArXiv*, 2022.

[70] Y. Zhang, J. Zhang, Z. Cui, S. Wu, and L. Wang, "A graph-based relevance matching model for ad-hoc retrieval," in *AAAI*, 2021.

[71] X. Yu, W. Xu, Z. Cui, S. Wu, and L. Wang, "Graph-based hierarchical relevance matching signals for ad-hoc retrieval," in *WWW*, 2021.

[72] C. Xiong, Z. Dai, J. Callan, Z. Liu, and R. Power, "End-to-end neural ad-hoc ranking with kernel pooling," in *SIGIR*, 2017, pp. 55–64.

[73] K. Popat, S. Mukherjee, J. Strötgen, and G. Weikum, "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.

[74] H. Rashkin, E. Choi, J. Y. Jang, S. Volkova, and Y. Choi, "Truth of varying shades: Analyzing language in fake news and political fact-checking," in *EMNLP*, 2017.

[75] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[76] W. Y. Wang, ""liar, liar pants on fire": A new benchmark dataset for fake news detection," in *ACL*, 2017.

[77] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *NAACL*, 2019.

[78] Q. Li, Z. Han, and X.-M. Wu, "Deeper insights into graph convolutional networks for semi-supervised learning," *ArXiv*, vol. abs/1801.07606, 2018.

[79] J. Zhang, Y. Zhu, Q. Liu, M. Zhang, S. Wu, and L. Wang, "Latent structures mining with contrastive modality fusion for multimedia recommendation," *Arxiv*, 2021.

[80] J. Yu, H. Yin, J. Li, Q. Wang, N. Q. V. Hung, and X. Zhang, "Self-supervised multi-channel hypergraph convolutional network for social recommendation," in *WWW*, 2021, pp. 413–424.
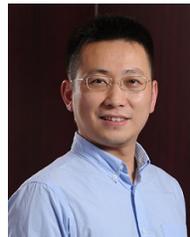
**Weizhi Xu** is currently pursuing his master's degree of Computer Science at the Center for Research on Intelligent Perception and Computing (CRIPAC) at National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA). His current research interests mainly include graph representation learning, fake news detection and recommender systems.

**Qiang Liu** is an Associate Professor with the Center for Research on Intelligent Perception and Computing (CRIPAC), Institute of Automation, Chinese Academy of Sciences (CASIA). He received his PhD degree from CASIA. Currently, his research interests include data mining, recommender systems, text mining, knowledge graph, and graph representation learning. He has published papers in top-tier journals and conferences, such as IEEE TKDE, AAAI, IJCAI, NeurIPS, WWW, SIGIR, CIKM and ICDM.

**Shu Wu** received his B.S. degree from Hunan University, China, in 2004, M.S. degree from Xiamen University, China, in 2007, and Ph.D. degree from Department of Computer Science, University of Sherbrooke, Quebec, Canada, all in computer science. He is an Associate Professor with the Center for Research on Intelligent Perception and Computing (CRIPAC) at National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA). He has published more than 50 papers in the areas of data mining and information retrieval in international journals and conferences, such as IEEE TKDE, IEEE THMS, AAAI, ICDM, SIGIR, and CIKM. His research interests include data mining, information retrieval, and recommendation systems.

**Junfei Wu** is currently pursuing his Ph.D. degree of Computer Science at the Center for Research on Intelligent Perception and Computing (CRIPAC) at National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA). His current research interests mainly include fake news detection.

**Liang Wang** received both the BEng and MEng degrees from Anhui University in 1997 and 2000, respectively, and the PhD degree from the Institute of Automation, Chinese Academy of Sciences (CASIA) in 2004. From 2004 to 2010, he was a research assistant at Imperial College London, United Kingdom, and Monash University, Australia, a research fellow at the University of Melbourne, Australia, and a lecturer at the University of Bath, United Kingdom, respectively. Currently, he is a full professor of the Hundred Talents Program at the National Lab of Pattern Recognition, CASIA. His major research interests include machine learning, pattern recognition, and computer vision. He has widely published in highly ranked international journals such as IEEE TPAMI and IEEE TIP, and leading international conferences such as CVPR, ICCV, and ECCV. He has served as an Associate Editor of IEEE TPAMI, IEEE TIP, and PR. He is an IEEE Fellow and an IAPR Fellow.