

# Short Video is not only Video: Multimodal Unified Social Hypergraph Contrastive Enhancement for Fake News Video Detection

Anonymous ACL submission

## Abstract

Nowadays, fake short videos have seriously affected people’s perception of news and situational awareness of event development. Previous work mainly focuses on the characteristics and dissemination of the news, and there is no in-depth mining of the social relationships and feature relationships of videos. This paper proposes a **Multimodal Unified Social Hypergraph Contrastive Enhancement** method (**MUHC**) for fake news videos detection. First, a unified social hypergraph is innovatively established for the representation of potential relationships in short videos. Meanwhile, a multimodal contrastive learning method for intra-modal and inter-modal relationships are designed to integrate different modalities. The above approach enhances data scalability while learning deeper about the potential relationships of the videos. Extensive experiments demonstrate that the method outperforms state-of-the-art on benchmark dataset.<sup>1</sup>

## 1 Introduction

The emergence of social networks has not only brought about the convenience of communication and access to knowledge, but also resulted in an abundance of fake news, which has caused significant impacts in politics(Fisher et al., 2016), economy(ElBoghdady, 2013), culture(Olan et al., 2024) and health(Chen et al., 2023). Nowadays, with the rise of short video platforms, the carrier of information exists not only in simple text, image, and audio, but in social relationships and user interactions, making the detection of fake news increasingly difficult(Comito et al., 2023). Therefore, it is urgent to design an automated fake news video detection method to decrease the negative impact.

Existing multimodal fake news detection methods are mainly classified into two aspects, social relationship-based(Shu et al., 2020) and semantic

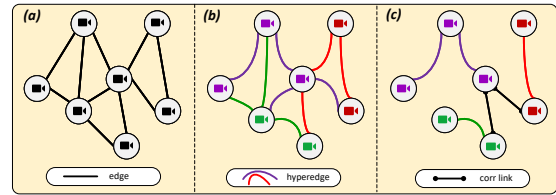


Figure 1: Illustration of a) video potential relationship graph, b) video potential relationship hypergraph and c) unified social hypergraph. **HINT**: where corr link connects event relationship. It can be clearly seen that the unified social hypergraph can represent the same video correlation using fewer edges.

feature-based(Wang et al., 2018). For the first stage, social relationship is mainly represented between different news through graph. In recent years, the methods are mainly used to construct the propagation relationship graph through commenting and retweeting mechanism(Cheng et al., 2021; Zhang et al., 2023). Graph structure mining is performed through GNN-based methods to accomplish the node classification or graph classification tasks to predict fake news. Further, related works focus on the learning and mining of attribute graphs, fusing graph and modal information for classification tasks(Nguyen et al., 2020; Phan et al., 2023).For the other stage, existing methods mainly rely on extracting, fusing(Zhou et al., 2020; Nan et al., 2021), and enhancing different modal data(Zhu et al., 2022; Qi et al., 2021). To avoid simple feature fusion engineering, some researchers perform information integrity through fact checking(Vo and Lee, 2018), reading interest(Wu et al., 2023), external knowledge(Hu et al., 2021), etc. These methods have been confirmed to achieve good results.

However, there are certain challenges in the above works. For social relations, due to the different recommendation mechanisms, short video data social relations are more complex. It is hard to obtain the propagation relations of the videos

<sup>1</sup>Code will release if manuscript acceptable.

in the same way as Twitter and Weibo. Therefore, the traditional cascade mining is difficult to obtain the complete short video relationships. Meanwhile, the more complex network structure provides challenges for short video potential relationship mining, and it is difficult for traditional graph to mine the global information of the nodes, which restricts the learning of video potential relationships and styles. For the feature aspect, the above methods are overly dependent on external features, and the features of the video itself only remain in the work of extraction and splicing using pre-trained models, without aligning the relationships of the modalities. In addition, due to the specificity of the fake news detection, there is not enough labeled data for experimental analysis. Overall, the **motivation** lies in design a model that can both deeply mine large-scale video content features and extract video social relationship features.

To solve the above issues, this paper designs a **Multimodal Unified Social Hypergraph Contrastive Enhancement** method (**MUHC**) for Fake News Video Detection. To address the first challenge, this paper designs a unified social hypergraph establishment method for short videos. Specifically, unlike propagation graph, this paper combining ordinary graph and hypergraph based on the similarity attributes and event attributes between videos for representing the relationships between short videos. This approach can represent social relationships with different graph structures more effectively and reduces the effect of redundant edges. For the second challenge, this paper designs a multimodal contrastive learning method. Combining the unified social hypergraph constructed above, intra-modal is performed by data augmentation of different modal features, and then multimodal consistency alignment is performed for inter-modal features for inter-modal contrastive learning and network tuning. Finally, fine-tuning is performed to learn the video semantics based on the existing tags. Overall, the main contributions of this work are as follows.

- **Unified Social Potential Hypergraph.** Based on video social relationships, a novelty unified social potential hypergraph is built for addressing the shortcomings of graph structure, and a more effective short video potential hypergraph representation is obtained.
- **Multimodal Contrastive Learning with Hypergraph.** Combined with the above unified

hypergraph, a multimodal contrastive learning method is designed for intra-modal and inter-modal feature alignment and enhancement, and the self-supervised learning method lays the foundation for the extension of unlabeled datasets.

- **Better Performance.** Experimental results show that the method proposed in this paper outperform the existing fake news detection methods on the benchmark dataset and achieves state-of-the-art.

## 2 Related Works

To understand the existing works, this section investigates the application of graph representation learning and multimodal contrastive learning in fake news detection.

### 2.1 Graph Representation Learning

Fake news detection methods based on graph representation learning mainly focus on the propagation structure. Vosoughi et al. explored the diffusion patterns of information on Twitter, revealing that misinformation tends to spread more extensively than real news(Vosoughi et al., 2018). Zhou et al. subsequently detailed the characteristics of how fake news propagates(Zhou and Zafarani, 2019). Moreover, the GCNFN model enhances comment embedding by incorporating users' profiles(Monti et al., 2019), while the UPFD model considers users' past posts to reflect their inherent biases(Dou et al., 2021). These approaches demonstrate robust detection capabilities but require high-quality social network data. Furthermore, some researchers have focused on leveraging graph-structured data. Wang et al. introduced KMGCN, a framework that employs graph convolution networks for extracting textual features(Wang et al., 2020). Rosenfeld et al. developed a Weisfeiler-Lehman graph kernel that is independent of text, user, and time inputs, illustrating that structural encoding of information cascades can significantly aid in assessing the credibility of news(Rosenfeld et al., 2020). However, the above methods can rely on relatively few features, and the detection efficiency is lower than that of multi-dimension feature methods. To mine the potential social relationships of news, Qi et al. designed a short video social graph characterization method based on correlating with neighbors, which achieved good results.(Qi et al., 2023b) But the method is augmented by debunk videos, which is

not able to achieve the purpose of early detection. In addition, inspired by hypergraphs(Feng et al., 2019; Gao et al., 2022), Sun et al. designed HG-SL for learning user spreading behavior(Sun et al., 2023), but still have problems.

