

Exploring Faithful Rationale for Multi-hop Fact Verification via Saliency-Aware Graph Learning

Jiasheng Si, Yingjie Zhu, Deyu Zhou*

School of Computer Science and Engineering, Key Laboratory of Computer Network and Information Integration, Ministry of Education, Southeast University, China
 {jasenchn, yj_zhu, d.zhou}@seu.edu.cn

Abstract

The opaqueness of the multi-hop fact verification model imposes imperative requirements for explainability. One feasible way is to extract *rationales*, a subset of inputs, where the performance of prediction drops dramatically when being removed. Though being explainable, most rationale extraction methods for multi-hop fact verification explore the semantic information within each piece of evidence individually, while ignoring the topological information interaction among different pieces of evidence. Intuitively, a faithful rationale bears complementary information being able to extract other rationales through the multi-hop reasoning process. To tackle such disadvantages, we cast explainable multi-hop fact verification as subgraph extraction, which can be solved based on graph convolutional network (GCN) with saliency-aware graph learning. In specific, GCN is utilized to incorporate the topological interaction information among multiple pieces of evidence for learning evidence representation. Meanwhile, to alleviate the influence of noisy evidence, the saliency-aware graph perturbation is induced into the message passing of GCN. Moreover, the multi-task model with three diagnostic properties of rationale is elaborately designed to improve the quality of an explanation without any explicit annotations. Experimental results on the FEVEROUS benchmark show significant gains over previous state-of-the-art methods for both rationale extraction and fact verification.

1 Introduction

The wide availability of user-provided content on online social media facilitates the rapid dissemination of unfounded rumors and misinformation. Fact verification, automatically assessing the veracity of a textual claim against multiple pieces of evidence retrieved from external sources, has gained intense attention to combat misinformation spread on the internet (Zhou and Zafarani 2020; Guo, Schlichtkrull, and Vlachos 2022; Si et al. 2021). However, as the opaqueness of the model diminishes user confidence and impedes the discovery of harmful biases (Kotonya and Toni 2020a), it is essential to understand the “reasoning” behind the model prediction, i.e., the explainability of the fact verification approaches.

*Corresponding author.

Claim: Olympic athlete **May Wafic Sardouk** represented **Lebanon at the 1988 Summer Olympics** in **Seoul, Korea**, landing in the **6th position in the Heat 4 event**.

Evidence:

S₁(wiki/May_Sardouk): *May Wafic Sardouk* (Arabic: *وفيق مي ساردوك*; born June 4, 1963) is a *Lebanese* Olympic athlete.

S₂(wiki/May_Sardouk): She *represented Lebanon in 1988 Summer Olympics in Seoul*.

S₃(wiki/May_Sardouk): Sardouk and Nancy Khalaf were the only female participants for Lebanon in that tournament among a total of 21 participant for Lebanon.

S₄(wiki/Seoul): *Seoul*, officially the Seoul Special City, is the capital and largest metropolis of *South Korea*.

S₅(wiki/1988 Summer Olympics): The 1988 Summer Olympics, ..., was an international multi-sport event held from 17 September to 2 October 1988 in Seoul, South Korea.

T₆(wiki/May_Sardouk):

Heat 4		
Rank	Athlete	Time
1	Diane Dixon (USA)	52.45
2	Ute Thimm (FRG)	52.79
...
6	May Sardouk (LIB)	1:00.01

Label: SUPPORTS

Figure 1: An example from FEVEROUS dataset, where S_1 , S_2 , S_4 and two table cells in T_6 are considered as rationales.

A straightforward way to generate explanations for fact verification is to use **rationale extraction** (Zaidan, Eisner, and Piatko 2007; DeYoung et al. 2020; Atanasova et al. 2020; Glockner, Habernal, and Gurevych 2020; Atanasova et al. 2022), a *post-hoc technique* searching for a minimal portion of input (i.e., rationales) that can be sufficient (i.e., solely based on the rationales) to derive a veracity prediction. The intuition is that the retrieved evidence for verifying the claim comprises noisy evidence inevitably, and the interaction of true evidence¹ is adequate for the multi-hop fact verification model to reach a disposition accurately. It

¹Note that we make a distinction between *true evidence* and *noisy evidence* conceptually, and define the *true evidence* as the *rationale* since the term “rationale” implies human-like intent (Wiegrefe and Marasovic 2021).

contrasts with the methods heuristically exploring the importance of input features, such as attention-based methods (Chen et al. 2022; Wu et al. 2021) or gradient-based methods (Sundararajan, Taly, and Yan 2017), which inevitably induce low-scoring features, drawing criticism recently (Jain and Wallace 2019). In this work, we focus on how to extract valid rationales for explaining multi-hop fact verification model.

Existing rationale extraction methods for multi-hop fact verification model usually rely on the FEVER dataset (Thorne et al. 2018; DeYoung et al. 2020; Bekoulis, Papagiannopoulou, and Deligiannis 2021). These methods typically decompose the model into the extractor module and the predictor module independently (Lei, Barzilay, and Jaakkola 2016; Jain et al. 2020; Kotonya and Toni 2020b), where the former is trained to assign a mask score over the subset (e.g., sentences or tokens) of inputs to capture the rationale, then the latter makes predictions exclusively on the rationales provided by the extractor. The quality of the explanation depends upon the strategy of the mask vector training.

However, despite the salient progress, there are still limitations required to be addressed for explainable multi-hop fact verification (Jiang et al. 2020; Ostrowski et al. 2021; Aly et al. 2021; Atanasova et al. 2020; Jain et al. 2020; Atanasova et al. 2022). Inherently, in multi-hop fact verification, the claims may be verified by aggregating and reasoning over multiple pieces of rationales. For example in Fig.1, the truthfulness of the claim can be assessed by aggregating four pieces of true evidence (i.e., rationales), including sentences and table cells, surrounded by multiple pieces of noisy evidence. Intuitively, the rationale $S1$ carries the complementary information capable of extracting the rationale $S2$ through the multi-hop reasoning process. However, in the mask vector learning process, current rationale extraction methods are mainly based on the semantic information of individual semantic units (sentence or token) within the input, failing to capture the topological information interaction among multiple pieces of the semantic unit in the multi-hop reasoning process for rationale extraction, which we argue is crucial for the explainable multi-hop fact verification.