## 2.2 Multimodal Contrastive Learning

Contrastive learning has been widely used in multimodal representation and learning(Qian et al., 2022). Radford et al. designed a text-image pre-training model using contrast learning, CLIP(Radford et al., 2021). Based on this, Zhou et al. proposed a CLIP-guided learning multimodal fake news detection method for image-text consistency detection(Zhou et al., 2023). Moreover, leveraging different contrastive properties, ALBEF(Li et al., 2021) introduced a strategy for mining hard negatives based on the distribution of contrastive similarities, while BLIP(Li et al., 2022) and VLMO(Bao et al., 2022) implemented a comparable approach. To deepen the exploration of contrastive information, Chen et al. addressed the cross-modal ambiguity learning challenge from an information theory standpoint, adopting specific methods for detection depending on the ambiguities present across various modalities(Chen et al., 2022). Based on this, Wang et al. established a cross-modal contrastive learning fake news detection method, which realized more accurate image-text alignment(Wang et al., 2023). However, the above method mainly focuses on the graphic domain and mainly compares the traditional modalities with limited learnable features. To solve above problems, Yin et al. designed a contrast loss and graph autoencoder to effectively learn the potential features of the propagation graph(Yin et al., 2024). However, it still did not solve the problem of contrast relationship between traditional modalities and social relation modalities.

## 3 Paramilitary

### 3.1 Problem Definition

**Definition 1 (Unified Hypergraph)** We define the unified hypergraph  $U = (\mathcal{V}, \mathcal{E})$  as an extended graph structure where  $\mathcal{V}$  is the set of vertices and  $\mathcal{E}$  is the set of edges. Each edge  $e \in \mathcal{E}$  is a non-empty subset  $e \subseteq \mathcal{V}$  such that  $|e| \geq 2$ .

**Problem 1 (Unified Social Potential Hypergraph)** Define the Unified Potential Relationship Hypergraph as  $\mathcal{H} = (\mathcal{V}, \mathcal{E}, \mathbf{W}, \mathbf{H})$ , where  $\mathcal{V}$  represents the nodes of the hypergraph. The hyperedges,

$\mathcal{E} \subseteq \mathcal{P}(\mathcal{V})$ , represent the complex relationships among videos and events, where  $\mathcal{P}(\mathcal{V})$  denotes the power set of  $\mathcal{V}$ . The node attributes  $\mathbf{X} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathbb{R}^{N \times D}$ , describe the features of the videos, with  $N$  being the number of videos/nodes and  $D$  the dimension of the attribute features. The adjacency matrix  $\mathbf{H} \in \mathbb{N}^{|\mathcal{V}| \times |\mathcal{E}|}$  represents the connectivity within the hypergraph, reflecting the relationships among nodes. Further details see section 4.

**Problem 2 (Fake News Video Detection)** The video dataset is defined as  $V = \{V_1, V_2, \dots, V_n\}$ , where each  $V_i$  is a video instance. Each instance  $V_i$  is modeled as a video potential relationship processes. We aim to learn a self-supervised function,  $g$ , defined as:

$$g : V \rightarrow Y,$$

where  $V$  represents video instances with their potential relationship processes and  $Y \in \{\text{Fake}, \text{Real}\}$  denotes fake or real videos.

### 3.2 Video Unimodal Features Representations

Short videos have rich modalities, such as text, audio, image and frame et al. In social networks, short videos also have social contexts (which are also named metadata), such as user information and comments.

For modal extractions, this paper use pre-trained models, and the specific extraction model is shown in Figure 2(a).

## 4 Methodology

This section introduce the proposed method MUGC in detail. As shown in Figure 2, we first perform feature extraction and aggregation of different modalities of the video. After that, the unified social potential hypergraph is built based on the video **Social Link** and **Corr Link**, and the hypergraph structure is embedded with the representation. Furthermore, intra-modal and inter-modal contrastive learning is performed on the representations to perform relevance enhancement of the different features, with the pruning optimization of network parameters. Finally, the network is tuned by labeled data to complete the supervised fake news video detection. The specific details and computational process are shown below.

### 4.1 Video Feature Integration

There may be correlations between different modals of short video. This subsection follows

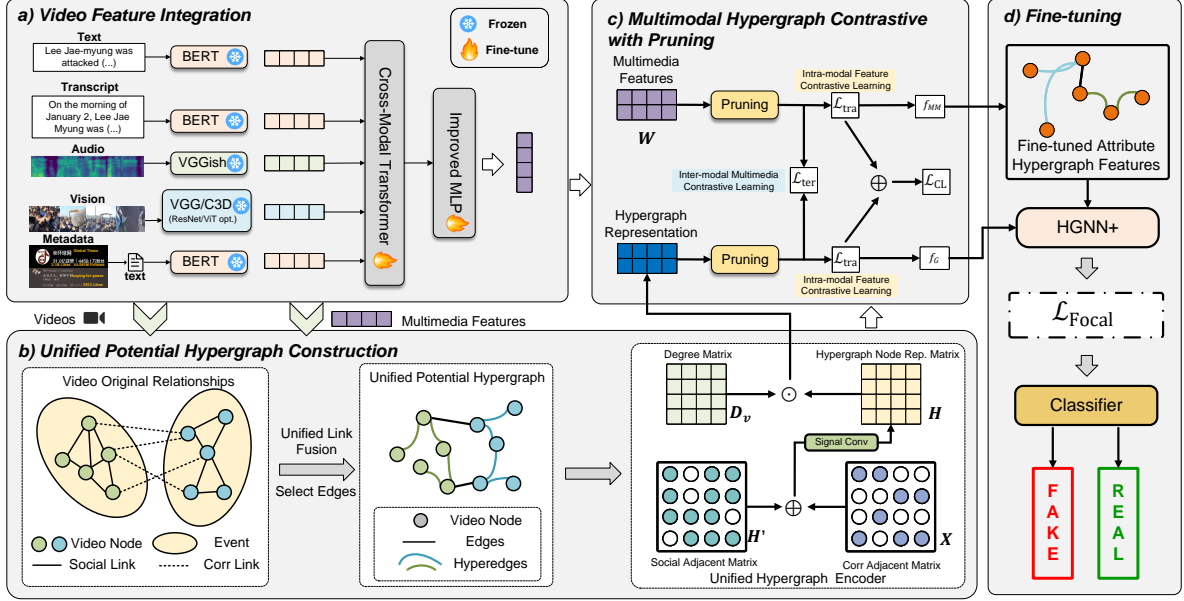


Figure 2: Framework of MUHC. The whole framework is divided into four modules.

the feature extraction method of Section 3.2, and expanded from co-attention(Lu et al., 2016) to a **cross-modal transformer** method. For consistency modeling for different modalities, which is used to learn the differences and correlations between different modalities. The specific equation is eq(1).

$$\widetilde{h}_{x1}, \dots, \widetilde{h}_{xn} = \text{cross-modal-transformer}\{h_{x1}, \dots, h_{xn}\}, \quad (1)$$

where  $h_{x1}, \dots, h_{xn}$  is the original representation of the different modes,  $\widetilde{h}_{x1}, \dots, \widetilde{h}_{xn}$  is the computed new representation with modal alignment. The specific combining method is shown in Fig. 2.

After computation, the corresponding new modal information is aggregated using the transformer to obtain the corresponding modal aggregation representation as shown in equation 2.

$$\mathbf{X} = \text{MEAN}\{\widetilde{h}_T, \widetilde{h}_O, \widetilde{h}_A, \widetilde{h}_I, \widetilde{h}_F, \widetilde{h}_M\}, \quad (2)$$

where MEAN is an average computation based on the modal representation, which is computed and then passed into a modified MLP for vector representation and parameter augmentation.

## 4.2 Unified Social Potential Hypergraph Construction

In this section, this paper builds a special kind of hypergraph, Unified Hypergraph, for representation of different relations. For Social Link, hyperedge

is used for the representation, and for Corr Link, the ordinary graph edge is used. Finally the potential hypergraph representation of short videos with stronger performance combining the hyper-edge and the ordinary edge is built.

### 4.2.1 Social Link

When a outbreak news appears, short video platforms will have amount of videos of the same type in a very short period of time, which is defined as an event. Different videos in the same event are fake and real, but generally there will be a strong correlation. The correlation of the same event is defined as a Social Link. We build a hypergraph of the Social Link with hyperedges, where all the video nodes of the same event are connected by a single hyper edge. The set of these hyperedges is defined as  $\mathcal{E}_s$ . Equation (3) indicates the method of hypergraph building.

$$\begin{cases} \mathbf{w}_c = \text{copy}(\text{sigmoid}(w_c), M_c) \\ \mathbf{W} = \text{diag}(\mathbf{w}_1^1, \dots, \mathbf{w}_1^{M_1}, \dots, \mathbf{w}_C^1, \dots, \mathbf{w}_C^{M_C}), \\ \mathbf{H}' = \mathbf{H}'_1 \|\mathbf{H}'_2\| \dots \|\mathbf{H}'_C \end{cases} \quad (3)$$

where  $w_c \in R$  represent a hyperparameter and all hyperedges within a specific hyperedge group  $c$ . The  $\text{sigmoid}(\cdot)$  function is used for element-wise normalization. The representation  $\mathbf{w}_j = (\mathbf{w}_C^1, \dots, \mathbf{w}_C^{M_C}) \in R^{M_C}$  denotes the weight embedding for hyperedge group  $c$ . The  $\text{copy}(a, b)$  function generates a embedding of size  $b$ , with the

value  $a$  replicated  $b$  times. The total number of hyperedges,  $M$ , is given by the sum  $M_1 + M_2 + \dots + M_c$ . The weight matrix  $\mathbf{W} \in R^{M \times M}$  is diagonal, with each diagonal entry  $\mathbf{W}^{ii}$  representing the weight of hyperedge  $e_i$ . The incidence matrix  $\mathbf{H}' \in \mathbb{N}^{N \times M}$  is formed by concatenating the incidence matrices of various hyperedge groups.

#### 4.2.2 Corr Link

However, for the above hypergraph building method, events cannot be linked to each other. Therefore, we built Corr Link, which establishes a link based on the similarity between videos and videos, and is calculated using the Minkowski distance formula, which is shown below:

$$\text{sim}(\mathbf{v}_i, \mathbf{v}_j) = \left( \sum_{k=1}^D |\mathbf{x}_{i,k} - \mathbf{x}_{j,k}|^D \right)^{1/D}, \quad (4)$$

where  $x_{i,k}, x_{j,k}$  are representations of different videos.

However, a huge number of correlation edges can be obtained through the equation. To ensure the stronger correlation between different videos, this paper ensures that each node can select at most  $K$  edges. The specific equation is as follows.

$$\mathcal{E}_w = \left\{ (\mathbf{v}_i, \mathbf{v}_j) \mid \text{sim}(\mathbf{v}_i, \mathbf{v}_j) \text{ in top-}K \text{ of } [\text{sim}(\mathbf{v}_i, \mathbf{v}_j)]_{p=1}^N \right\}, \quad (5)$$

where  $K$  is a trainable hyperparameter. Note that the following relationship  $\forall \mathbf{v}_i \in \mathcal{V}$  holds.

#### 4.2.3 Unified Link Fusion

Combining the hyperedges created above with ordinary edges results in a newly defined Unified Hypergraph, and the set of edges can be denoted by  $\mathbf{H}$ .

$$\mathbf{H} = (\mathbf{H}'_1 \parallel \mathbf{H}'_2 \parallel \dots \parallel \mathbf{H}'_c) \oplus (\mathbf{A}_1 \parallel \mathbf{A}_2 \parallel \dots \parallel \mathbf{A}_c), \quad (6)$$

where  $\mathbf{H}'_1 \parallel \mathbf{H}'_2 \parallel \dots \parallel \mathbf{H}'_c$  denotes hyperedges set.  $\mathbf{A}_1 \parallel \mathbf{A}_2 \parallel \dots \parallel \mathbf{A}_c$  denotes corr edges set.

The graph is a special form of hypergraph, to better characterize the above graph structure, this paper introduces an iterative version of the hypergraph convolution method, HGNN+, which can be efficiently fused for different nodes and hyperedges, and is calculated as follows.

$$\mathbf{X}^{t+1} = \sigma \left( \mathbf{D}_v^{-1/2} \mathbf{H} \mathbf{W} \mathbf{D}_e^{-1} \mathbf{H}^\top \mathbf{D}_v^{-1/2} \mathbf{X}^t \Theta \right), \quad (7)$$

where parameter  $\Theta$  is learnable. This filter  $\Theta$  is used on the nodes of the hypergraph to extract features. and  $\mathbf{D}_v$  is the diagonal matrix of degrees.

After the above equations are calculated, a unified social potential hypergraph representation, named  $h_{UG}$ , can be obtained and used as input for subsequent contrastive learning task.

### 4.3 Multimodal Hypergraph Contrastive Learning with Pruning

#### 4.3.1 Pruning

Data imbalance is a very important issue that affects the performance of comparative learning. Nowadays, there have been many works to solve the problem of data imbalance (Jiang et al., 2021; Frankle and Carbin, 2018). In this paper, we follow a pruning method for data augmentation on unbalanced data (Frankle and Carbin, 2018).