To address such disadvantage, we introduce a GCN-based model (Kipf and Welling 2017) with **S**aliency-aware **G**raph **P**erturbation, namely SaGP, where *multi-hop fact verification* and *sentence-level rationale extraction* are optimized jointly. The **core novelty** here is that we frame the rationale extraction of multi-hop fact verification as *subgraph extraction* via searching for the rationale subgraph with minimal nodes while sufficiently maintaining the prediction accuracy. Specifically, we use GCN to integrate the topological interaction of information among different pieces of evidence to update evidence representation. Meanwhile, to alleviate the influence of the noisy evidence in this process, we induce a learnable saliency-aware perturbation (edge mask, node mask) into the message passing process of GCN to approximate the deletion of the superfluous edges or nodes in the input graph. It guarantees that the information masked out from the graph is not propagated for evidence representation learning. Then the assignment vector over each

node is learned to indicate whether the evidence could be contained in the rationale subgraph, which approximates the mask vector learning following prior works. Moreover, we incorporate the multi-task learning paradigm and define three diagnostic properties (i.e., *Fidelity*, *Compact*, *Topology*) as additional optimizing signals to guide the learning of rationale subgraph.

The main contributions are listed as follows: (I) We frame the explainable multi-hop fact verification as subgraph extraction, where a GCN-based model with saliency-aware graph learning is proposed. (II) The multi-task model with three diagnostic properties is designed and optimized to improve the quality of extracted explanations without accessing the rationale supervision. (III) Experimental results on the FEVEROUS dataset show the superior performance of the proposed approach.

2 Related Works

Fact verification is the task of assessing the veracity of the claim backed by multiple pieces of validated evidence, which can be decomposed into two stages: evidence retrieval and claim verification (Si et al. 2021). Thus, two aspects can be explored for explaining the fact verification model: (I) retrieving faithful evidence as precisely as possible in the evidence retrieval stage (Wan et al. 2021); (II) extracting rationales with the explainable techniques in the claim verification stage. Researchers mainly focus on the latter by developing explainability techniques for fact verification, which broadly fall into the scope of attention-based methods by treating the attention weights over input representation as a measure of credibility score (Kotonya and Toni 2020a; Luo et al. 2021). Attention mechanisms vary in the information type of input, including self-attention (Popat et al. 2018), co-attention (Shu et al. 2019; Wu et al. 2021; Yang et al. 2019; Wu et al. 2020). It provides an optional way for multi-hop fact verification, while recently some works argued that attention weights cannot guarantee the inattention of low-confidence features, and are not valid explanations for model prediction (Jain and Wallace 2019; Serrano and Smith 2019; Meister et al. 2021).

Other researches focus on exploring the rationale extraction for explaining the fact verification via perturbing the input (Luo et al. 2021; Kotonya and Toni 2020a), which usually consists of two modules, i.e., extractor and predictor (Atanasova et al. 2020; Paranjape et al. 2020; Sha, Camburu, and Lukasiewicz 2021). Generally, the predictor is used to devise the decision based on the rationales generated from the extractor rather than the whole input. However, it is tricky to jointly train the two separate modules because of the intractable rationale sampling. It is partially solved via reinforcement learning (Lei, Barzilay, and Jaakkola 2016; Yoon, Jordon, and van der Schaar 2019), or reparameterization techniques (Bastings, Aziz, and Titov 2019). An alternative way is to adopt multi-task learning (Atanasova et al. 2022; Wiegrefe and Marasovic 2021), which provides more label-informed rationales than prior methods. However, they mainly extract the rationales with the semantic information of individual semantic units of the input, which is not able to

capture the topological information interaction among different pieces of evidence in the mask vector learning process. It may not be suitable to explain multi-hop fact verification (Glockner, Habernal, and Gurevych 2020).

Different from the approaches mentioned above, our work is the first to formulate the explainable multi-hop fact verification as subgraph extraction. Moreover, we utilize three properties to guide the extraction of the rationale subgraph in multi-task learning.

3 Methodology

3.1 Problem Setting

Assume that a claim may be verified with multi-hop evidence, including sentences, table cells, or a combination of multiple sentences or table cells (Aly et al. 2021). Given a claim c with associated textual evidence (i.e., sentences, table captions) $\{t_1, t_2, \dots, t_S\}$ or tabular evidence (i.e., table cells) $\{c_1, c_2, \dots, c_C\}$, a heterogeneous evidence graph $G = (X, A)$ is constructed to model how the claim is associated with the textual evidence and the tabular evidence. Each node $x_i \in X$ represents an evidence sequence by concatenating the claim and the textual evidence (i.e., sentences or table captions) or tabular evidence (i.e., cell sequences). A denotes an adjacency matrix for the undirected fully connected graph with the edge weight equal to 1. Given a trained GNN-based multi-hop fact verification model, we aim to extract the rationale subgraph (*RA-subgraph*) for this model by incorporating the topological information interaction among different pieces of evidence. Nodes within the extracted subgraph are treated as the rationale. Therefore, the aim of explainable multi-hop fact verification is to infer the claim verification label \hat{y}^c as *SUPPORTS* or *REFUTES* and to assign each sentence or table cell with a Boolean label denoted as $\hat{y}_i^c \in \{0, 1\}$, where $\hat{y}_i^c = 1$ denotes that the sequence i is in the RA-subgraph which actually benefits the predictions.