#### 4.3.2 Intra-modal Contrastive Learning

Contrastive learning is an effective self-supervised learning method, and this and the next subsection will be devoted to the contrastive learning method used. Intra-modal contrastive learning is mainly applied to the analysis of correlations within similar modalities, where the original modal features are compared with the pruning enhanced features to calculate the contrast loss with the following equation:

$$\mathcal{L}_{\text{tra}} = -\log \sum_{\mathbf{v}_i \in \mathcal{V}} \frac{\exp [\text{sim}(\tilde{z}^i, \tilde{z}_p^j) / \tau_{\text{tra}}]}{\sum_{p=1}^{2n} \mathbb{1}_{[i \neq p]} \exp [\text{sim}(\tilde{z}^i, \tilde{z}_p^p) / \tau_{\text{tra}}]}, \quad (8)$$

where  $\tilde{z}, \tilde{z}_p$  are the features before and after pruning.  $\tau$  is temperature which is a hyperparameter.

#### 4.3.3 Inter-modal Contrastive Learning

Inter-modal contrastive learning is mainly used for the comparison between the graph structure and the features of the video itself for the alignment between different modal data. All contrast losses are calculated and summed to get the inter-modal contrast loss.

$$\mathcal{L}_{\text{ter}} = -\sum_i \log \sum_{\mathbf{v}_i \in \mathcal{V}} \frac{\exp [\text{sim}(\tilde{z}^i, \tilde{z}^j) / \tau_{\text{ter}}]}{\sum_{p=1}^{2n} \mathbb{1}_{[i \neq p]} \exp [\text{sim}(\tilde{z}^i, \tilde{z}_p^p) / \tau_{\text{ter}}]}. \quad (9)$$

#### 4.3.4 Contrastive Loss Integration

To integrate contrast learning loss, a hyperparameter  $\lambda$  is set for joint contrast loss to use for tuning the network.

$$\mathcal{L}_{CL} = \lambda \mathcal{L}_{\text{tra}} + (1 - \lambda) \mathcal{L}_{\text{ter}}. \quad (10)$$

## 4.4 Fine-tuning

To obtain the final model prediction structure, a fine-tune attribute graph classification method is designed in this paper. Both graph structure and feature representation are derived from the comparison learning trained structure, and a more robust loss function Focal is used for model training inference and prediction.

$$\mathcal{L}_{\text{focal}} = -\frac{1}{|\mathcal{V}_i|} \sum_{i \in \mathcal{V}_i} \sum_{c=0}^C \alpha_c y_{ic} (1 - \hat{y}_{ic})^\gamma \log(\hat{y}_{ic}), \quad (11)$$

where  $\alpha, \gamma$  are learnable hyperparameters.

## 5 Experiments and Analysis

In this section, we will design extensive experiments to measure the effectiveness of our proposed method MUGC. The following 3 research questions are proposed for subsequent experimental investigations.

- **RQ1.** How does MUGC perform in the fake news video detection?
- **RQ2.** What is the performance of the proposed unified hypergraph and contrastive learning method?
- **RQ3.** What is the interpretability significance of the method proposed and why not use other modules to solve the task?

### 5.1 Dataset

Due to the lack of datasets, this paper adopts the **only** available benchmark of short social network videos, FakeSV(Qi et al., 2023a). The dataset is obtained from Douyin and Kuaishou, which has a total of 738 events, 1827 real videos and 1827 videos are fake. The video has semantic features such as original text, audio, and visual, and social features such as comment and user.

### 5.2 Experiment Settings

#### 5.2.1 Metrics

The experiments follow the criteria of benchmark dataset, also using 5-fold cross-validation with interval estimation. For the metrics, for fake news detection which is essentially a classification task, Accuracy, Recall, Precision, F1-score are used.

#### 5.2.2 Parameter Setups

The experiments in this paper use AMD Ryzen Threadripper PRO 5995WX 64-Cores CPU and NVIDIA RTX A6000 as the experimental environment, and the maximum memory usage for data loading is 29G. textual feature extraction is performed using BERT-base-uncased, and the batch size is 64. Indicating that the features are 768, and the K parameter of Corr Link is set to 10. Other hyperparameters are analyzed in Appendix.

#### 5.3 Baselines

In order to prove the superiority of the method, some of the more advanced algorithms are compared in terms of unimodal, multimodal and large language models as follows:

##### 5.3.1 Uni-modal

Traditional analytical methods mainly explore the expressive unimodal features. This paper mainly use BERT(Devlin et al., 2019), VGGish(Hershey et al., 2017), VGG19(Simonyan and Zisserman, 2015), and C3D(Ji et al., 2013) to analyze the characteristics of their respective modalities. In addition, Zhang et al. designed a news text detection method for news content and comments unified about emotions(Zhang et al., 2021).

##### 5.3.2 Multi-modal

Existing multimodal methods mainly focus on text-image methods, such as EANN(Wang et al., 2018), and video semantic modal methods, such as Serrano et al.(Serrano et al., 2020) and FANVM(Choi and Ko, 2021). But there are fewer analysis that have both semantic features and social features. In this paper, we use the benchmark of FakeSV dataset, SV-FEND, as a baseline to compare with our proposed method(Qi et al., 2023a).

##### 5.3.3 Large Language Model

With the development of the large language model, it is more capable of giving more positively correlated responses to linguistic information. To explore the application ability of GPT for the task of fake news detection, this paper designs the prompt method and conducts experiments through the APIs of GPT3.5-turbo (OpenAI, 2022) and GPT4.0-turbo (OpenAI, 2023), respectively. For GPT3.5, the text is inputted. For the corresponding multimodal large model of GPT4.0, the content of text and video pumping frames are inputted to obtain positive incentive feedback.

Table 1: Experimental results of different methods. (Acc.: accuracy, Prec.: precision, Rec.: recall)

	Method	Acc.	Prec.	Rec.	F1
Uni-modal	Text(BERT)	77.14 $\pm$ 2.75	77.18 $\pm$ 2.73	77.12 $\pm$ 2.72	77.10 $\pm$ 2.76
	Audio(VGGish)	66.78 $\pm$ 1.29	67.10 $\pm$ 1.20	66.74 $\pm$ 1.26	66.58 $\pm$ 1.34
	Image(VGG19)	69.72 $\pm$ 3.21	69.80 $\pm$ 3.30	69.82 $\pm$ 3.29	69.65 $\pm$ 3.21
	Video(C3D)	69.59 $\pm$ 1.77	70.07 $\pm$ 1.81	69.56 $\pm$ 1.75	69.38 $\pm$ 1.77
	Text(Dual Emotion)	75.68 $\pm$ 2.04	77.18 $\pm$ 2.74	77.09 $\pm$ 2.74	77.05 $\pm$ 2.77
Multi-modal	Serrano et al.	68.72 $\pm$ 2.30	70.75 $\pm$ 2.11	68.73 $\pm$ 2.30	67.92 $\pm$ 2.53
	FANVM	76.39 $\pm$ 1.14	75.41 $\pm$ 1.54	73.73 $\pm$ 1.24	74.19 $\pm$ 1.18
	EANN	77.87 $\pm$ 2.03	77.91 $\pm$ 2.02	77.87 $\pm$ 2.02	77.86 $\pm$ 2.03
	SV-FEND	78.88 $\pm$ 1.98	79.41 $\pm$ 1.85	78.89 $\pm$ 1.93	78.79 $\pm$ 2.01
LLM	GPT3.5-turbo	57.49 $\pm$ 3.39	60.82 $\pm$ 4.12	57.48 $\pm$ 2.94	58.42 $\pm$ 3.24
	GPT4-turbo	65.05 $\pm$ 1.98	66.02 $\pm$ 2.13	65.05 $\pm$ 2.16	65.31 $\pm$ 2.06
	<b>MUHC(Ours)</b>	<b>89.82 <math>\pm</math>1.33</b>	<b>89.96 <math>\pm</math>1.91</b>	<b>89.82 <math>\pm</math>1.81</b>	<b>89.80 <math>\pm</math>1.27</b>
	Improve(%)	10.94 $\uparrow$	10.55 $\uparrow$	10.93 $\uparrow$	11.01 $\uparrow$

## 5.4 Experiment Performance (RQ1)

We completed the experiments with the above baselines and our method separately and performed the interval estimation. The experimental results are shown in Table 1. In the following, we will analyze the experimental results in depth.

**Performance Analysis:** The experimental results show that the experimental metrics of our method exceeds 89.5%, which is an improvement of 10% compared to the benchmark method. This result is compared with other unimodal, multimodal methods. In terms of unimodal methods, it can be seen that the modality that determines whether a short video is real or fake is mainly text, which can express more semantic than other features. In addition, although the method Zhang et al. (2021) proposed combines the news text features with the sentiment of the comments, it still does not have better results due to the semantic features. In terms of multi-modal methods, such as FANVM and EANN have a greater performance improvement than unimodal methods, but they still have a certain gap with our benchmark. For the method proposed in this paper, the main reason is that it is able to enhance the inter-modal features using contrastive learning while obtaining the features within the video modalities. Moreover, the short video social relationship mining innovatively proposed is able to better obtain the potential relationship between videos, which improves the feature mining ability.

**Compare with LLM:** Due to the limitation of data and resources, the common method of applying LLM at this stage is mainly prompt without fine-tuning. However, the LLMs are mainly generative models, which cannot achieve better results in the classification task of learning features mentioned above. As a result, the general metrics of the LLM are very low. The future work can consider the en-

hancement and expansion of the data to accomplish the fine-tuning for large language models.

## 5.5 Ablation Studies (RQ2)

In this subsection, the Unified Social Potential Hypergraph and the contrastive learning approach are ablated and analyzed, respectively. Besides, the practical implications of the above modules are explained below.

### 5.5.1 Study on Unified Social Potential Hypergraph

- **w/o PH.** the proposed unified social potential hypergraph is removed altogether, leaving only the basic features of the video and comparative learning.
- **w/o Social Link.** remove the hyperedge module and leave the PH module with only Corr Link aka similarity module.
- **w/o Corr Link.** remove the similarity module and leave only the hyperedge module.

**Analysis:** Experimental results can clearly prove that the unified social potential hypergraph has very important performance in fake news video detection based on social relationships. When all the potential hypergraphs are removed, it can be seen that the experimental metrics are reduced by about 1%, which is a significant difference. Additionally, Social Link and Corr Link both show an overall performance decrease of more than 1% when removed. Furthermore, to demonstrate the advantage of unified hypergraph performance, this experiment compares the inference speed and performance of general graph modeling. The results indicate that the unified hypergraph has a slight improvement in both inference speed and performance compared to graph.

Table 2: Ablation study on unified social potential hypergraph(PH) and its various components. HINT: Time is inference time.

Method	Acc.	Prec.	Rec.	F1	Time
<b>Ours</b>	<b>89.83</b>	<b>89.97</b>	<b>89.91</b>	<b>89.91</b>	1.0X
w/o PH	82.44	82.16	82.34	82.17	-
w/o Social Link	87.99	88.07	88.00	88.01	-
w/o Corr Link	88.84	88.80	88.88	88.85	-
w/o Hyperedge	89.14	89.36	89.22	89.22	0.8X

Table 3: Ablation result of Multimodal Hypergraph Contrastive Enhancement

Method	Acc.	Prec.	Rec.	F1
<b>Ours</b>	<b>89.82</b>	<b>89.96</b>	<b>89.82</b>	<b>89.80</b>
w/o intra-modal	88.83	88.80	88.89	88.82
w/o inter-modal	87.84	87.89	87.86	87.87

### 5.5.2 Study on Contrastive Learning

In this section we will analyze the multimodal contrastive learning methods. To prove that inter-modal and intra-modal methods are both effective, above modules are removed respectively, with the following details:

- **w/o intra-modal** remove intra-modal contrastive learning module, compute only the relationship between video features and hypergraph features.
- **w/o inter-modal** remove inter-modal contrastive learning module, only intra-modal relations are computed after data augmentation.

**Analysis:** Two contrastive learning methods are removed separately and it can be seen that there is a certain decrease about 2% in metrics after removed. This indicates that contrastive learning can better learn the relationship and modalities. However, it can be seen that in fact the metrics of contrastive learning are not particularly significantly improved compared to the potential hypergraph. The main reason for proposing the method is to propose a self-supervised learning short video social relationship detection mechanism. For unlabeled data, our proposed method can also enhance the performance through augmentation and comparison. The method is not limited to the FakeSV dataset, but can be extended to other related fields as well in the future.

### 5.6 Interpretability Studies (RQ3)

The above experiments show that some of the feature extraction techniques and methods used are might outdated, but the method is the best strategy

Table 4: Result of module feasibility and interpretability study. SVA. is F1-score of SV-FEND.

Method	SVA.	F1
<b>Ours</b>	<b>78.88</b>	<b>89.80</b>
w/ ResNet50	78.45	89.63
w/ ViT	78.46	89.74
w/ VST	78.12	89.06

obtained by combining the ROI analysis. In this subsection, we will take the vision module as an example to briefly analyze the performance metrics. The experimental results are shown in Table 4.

For image, this paper selects ResNet(He et al., 2016) and ViT(Dosovitskiy et al., 2021) as feature extractors. Additionally, to demonstrate the method’s validity, the video processing unit’s C3D extracted content is replaced with Video Swin Transformer(Liu et al., 2022). It can be seen that the performance of the method in this paper slightly surpasses newer methods. Moreover, methods such as ViT and Video Swin Transformer have significantly more parameters compared to VGG19. Considering the return on investment, the aforementioned methods are not suitable for this framework. This proves that newer methods do not necessarily lead to better results.