3.2 The Architecture

We propose the method **SaGP**, which consists of three key components as shown in Fig. 2: (I) an embedding layer, where a pre-trained language model is employed to obtain the initial node representations of the input graph. (II) a graph perturbation layer, where evidence representation is updated by inducing the salience-aware graph perturbation into the message passing of GCN. (III) a rationale extraction layer, where the RA-subgraph is extracted by directly optimizing for three diagnostic properties (i.e., Fidelity, Compact, Topology) of rationales in the multi-task model. The details of each component are provided in the following sections.

3.3 Embedding Layer

It is challenging to manipulate the tabular evidence at cell-level. Following Chen et al. (2020); Kotonya et al. (2021), we employ the simple table linearization template to generate contextualized per-cell sequence representations to form the cell sequence of the table, where each table cell is linearized as “*In wikipedia, the header is tableheader1 (and tableheader2), the value is tablecell.*”. Different from

TAPAS (Herzig et al. 2020), we do not consider the table structure since there are some noisy cells contained in the table, which might confuse the model. It should be pointed out that in our experiments, there is no substantial improvement when employing elaborated templates such as Kotonya et al. (2021).

For the initial node representation in the input evidence graph, a RoBERTa (Liu et al. 2019) model is employed to encode each evidence sequence x_i to its contextual embedding $h_i = \text{encoder}(x_i)$ by selecting the hidden representation of the *CLS* token from the final layer.

3.4 Graph Perturbation Layer

To identify rationales by incorporating the topological interaction of different pieces of evidence while mitigating the influence of noisy evidence, we induce a learnable salience-aware graph perturbation into the messaging passing of GCN. Specifically, given a trained GCN-based model, taking an input graph as input, we probe the GCN layer using a learnable salience-aware graph perturbation matrix (a.k.a., edge mask or node mask) (Ying et al. 2019), where the gates in the perturbation matrix indicate which edges or nodes are necessary and which can be disregarded. The assumption behind this is that the nodes within the neighborhood might contain noise or even conflicting information, while the *zero* edges or nodes bear no information and can be removed from the message aggregation for node updating, which is analogous to the perturbation of input. The masking may also be equivalently seen as adding a certain type of noise to the message passing process in the GCN.

Let f be a trained GCN layer for node representation learning,

$$f(A, H; W) = \text{relu}(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} H W), \quad (1)$$

where $\tilde{A} = A + I$, I is the identity matrix, \tilde{D} is the degree matrix, H denotes the evidence embeddings and W denotes the parameters of GCN.

Given the trained GCN layer with parameters W frozen, we seek to introduce the learnable parameters to mitigate the impact of noise evidence in the GCN message passing process. For the edge mask, we introduce a learnable perturbation matrix P with the same size as the adjacency matrix A to approximate zero out the superfluous entries in the adjacency matrix. Each element p_{ij} in matrix P indicates the importance score for message aggregation from node i to node j , where the sigmoid transformation is utilized to restrict the matrix with entries in $[0, 1]$. Specifically, if element $p_{ij} \rightarrow 0$, it results in the deletion of the edge from node i to node j . We populate P in an asymmetrical manner since we argue the information discrepancy for different nodes within a node pair. In other words, the node i is useful for updating the node j cannot guarantee that the opposite is valid. The calculation in the GCN can be rewritten into

$$\tilde{f}(A, H; W) = \text{relu}(((\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2}) \odot \sigma(P)) H W), \quad (2)$$

where \odot denotes the element product.

For node mask, similarly, we introduce a learnable perturbation matrix M to assign the attribution score for each node

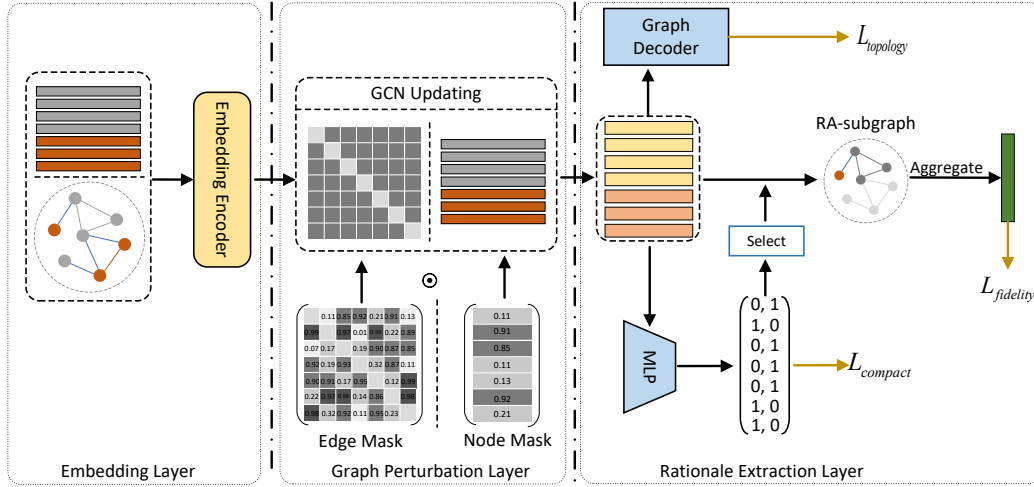


Figure 2: The overall framework of the proposed SaGP.

with the sigmoid transformation, which indicates the degree of node features used in the message passing in the graph.

$$\tilde{f}(A, H; W) = \text{relu}(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} (H \odot \sigma(M)) W) \quad (3)$$

There are no supervision signals on the graph to update the perturbation matrix of the probe. Therefore, two regulations are applied to avoid all entries in the matrix approximating 1: (I) the sum regulation \mathcal{L}_{sum} of all entries in the perturbation matrix to constrain the size of the perturbation; (II) the information entropy regulation $\mathcal{L}_{entropy}$ to reduce the uncertainty of the perturbation matrix.

3.5 Rationale Extraction Layer

Connecting the explanation extractor module and predictor module with the learnable mask vector is indispensable for valid rationale extraction. Previous methods generally devise complicated strategies for linking the two modules (e.g., REINFORCE-based (Lei, Barzilay, and Jaakkola 2016), heuristic-based (Jain et al. 2020)), which is inefficient. Following Atanasova et al. (2022), we employ the multi-task learning paradigm in which explainable multi-hop fact verification is factorized into first extracting the RA-subgraph from the input graph and then conducting the claim verification conditioned on RA-subgraph. Moreover, three diagnostic properties of the rationale subgraph are optimized jointly.

Rationale Subgraph Extraction Given the node representation $\tilde{U} \in \mathbb{R}^{N \times d}$ learned by the **perturbed GCN** \tilde{f} (in Sec. 3.4), we first extract the subgraph via the node assignment $S \in \mathbb{R}^{N \times 2}$, where the first dimension indicates the probability of nodes is the rationale and could be included in the rationale subgraph G_{sub} (Yu et al. 2021), and the second dimension denotes the non-rationale and the complementary subgraph of G_{sub} in input graph, namely \bar{G}_{sub} .

$$S = \text{softmax}(MLP(\tilde{U}; W_{sub})), \quad (4)$$

where MLP is a linear layer, W_{sub} is the learnable parameters. Then, the representation R_{sub} of the RA-subgraph G_{sub} can be aggregated by taking the first column of S , i.e., $R_{sub} = S_1^T \tilde{U}$. Consequently, a linear layer with softmax operation on R_{sub} predicts the target label $\hat{y}_{sub}^c \in \mathbb{R}^2$.

Property We propose three diagnostic properties to regularize the rationale subgraph extraction.

Fidelity. The faithful RA-subgraph extracted from an input graph ought to exhibit meaningful influence on its prediction. Thus, we define the cross entropy function to measure the *fidelity* for faithfulness of RA-subgraph (DeYoung et al. 2020). In particular, given the node representation U learned from the **standard GCN** f , a linear layer with softmax operation is employed to map the representation of the full graph to its prediction \hat{y}_{full}^c by aggregating U . We assume that the rationales within the subgraph are adequate to match the original prediction.

$$L_{fidelity} = \text{CrossEntropy}(\hat{y}_{sub}^c(\tilde{U}), \hat{y}_{full}^c(U)), \quad (5)$$

intuitively, a low score here implies that the rationale contained in the subgraph is indeed susceptible on the prediction.

Compact. The compact measures how distinctive the node assignment within the graph is. Including it as an objective can serve as an additional regularization on S for the model to be compact in the extracted subgraph. The poor perturbation, on the one hand, will enable $p(x_i) \in G_{sub}$ and $p(x_i) \in \bar{G}_{sub}$ to be close. Additionally, the model is prone to assigning all the nodes to $p(x_i) \in G_{sub}$, resulting in redundancy for the RA-subgraph. To do so, we employed the compact loss to distinguish node assignments and urge the RA-subgraph to be compact,

$$L_{compact} = \|\text{norm}(S^T A S) - I_2\|_F, \quad (6)$$

where $\text{norm}(\cdot)$ denotes the normalization, $\|\cdot\|_F$ denotes the Frobenius norm, I_2 is the identity matrix.

Topology. The topology measures whether the representation learned from the perturbed GCN \tilde{f} preserves the original topological information. We consider this a valuable objective to re-calibrate due to the asymmetrical information passing in the perturbed graph. Motivated by graph auto-encoder (Kipf and Welling 2016), we define a decoder to reconstruct the original graph adjacency A using the perturbed node representations \tilde{U} ,

$$\begin{aligned} L_{topology} &= CrossEntropy(\hat{A}, A), \\ \hat{A} &= \sigma(\tilde{U}\tilde{U}^T) \end{aligned} \quad (7)$$

where σ denotes the sigmoid function.

3.6 Objective for Learning

The overall objective function \mathcal{L} is minimized over the above modules:

$$\begin{aligned} \mathcal{L} &= \lambda_1 \mathcal{L}_{fidelity} + \lambda_2 \mathcal{L}_{compact} + \lambda_3 \mathcal{L}_{topology} \\ &+ \lambda_4 \mathcal{L}_{sum} + \lambda_5 \mathcal{L}_{entropy}, \end{aligned} \quad (8)$$

where λ_{1-5} are hyperparameters. In addition, we elucidate how the rationale supervision can be used for rationale extraction, where the $\mathcal{L}_{compact}$ is replaced by the rationale loss, which treats S as a standard multi-label problem via a sigmoid layer and binary cross entropy function. This is a compromised way to use supervision at the node-level, as it is difficult to obtain the real edge tag.

4 Experimental Setup

This section describes the datasets, evaluation metrics and the baselines in the experiments.

Dataset. We conduct our experiments on the large scale dataset FEVEROUS (Aly et al. 2021), which is a multi-hop dataset with different types of evidence. To construct the explanatory dataset with rationale annotation, the claims with the *NOT ENOUGH INFO* label are deleted because there are no gold standard rationales corresponding to that label. Moreover, apart from retaining the rationales in the dataset, we retrieved the complementary evidence relevant to the claim as the *noisy evidence* using the method described in Aly et al. (2021). We construct the dataset where each claim associates with 20 pieces of evidence, including *true evidence* (i.e., rationale) and *noisy evidence*, wherein the verdict of the claim requires reasoning over the aggregation of the sentences and table cells (1.43 sentences and 3.42 table cells in the training dataset, 1.43 sentences and 2.83 table cells in the test dataset on average). The statistics of the dataset are shown in Tab.1.