## 6 Conclusion

This paper explores the establishment of potential relationships and the implementation of modal enhancement methods for short videos. Specifically, this paper designs a multimodal fake news detection algorithm based on unified social hypergraph enhancement. Separately, this paper proposes a potential social relationship hypergraph establishment without retweet information for short video social relationship mining. Among them, a Social Link building method based on hyperedge and a Corr Link building method based on ordinary edge are designed. The combination of both can be better for the representation of the proposed relationship hypergraph. Additionally, for multimodal data, we propose an intra-modal and inter-modal contrastive learning method for enhancing the relationships between different modal features and facilitating the extension to large-scale datasets. Experimental results show that the detection method proposed in this paper achieves sota results on the benchmark dataset with good efficiency.



## 626 Limitations

627 The dataset available for this method is currently  
628 limited; When performing the design of potential  
629 hypergraphs, the data needs to be labeled and de-  
630 signed according to the method proposed. Besides,  
631 the method may not work best in the multimodal  
632 data aggregation module. This paper only provides  
633 a benchmark method for data mining in the corre-  
634 sponding scenario, and its metrics are not guaran-  
635 teed to be the highest. It should be noted that the  
636 feature extraction module of the method proposed  
637 in this paper can be discussed in depth, but it is not  
638 specifically built in this paper based on such task.

## 639 References

640 Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu,  
641 Owais Khan Mohammed, Kriti Aggarwal, Subho-  
642 jit Som, Songhao Piao, and Furu Wei. 2022. Vlmo:  
643 Unified vision-language pre-training with mixture-of-  
644 modality-experts. *Advances in Neural Information  
645 Processing Systems*, 35:32897–32912.

646 Mu-Yen Chen, Yi-Wei Lai, and Jiunn-Woei Lian. 2023.  
647 Using deep learning models to detect fake news about  
648 covid-19. *ACM Transactions on Internet Technology*,  
649 23(2):1–23.

650 Yixuan Chen, Dongsheng Li, Peng Zhang, Jie Sui, Qin  
651 Lv, Lu Tun, and Li Shang. 2022. Cross-modal am-  
652 biguity learning for multimodal fake news detection.  
653 In *Proceedings of the ACM Web Conference 2022*,  
654 pages 2897–2905.

655 Lu Cheng, Ruocheng Guo, Kai Shu, and Huan Liu. 2021.  
656 Causal understanding of fake news dissemination  
657 on social media. In *Proceedings of the 27th ACM  
658 SIGKDD Conference on Knowledge Discovery &  
659 Data Mining*, pages 148–157.

660 Hyewon Choi and Youngjoong Ko. 2021. Using topic  
661 modeling and adversarial neural networks for fake  
662 news video detection. In *Proceedings of the 30th  
663 ACM International Conference on Information &  
664 Knowledge Management*, pages 2950–2954.

665 Carmela Comito, Luciano Caroprese, and Ester  
666 Zumpano. 2023. Multimodal fake news detection on  
667 social media: a survey of deep learning techniques.  
668 *Social Network Analysis and Mining*, 13(1):101.

669 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
670 Kristina Toutanova. 2019. BERT: Pre-training of  
671 deep bidirectional transformers for language under-  
672 standing. In *Proceedings of the 2019 Conference  
673 of the North American Chapter of the Association  
674 for Computational Linguistics: Human Language  
675 Technologies*, pages 4171–4186.

676 Alexey Dosovitskiy, Lucas Beyer, Alexander  
677 Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,

Thomas Unterthiner, Mostafa Dehghani, Matthias  
678 Minderer, Georg Heigold, Sylvain Gelly, Jakob  
679 Uszkoreit, and Neil Houlsby. 2021. An image  
680 is worth 16x16 words: Transformers for image  
681 recognition at scale. In *International Conference on  
682 Learning Representations*. 683

Yingtong Dou, Kai Shu, Congying Xia, Philip S Yu, and  
684 Lichao Sun. 2021. User preference-aware fake news  
685 detection. In *Proceedings of the 44th international  
686 ACM SIGIR Conference on Research and Develop-  
687 ment in Information Retrieval*, pages 2051–2055. 688

Dina ElBoghdady. 2013. Market quavers after fake ap  
689 tweet says obama was hurt in white house explosions.  
690 *The Washington Post*, 23. 691

Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong  
692 Ji, and Yue Gao. 2019. Hypergraph neural networks.  
693 In *Proceedings of the AAAI conference on Artificial  
694 Intelligence*, volume 33, pages 3558–3565. 695

Marc Fisher, John Woodrow Cox, and Peter Hermann.  
696 2016. Pizzagate: From rumor, to hashtag, to gunfire  
697 in dc. *Washington Post*, 6:8410–8415. 698

Jonathan Frankle and Michael Carbin. 2018. The lottery  
699 ticket hypothesis: Finding sparse, trainable neural  
700 networks. In *International Conference on Learning  
701 Representations*. 702

Yue Gao, Yifan Feng, Shuyi Ji, and Rongrong Ji. 2022.  
703 Hgnn+: General hypergraph neural networks. *IEEE  
704 Transactions on Pattern Analysis and Machine Intel-  
705 ligence*, 45(3):3181–3199. 706

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian  
707 Sun. 2016. Deep residual learning for image recog-  
708 nition. In *Proceedings of the IEEE Conference on  
709 Computer Vision and Pattern Recognition*, pages 770–  
710 778. 711

Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis,  
712 Jort F Gemmeke, Aren Jansen, R Channing Moore,  
713 Manoj Plakal, Devin Platt, Rif A Saurous, Bryan  
714 Seybold, et al. 2017. CNN architectures for large-  
715 scale audio classification. In *Proceedings of the 2017  
716 IEEE International Conference on Acoustics, Speech  
717 and Signal Processing*, pages 131–135. 718

Linmei Hu, Tianchi Yang, Luhao Zhang, Wanjun Zhong,  
719 Duyu Tang, Chuan Shi, Nan Duan, and Ming Zhou.  
720 2021. Compare to the knowledge: Graph neural fake  
721 news detection with external knowledge. In *Proceed-  
722 ings of the 59th Annual Meeting of the Association for  
723 Computational Linguistics and the 11th International  
724 Joint Conference on Natural Language Processing*,  
725 pages 754–763. 726

Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 2013.  
727 3D convolutional neural networks for human action  
728 recognition. *IEEE Transactions on Pattern Analysis  
729 and Machine Intelligence*, 35(1):221–231. 730