Metrics. We first adopt the metrics proposed for the ERASER benchmark (DeYoung et al. 2020) to measure the agreement with the human-annotated rationales by evaluating the macro F1, Precision, and Recall metrics, where we choose the sentences as the basic unit for rationales. We also report the Exact Match Accuracy (Ext.acc) for strict measuring rationale comparison. Secondly, We adopt the macro F1 and Accuracy metrics for claim verification evaluation. In addition, we report the joint accuracy of the claim verification and the rationale extraction, where we consider the

FEVEROUS	Num.Supp	Num.Ref	Avg.Ra	Avg.S	Avg.C
Train	41,835	27,215	4.85	1.43	3.42
Test	3,908	3,481	4.26	1.43	2.83

Table 1: Statistics of the FEVEROUS dataset. *Num.Supp* and *Num.Ref* are the number of claims with *SUPPORT* label and *REFUTE* label. *Avg.Ra*, *Avg.S*, and *Avg.C* denote the average number of *rationales*, *sentence rationales*, *table cell rationales* per claim, respectively.

prediction is correct if the predictions are correct for the two tasks, where Acc.Full denotes the correct prediction of the claim with all rationales, Acc.Part denotes the correct prediction of the claim with one piece of rationale.

Baselines. We compare the proposed SaGP model with the following baselines for claim verification and rationale extraction under unsupervised and supervised settings, including (I) the pipeline method from ERASER (DeYoung et al. 2020), which verdicts the claim using one rationale with the highest score from the extractor. (II) the information bottleneck (IB) method (Paranjape et al. 2020), which extracts sentence-level rationales by measuring the mutual information with the label. (III) the two-sentence selecting (TSS) method (Glockner, Habernal, and Gurevych 2020), which extracts the rationales by utilizing the loss logits of the rationales. (IV) the DeClarE (Popat et al. 2018) and Transformer-XH (Zhao et al. 2020), which extract the rationales with the attention score.

Setup. Our trained GNN-based model adopts the base version of the pre-trained RoBERTa model followed by the GCN with 2 layers, where each node corresponds to the concatenation of the claim sequence and the textual sequence or cell sequence. We also insert the *WikiTitle* into the two sequences as the bridge information. We take the *CLS* token representation as the initial node representations. The maximum number of input tokens to RoBERTa is 140. The original model has 85.6% on label accuracy of claim verification. The hyperparameters λ_1 , λ_2 , λ_3 are set to 1, 1, 1 respectively, and λ_4 , λ_5 are defined as $5e-3$, 0.1 for edge mask, and 0.1, 1 for node mask. We adopt the instance-level explanation, where each instance is trained 100 epochs, with the learning rate being settled to $1e-2$. For comparison with baselines, we choose the attention score with the threshold of 0.5 for attention-based methods (DeClarE and Transformer-XH), the threshold is set to 0.2 for IB since the portion of the gold standard with input is nearly 0.2, we choose the first two sentences for TSS baseline because of the expensive cost of computing. All experimental setups of the baselines are followed from the original papers.

5 Experimental Results

In this section, we evaluate the SaGP model in different aspects. Firstly, we compare the overall performance with the baselines using different types of masks under unsupervised and supervised settings. Then, we evaluate the effect of using different diagnostic properties as additional signals for rationale extraction. Finally, following (Lucic et al. 2022), we explore various desirable properties of the edge mask.

Model	Claim		Rationale			Claim & Rationale			
	F1.c	Acc.c	F1.r	Ext.acc.r	P.r	R.r	Acc.Part	Acc.Full	
Unsupervised									
TSS-U	34.61	52.93	18.75	16.83	36.57	14.59	23.77	1.13	
DeClarE	68.23	69.18	27.59	13.63	31.46	31.71	43.85	9.81	
IB-U	77.30	77.30	65.28	20.08	78.01	67.30	75.36	15.76	
Edge	SaGP	85.05±0.02	85.15±0.02	80.08±0.01	45.33±0.05	79.15±0.03	88.30±0.01	82.92±0.03	41.17±0.05
	-T.	85.04±0.02	85.15±0.02	80.01±0.01	45.30±0.06	79.14±0.01	88.30±0.01	82.82±0.03	40.11±0.05
Mask	-C.	85.04±0.05	85.15±0.07	80.25±0.16	46.22±1.41	79.80±0.97	87.68±1.09	82.85±0.06	41.14±1.57
	-T.&C.	85.01±0.04	85.11±0.04	80.15±0.01	45.23±0.03	79.14±0.01	88.46±0.01	82.92±0.05	40.01±0.01
Node	SaGP	82.24±0.13	82.26±0.13	70.47±0.08	38.56±0.13	75.19±0.12	76.40±0.05	75.03±0.08	33.61±0.01
	-T.	82.25±0.12	82.25±0.12	70.50±0.09	38.60±0.10	75.19±0.12	76.37±0.07	75.04±0.06	33.65±0.04
Mask	-C.	81.80±0.19	81.81±0.19	70.34±0.26	36.97±0.55	73.60±1.05	78.28±1.72	75.36±0.54	32.18±0.56
	-T.&C.	81.85±0.15	81.85±0.16	70.17±0.12	37.50±0.18	74.27±0.05	77.01±0.18	74.78±0.22	32.64±0.04
All	SaGP	82.06±0.12	82.08±0.12	70.40±0.21	38.66±0.27	74.99±0.27	76.27±0.14	75.27±0.81	33.90±0.25
	-T.	81.77±0.11	81.78±0.11	70.14±0.20	37.40±0.36	74.23±0.21	76.95±0.16	74.67±0.15	32.66±0.33
	-C.	81.89±0.09	81.90±0.09	73.64±4.80	40.17±3.60	75.81±2.09	81.11±5.70	76.59±2.54	34.99±3.03
	-T.&C.	82.03±0.11	82.05±0.11	70.38±0.23	38.60±0.26	74.93±0.28	76.30±0.15	74.64±0.05	33.84±0.24
Supervised									
BERT Blackbox	64.72	65.20	-	-	-	-	-	-	
Pipeline	69.76	69.80	77.56	44.83	76.87	86.75	62.77	31.23	
TSS-S	72.99	74.36	44.15	19.42	85.67	34.12	67.75	11.76	
IB-S	79.14	79.17	65.68	20.08	78.91	67.31	76.70	16.37	
Transformer-XH	74.05	74.33	76.70	49.10	79.43	80.47	69.17	40.22	
Edge Mask	85.12±0.01	85.25±0.01	80.49±0.02	48.22±0.01	81.18±0.02	86.14±0.02	82.77±0.01	43.36±0.01	
SaGP Node Mask	81.53±0.06	81.54±0.06	84.50±0.66	56.23±0.23	85.51±0.06	86.48±0.02	78.10±0.11	47.67±0.29	
All	82.10±0.04	82.15±0.03	85.80±0.07	61.94±0.26	87.89±0.07	87.05±0.06	78.76±0.11	53.19±0.30	

Table 2: The performances of different approaches for claim verification and rationale extraction on the FEVEROUS dataset under two settings (mean and standard deviation over three random seed runs), where $-U$, $-S$ denote the unsupervised and supervised version of the model, respectively. $T.$, $C.$ denote the *Topology* and *Compact* properties. *All* denotes the combination of edge mask and node mask. The best results are marked in bold.

5.1 Overall Results

Unsupervised Setting We mainly explore the rationale extraction under the unsupervised setting with different diagnostic properties as the additional signals. The top of Tab. 2 reports the overall results of our model against the baselines. As shown in Tab. 2, our model, especially for the edge mask, significantly outperforms the baselines on both the claim verification task and rationale extraction task. It is worth pointing out that our model outperforms the baselines by over 20% on the Ext.acc, which demonstrates the effectiveness of the graph network for extracting rationales. However, compared with using the edge mask, the performance of using the node mask decreases on all evaluation metrics in varying degrees, and nearly 7% on ACC.Full metrics and 8% on Ext.acc. We conjecture that this might attribute to the deficiency of indispensable information in the reasoning process caused by the node mask. Inherently, the node mask can be considered as the noise signal added to the feature of rationales directly.

Effect of Property. We explore the effectiveness of us-

ing different diagnostic properties for rationale extraction by conducting the ablation study². As shown in Tab. 2, using **topology** as an additional objective aims to regularize the graph structure in the multi-hop reasoning process. We note a slight decrease when removing the topology from edge mask while no effect occurs for node mask. This can be explained by the fact that the topology focuses primarily on the structure other than the input feature. Furthermore, for **compact**, we observe that the standard deviation is relatively higher than other properties on rationale extraction metrics. It might be explained that the extraction process is unstable when perturbing the input. Although the performance decreases, we still include the compact as a training objective to make the RA-subgraph compact.

Supervised Setting To explore the impact of rationale supervision, we replace the compact loss with the rationale

²Note that we take the *fidelity* as the primary property for guiding rationale extraction and only conduct the ablation study on the *compact* and *topology* properties.

Model	FEVEROUS		
	Fidelity (\downarrow)	Size (\uparrow)	Sparsity (\downarrow)
Edge Mask	1.95 \pm 0.59	367.40 \pm 0.89	3.31 \pm 0.23
SaGP	-C.	1.53 \pm 0.02	361.12 \pm 0.01
	-T.	1.42 \pm 0.01	361.44 \pm 1.24
	-C. & T.	1.42 \pm 0.00	361.45 \pm 1.21

Table 3: Evaluation of the edge mask matrix. \downarrow denotes the lower is better.

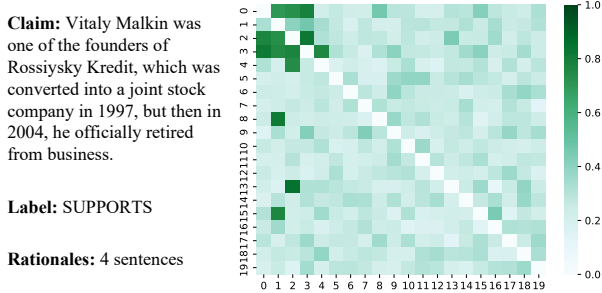


Figure 3: An example with visualization of the edge mask.

label by formulating it as a multi-label problem. Such rationale supervision can also be considered as the *plausible* property (Jain et al. 2020). From the bottom of Tab. 2 we can see, our model outperforms all baselines on different metrics of both tasks. The ceiling performance of the claim verification task remains the same with the unsupervised setting, while the rationale supervision brings an improvement on Acc.Full metrics. Moreover, compared with the unsupervised setting, we observe that there is performance improvement in rationale extraction metrics with different types of masks. It is worth pointing out that the rationale supervision provides more improvement for the node mask compared with the edge mask. The reason is that the rationale supervision directly forces the model to learn from the rationales while it is not available for edge mask.

5.2 Analysis

We explore how well the edge mask satisfies the objective of approximating the deletion of the redundant edges. The node mask is not considered as it mainly focuses on the input features as prior works have done. We evaluate the edge mask in terms of three metrics: (I) Fidelity measures the proportion of claims where the prediction is retained after the perturbation of the edge. (II) Size denotes the number of removed edges. (III) Sparsity measures the proportion of edges that are retained. As shown in Tab. 3, there are nearly 367.4 edges removed from the total of 420 edges, i.e., 3.31% edges are retained in the graph, while with only 1.95% decreases in the accuracy of claim. This reveals that a large portion of edges in the input graph is not helpful for fact verification, and the model relies heavily on a few significant edges for updating the node representation. Moreover, we explore the effectiveness of the indirect diagnostic properties of edge mask. It