731	Ziyu Jiang, Tianlong Chen, Bobak J Mortazavi, and	Peng Qi, Yuyan Bu, Juan Cao, Wei Ji, Ruihao Shui,	783
732	Zhangyang Wang. 2021. Self-damaging contrastive	Junbin Xiao, Danding Wang, and Tat-Seng Chua.	784
733	learning. In <i>International Conference on Machine</i>	2023a. FakeSV: A multimodal benchmark with rich	785
734	<i>Learning</i> , pages 4927–4939.	social context for fake news detection on short video	786
		platforms. In <i>Proceedings of the AAAI Conference</i>	787
735	Junnan Li, Dongxu Li, Caiming Xiong, and Steven	<i>on Artificial Intelligence</i> , pages 14444–14452.	788
736	Hoi. 2022. BLIP: Bootstrapping language-image		
737	pre-training for unified vision-language understand-	Peng Qi, Juan Cao, Xirong Li, Huan Liu, Qiang Sheng,	789
738	ing and generation. In <i>International Conference on</i>	Xiaoyue Mi, Qin He, Yongbiao Lv, Chenyang Guo,	790
739	<i>Machine Learning</i> , pages 12888–12900.	and Yingchao Yu. 2021. Improving fake news detec-	791
		tion by using an entity-enhanced framework to fuse	792
740	Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare,	diverse multimodal clues. In <i>Proceedings of the 29th</i>	793
741	Shafiq Joty, Caiming Xiong, and Steven Chu Hong	<i>ACM International Conference on Multimedia</i> , pages	794
742	Hoi. 2021. Align before fuse: Vision and language	1212–1220.	795
743	representation learning with momentum distillation.		
744	<i>Advances in Neural Information Processing Systems</i> ,	Peng Qi, Yuyang Zhao, Yufeng Shen, Wei Ji, Juan Cao,	796
745	34:9694–9705.	and Tat-Seng Chua. 2023b. Two heads are better	797
		than one: Improving fake news video detection by	798
746	Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang,	correlating with neighbors. In <i>Findings of the Asso-</i>	799
747	Stephen Lin, and Han Hu. 2022. Video swin trans-	<i>ciation for Computational Linguistics 2023</i> , pages	800
748	former. In <i>Proceedings of the IEEE/CVF Conference</i>	11947–11959.	801
749	<i>on Computer Vision and Pattern Recognition (CVPR)</i> ,		
750	pages 3202–3211.	Yiyue Qian, Chunhui Zhang, Yiming Zhang, Qianlong	802
		Wen, Yanfang Ye, and Chuxu Zhang. 2022. Co-	803
751	Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh.	modality graph contrastive learning for imbalanced	804
752	2016. Hierarchical question-image co-attention for	node classification. <i>Advances in Neural Information</i>	805
753	visual question answering. In <i>Advances in Neural</i>	<i>Processing Systems</i> , 35:15862–15874.	806
754	<i>Information Processing Systems</i> .		
755	Federico Monti, Fabrizio Frasca, Davide Eynard, Da-	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	807
756	mon Mannion, and Michael M Bronstein. 2019. Fake	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	808
757	news detection on social media using geometric deep	try, Amanda Askell, Pamela Mishkin, Jack Clark,	809
758	learning. <i>arXiv preprint arXiv:1902.06673</i> .	et al. 2021. Learning transferable visual models from	810
		natural language supervision. In <i>International Con-</i>	811
759	Qiong Nan, Juan Cao, Yongchun Zhu, Yanyan Wang,	<i>ference on Machine Learning</i> , pages 8748–8763.	812
760	and Jintao Li. 2021. MDFEND: Multi-domain fake		
761	news detection. In <i>Proceedings of the 30th ACM In-</i>	Nir Rosenfeld, Aron Szanto, and David C Parkes. 2020.	813
762	<i>ternational Conference on Information &amp; Knowledge</i>	A kernel of truth: Determining rumor veracity on	814
763	<i>Management</i> , pages 3343–3347.	twitter by diffusion pattern alone. In <i>Proceedings of</i>	815
		<i>The Web Conference 2020</i> , pages 1018–1028.	816
764	Van-Hoang Nguyen, Kazunari Sugiyama, Preslav		
765	Nakov, and Min-Yen Kan. 2020. Fang: Leveraging	Juan Carlos Medina Serrano, Orestis Papakyriakopou-	817
766	social context for fake news detection using graph	los, and Simon Hegelich. 2020. NLP-based feature	818
767	representation. In <i>Proceedings of the 29th ACM In-</i>	extraction for the detection of COVID-19 misinfor-	819
768	<i>ternational Conference on Information &amp; Knowledge</i>	mation videos on YouTube. In <i>Proceedings of the 1st</i>	820
769	<i>Management</i> , pages 1165–1174.	<i>Workshop on NLP for COVID-19 at ACL 2020</i> .	821
770	Femi Olan, Uchitha Jayawickrama, Emmanuel Ogiem-	Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dong-	822
771	wonyi Arakpogun, Jana Suklan, and Shaofeng Liu.	won Lee, and Huan Liu. 2020. Fakenewsnet: A data	823
772	2024. Fake news on social media: the impact on soci-	repository with news content, social context, and spa-	824
773	ety. <i>Information Systems Frontiers</i> , 26(2):443–458.	tiotemporal information for studying fake news on	825
		social media. <i>Big data</i> , 8(3):171–188.	826
774	OpenAI. 2022. Chatgpt: Optimizing language mod-	Karen Simonyan and Andrew Zisserman. 2015. Very	827
775	els for dialogue. <a href="https://openai.com/blog/chatgpt/">https://openai.com/blog/</a>	deep convolutional networks for large-scale image	828
776	<a href="https://openai.com/blog/chatgpt/">chatgpt/</a> .	recognition. In <i>Proceedings of the 3rd International</i>	829
		<i>Conference on Learning Representations</i> .	830
777	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> ,	Ling Sun, Yuan Rao, Yuqian Lan, Bingcan Xia, and	831
778	<a href="#">arXiv:2303.08774</a> .	Yangyang Li. 2023. HG-SL: Jointly learning of	832
		global and local user spreading behavior for fake	833
779	Huyen Trang Phan, Ngoc Thanh Nguyen, and Dosam	news early detection. In <i>Proceedings of the AAAI</i>	834
780	Hwang. 2023. Fake news detection: A survey of	<i>Conference on Artificial Intelligence</i> , volume 37,	835
781	graph neural network methods. <i>Applied Soft Comput-</i>	pages 5248–5256.	836
782	<i>ing</i> , page 110235.	Nguyen Vo and Kyumin Lee. 2018. The rise of	837
		guardians: Fact-checking url recommendation to	838