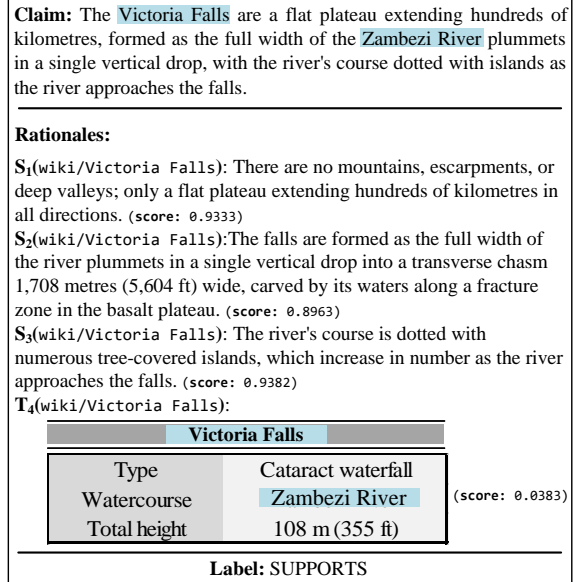


Figure 4: A case with failing to identify rationales within T₄.

can be observed from Tab. 3 that compared with the compact constrained on the rationale subgraph, the topology property affects the learning process of edge mask to some extent. In Fig. 3, we visually present the edge mask logits of the example from the test dataset. As expected, it clearly shows the significant edges within the adjacency matrix.

5.3 Error Analysis

We conduct an error analysis on rationales extracted by our model on 50 randomly chosen examples from the test set of the FEVEROUS dataset. The main errors are summarized as follows: (I) the noisy evidence contains bags of tokens that highly overlap with the claim, which brings difficulty in accurately understanding the semantic information for the model. (II) the model fails to capture the implicit correlation between different entities, especially for the table cells. For example, in Fig. 4, the relation between “Victoria Falls” and “Zambezi River” present in the table while not mentioned in the sentence. This scenario requires the model with a higher-level understanding to parse the structure of tables.

6 Conclusion

By framing explainable multi-hop fact verification as the subgraph extraction, we propose a novel GCN-based method with salience-aware graph learning to jointly model the multi-hop fact verification and the rationale extraction. Moreover, we introduce three diagnostic properties as additional training objectives to improve the quality of the extracted rationale in the multi-task model. The results on the FEVEROUS benchmark dataset demonstrate the effectiveness of our model. In the future, we will explore how to adapt the model to other domains.

7 Acknowledgments

We would like to thank anonymous reviewers for their valuable comments and helpful suggestions. We would also like to thank Professor Yulan He for her valuable suggestions for our paper. This work was funded by the National Natural Science Foundation of China (62176053).

References

- Aly, R.; Guo, Z.; Schlichtkrull, M. S.; Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; Cocarascu, O.; and Mittal, A. 2021. FEVEROUS: Fact Extraction and VERification Over Unstructured and Structured information. In *proceedings of NeurIPS*.
- Atanasova, P.; Simonsen, J. G.; Lioma, C.; and Augenstein, I. 2020. Generating Fact Checking Explanations. In *proceedings of ACL*, 7352–7364.
- Atanasova, P.; Simonsen, J. G.; Lioma, C.; and Augenstein, I. 2022. Diagnostics-Guided Explanation Generation. In *proceedings of AAAI*, 10445–10453.
- Bastings, J.; Aziz, W.; and Titov, I. 2019. Interpretable Neural Predictions with Differentiable Binary Variables. In *proceedings of ACL*, 2963–2977.
- Bekoulis, G.; Papagiannopoulou, C.; and Deligiannis, N. 2021. A Review on Fact Extraction and Verification. *ACM Computing Surveys*, 55(1): 1–35.
- Chen, J.; Bao, Q.; Sun, C.; Zhang, X.; Chen, J.; Zhou, H.; Xiao, Y.; and Li, L. 2022. LOREN: Logic-Regularized Reasoning for Interpretable Fact Verification. In *proceedings of AAAI*, 10482–10491.
- Chen, W.; Wang, H.; Chen, J.; Zhang, Y.; Wang, H.; Li, S.; Zhou, X.; and Wang, W. Y. 2020. TabFact: A Large-scale Dataset for Table-based Fact Verification. In *proceedings of ICLR*.
- DeYoung, J.; Jain, S.; Rajani, N. F.; Lehman, E.; Xiong, C.; Socher, R.; and Wallace, B. C. 2020. ERASER: A Benchmark to Evaluate Rationalized NLP Models. In *proceedings of ACL*, 4443–4458.
- Glockner, M.; Habernal, I.; and Gurevych, I. 2020. Why Do You Think That? Exploring Faithful Sentence-level Rationales without Supervision. In *proceedings of Findings of EMNLP*, 1080–1095.
- Guo, Z.; Schlichtkrull, M. S.; and Vlachos, A. 2022. A Survey on Automated Fact-Checking. *Transactions of the Association for Computational Linguistics*, 10: 178–206.
- Herzig, J.; Nowak, P. K.; Müller, T.; Piccinno, F.; and Eisen-schlos, J. M. 2020. TaPas: Weakly Supervised Table Parsing via Pre-training. In *proceedings of ACL*, 4320–4333.
- Jain, S.; and Wallace, B. C. 2019. Attention is not Explanation. In *proceedings of NAACL-HLT*, 3543–3556.
- Jain, S.; Wiegrefe, S.; Pinter, Y.; and Wallace, B. C. 2020. Learning to Faithfully Rationalize by Construction. In *proceedings of ACL*, 4459–4473.
- Jiang, Y.; Bordia, S.; Zhong, Z.; Dognin, C.; Singh, M. K.; and Bansal, M. 2020. HoVer: A Dataset for Many-Hop Fact Extraction And Claim Verification. In *proceedings of Findings of the EMNLP*, 3441–3460.
- Kipf, T. N.; and Welling, M. 2016. Variational Graph Auto-Encoders. *CoRR*, abs/1611.07308.
- Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *proceedings of ICLR*.
- Kotonya, N.; Spooner, T.; Magazzeni, D.; and Toni, F. 2021. Graph Reasoning with Context-Aware Linearization for Interpretable Fact Extraction and Verification. In *proceedings of the Fourth Workshop on Fact Extraction and VERification (FEVER)*, 21–30.
- Kotonya, N.; and Toni, F. 2020a. Explainable Automated Fact-Checking: A Survey. In *proceedings of COLING*, 5430–5443.
- Kotonya, N.; and Toni, F. 2020b. Explainable Automated Fact-Checking for Public Health Claims. In *proceedings of EMNLP*, 7740–7754.
- Lei, T.; Barzilay, R.; and Jaakkola, T. S. 2016. Rationalizing Neural Predictions. In *proceedings of EMNLP*, 107–117.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR*, abs/1907.11692.
- Lucic, A.; ter Hoeve, M. A.; Tolomei, G.; de Rijke, M.; and Silvestri, F. 2022. CF-GNNExplainer: Counterfactual Explanations for Graph Neural Networks. In *proceedings of AISTATS*, volume 151, 4499–4511.
- Luo, S.; Ivison, H.; Han, S. C.; and Poon, J. 2021. Local Interpretations for Explainable Natural Language Processing: A Survey. *CoRR*, abs/2103.11072.
- Meister, C.; Lazov, S.; Augenstein, I.; and Cotterell, R. 2021. Is Sparse Attention more Interpretable? In *proceedings of ACL/IJCNLP*, 122–129.
- Ostrowski, W.; Arora, A.; Atanasova, P.; and Augenstein, I. 2021. Multi-Hop Fact Checking of Political Claims. In *proceedings of IJCAI*, 3892–3898.
- Paranjape, B.; Joshi, M.; Thickstun, J.; Hajishirzi, H.; and Zettlemoyer, L. 2020. An Information Bottleneck Approach for Controlling Conciseness in Rationale Extraction. In *proceedings of EMNLP*, 1938–1952.
- Popat, K.; Mukherjee, S.; Yates, A.; and Weikum, G. 2018. DeClarE: Debunking Fake News and False Claims using Evidence-Aware Deep Learning. In *proceedings of EMNLP*, 22–32.
- Serrano, S.; and Smith, N. A. 2019. Is Attention Interpretable? In *proceedings of ACL*, 2931–2951.
- Sha, L.; Camburu, O.; and Lukasiewicz, T. 2021. Learning from the Best: Rationalizing Predictions by Adversarial Information Calibration. In *proceedings of AAAI*, 13771–13779.
- Shu, K.; Cui, L.; Wang, S.; Lee, D.; and Liu, H. 2019. dEFEND: Explainable Fake News Detection. In *proceedings of SIGKDD*, 395–405.
- Si, J.; Zhou, D.; Li, T.; Shi, X.; and He, Y. 2021. Topic-Aware Evidence Reasoning and Stance-Aware Aggregation for Fact Verification. In *proceedings of ACL/IJCNLP*, 1612–1622.