839	combat fake news. In <i>The 41st international ACM SIGIR Conference on Research &amp; Development in Information Retrieval</i> , pages 275–284.	Yongchun Zhu, Qiang Sheng, Juan Cao, Qiong Nan, Kai Shu, Minghui Wu, Jindong Wang, and Fuzhen Zhuang. 2022. Memory-guided multi-view multi-domain fake news detection. <i>IEEE Transactions on Knowledge and Data Engineering</i> .	893
840			894
841			895
842	Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. <i>Science</i> , 359(6380):1146–1151.		896
843			897
844			
845	Longzheng Wang, Chuang Zhang, Hongbo Xu, Yongxiu Xu, Xiaohan Xu, and Siqi Wang. 2023. Cross-modal contrastive learning for multimodal fake news detection. In <i>Proceedings of the 31st ACM International Conference on Multimedia</i> , pages 5696–5704.		
846			
847			
848			
849			
850	Yaqing Wang, Fenglong Ma, Zhiwei Jin, Ye Yuan, Guangxu Xun, Kishlay Jha, Lu Su, and Jing Gao. 2018. EANN: Event adversarial neural networks for multi-modal fake news detection. In <i>Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery &amp; Data Mining</i> , pages 849–857.		
851			
852			
853			
854			
855			
856			
857	Youze Wang, Shengsheng Qian, Jun Hu, Quan Fang, and Changsheng Xu. 2020. Fake news detection via knowledge-driven multimodal graph convolutional networks. In <i>Proceedings of the 2020 International Conference on Multimedia Retrieval</i> , pages 540–547.		
858			
859			
860			
861			
862	Lianwei Wu, Pusheng Liu, and Yanning Zhang. 2023. See how you read? multi-reading habits fusion reasoning for multi-modal fake news detection. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 37, pages 13736–13744.		
863			
864			
865			
866			
867	Shu Yin, Peican Zhu, Lianwei Wu, Chao Gao, and Zhen Wang. 2024. Gamc: an unsupervised method for fake news detection using graph autoencoder with masking. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 347–355.		
868			
869			
870			
871			
872	Kaiwei Zhang, Junchi Yu, Haichao Shi, Jian Liang, and Xiao-Yu Zhang. 2023. Rumor detection with diverse counterfactual evidence. In <i>Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pages 3321–3331.		
873			
874			
875			
876			
877	Xueyao Zhang, Juan Cao, Xirong Li, Qiang Sheng, Lei Zhong, and Kai Shu. 2021. Mining dual emotion for fake news detection. In <i>Proceedings of the Web Conference 2021</i> , pages 3465–3476.		
878			
879			
880			
881	Xinyi Zhou, Jindi Wu, and Reza Zafarani. 2020. SAFE: Similarity-aware multi-modal fake news detection. In <i>Pacific-Asia Conference on Knowledge Discovery and Data Mining</i> , pages 354–367.		
882			
883			
884			
885	Xinyi Zhou and Reza Zafarani. 2019. Network-based fake news detection: A pattern-driven approach. <i>ACM SIGKDD explorations newsletter</i> , 21(2):48–60.		
886			
887			
888	Yangming Zhou, Yuzhou Yang, Qichao Ying, Zhenxing Qian, and Xinpeng Zhang. 2023. Multimodal fake news detection via clip-guided learning. In <i>2023 IEEE International Conference on Multimedia and Expo</i> , pages 2825–2830.		
889			
890			
891			
892			

898	<b>A Baseline Details</b>		
899	<b>A.1 Uni-Modal (5 Methods):</b>		
900	• <b>BERT</b> (Devlin et al., 2019): applied for analyzing textual data extracted from video titles and transcripts.		944
901			945
902			946
903	• <b>VGGish</b> (Hershey et al., 2017): processes audio segments from videos to compute spectrograms, crucial inputs for the VGGish model.		947
904			948
905			949
906	• <b>VGG19</b> (Simonyan and Zisserman, 2015): used to encode individual video frames, capturing intricate visual details essential for precise video content analysis.		950
907			951
908			
909			
910	• <b>C3D</b> (Ji et al., 2013): combined with an MLP layer for generating predictions, particularly effective in recognizing temporal patterns in video sequences.		
911			
912			
913			
914	• <b>Dual Emotion</b> (Zhang et al., 2021): introduces dual emotion features to improve existing fake news detection systems by distinguishing between publisher and social emotions.		
915			
916			
917			
918			
919	<b>A.2 Multi-Modal (4 Methods):</b>		
920	• <b>Serrano et al. 2020</b> extract TF-IDF features from video titles and the first 100 comments, employed in logistic regression and SVM classifiers.		
921			
922			
923			
924	• <b>FANVM</b> (Choi and Ko, 2021) aims to detect fake news videos across various topics, estimating topic distributions from video descriptions and comments.		
925			
926			
927			
928	• <b>EANN</b> (Wang et al., 2018) detects fake news using both image and text data, leveraging an event adversarial neural network to learn invariant features across different events.		
929			
930			
931			
932	• <b>SV-FEND</b> (Qi et al., 2023a) integrates multi-modal data using cross-modal transformers to capture correlations among text, audio, visual data, and social contextual features.		
933			
934			
935			
936	<b>A.3 Large Language Models (2 Methods):</b>		
937	• <b>GPT3.5-turbo</b> (OpenAI, 2022) is an optimized iteration of the GPT-3.5 model, engineered to deliver rapid responses without compromising text quality. Tailored for applications needing real-time natural language processing, it reduces latency and computational demands through fine-tuned efficiencies.		
938			
939			
940			
941			
942			
943			
		• <b>GPT-4-turbo</b> (OpenAI, 2023) expands on GPT-3.5-turbo’s foundation by incorporating advanced training methodologies and a broader dataset. This evolution enhances its ability to produce nuanced and contextually precise outputs, making it ideal for addressing intricate conversational contexts and diverse linguistic tasks.	
			944
			945
			946
			947
			948
			949
			950
			951
	<b>B Hyperparameter Details</b>		952
	To identify the optimal combination of model hyperparameters, the hyperparameters of corr link, feature dimension, contrastive learning, and pruning are thoroughly examined. The specific results are presented in Table 5.		953
			954
			955
			956
			957
			958
			959
			960
			961
			962
			963
			964
			965
			966
			967
			968
			969
			970
			971
			972
			973
			974
			975
			976
	<b>C Interpretability Note and Case Note</b>		977
	For relationships between videos, Qi et al. 2023b have explored that establishing potential relationships in the use of disinformation videos can help in the detection of false short videos. However, in daily life, it does not make sense to perform the detection of fake short videos after obtaining a debunk video. The importance of this work is to quickly perform the detection of whether a short video is fake or not based on the video features and potential relationships through the key nodes of the event outbreak without debunk.		978
			979
			980
			981
			982
			983
			984
			985
			986
			987
			988
			989
			990
			991
			992

Table 5: Hyperparameter Analysis Results

Hyperparameter	Range	Optimal Value
Corr Link $K$	{1, 5, 10, 15, 20, 25}	10
Feature Dimension $D$	{128, 256, 512, 768, 1024}	768
Contrastive Learning $\lambda$	{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0}	0.4
Pruning Ratio $e$	{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0}	0.2

993 here. Therefore, discovering the nature of short  
994 videos based on events is more representative than  
995 learning features based on single video relation-  
996 ships.