Sundararajan, M.; Taly, A.; and Yan, Q. 2017. Axiomatic Attribution for Deep Networks. In *proceedings of ICML*, volume 70, 3319–3328.

Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; and Mittal, A. 2018. FEVER: a Large-scale Dataset for Fact Extraction and VERification. In *proceedings of NAACL-HLT*, 809–819.

Wan, H.; Chen, H.; Du, J.; Luo, W.; and Ye, R. 2021. A DQN-based Approach to Finding Precise Evidences for Fact Verification. In *proceedings of ACL/IJCNLP*, 1030–1039.

Wiegrefe, S.; and Marasovic, A. 2021. Teach Me to Explain: A Review of Datasets for Explainable Natural Language Processing. In *proceedings of NeurIPS*.

Wu, L.; Rao, Y.; Lan, Y.; Sun, L.; and Qi, Z. 2021. Unified Dual-view Cognitive Model for Interpretable Claim Verification. In *proceedings of ACL/IJCNLP*, 59–68.

Wu, L.; Rao, Y.; Zhao, Y.; Liang, H.; and Nazir, A. 2020. DTCA: Decision Tree-based Co-Attention Networks for Explainable Claim Verification. In *proceedings of ACL*, 1024–1035.

Yang, F.; Pentyala, S. K.; Mohseni, S.; Du, M.; Yuan, H.; Linder, R.; Ragan, E. D.; Ji, S.; and Hu, X. B. 2019. XFake: Explainable Fake News Detector with Visualizations. In *proceedings of WWW*, 3600–3604.

Ying, Z.; Bourgeois, D.; You, J.; Zitnik, M.; and Leskovec, J. 2019. GNNExplainer: Generating Explanations for Graph Neural Networks. In *proceedings of NeurIPS*, 9240–9251.

Yoon, J.; Jordon, J.; and van der Schaar, M. 2019. INVASE: Instance-wise Variable Selection using Neural Networks. In *proceedings of ICLR*.

Yu, J.; Xu, T.; Rong, Y.; Bian, Y.; Huang, J.; and He, R. 2021. Graph Information Bottleneck for Subgraph Recognition. In *proceedings of ICLR*.

Zaidan, O.; Eisner, J.; and Piatko, C. 2007. Using “Annotator Rationales” to Improve Machine Learning for Text Categorization. In *proceedings of NAACL*, 260–267.

Zhao, C.; Xiong, C.; Rosset, C.; Song, X.; Bennett, P. N.; and Tiwary, S. 2020. Transformer-XH: Multi-Evidence Reasoning with eXtra Hop Attention. In *proceedings of ICLR*.

Zhou, X.; and Zafarani, R. 2020. A Survey of Fake News: Fundamental Theories, Detection Methods, and Opportunities. *ACM Computing Surveys*, 53(5): 109:1–109:40.